

# Efficient and Adaptive Estimation for Semiparametric Models

**Springer**

*New York*

*Berlin*

*Heidelberg*

*Barcelona*

*Budapest*

*Hong Kong*

*London*

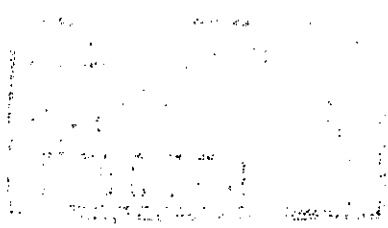
*Milan*

*Paris*

*Santa Clara*

*Singapore*

*Tokyo*



51.127

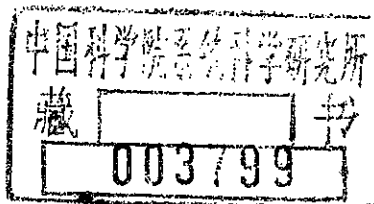
B 583

Peter J. Bickel    Chris A. J. Klaassen  
Ya'acov Ritov    John A. Wellner

# Efficient and Adaptive Estimation for Semiparametric Models



Springer



Peter J. Bickel  
Department of Statistics  
University of California  
at Berkeley  
Berkeley, CA 94720  
USA

Chris A.J. Klaassen  
Korteweg – de Vries Institute  
for Mathematics  
University of Amsterdam  
Plantage Muidersgracht 24  
1018 TV Amsterdam  
The Netherlands

Ya'acov Ritov  
Department of Statistics  
The Hebrew University of Jerusalem  
Jerusalem 91905  
Israel

Jon A. Wellner  
Department of Statistics  
University of Washington  
Seattle, WA 98195  
USA

©1993 The Johns Hopkins University Press

Paperbound edition published in 1998 by Springer-Verlag New York, Inc., by arrangement with The Johns Hopkins University Press.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the publisher.

Published in the United States by Springer-Verlag New York, Inc.

Springer-Verlag New York, Inc.  
175 Fifth Avenue  
New York, NY 10010  
USA

Library of Congress Cataloging-in-Publication Data

Efficient and adaptive estimation for semiparametric models / by Peter J. Bickel . . . [et al].  
p. cm.—(Johns Hopkins series in the mathematical sciences : [no. not given])

Includes bibliographical references and index.

ISBN 0-387-98473-9

1. Estimation theory. 2. Asymptotic expansions. I. Bickel, Peter J. II. Series.

QA276.8.E374 1998

519.5'44—dc21

98-17538

Production managed by Allan Abrams; manufacturing supervised by Thomas King.  
Printed and bound by Maple-Vale Book Manufacturing Group, York, PA.  
Printed in the United States of America.

9 8 7 6 5 4 3 2 1

ISBN 0-387-98473-9 Springer-Verlag New York Berlin Heidelberg SPIN 10667414

To Nancy	P.J.B.
To Angeli and in memory of my Mother, G.H. Claassen	C.A.J.K.
To Ilana and in memory of my Father A.R.	Y.R.
To my Father, Charles, and in memory of my Mother, Ethel	J.A.W.



# Contents

**List of Examples** ix

**List of Figures** xv

**Preface** xvii

- 1 **Introduction**
  - 1.1 Preliminaries 1
  - 1.2 Problems and scope 2
  - 1.3 Examples 5
- 2 **Asymptotic Inference for (Finite-Dimensional) Parametric Models**
  - 2.1 Regular parametric models in the I.I.D. case 11
  - 2.2 Regular estimates of Euclidean parameters 17
  - 2.3 The information bound and the Hájek–Le Cam convolution and asymptotic minimax theorems 23
  - 2.4 Nuisance parameters, adaptation, and some geometry 27
  - 2.5 Construction of  $\sqrt{n}$ -consistent and efficient estimates 40
- 3 **Information Bounds for Euclidean Parameters in Infinite-Dimensional Models**
  - 3.1 Introduction and overview 46
  - 3.2 Tangent spaces 48
  - 3.3 Information bounds via derivatives of functions: the nonparametric approach 57
  - 3.4 Information bound calculations via scores: the semiparametric approach 70
- 4 **Euclidean Parameters: Further Examples**
  - 4.1 Introduction: models 83
  - 4.2 Semiparametric group models 88
  - 4.3 Regression models 103
  - 4.4 Biased sampling models 113
  - 4.5 Mixture models 125
  - 4.6 Missing data models 143
  - 4.7 Transformation models 153

<b>5</b>	<b>Information Bounds for Infinite-Dimensional Parameters</b>	
5.1	Introduction	176
5.2	Convolution theorems for regular estimates of infinite-dimensional parameters	177
5.3	Examples	191
5.4	Differentiability of functions	201
5.5	The “calculus” of efficient score and influence operators	210
<b>6</b>	<b>Infinite-Dimensional Parameters: Further Examples</b>	
6.1	Introduction	221
6.2	Constrained families	221
6.3	Group models	229
6.4	Biased sampling models	240
6.5	Mixture models and models with monotonicity constraints	261
6.6	Missing data and censoring	271
6.7	Transformation models	292
<b>7</b>	<b>Construction of Estimates</b>	
7.1	Introduction	298
7.2	$M$ -estimates for Euclidean parameters	301
7.3	Generalized $M$ -estimates for Euclidean parameters	309
7.4	$GMC$ - and $GM$ -estimates corresponding to convex $D$	325
7.5	Estimation of $P$ and other infinite-dimensional parameters: methods, consistency, and rates of convergence	335
7.6	Estimation of infinite-dimensional parameters: asymptotics and applications	356
7.7	Joint estimation of Euclidean and infinite-dimensional parameters	382
7.8	Efficient estimation	391
	<b>Appendix</b>	
A.1	Vector spaces; linear functionals and dual spaces	414
A.2	Orthogonality and projection formulas	425
A.3	Conditional expectation formulas	430
A.4	Projection on sumspaces and ACE	436
A.5	Derivatives	453
A.6	Metrics on classes of probability measures and probability inequalities	464
A.7	Limit theorems, weak convergence, and tightness	468
A.8	Hoffmann-Jørgensen-Dudley weak convergence theory	475
A.9	Contiguity	498
A.10	The master theorem for asymptotic generalized $M$ -estimates	514
	<b>List of Symbols</b>	521
	<b>Bibliography</b>	527
	<b>Author Index</b>	549
	<b>Subject Index</b>	553



# List of Examples

## CHAPTER 1

### SECTION 3

1. Location models 5
2. Linear regression with nonparametric errors 6
3. Mixture models 7
4. Biased sampling models 8
5. Censored linear regression 9
6. Constraint defined models 9
7. Cox's proportional hazards model 9

## CHAPTER 2

### SECTION 1

1. Exponential family 14
2. Translation model with known shape  $f$  15
3. Weibull translation model 15
4. Three-parameter lognormal model 15

### SECTION 2

1. Stein's estimator of a normal mean 22
2. Minimum Kolmogorov distance estimator of center of symmetry 22
3. Bickel-Hodges estimate of center of symmetry 23

### SECTION 3

1. Quadratic loss 26
2. Zero-one loss 27

### SECTION 4

1. Normal location-scale 29
2. Reparametrization of normal location-scale 29
3. The bivariate normal distribution 32
4. The multinomial distribution 33
5. The Gaussian linear regression model 35
6. The bivariate normal distribution, continued 36

**CHAPTER 3****SECTION 1**

1. Density estimation 48

**SECTION 2**

1. All probabilities dominated by  $\mu$ ,  $\mathbf{P} = \mathbf{M}_\mu$  52
2. All probabilities symmetric about  $\theta_0$  fixed 53
3. Constraint defined models 53
4. The symmetric location model 55

**SECTION 3**

1. The symmetric location model 58
2. Estimation of the mean,  $\mathbf{P}$  unconstrained 67
3. Estimation of the mean,  $\mathbf{P}$  constrained 68
4.  $d_K$ -differentiable parameters on  $\mathbf{M} = \{\text{all } P \text{ on } \mathbf{X}\}$  69

**SECTION 4**

1. The symmetric location model 75
2. Cox model without censoring 77

**CHAPTER 4****SECTION 2**

1. The Gaussian linear model 89
2. The general linear model 89, 95
3. The elliptic model 89, 96
4. Two-sample location and scale 100
5. Two-sample scale mixture model 101

**SECTION 3**

1. Nonlinear regression 104, 105
2. Heteroscedasticity 104, 105
3. Projection pursuit regression 107, 109
4. Periodic function regression 107, 110
5. "Partial spline" regression models 107, 110
6. Logistic partial spline model 111

**SECTION 4**

1. Case control studies 113, 118, 121
2. Choice-based sampling 114, 118
3. Truncated regression and extensions 115, 119
4. Vardi's model and stratified sampling 115, 122

**SECTION 5**

1. The Neyman-Scott models 127, 133
2. Errors in variables 127, 135, 138
3. Bernoulli pairs with common odds ratio 128
4. Paired exponential mixture model 134

## SECTION 6

1. Nonparametric core model 144
2. Missing observations on a component of  $X^0$  144
3. Regression with missing observations on the covariates 146
4. Linear regression with right censoring 147

## SECTION 7

1. Linear regression-transformation model 153
2. Regression-transformation model 154
3. Joint distribution-transformation model 154, 170, 171, 172
4. Copula models 155
  - 4.A. Copula model with one unknown marginal df 155
  - 4.B. Copula model with two unknown marginals 156
  - 4.I. Archimedean copulas 156
    - 4.I.1. Clayton-Oakes Archimedean copula 156,173
    - 4.I.2. Stable frailty Archimedean copula 157
    - 4.I.3. Frank's Archimedean copula 157, 173
  - 4.II. Morgenstern distributions 157, 174
  - 4.III. Bivariate normal copula function 157, 174
  - 4.IV. Plackett's constant odds model 157
5. Regression-copula models 158
6. Linear regression-transformation model with unknown error distribution 158

## CHAPTER 5

## SECTION 3

1. Estimation of a df  $F$  on  $R$  191
2. Estimation of a df  $F$  with  $\gamma = \gamma(F)$  known 193
3. Estimation of a symmetric df 193
4. Estimation of a df  $F$  in a regular parametric model 195
5. Estimation of a cumulative hazard function 196
6. Estimation of mean residual life 197
7. Estimation of mean residual life in a family  $P$  with known mean 198
8. Estimation of a probability measure  $P$  199

## SECTION 4

1. Indicator censoring model 207

## SECTION 5

1. Estimation of a distribution function up to its mean 216
2. Estimation of the distribution function  $G$  corresponding to the baseline cumulative hazard function  $\Lambda$  in the Cox model 217

## CHAPTER 6

## SECTION 2

1. Estimation of a constrained distribution  $P$  222

2. Estimation of a bivariate distribution  $P$  with one known marginal 223
3. Estimation of a bivariate distribution  $P$  with two known marginals 225
4. Estimation of a distribution with independence 227

## SECTION 3

1. Spherically symmetric distributions in  $R^d$  229, 232
2. Coordinatewise symmetric distributions 229, 233
3. Exchangeable distributions 229, 233
4. Cyclically symmetric distributions 230, 234
5. Dihedrally symmetric distributions 230, 234
6. Rotationally symmetric distributions on the sphere 230
7.  $A = \{\text{Location group}\}$ ,  $T = \{\text{identity}\}$  235
8. Symmetric location model 235
9. Elliptic distributions 237
10. Partial splines, projection pursuit, and period regression, generalized 240

## SECTION 4

1. Random truncation model; nonparametric view 240
2. Random truncation model; score operator approach 247
3. Truncated regression 253

## SECTION 5

1. Mixtures of uniforms: distributions with monotone density 261, 263
2. Scale mixtures of exponentials: distributions with completely monotone density 261, 264, 266
3. Poisson mixture model 262, 264, 266
4. Location mixtures of Gaussians 262, 264, 266
5. Scale mixtures of centered Gaussians 262, 264, 266
6. Increasing failure rate (IFR) distributions 269, 270
7. Increasing failure rate average (IFRA) distributions 269
8. Schur-concave distributions 269, 271

## SECTION 6

1. Random censoring; nonparametric view 272
2. Random censoring; score operators and martingale theory 276
3. Random censoring; score operators and integral equation theory 283
4. Censored linear regression; example 4.6.4 continued 284
5. Double censoring 285
6. Bivariate censoring 289

## SECTION 7

1. Linear regression-transformation model 292
- 1.A. The Cox proportional hazards model 293
2. Joint distribution-transformation model 294
3. Copula model with one unknown marginal df 295

**CHAPTER 7****SECTION 2**

1. The median 303
2. Symmetric location in  $k$  dimensions 304
3. Regression models 305
4. Elliptic models 307
5. Mixture models and conditional likelihood 308
6. The Has'minskii-Ibragimov model 308

**SECTION 3**

1.  $M$ -estimates 310
2. The Hodges-Lehmann estimate 310, 315
3.  $M$ -estimates for  $v$  with  $\eta$  estimated 311
4. Estimation of location for the Cauchy distribution 317
5. Minimum Cramér-von Mises distance estimation for the exponential distribution 318
6. The regression-transformation model 319

**SECTION 4**

1. Exponential family 326
2.  $M$ -estimates in the linear model 328
3. Linear regression with right censoring 329
4. The Cox estimate 330

**SECTION 5**

1. Exponential mixture model 337
2. Absolutely continuous distributions on  $R$  338, 339
3. Absolutely continuous distributions symmetric about zero 340
4. Biased sampling 340
5. Random censoring 342
6. Symmetric location 343
7. Density estimation by the method of sieves 345, 350
8. Estimating a joint distribution with one or both marginals known 346
9. Regression estimation by penalized maximum likelihood 347, 353
10. Regularized MLE in the Has'minskii-Ibragimov model 348

**SECTION 6**

1. The Nelson-Aalen and Kaplan-Meier estimators 358
2. Estimating a joint distribution with one marginal known 360
3. Biased sampling 362
4. Density estimation by sieves 366
5. Regression by penalized maximum likelihood 369
6. The Has'minskii-Ibragimov model 372
7. The joint distribution-transformation model 380

## SECTION 7

1. Biased sampling regression 383
2. Censored regression 388

## SECTION 8

1. Symmetric location 398, 400, 403
2. Regression 399, 401
3. Elliptic model 399, 401
4. Paired exponential mixture model 399, 401

# List of Figures

## CHAPTER 2

1. Projection of score functions 31
2. Projection of influence functions 31
3. Score and influence function projections 31

## CHAPTER 3

1. Projection of pathwise derivatives 59
2. Projection of influence functions 66
3. Projection of score functions 71

## CHAPTER 5

1. The function diagram for theorem 1 203

## CHAPTER 6

1. Clayton-Oakes copula model:  
ARE of empirical df to efficient estimator 295
2. Frank copula model:  
ARE of empirical df to efficient estimator 296

## APPENDIX

1. Alternating projections 437
2. Rényi's distribution with  $\rho = 1$  442





# Preface

This book is about estimation in situations when we believe we have enough knowledge to model some features of the data parametrically, but are unwilling to assume anything for other features. For example in the two sample case of the famous Cox proportional hazards model, we assume the treatment effect can be modeled multiplicatively (parametrically) on some completely unknown scale. Such models have arisen in a wide variety of contexts in recent years, particularly in economics, epidemiology, and astronomy. The complicated structure of these models typically requires us to consider nonlinear estimation procedures which often can only be implemented algorithmically. The theory of these procedures is necessarily based on asymptotic approximations, while actual performance for finite sample sizes is often gauged best by simulations.

Therefore our focus is on asymptotic theory. We limit ourselves to models for independent, identically distributed observations, the basic building blocks of most models for data. Easily understandable general results are available in this context, and, as usual, the methods and theory for these cases guide extensions to more complicated models.

Our goals are:

1. To show how the information bounds and methods of estimation developed in the contexts of non- and semiparametric models can be viewed as natural extensions and developments of the corresponding bounds and methods in the classical parametric model context.
2. To apply these techniques in as broad a range of models as possible, illustrating the ease with which heuristic calculations of "optimal behavior" can be carried out.
3. To develop the theory of information bounds for estimation of infinite-dimensional parameters.
4. To develop a coherent heuristic view of the methods of estimation actually used in semiparametric models.

A companion goal of giving simple necessary and sufficient conditions for particular methods to work as expected asymptotically has met with only partial success.

A 600-page book poses a challenge for any reader, and perhaps some introductory guidelines and suggestions will prove helpful:

Chapters 1–3 are fundamental and should be studied, at least to the point of understanding theorem statements, by all readers. Chapter 4 should be dipped into by all readers interested in Euclidean parameter examples. The sections of chapter 4 are relatively independent of each other, but rely heavily on the development in chapters 2 and 3 together with some material from appendix sections A.1–A.3 and occasionally section A.4. Sections 7.1–7.4 contain an important, but relatively elementary, part of the theory of estimate construction for Euclidean parameters. They, along with chapter 2 and the complete proofs for sections 7.2 through 7.4 given in section A.10, are essentially free-standing. Most material in chapters 1–4 and 7.1–7.4 require only a minimum of functional analysis.

Chapter 5, which parallels chapter 3, presents the basic theory of information bounds for infinite-dimensional parameters. Similarly, chapter 6 parallels chapter 4. Although these chapters are not a prerequisite for reading sections 7.5–7.8, which still deal extensively with Euclidean parameters, these sections, and chapters 5 and 6, require much more familiarity with functional analysis than the earlier chapters and sections of chapter 7. We have collected and summarized the necessary functional analysis and other material required for reading this book in the appendix. Most sections of the appendix can be read independently, but there are strong interconnections between several sections of the appendix, especially section A.2 through A.4. Extensive references to standard texts are given as needed. We refer frequently to the appendices for supporting material and results. If the reader finds these immediate results unfamiliar, then the corresponding entire appendix section should be used as background reading.

This book began with the 1983 Mathematical Sciences Lectures at Johns Hopkins University given by Bickel and Wellner. We (Bickel and Wellner) are grateful to the Department of Mathematical Sciences at Johns Hopkins for the opportunity to present our ideas on the subject, and for spurring us on to develop them further. We (Bickel and Wellner) realized quickly that we knew fairly little when we presented these lectures. Even though Chris Klaassen and Ya'acov Ritov soon joined us in our study of the area, learning enough to write this book took seven years. On the whole it has been fun.

We would like to particularly acknowledge our intellectual debt to Lucien Le Cam, Charles Stein, and Johann Pfanzagl. Le Cam could have written this book in 1956, and has published a treatise (Le Cam (1986)) which probably includes ours—but the language of our books is somewhat different! Stein (1956) introduced the heuristic underlying the development by Levit, Ibragimov and Has'minskii, and Pfanzagl of information-bound theory in semiparametric models. Pfanzagl and Wefelmeyer (1982) wrote a precursor of our book which, along with Begun, Hall, Huang, and Wellner (1983) and Bickel (1982), formed the basis of our lectures. Our book, in part, is an effort to translate the ideas of these authors into language which is current in mainstream statistics and to demonstrate their wide applicability.

We also wish to acknowledge comments and constructive criticism by a number of colleagues: Jack Cuzick, Richard Gill, Friedrich Götze, Piet Groeneboom, Jack Hall, Nicholas Jewell, Whitney Newey, Aad van der Vaart, and Willem van Zwet. Section A.8 of the appendix was developed in joint work by Aad van der Vaart and Wellner, and has benefitted from comments and corrections by Peter Gänssler and Michael Wichura. Various preliminary versions of this book have been presented in graduate courses at Berkeley, Seattle, Leiden, and Amsterdam. We owe thanks to the students in these courses for comments and lists of errata. We wish to single out, in particular, Kun Jin, Panos Lorentziadis, and Peter Sasieni, Jens Praestgaard, and Mary Emond.

Despite our serious effort to catch them all, errors undoubtedly remain. We take full responsibility for these, and will greatly appreciate receiving information concerning any errors which readers notice.

Our research and writing on this book has been supported in part by a number of government agencies and foundations: Bickel and Ritov have received research support from the Office of Naval Research. Wellner has received support from the National Science Foundation and the John Simon Guggenheim Foundation. Grants from the Israel-Netherlands Cultural Exchange program funded visits of Ritov to Klaassen and Wellner in Leiden and of Klaassen to Ritov in Jerusalem. A month of editing work at the Mathematical Sciences Research Institute in Berkeley by Klaassen and Wellner in the fall of 1991 was supported by the National Science Foundation.

We thank Richard O'Grady, our last editor at Johns Hopkins Press, for spurring us on to finish.

Last, but not least, we wish to thank Chris Bush. Her fabulously quick and accurate mathematical typing (and e-mails to authors scattered over the globe) made this four-author collaboration feasible.



EFFICIENT AND ADAPTIVE ESTIMATION FOR  
SEMPARAMETRIC MODELS



# 1 | Introduction

## 1.1 PRELIMINARIES

Let  $X_1, \dots, X_n$  be a sample from the probability distribution  $P$  on  $(\mathbf{X}, \mathcal{B})$ , where  $\mathbf{X}$  is some Euclidean sample space and  $\mathcal{B}$  denotes its Borel  $\sigma$ -field. That is,  $X_1, \dots, X_n$  are independent and identically distributed with common unknown distribution  $P$  on  $(\mathbf{X}, \mathcal{B})$ . The problem of statistical inference in this context can be viewed as using the sample to gain information about some features of  $P$ . We often think of these features as a vector parameter  $v(P) = (v_1(P), \dots, v_m(P))$ . Classically we proceed by assuming that (at least approximately)  $P$  is determined by  $v$  and perhaps an additional nuisance parameter  $\eta$ . If  $\theta = (v, \eta)$ , we think of  $P$  as ranging over  $\mathbf{P}$ , a subset of the collection  $\mathbf{M}$  of all probability measures on  $(\mathbf{X}, \mathcal{B})$ , which is describable by a map  $\theta \rightarrow P_\theta$  with  $\theta$  ranging over a Euclidean parameter space  $\Theta$ . This is the usual notion of a parametric model. It has as elements the sets  $\mathbf{P}$ ,  $\Theta$ , and the parametrization  $\theta \rightarrow P_\theta$ . We will reserve the term *model* without qualification for a subset  $\mathbf{P}$  of  $\mathbf{M}$ .

The “nuisance parameters”  $\eta$  necessary to describe  $P$  may influence the variability of estimates of  $v$ , but do not need to be known as closely as  $v$ . This distinction is certainly not hard and fast, but such hierarchies clearly exist in practice. For example, in one-sample problems, location parameters are of greatest interest, then scale parameters, then skewness and kurtosis; in linear models the hierarchy of main effects, interactions, and then scale parameters is common; in the Cox regression model (to be described in example 7 in section 3) the regression parameters are usually of primary interest while the common hazard function is a nuisance parameter of secondary interest.

In recent years, realism and larger data sets have led statisticians to focus increasingly on *nonparametric* or *semiparametric* models in which  $\mathbf{P}$  is a “big” subset of the collection of all probability distributions. Roughly speaking, we will use the term *nonparametric* for models  $\mathbf{P}$  which are very large subsets (or all) of  $\mathbf{M}$ , and the word *semiparametric* will be used for models  $\mathbf{P}$  which are intermediate between  $\mathbf{M}$  and some classical parametric model  $\mathbf{P}_0$ . These models are often described in terms of a Euclidean parameter  $\theta$ , as above, together with a function  $g$  or  $G$  in some set of functions  $\mathbf{G}$ . (Often  $g^{1/2}$  is an element of some

Hilbert space  $\mathbf{H}$ , or at least  $G$  is an element of some Banach space  $\mathbf{B}$ .) Then the model is specified by the sets  $\mathbf{P}$ ,  $\Theta$ ,  $\mathbf{G}$ , and the parametrization is given by the function  $(\theta, g) \rightarrow P_{(\theta, g)}$  for  $(\theta, g) \in \Theta \times \mathbf{G}$ .

Another useful way of thinking about semiparametric models is in terms of the important notion of the *tangent space* of any given model, a notion which we will define carefully in section 3.2. Informally however, the tangent space at a point  $P_0$  of the model  $\mathbf{P}$  is simply the subspace  $\dot{\mathbf{P}}$  of  $L_2(P_0)$  spanned by the score functions of all regular parametric submodels in  $\mathbf{P}$  passing through  $P_0$ . For a *nonparametric* model,  $\dot{\mathbf{P}}$  consists of the entire subspace  $L_2^0(P_0)$  of  $L_2(P_0)$  which is orthogonal to 1 (i.e., with  $P_0$  - mean zero). In our terminology, a model  $\mathbf{P}$  is *semiparametric* at  $P_0 \in \mathbf{P}$  if  $\dot{\mathbf{P}}$  is infinite-dimensional, and yet  $\dot{\mathbf{P}}$  is a proper subspace of  $L_2^0(P_0)$ , the tangent space of a completely nonparametric model. Note that any regular parametric model  $\mathbf{P}_0$  has a finite-dimensional tangent space  $\dot{\mathbf{P}}_0$  at  $P_0 \in \mathbf{P}_0$ , namely the subspace  $\dot{\mathbf{P}}_0$  spanned by the finite collection of score functions of the model.

Of course, any subset  $\mathbf{P}$  of  $\mathbf{M}$  can be “parametrized” by a map  $\theta \rightarrow P_\theta$ ,  $\theta \in \Theta$ . We can even always take  $\Theta$  to be finite-dimensional. However, regular parametric models (which we define in chapter 2) put smoothness conditions on the map which make the distinction clear.

As we will see in the following, semiparametric models  $\mathbf{P}$  can frequently be described or parametrized in more than one way. One common way is in terms of some “mixed” parametrization in which the model  $\mathbf{P}$  is built up from elementary pieces, such as in the usual symmetric location model or the Cox regression model. Another useful method of describing a semiparametric model is via “side conditions” or constraints, a method which we will elaborate upon in example 6 in section 3.

In general, we shall focus on functionals  $v: \mathbf{P} \rightarrow R$ , or  $R^m$ , or even an abstract space, and their extensions to functionals defined on  $\mathbf{M}$ . Such a functional evaluated at  $\mathbf{P}$  is a *parameter*, while an extension evaluated at the empirical distribution is an *estimate* of the parameter.

## 1.2 PROBLEMS AND SCOPE

Our goals in this monograph are to address the following questions:

- A. How well can Euclidean or finite-dimensional parameters of a nonparametric or semiparametric model be estimated? (What structure must efficient estimates have if they exist?)
- B. What is the price in efficiency which must be paid when we extend a classical parametric model  $\mathbf{P}_0$  to a supermodel  $\mathbf{P}$ ? In what situations is such an extension “free?” Or, equivalently, when can we adapt?
- C. How well can the infinite-dimensional or “nuisance” parameters of a semiparametric model be estimated?
- D. What are general methods and techniques for constructing asymptotically efficient estimates for such models?



As the reader will see, we have made considerable headway on answers to questions A–C. Although we have made some progress on question D, the overall picture is somewhat disappointing. There are a number of methods that heuristically should yield procedures with the properties we want. But which approach works best or can most easily be proved to work depends on the example or class of examples.

We begin in chapter 2 with a review of lower bound theory and estimation for regular parametric models. Section 2.4, where we develop some basic geometric results in the parametric case, is especially important for later chapters.

Chapter 3 contains our basic answers to questions A and B. The material in section 3.2 on tangent spaces is basic to the rest of the monograph. In section 3.3 we give an approach to information bounds for semiparametric models which derives from a more nonparametric view involving pathwise (Hadamard) differentiation of Euclidean-valued functions  $v$ . It stems from figure 2.4.2 and propositions 2.4.1.B, 2.4.2, and 2.4.3. This approach began in the work of Levit (1978) and Koshevnik and Levit (1976), and has been further developed by Pfanzagl and Wefelmeyer (1982), (1985), Pfanzagl (1990), and Van der Vaart (1988a). A key element of our approach is the emphasis on the asymptotic linearity of efficient procedures and identification of the linear approximation—what we call the *efficient influence function*. An alternative approach to information bounds for semiparametric models based on scores and score operators developed by Begun, Hall, Huang, and Wellner (1983) is outlined in section 3.4. This approach, which is based on what we call the *efficient score function*, continues the theme begun in proposition 2.4.1.A and figure 2.4.1 of section 2.4.

These same themes, but with different emphases and with greater generality, have been developed in the fundamental works of Le Cam (1986), Le Cam and Yang (1990), and Ibragimov and Has'minskii (1981).

Chapter 4 illustrates the methods of chapter 3 by applying them to a wide range of examples, including group models, regression models, biased sampling models, mixture models, missing data models, and transformation models. The prime examples of models for which the answer to question B is affirmative are given in section 4.1.

In chapter 5 we develop information bound theory for infinite-dimensional or “nuisance” parameters of nonparametric or semiparametric models. As we will see, estimation of infinite-dimensional nuisance parameters at the usual  $n^{-1/2}$  rate depends crucially on the continuity (or, equivalently, boundedness) of the inverse of a certain linear operator. More generally, when parameters are defined implicitly, as they are in many models of interest, which parameters are pathwise differentiable? For these parameters we will be able to state convolution theorems which suggest that they are estimable at the usual  $n^{-1/2}$  rate. The basic theorem of Van der Vaart—concerning necessary and sufficient conditions for pathwise differentiability of implicitly defined functions—is given in section 5.4, and various special cases and examples thereof are worked out in section 5.5. Chapter 6 contains further examples of the theory developed in chapter 5.

Finally, chapter 7 is concerned with the construction of estimators, both efficient and inefficient, for essentially all the examples we consider in chapters 2–6. We begin by extending the approach to  $M$ -estimates of Euclidean parameters due to Huber (1967) to the more general types of implicitly defined estimates introduced by Filippova (1962). We show how these methods can be applied to a number of important examples. We then consider estimation of function-valued parameters using methods motivated by parametric maximum likelihood and give conditions under which these methods behave as we heuristically expect them to. Unfortunately, the general conditions are easy to state, but verification in particular examples is formidable and requires taking full advantage of the individual structure of the models we consider. Finally we give a heuristic discussion of the potential for producing efficient estimates possessed by our various approaches, and show how this potential is realized in a number of our examples.

The appendix contains much of the basic background theory: linear operators and their adjoints and inverses; projection theorems and formulas, and conditional expectation formulas useful in calculation of projections; theory of projections on nonorthogonal subspaces; basic differentiation theory—Gâteaux, Hadamard, and Fréchet; weak convergence theory—both standard and the newer theory of Hoffmann-Jørgensen and Dudley; the contiguity theory of Le Cam; and a summary of finite-dimensional  $M$ -estimation theory.

### *What Remains to Be Done?*

The easy answer to this question is: A lot! We have not formally looked at testing or confidence regions. Admittedly asymptotic inference for Euclidean parameters is relatively straightforward: Efficient estimates coupled with consistent estimates of their variances yield efficient tests and confidence bands, as discussed for instance in Pfanzagl and Wefelmeyer (1982). Again, heuristic construction of consistent variance estimates using estimates of the efficient influence function is straightforward. But we have not investigated such constructions in our examples and a fortiori have not considered alternative approaches, such as the jackknife or the bootstrap. We have also not addressed the more difficult questions of testing goodness of fit to semiparametric models, or the related important questions of model selection and diagnostic construction. Nor have we considered robustness properties of the procedures we discuss, although the criteria developed by Hampel (see Hampel, Ronchetti, Rousseeuw, and Stahel (1986)) should obviously be taken into account and the criterion of efficiency challenged.

We complete the list of our omissions with the biggest. The entire development in this monograph and all examples treated are for i.i.d. data. The Local Asymptotic Normality (LAN) theory of Le Cam (1960), on which we based our development, is, of course, applicable quite generally to experiments which can be approximated locally by Gaussian shift experiments. The general treatment can be found in Le Cam (1986) and Le Cam and Yang (1990), and applications to situations in which data are independent but not identically distributed or are dependent have been given by many authors. Two recent very general treat-

ments, with examples, are by Has'minskii and Ibragimov (1991) and Van der Vaart and Wellner (1990). Unfortunately, in general, the theory loses some elegance and simplicity. In the i.i.d. case parameters are functions of the (marginal)  $P$ , which itself can be estimated without conditions by the empirical distribution  $\hat{P}_n$ . In general, the joint distributions  $P^{(n)}$  of  $(X_1, \dots, X_n)$  have to be considered, and there is no estimate of  $P^{(n)}$  which is satisfactory if we assume  $P^{(n)}$  can range freely. It is, however, possible, as suggested by Levit (1978), to develop a simple treatment of information bounds in a number of important non-i.i.d. cases.

Why don't we look at all these questions here? This book started out as a 150-page account of lectures given by Bickel and Wellner at Johns Hopkins University in 1983. As we proceeded, starting in 1984, we realized how little we knew when we gave the lectures. We began to learn more from the existing relevant literature and to answer some of the new questions that arose as we progressed. We began to appreciate in how many ways this work relied on functional analysis and abstract probability theory and to collect the relevant material in an appendix. On crossing the 5-year and 450-page mark we realized that enough was enough! Although we ourselves intend to look in some of the directions we mentioned, we are only too glad to be joined by some of our readers.

### 1.3 EXAMPLES

The examples which we now give serve to illustrate the wide range of models which will concern us throughout the rest of the book. More examples will be introduced and discussed (primarily) in chapters 4 and 6. Constructions of estimates in these examples are given in chapter 7. A complete list of all examples treated in the text may be found in the table of contents. There are many more examples which we have *not* treated, and new examples seem to be appearing at a rapid rate in econometrics, biostatistics, and other fields.

#### Example 1. Location models.

The familiar location model is

$$(1) \quad X = v + \varepsilon, \quad v \in R,$$

where  $\varepsilon$  is an "error." The usual parametric model has  $\varepsilon \sim N(0, \eta^2)$ ,  $\eta > 0$ , where  $N(\mu, \sigma^2)$  is the normal distribution with mean  $\mu$  and variance  $\sigma^2$ , and  $\sim$  denotes "distributed according to." Thus  $\mathbf{P}_0$  is just the set of normal distributions, and the parameter of interest  $v$  can be described as the mean, median, center of symmetry, etc.

Interest in  $\mathbf{P}_0$  arises in at least two different situations:

- (i) Measurement model:  $N(0, \eta^2)$  errors are being made in measuring an unknown constant  $v$ .
- (ii) Paired comparisons model: The  $X_i$ 's are differences in some measured response between treated and control members of pairs of subjects. The treatment effect  $v$  is assumed additive and the joint distribution of responses before treatment is bivariate normal with equal means.

The paired comparisons model has a natural nonparametric extension if we continue to take the treatment effect to be additive, namely to assume only that  $\varepsilon$  is distributed symmetrically about zero so that

$$(2) \quad \mathbf{P} = \{ P : P \text{ is symmetric} \} .$$

In this case we can think of  $\nu$  as the center of symmetry of  $P$  and take the nuisance parameter  $G$  to be the distribution of  $\varepsilon$ :

$$(3) \quad \mathbf{G} = \{ \text{all distributions symmetric about } 0 \} .$$

Natural extensions of the measurement model suppose  $\varepsilon$  to be arbitrarily distributed or at least (Huber (1964)) to be an  $N(\nu, \eta^2)$  distribution "contaminated" by arbitrary "gross errors":

$$(4) \quad \mathbf{P} = \{ P : P = (1 - p)N(\nu, \eta^2) + pG, \ G \text{ arbitrary} \} ,$$

where  $0 < p < 1$  is fixed. This extension raises an important conceptual issue. In both models (2) and (4), the map  $(\nu, G) \rightarrow P \in \mathbf{P}$  is well defined. However, in model (2) this map is one-to-one, and hence invertible, so that the parameter  $\nu$  is well defined as a function of  $P$ ; or, equivalently,  $\nu$  is identifiable. This is clearly not the case in model (4). It can be argued that in the measurement model  $\nu$  has an existence independent of  $P$ ; see Huber (1981) or Bickel (1981). However, in most situations it is conceptually more satisfactory to limit oneself to parameters which can be identified operationally—i.e., as functions of  $P$ . In this book we will limit ourselves to situations in which the parameters of interest  $\nu$  are functions on  $\mathbf{P}$  (perhaps defined implicitly), usually to some Euclidean space. Alternatively, we insist that the parameter  $\nu$  defined as a function on  $\mathbf{P}_0$  have an extension to  $\mathbf{P}$  which we specify and is to be our object of study. These questions of identifiability can be nontrivial when the model is defined structurally. For convenience, we may sometimes describe models in terms of unidentifiable parameters  $\eta$ , but with the understanding that, within the model so defined,  $\nu$  is definable as a function on  $\mathbf{P}$ .

As a third nonparametric extension of  $\mathbf{P}_0$  which illustrates these points of view, we can take

$$(5) \quad \mathbf{P} = \left\{ P : \int x^2 dP(x) < \infty \right\} ,$$

and  $\nu(P) \equiv E_P X = \int x dP(x)$ . Here we are really interested, as is typically the case in sample surveys, in the population mean  $\nu$ , and  $\mathbf{P}$  is a very natural nonparametric extension of  $\mathbf{P}_0$ .  $\square$

Multivariate location (and scale) models will be discussed in chapter 4. These location models and the following example of linear regression can all be viewed as special cases of the general class of "transformation" or "group" models, a class of models which will be defined and further studied in chapter 4.

**Example 2. Linear regression with nonparametric errors.**

The Gaussian model here is  $X = (Z, Y)$ , where

$$(6) \quad Y = Z^T v + \varepsilon$$

where  $\varepsilon \sim N(\mu, \sigma^2)$ ,  $Z \sim H$  is an  $m \times 1$  column vector independent of  $\varepsilon$ ,  $v$  is an  $m \times 1$  column vector of unknown parameters,  $E|Z|^2 < \infty$ , and

$$(7) \quad E(Z - EZ)(Z - EZ)^T \equiv \Sigma \text{ is nonsingular.}$$

This is a stochastic version of the ordinary linear model which is used here only to accommodate our restriction to one-sample models. In this parametric model  $P_0 \equiv \{P_\theta : \theta = (v, \mu, \sigma^2) \in R^m \times R \times R^+\}$  and the parameter  $\theta$  is identifiable since, from the usual least squares formulae,

$$(8) \quad \begin{aligned} v &= \Sigma^{-1} E((Z - EZ)Y), \\ \mu &= EY - EZ^T v, \\ \sigma^2 &= \text{Var}(Y|Z). \end{aligned}$$

One natural extension of  $P_0$  is to assume an arbitrary distribution  $H$  of  $Z$ ; a further extension is to let  $\varepsilon$  have an arbitrary distribution  $G$ . Then

$$(9) \quad P = \{P : Z \sim H \text{ satisfying (7), } v \in R^m, \varepsilon \sim G \text{ arbitrary}\}.$$

It will be seen that  $v$  is identifiable in model (9) and is even "adaptively estimable."

Other extensions are also possible along the lines of example 1. For instance,  $(v, \mu)$  can be defined by (8) for all joint distributions of  $X = (Z, Y)$  such that (8) makes sense; in other words, in this extension  $(v, \mu)$  are defined as the coefficients of the linear (population) regression of  $Y$  on  $(Z^T, 1)$ .  $\square$

#### Example 3. Mixture models.

A large class of interesting models consists of nonparametric mixtures of parametric models, or vice versa. In the first case, if  $f(x; v, \eta)$  is indexed by  $\theta = (v, \eta) \in \Theta \subset R^k$  and  $G \in \mathbf{G}$  is a collection of mixing distributions, then

$$(10) \quad f(x; v, G) \equiv \int f(x; v, \eta) dG(\eta) \quad \text{for } x \in \mathbf{X},$$

with  $v$  ranging over some subset of a Euclidean space, and  $G \in \mathbf{G}$  is the corresponding mixture family. Here

$$(11) \quad P = \{P : P \text{ has density } f \text{ given in (10) for some } v \text{ and } G \in \mathbf{G}\}.$$

Of course,  $\mathbf{G}$  must be specified so that  $v$  is identifiable.

There are many interesting special cases of the model (11); see, e.g., Kiefer and Wolfowitz (1956) or Lindsay (1980), (1983a,b,c). An especially interesting case which we will treat in chapters 4 and 7 is the classical "errors in variables" model described as follows: suppose that  $Z' \sim G$  where  $G$  has no Gaussian component (that is,  $G$  cannot be represented as a convolution of a nondegenerate Gaussian distribution with some other distribution), that  $(\varepsilon_1, \varepsilon_2) \sim N_2(0, \Sigma)$  with  $\Sigma$  unknown is independent of  $Z'$ , and that we observe  $X = (Z, Y)$  where

$$(12) \quad Z = Z' + \varepsilon_1,$$

$$(13) \quad Y = \alpha + \beta Z' + \varepsilon_2.$$

Thus,

$$(14) \quad \mathbf{P} = \{ P : P \text{ is the distribution of } X \text{ determined by (12) and (13) for some } \alpha, \beta, \Sigma, G \text{ with no Gaussian component} \}.$$

Note that  $Z'$  is not observed and plays the role of  $\eta$  in (10). The parameter of primary interest is  $\beta$  or the pair  $(\alpha, \beta)$ .

An interesting model which can be viewed as a parametric mixture of non-parametric distributions is the following bivariate model introduced by Clayton (1978). Let  $G$  and  $H$  be univariate survival functions and let  $W \sim \text{Gamma}(\nu, 1)$ . Then, conditional on  $W$ , suppose that  $X = (S, T)$ , where  $S, T$  are independent with survival functions  $G^W$  and  $H^W$  respectively. Hence,

$$(15) \quad P(S \geq s, T \geq t) = EG(s)^W H(t)^W \\ = \int_0^\infty G(s)^w H(t)^w \frac{w^{\nu-1}}{\Gamma(\nu)} e^{-w} dw$$

and

$$(16) \quad \mathbf{P} = \{ P : P \text{ has joint survival function given by (15) for some } \nu > 0 \text{ and some } G, H \}.$$

This model has also been discussed by Oakes (1982), and generalizations have been explored by Marshall and Olkin (1988).  $\square$

#### Example 4. Biased sampling models.

Suppose that given a covariate vector  $Z$  with values in  $R^d$ , the random variable  $Y$  takes on values in the set  $\{1, \dots, M\}$  with probabilities determined by the (parametric) function  $P(Y = y | Z = z, \nu) \equiv p(y | z, \nu)$  for  $y = 1, \dots, M$ . Suppose that  $Z \sim H$ . Then, conditional on  $Y$  taking values in a subset  $J$  of  $\{1, \dots, M\}$ , the distribution of  $X = (Z, Y)$  is given by

$$(17) \quad P(Y = y, Z \in A | Y \in J) = \frac{\int_A p(y | z, \nu) dH(z)}{\int_{R^d} \sum_{y' \in J} p(y' | z, \nu) dH(z)}$$

for  $y \in J$  and  $A \in \mathcal{B}^d$ . Here  $\theta = (\nu, H)$  and

$$(18) \quad \mathbf{P} = \{ P : P \text{ is given by (17) for some } \nu \text{ and } H \}.$$

The "choice-based" sampling model studied by Manski and Lerman (1977) and Cosslett (1981) involves selection of "sampling strata"  $J_1, \dots, J_s$  where each  $J_i$  is a subset of  $\{1, \dots, M\}$ , and then sampling "at random" from each "stratum" with distribution determined by (17) to estimate the parameters  $\nu$  and the distribution  $H$ .

A variety of models similar to this one, including the truncated regression model of Bhattacharya, Chernoff, and Yang (1983), related regression models

with stratification on the dependent variable studied by Hausman and Wise (1982) and Jewell (1985), and the selection bias models of Vardi (1985), can all be treated as special cases of a general class of "biased" or "stratified" sampling models. This will be done in chapters 4 and 6.  $\square$

**Example 5. Censored linear regression.**

In this model, introduced by Miller (1976) for applications in survival analysis, the pair  $(Z, Y)$  is as in the linear regression model of example 2, but we observe only

$$(19) \quad X = (Z, T, \Delta)$$

where  $T \equiv \min\{Y, C\}$ ,  $\Delta \equiv 1_{\{Y \leq C\}}$ , and  $C$  is a censoring time which, conditional on  $Z$ , is independent of  $Y$ . Thus the model  $\mathbf{P}$  is now the collection of distributions of  $X$  in (19):

$$(20) \quad \mathbf{P} = \{P: Z \sim H \text{ satisfying (7), } C \sim K(\cdot | z), \\ \varepsilon \sim G \text{ arbitrary, } v \in R^m \}.$$

Further conditions on the conditional distributions of  $C$  given  $Z$  are necessary for  $v$  to be identifiable; these will be discussed in chapter 4 along with information bounds for estimation of  $v$ .  $\square$

**Example 6. Constraint defined models.**

Suppose that  $\{\gamma_1, \dots, \gamma_r\}$  is a collection of functionals,  $\gamma_i: \mathbf{M} \rightarrow R$  (which are "regular" on some large subset  $\mathbf{M}_0$  of  $\mathbf{M}$ ), and let  $a_1, \dots, a_r$  be  $r$  fixed constants. Then consider

$$(21) \quad \mathbf{P} = \{P \in \mathbf{M} : \gamma_i(P) = a_i, \quad i = 1, \dots, r \}.$$

For example, if  $\mathbf{M}_0$  is the collection of all probability distributions on  $R^+$  with finite variances, and  $\gamma(P) \equiv \gamma_1(P) \equiv \sigma(P)/\mu(P)$  is the coefficient of variation, then the family  $\mathbf{P}$  consists of all distributions on  $R^+$  with coefficient of variation equal to  $a$ , and we can ask how well any other functional of  $P$ , such as the mean, or median, or  $P$  itself, can be estimated.

These kinds of families have been investigated by Koshevnik and Levit (1976), and Levit (1975). Information bounds for models such as these will be given in chapter 4. Note that many other families can be described in this way if  $\gamma$  is permitted to be vector- or function-valued: if  $\gamma(P)(x) \equiv P(-\infty, -x] - P[x, \infty)$ , then  $\gamma(P) \equiv 0$  defines the set of distributions which are symmetric about zero.  $\square$

**Example 7. Cox's proportional hazards model.**

As its name suggests, this model is most naturally described in terms of hazard functions. Suppose that a covariate vector  $Z \sim H$  and given  $Z = z$  the survival time  $Y$  has hazard function

$$(22) \quad \lambda(y | z, v) = r(z^T v) \lambda(y),$$

where  $\lambda(y) \equiv g(y)/(1 - G(y))$  is a common (baseline) hazard function,  $r$  is a fixed nonnegative function, and  $v \in R^m$ . The exponential function  $r(x) = e^x$  is a frequent choice for  $r$  in practice. Usually in survival analysis applications

one actually has available only censored observations from (22): if  $C$  is a censoring time with conditional distribution  $K(\cdot|z)$  which is conditionally independent of  $Y$ , we observe  $X \equiv (Z, T, \Delta)$  with  $T \equiv Y \wedge C$  and  $\Delta \equiv I_{\{Y \leq C\}}$ . Thus,  $\theta = (\nu, G, K, H)$ , and a full analysis of this model involves consideration of  $K$  and  $H$  as nuisance parameters in addition to  $G$ . In fact it can be shown that lack of knowledge of both  $K$  and  $H$  does not make the estimation of  $\nu$  more difficult asymptotically. Thus, we will consider  $K$  and  $H$  as known, and consider the model to be

$$(23) \quad \mathbf{P} = \{P: P \text{ has conditional hazard function given by (22)} \\ \text{for some } \nu \text{ and some } \lambda \}.$$

Cox (1972) introduced the model in essentially this generality, and proposed estimation of  $\nu$  by maximization of what has come to be called Cox's partial likelihood.  $\square$

Censoring has entered into examples 5 and 7 and can be considered as a complicating factor in almost any model. The type of censoring mechanism considered in these two examples, namely arbitrary right censorship in which the minimum of two variables and an indicator variable are observed, can also be generalized. We will not be able to give a full treatment here, but many real problems seem to involve censoring of some kind, and the methods presented in the following chapters should prove useful in dealing with many of these problems.



# 2 | Asymptotic Inference for (Finite-Dimensional) Parametric Models

## 2.1 REGULAR PARAMETRIC MODELS IN THE I.I.D. CASE

We shall review some of the basic results in asymptotic inference, particularly estimation, for regular parametric models. The statements and conditions are essentially those of Le Cam (1956), (1969), (1970), and Hájek (1970), (1972), but the basic heuristic goes back to Fisher (1922), (1925).

Let  $\mu$  be a fixed  $\sigma$ -finite measure on  $(X, \mathcal{B})$ , and let  $M_\mu$  be all probability measures  $P$  on  $(X, \mathcal{B})$  dominated by  $\mu$ ; i.e.,  $M_\mu \equiv \{P \in \mathcal{M} : P \ll \mu\}$ . Then, as usual, suppose that  $X_1, \dots, X_n$  are i.i.d. with common distribution  $P \in \mathcal{P}$ , where  $\mathcal{P}$  is dominated by  $\mu$ . Recall that we loosely defined  $\mathcal{P}$  to be parametric or finite-dimensional if we could write

$$\mathcal{P} = \{P_\theta : \theta \in \Theta\},$$

where

- (i)  $\Theta$  is a “nice” subset of  $R^k$ .
- (ii) The parametrization  $\theta \rightarrow P_\theta$  is “smooth.”

Let

$$(1) \quad p(\theta) = p(\cdot, \theta) = \frac{dP_\theta}{d\mu}(\cdot), \quad l(\theta) = \log p(\theta),$$

be the *density* and *log-likelihood* of  $P_\theta$  respectively.

**Convention.** If  $h(\theta)$  is a function on  $X$  for fixed  $\theta$ , then  $h(x, \theta)$  denotes its value at  $x$ .

The facts according to Fisher (for  $k = 1$ ) are:

- (iii) If some estimate  $T_n$  satisfies  $L_\theta(\sqrt{n}(T_n - \theta)) \rightarrow N(0, \sigma^2(\theta))$  as  $n \rightarrow \infty$  for all  $\theta$ , then

$$\sigma^2(\theta) \geq I^{-1}(\theta) \quad \text{where } I(\theta) = E_\theta \left( \frac{\partial}{\partial \theta} l(\theta) \right)^2.$$

- (iv) If  $\theta$  is identifiable, the maximum likelihood estimate  $\hat{\theta}$  solving

$$\sum_{i=1}^n \frac{\partial}{\partial \theta} l(X_i, \hat{\theta}) = 0$$

is efficient. That is,

$$L_{\theta}(\sqrt{n}(\hat{\theta} - \theta)) \rightarrow N(0, I^{-1}(\theta)).$$

Of course, these “facts” are not quite right. Some uniformity in the convergence in law in (iii) has to be required to avoid “superefficiency,” and some modification of maximum likelihood is often needed to obtain efficiency. Note that the inequality in (iii) can be viewed as an asymptotic version of the well-known Cramér-Rao (information) inequality for unbiased estimation of  $\theta$ . Fisher (1925, section 7, pages 710, 711) gave a heuristic argument for (iii) and (iv) apparently based on multinomial distributions. A rephrasing of that argument appears in Fisher (1956, section VI.3, pages 147–150); see Savage (1976) for further discussion.

To obtain the elegant statement of correct versions of (iii) and (iv) that we give in sections 3 and 5, it is convenient to view  $\mathbf{P}$  as a subset of  $L_2(\mu)$  via the embedding  $P \rightarrow s$ , where

$$p \equiv \frac{dP}{d\mu}, \quad s \equiv \sqrt{p}, \quad p(\theta) \rightarrow s(\theta).$$

Both this embedding and that of  $\mathbf{P}$  into  $L_1(\mu)$  via  $P \rightarrow p$  endow  $\mathbf{P}$  with the same topology, convergence in total variation; see (A.6.3). One advantage of this embedding is that we can replace awkward conditions on pointwise differentiability of  $l$  and integrability assumptions by the natural condition of Fréchet (Hellinger) differentiability of the map  $\theta \rightarrow s(\theta)$ . More significantly it enables us (in section 2.4) to give a geometric formulation of the theory, which then extends fairly readily to the semiparametric case. In the following, elements of  $R^k$  are written as column vectors,  $|\cdot|$  is the Euclidean norm, and  $\|\cdot\|$  is the Hilbert norm in  $L_2(\mu)$ :

$$\|f\|^2 \equiv \int f^2 d\mu.$$

**Definition 1.**  $\theta_0$  is a *regular point* of the parametrization  $\theta \rightarrow P_{\theta}$  if  $\theta_0$  is an interior point of  $\Theta$ , and

- (i) The map  $\theta \rightarrow s(\theta)$  from  $\Theta$  to  $L_2(\mu)$  is Fréchet differentiable at  $\theta_0$ : there exists a vector  $\dot{s}(\theta_0) = (\dot{s}_1(\theta_0), \dots, \dot{s}_k(\theta_0))^T$  of elements of  $L_2(\mu)$  such that
- (2)  $\|s(\theta_0 + h) - s(\theta_0) - \dot{s}^T(\theta_0)h\| = o(|h|)$  as  $h \rightarrow 0$ .

- (ii) The  $k \times k$  matrix  $\int \dot{s}(\theta_0)\dot{s}^T(\theta_0) d\mu$  is nonsingular.

This is exactly as in example A.5.2; see section A.5 for more on Fréchet derivatives.

**Definition 2.** A parametrization  $\theta \rightarrow P_{\theta}$  is *regular* if:

- (i) Every point of  $\Theta$  is regular.
- (ii) The map  $\theta \rightarrow \dot{s}_i(\theta)$  is continuous from  $\Theta$  to  $L_2(\mu)$  for  $i = 1, \dots, k$ .

Note that (i) implies that  $\Theta$  is open. We also note that by proposition A.5.1 Fréchet differentiability of  $\theta \rightarrow s(\theta)$  implies continuity of this map, and hence, by (A.6.3), the continuity of the more familiar map  $\theta \rightarrow p(\theta)$  from  $\Theta$  to  $L_1(\mu)$ . Define the *Fisher information matrix* of  $\theta$  by

$$(3) \quad I(\theta) = 4 \int \dot{s}(\theta) \dot{s}^T(\theta) d\mu .$$

We call  $\mathbf{P}$  a *regular parametric model* if it has a regular parametrization. In such models the “niceness” of  $\Theta$  is evident and the “smoothness” of the parametrization is made precise. They are the objects we shall study.

Define the *score function*  $\dot{\mathbf{i}}$  of an observation by

$$(4) \quad \dot{\mathbf{i}}(\theta) = 2 \frac{\dot{s}(\theta)}{s(\theta)} 1_{[s(\theta) > 0]} = \frac{\dot{p}(\theta)}{p(\theta)} 1_{[p(\theta) > 0]} ,$$

where

$$(5) \quad \dot{p}(\theta) = 2s(\theta)\dot{s}(\theta) .$$

If  $\theta$  is a regular point, then  $|\dot{\mathbf{i}}(\theta)| \in L_2(P_\theta)$ , and the more usual definition of the Fisher information matrix for  $\theta$  is given by

$$(6) \quad I(\theta) = \int \dot{\mathbf{i}}(\theta) \dot{\mathbf{i}}^T(\theta) dP_\theta .$$

By proposition A.5.3.F the two definitions of the Fisher information matrix given in (3) and (6) agree.

The following proposition gives sufficient conditions for regularity of a parametric model in terms of ordinary differentiability of the likelihood.

**Proposition 1.** Suppose  $\Theta$  is open and for all  $\theta$ :

- (i)  $p(x, \theta)$  is continuously differentiable in  $\theta$  for  $(\mu)$  almost all  $x$  with gradient  $\dot{p}(\theta)$ .
- (ii)  $|\dot{\mathbf{i}}(\theta)| \in L_2(P_\theta)$  with  $\dot{\mathbf{i}}(\theta)$  as in (4).
- (iii)  $I(\theta)$  defined in (6) is nonsingular and continuous in  $\theta$ .

Then, if we define

$$(7) \quad \begin{aligned} \dot{s}(\theta) &= \frac{1}{2} p^{-1/2}(\theta) \dot{p}(\theta) 1_{[p(\theta) > 0]} \\ &= \frac{1}{2} s(\theta) \dot{\mathbf{i}}(\theta) 1_{[p(\theta) > 0]} , \end{aligned}$$

the parametrization  $\theta \rightarrow P_\theta$  is regular with  $\dot{s}(\theta)$  from (7) as Fréchet derivative of  $s(\theta)$ .

**Proof.** Note that  $\dot{p}(\theta)$  vanishes  $(\mu)$  almost everywhere outside  $A(\theta) \equiv [p(\theta) > 0]$  because of (i), and that, hence, the definition of  $\dot{p}(\theta)$  in (5) is consistent with (7). It can be verified by (i) that for  $(\mu)$  almost all  $x$ ,

$$(a) \quad s(x, \theta + h) - s(x, \theta) = \int_0^1 \frac{1}{2} p^{-1/2}(x, \theta + \lambda h) h^T \dot{p}(x, \theta + \lambda h) d\lambda ,$$

and hence

$$(b) \quad \begin{aligned} & (s(x, \theta + h) - s(x, \theta)) 1_{A(\theta)}(x) \\ &= \frac{1}{2} p^{-1/2}(x, \theta) h^T \dot{p}(x, \theta) 1_{A(\theta)}(x) + o(|h|) \end{aligned}$$

holds. By (a), (ii), and (iii) it follows that

$$(c) \quad \begin{aligned} & \int |s(x, \theta + h) - s(x, \theta)|^2 d\mu(x) \\ & \leq \frac{1}{4} \int_0^1 h^T \int \dot{p}(x, \theta + \lambda h) \dot{p}^T(x, \theta + \lambda h) p^{-1}(x, \theta + \lambda h) d\mu(x) h d\lambda \\ &= \frac{1}{4} h^T I(\theta) h + o(|h|^2) \\ &= \frac{1}{4} h^T \int_{A(\theta)} \dot{p}(\theta) \dot{p}^T(\theta) p^{-1}(\theta) d\mu(x) h + o(|h|^2). \end{aligned}$$

Assume without loss of generality that  $h/|h|$  converges. Using (b) and (c) and applying lemma A.7.5 to  $|h|^{-1}(s(\theta + h) - s(\theta)) 1_{A(\theta)}$ , we obtain

$$(d) \quad \begin{aligned} & \int_{A(\theta)} |s(x, \theta + h) - s(x, \theta) - \frac{1}{2} s^{-1}(x, \theta) h^T \dot{p}(x, \theta)|^2 d\mu(x) \\ &= o(|h|^2), \end{aligned}$$

which combined with (c) yields

$$(e) \quad \int_{X-A(\theta)} |s(x, \theta + h) - s(x, \theta)|^2 d\mu(x) = o(|h|^2).$$

This implies the Fréchet differentiability of  $s$  at  $\theta$  with derivative (7). This yields, in view of (iii), the nonsingularity of  $I(\theta)$ . Hence every  $\theta$  is regular. From (i) and (iii)

$$\lim_{h \rightarrow 0} \dot{s}_i(x, \theta + h) = \dot{s}_i(x, \theta), \quad (\mu) \text{ a.e. } x \in A(\theta),$$

and

$$\begin{aligned} & \limsup_{h \rightarrow 0} \int_{A(\theta)} \dot{s}_i^2(x, \theta + h) d\mu(x) \\ & \leq \limsup_{h \rightarrow 0} \int \dot{s}_i^2(x, \theta + h) d\mu(x) \\ &= \int \dot{s}_i^2(x, \theta) d\mu(x) = \int_{A(\theta)} \dot{s}_i^2(x, \theta) d\mu(x) \end{aligned}$$

follow. Continuity of  $\theta \rightarrow \dot{s}(\theta)$  is obtained from this by the same argument as before.  $\square$

### Example 1. Exponential family.

Suppose  $\Theta$  is open and  $\{P_\theta\}$  is a curved exponential family,  $p(x, \theta) = \exp(c(\theta)T(x) - d(\theta))$ . If  $c$  has a differential and  $c(\Theta)$  is contained in the interior of the natural parameter space of the exponential family, then equation (2) applies.  $\square$

**Example 2. Translation model with known shape  $f$ .**

Suppose  $P_\theta$  is a one-dimensional translation parameter family  $p(x, \theta) = f(x - \theta)$ ,  $\theta$  real. Then hypothesis (i) of proposition 1 is equivalent to continuity of the derivative  $f'$ , ruling out the double exponential for instance. However, regularity of the model holds if and only if  $f$  is absolutely continuous with Radon-Nikodym derivative  $f'$  and  $\int [((f')^2/f)(x)] dx < \infty$  (Hájek and Šidák (1967, page 211), or corollary A.5.1), and hence the double exponential is regular. This example shows that only the requirement of *continuous differentiability* in (i) of proposition 1 is too strong.  $\square$

**Example 3. Weibull translation model.**

Suppose that  $\theta = (\alpha, \beta, \gamma) \in \Theta \equiv \{\theta : \alpha > 0, \beta > 0, \gamma \in R\}$  and

$$p(x, \theta) = \frac{\beta}{\alpha} \left(\frac{x - \gamma}{\alpha}\right)^{\beta-1} \exp\left(-\left(\frac{x - \gamma}{\alpha}\right)^\beta\right) 1_{(\gamma, \infty)}(x),$$

the Weibull translation model. This model is *not* regular, but the restricted model  $\mathbf{P} = \{P_\theta : \theta \in \Theta_0\}$ , where  $\Theta_0 \equiv \{\theta : \alpha > 0, \beta > 2, \gamma \in R\} \subset \Theta$ , is regular.  $\square$

**Example 4. Three-parameter lognormal model.**

Suppose that  $Y \sim N(\mu, \sigma^2)$  and  $X \equiv \gamma + \exp(Y)$ . This model, with  $\theta = (\mu, \sigma^2, \gamma) \in \Theta = R \times R^+ \times R$ ,

$$p(x, \theta) = \frac{1}{\sigma(x - \gamma)} \phi\left(\frac{\log(x - \gamma) - \mu}{\sigma}\right) 1_{(\gamma, \infty)}(x)$$

where  $\phi$  is the standard normal density function, and  $\mathbf{P} = \{P_\theta : \theta \in \Theta\}$ , is a lognormal translation family. It is a regular model, but yields unbounded likelihood functions, and hence the method of maximum likelihood fails; see, e.g., Hill (1963), Cohen and Whitten (1980), and Griffiths (1980) for solutions. In spite of this, the methods to be discussed in section 2.5 based on scores yield efficient estimates.  $\square$

The first equality in (4) suggests another useful local embedding of  $\mathbf{P}$  into  $L_2(P_{\theta_0})$  (for  $\theta_0$  fixed) given by

$$P_\theta \longleftrightarrow r(\theta) \equiv 2 \left(\frac{s(\theta)}{s(\theta_0)} - 1\right) 1_{[s(\theta) > 0]}.$$

Regularity at  $\theta_0$  is equivalent to (see proposition A.5.3.E and F):

- (ia)  $\theta \rightarrow r(\theta)$  is Fréchet differentiable at  $\theta_0$  (in  $L_2(P_{\theta_0})$ ) with derivative  $\dot{\mathbf{l}}(\theta_0)$ .
- (ib)  $P_{\theta_0+h}(s(\theta_0) = 0) = o(|h|^2)$ .
- (ii)  $I(\theta_0) = \int \dot{\mathbf{l}}(\theta_0) \dot{\mathbf{l}}^T(\theta_0) dP_{\theta_0}$  is nonsingular.

The function  $r(\theta)$ , which belongs to  $L_2(P_{\theta_0})$ , is a useful proxy for  $\mathbf{l}(\theta)$ , which may not belong to  $L_2(P_{\theta_0})$ .

Regularity of  $\theta$  is enough to guarantee a score function identity which is basic to the Cramér-Rao information bound calculation.

$$(8) \quad \int \dot{\mathbf{l}}(\theta) dP_\theta = 0$$

or, equivalently, by (4) and (5),

$$(9) \quad \langle \dot{s}_i(\theta), s(\theta) \rangle = 0 \quad \text{for } i = 1, \dots, k,$$

where  $\langle f, g \rangle \equiv \int fg \, d\mu$  is the inner product in  $L_2(\mu)$ .

**Proof of (9).** By hypothesis

$$(a) \quad \langle s(\theta), s(\theta) \rangle = 1 \quad \text{for all } \theta.$$

Fréchet-differentiating with respect to  $\theta_i$  yields  $\langle \dot{s}_i(\theta), s(\theta) \rangle = 0$ . □

Further,

$$(10) \quad \dot{s}(\theta) = \dot{s}(\theta) 1_{[s(\theta) > 0]} \quad \text{a.e. } \mu.$$

This follows since

$$\begin{aligned} & \int (s(\theta + h) - s(\theta) - \dot{s}^T(\theta)h 1_{[s(\theta) > 0]})^2 \, d\mu \\ &= \int (s(\theta + h) - s(\theta) - \dot{s}^T(\theta)h)^2 1_{[s(\theta) > 0]} \, d\mu \\ &+ \int s^2(\theta + h) 1_{[s(\theta) = 0]} \, d\mu \\ &= o(|h|^2) \end{aligned}$$

by (2) and (ib).

The fundamental consequence of regularity at a point  $\theta$  is the *local asymptotic normality* (LAN) of the model given by the following basic proposition. Define the log-likelihood of  $(X_1, \dots, X_n)$  by

$$L_n(\theta) = \sum_{i=1}^n \mathbf{l}(X_i, \theta)$$

and the score function by

$$(11) \quad S_n(\theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \dot{\mathbf{l}}(X_i, \theta).$$

**Proposition 2.** Suppose that  $\mathbf{P} = \{P_\theta : \theta \in \Theta\}$  is a regular parametric model, and write

$$(12) \quad L_n(\theta + \frac{t}{\sqrt{n}}) - L_n(\theta) = t^T S_n(\theta) - \frac{1}{2} t^T I(\theta) t + R_n(\theta, t).$$

Then  $R_n(\theta, t) \rightarrow 0$  in  $P_\theta$  probability uniformly for  $\theta \in K$  compact  $\subset \Theta$  and  $|t| \leq M$ ; i.e., for any compact set  $K \subset \Theta$ ,  $0 < M < \infty$ , and  $\varepsilon > 0$ ,

$$(13) \quad \sup_{|t| \leq M} \sup_{\theta \in K} P_\theta(|R_n(\theta, t)| > \varepsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Moreover,

$$(14) \quad L_\theta(S_n(\theta)) \rightarrow N(0, I(\theta))$$

uniformly in  $\theta \in K$  for compact  $K \subset \Theta$ , where  $N(\mu, \Sigma)$  is the multivariate normal distribution with mean  $\mu$  and covariance matrix  $\Sigma$ . Finally, uniformly in

$\theta \in K, |t| \leq M,$

$$(15) \quad S_n(\theta + \frac{t}{\sqrt{n}}) - S_n(\theta) + I(\theta)t \rightarrow_{P_\theta} 0 \quad \text{as } n \rightarrow \infty.$$

**Remark.** Note that the left side of (12) is a well-defined extended real-valued random variable under  $P_\theta$ . For a careful definition of the uniform convergence in (14) see definition 2.2.3.

**Proof.** The proof of this uniform version of the LAN property is given in appendix 9 along with other contiguity theory facts and proofs. The proofs of (13) and (14) are very similar to the proof of Le Cam's second lemma, and are based on the proof given by Ibragimov and Has'minskii (1981, theorem II.1.2, page 119); equation (15) has been noted in (6.43) of Bickel (1982). Under stronger conditions this proposition follows easily by Taylor expansion of  $L_n$  to two terms. See, e.g., Lehmann (1983, proof of theorem 6.2.3, page 415).  $\square$

An important property that follows from this basic proposition is contiguity of the product measures corresponding to  $\theta$  and  $\theta + t/\sqrt{n}$  at regular  $\theta$ . We recall the definition and basic properties of contiguity; also see section A.9.

**Definition 3.** Two sequences of probability measures  $\{P_n\}, \{Q_n\}$ , each pair defined on the same space,  $(X_n, \mathcal{B}_n)$  are called *contiguous*, and we write  $\{P_n\} \triangleleft \{Q_n\}$  if  $P_n(A_n) \rightarrow 0$  if and only if  $Q_n(A_n) \rightarrow 0$ , for  $A_n \in \mathcal{B}_n$ .

**Proposition 3.** If  $\theta \rightarrow P_\theta$  is a regular parametrization and  $\theta_n \rightarrow \theta$ , then  $\{P_{\theta_n + t_n/\sqrt{n}}^n\}$  and  $\{P_{\theta_n}^n\}$  are contiguous for any bounded sequence  $\{t_n\}$ .

**Proof.**  $\{P_{\theta_n + t_n/\sqrt{n}}^n\} \triangleleft \{P_{\theta_n}^n\}$  follows from (12), (14), and corollary A.9.1 of Le Cam's first lemma. To prove  $\{P_{\theta_n}^n\} \triangleleft \{P_{\theta_n + t_n/\sqrt{n}}^n\}$ , use (12) and (14) together with Le Cam's lemma A.9.3 and corollary A.9.1.  $\square$

Proposition 2 indicates that whatever be  $M$  finite, the log-likelihood of the local model  $\{P_{\theta+t/\sqrt{n}} : |t| \leq M\}$  is approximated in a weak sense, for  $n$  large, by that of the model  $\{Q_t : |t| \leq M\}$ , where we observe

$$S_n(\theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n i(X_i, \theta)$$

distributed under  $Q_t$  as  $N(I(\theta)t, I(\theta))$  as if  $\theta$  is known but the local deviation  $t$  is not. This weak type of approximation has been shown to have profound consequences by Le Cam in a series of papers starting in 1956, following the lead of Wald (1943). We explore some of these consequences concentrating on point estimation.

## 2.2 REGULAR ESTIMATES OF EUCLIDEAN PARAMETERS

Let  $v : \mathbf{P} \rightarrow R^m$  be a Euclidean parameter, where  $\mathbf{P}$  is a general (not necessarily parametric) model. An estimate of  $v$  is any (measurable)  $m$ -vector  $T_n(X_1, \dots, X_n)$  which depends only on  $(X_1, \dots, X_n)$ . We study the limiting behavior of sequences  $\{T_n(X_1, \dots, X_n)\}_{n \geq 1}$  whose members are related in some natural way in the expectation that the limiting behavior will be an approximation to that for fixed  $n$  and  $P$ . We call such sequences estimates also.

The first property a reasonable sequence of estimates should have is *consistency*:

$$P(|T_n - v(P)| \geq \varepsilon) \rightarrow 0 \quad \text{for all } P \in \mathbf{P} \text{ and for all } \varepsilon > 0.$$

In this work we focus on how close estimates can get to  $v$  on the  $n^{-1/2}$  scale. So we need to define  $\sqrt{n}$ -consistency:

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} P(n^{1/2} |T_n - v(P)| \geq M) = 0,$$

or, in the familiar  $O_p$ ,  $o_p$  notation,

$$T_n - v(P) = O_p(n^{-1/2}).$$

It is desirable to have properties such as consistency,  $\sqrt{n}$ -consistency, hold as uniformly in  $\mathbf{P}$  as possible. Otherwise “ $n$  large” depends on the unknown  $P$ , and we cannot even in principle specify what an adequate sample size is. Unfortunately, uniformity in consistency and other properties is often unachievable on all of  $\mathbf{P}$  and has to be replaced by uniformity on “small” subsets of  $\mathbf{P}$ .

In particular, we metrize  $\mathbf{P}$  with the variational distance (A.6.1) and require uniformity on compact subsets of  $\mathbf{P}$ .

**Definition 1.**  $T = \{T_n\}$  is *uniformly consistent* if, for every  $\varepsilon > 0$ ,

$$\sup\{P(|T_n - v(P)| \geq \varepsilon) : P \in \mathbf{K}\} \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

for all compact subsets  $\mathbf{K}$  of  $\mathbf{P}$ .

**Definition 2.**  $T = \{T_n\}$  is *uniformly  $\sqrt{n}$ -consistent* if

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} \sup\{P(n^{1/2} |T_n - v(P)| \geq M) : P \in \mathbf{K}\} = 0$$

for all compact subsets  $\mathbf{K}$  of  $\mathbf{P}$ .

Here are some desirable properties which are stronger than  $\sqrt{n}$ -consistency, but which are often possessed by procedures used in practice.

**Convention.** Write  $E_P$  and  $L_P$  for expectations and distributions of measurable functions of  $(X_1, \dots, X_n)$  under  $P$ . In parametric models shorten  $E_{P_\theta}$  and  $L_{P_\theta}$  to  $E_\theta$ , and  $L_\theta$ . Let  $\mathbf{K}$  denote an arbitrary (pre-) compact subset of  $\mathbf{P}$ .

**Definition 3.**  $T = \{T_n\}$  is a *uniformly regular* estimate of  $v(P) \in R^m$  if  $\{L_P(\sqrt{n}(T_n - v(P)))\} \equiv \{L_P(Z_n)\}$  converge uniformly on compact subsets  $\mathbf{K}$  of  $\mathbf{P}$ . That is,

$$\sup_{P \in \mathbf{K}} |E_P g(Z_n) - E_P g(Z)| \rightarrow 0$$

for all  $g$  continuous and bounded on  $R^m$ , all  $\mathbf{K}$ , and a fixed family  $\Lambda = \{L_P(Z)\} = \{L_P\}$  of probability distributions on  $R^m$ .

We usually ask for more.

**Definition 4.**  $T$  is a *uniformly Gaussian regular* estimate of  $v(P)$  if it is uniformly regular and for each  $P$  the limit  $L_P = L_P(Z)$  of  $L_P(Z_n) = L_P(\sqrt{n}(T_n - v(P)))$  is  $m$ -variate Gaussian with mean 0 and covariance matrix which we denote  $\Sigma(P, T)$ ; i.e.,  $Z \sim N(0, \Sigma(P, T))$ .



A further structural property which aids our understanding of how  $T$  acts and connects this theory with that of robust estimation is given by the next definition.

**Definition 5.**  $T$  is an asymptotically linear estimate of  $\nu$  if there exists

$$\psi : \mathbf{X} \times \mathbf{P} \rightarrow R^m$$

such that for all  $P \in \mathbf{P}$

$$(1) \quad \begin{aligned} &|\psi(\cdot, P)| \in L_2(P), \\ &\int \psi(x, P) dP = 0, \\ &T_n = \nu(P) + n^{-1} \sum_{i=1}^n \psi(X_i, P) + o_p(n^{-1/2}). \end{aligned}$$

Then

$$(2) \quad \Sigma(P, T) = \int \psi \psi^T(x, P) dP.$$

We call the function  $\psi(\cdot, P)$  (which is unique a.e.  $P$ ) the *influence function* of  $T$ . The observation  $X_i = x$ , to first order, contributes  $n^{-1}\psi(x, P)$  to the error  $T_n - \nu(P)$ , in essential agreement with the notions of Hampel (1974). If (1) holds for just  $P = P_0$ , then we call  $T$  asymptotically linear at  $P_0$ .

Suppose that  $\nu$  is a parameter defined on all of  $\mathbf{M} = \{ \text{all probability distributions on } \mathbf{X} \}$  satisfying the following regularity conditions:

- (i) For all  $P_0 \in \mathbf{M}$ ,  $\nu$  is continuously Fréchet differentiable at  $P_0$  with respect to  $d_K$  given by (A.6.8).
- (ii) For all  $P_0 \in \mathbf{P}$ , the derivative  $\dot{\nu}$  has the representation

$$\dot{\nu}(P_0)(P) = \int \psi(x, P_0) dP(x),$$

where  $\psi$  is continuous, bounded in  $x$ , continuous in  $P_0$  with respect to the total variation metric  $d_\nu$ , and

$$\int \psi(x, P_0) dP_0(x) = 0.$$

Then, if  $P_n$  is the empirical distribution, the estimator  $\nu(P_n)$  of  $\nu(P)$  is asymptotically linear with influence function  $\psi$ . Moreover,  $\psi$  is the Gâteaux derivative of  $\nu$  on  $\mathbf{M}$  at  $P$ :

$$(3) \quad \begin{aligned} \lim_{\varepsilon \rightarrow 0} \frac{\nu((1-\varepsilon)P + \varepsilon Q) - \nu(P)}{\varepsilon} &= \frac{\partial}{\partial \varepsilon} \nu((1-\varepsilon)P + \varepsilon Q) \Big|_{\varepsilon=0} \\ &= \int \psi(x, P) dQ. \end{aligned}$$

By taking  $Q$  to be point mass at  $x$  we recapture  $\psi(x, P)$ , making our definition agree with that of Hampel (1974). Proofs of these and stronger results may be found for instance in Huber (1981, section 2.5).

**Definition 6.**  $T$  is a *uniformly asymptotically linear* estimate of  $v(P)$  if (1) holds uniformly on compact subsets of  $\mathbf{P}$ ; i.e. the  $o_p(n^{-1/2})$  in (1) holds uniformly in  $P \in \mathbf{K}$ .

Existence of estimates of  $v$  possessing these properties places strong restrictions on  $v$  and the limit laws of the estimates, as we see in the following proposition. Metrize the set  $\Lambda$  of probability measures on  $R^m$  by the Prohorov metric  $d_{pr}$  given in (A.6.10) so that the topology is that of convergence in law, but retain the variational distance topology on  $\mathbf{P}$ .

**Proposition 1.** Suppose that  $T$  is an estimate of  $v$ .

- A. If  $T$  is uniformly consistent (on compact subsets of  $\mathbf{P}$ ), then the map  $P \rightarrow v(P)$  is continuous from  $\mathbf{P}$  to  $R^m$ .
- B. If  $T$  is uniformly regular with limit laws  $\{L_P\} \subset \Lambda$ , then the map  $P \rightarrow L_P$  is continuous from  $\mathbf{P}$  to  $\Lambda$ .
- C. If  $T$  is uniformly Gaussian regular, then the map  $P \rightarrow \Sigma(P, T)$  from  $\mathbf{P}$  to the set of nonnegative definite symmetric  $m \times m$  matrices is continuous.
- D. If  $T$  is uniformly asymptotically linear, and  $P \rightarrow \Sigma(P, T)$  (given by (2)) is continuous, then  $T$  is uniformly Gaussian regular with covariance matrix  $\Sigma(P, T)$ .

**Proof.** Without loss of generality, take  $\mathbf{P}$  to be compact.

A. Define maps  $g_n, g$  from  $\mathbf{P}$  to  $\Lambda$  by

$$g_n(P) = L_P(T_n),$$

$$g(P) = \delta_{v(P)} \equiv \text{point mass at } v(P).$$

Claim A is equivalent to continuity of  $g$ . The maps  $g_n$  are continuous for each  $n$ , since by (A.6.3) and (A.6.5)  $d_v(P_m, P) \rightarrow 0$  implies  $d_v(P_m^n, P^n) \rightarrow 0$  for the product measures  $P_m^n, P^n$ . Therefore A follows from

$$(a) \quad \sup\{d_{pr}(g_n(P), g(P)) : P \in \mathbf{P}\} \rightarrow 0.$$

But by a theorem of Strassen (1965), if  $P(|T_n - v(P)| \geq \varepsilon) \leq \varepsilon$ , then  $d_{pr}(g_n(P), g(P)) \leq \varepsilon$ ; see Theorem A.6.1. Since

$$\sup\{P(|T_n - v(P)| \geq \varepsilon) : P \in \mathbf{P}\} \rightarrow 0,$$

(a) and claim A follow.

B. Define

$$h_n(P) = L_P(\sqrt{n}(T_n - v(P))),$$

$$h(P) = L_P.$$

Since uniform regularity implies uniform consistency (on compacts),  $g$  is continuous and hence so is  $h_n$ . Since regularity corresponds to uniform convergence of  $h_n$  to  $h$ , claim B also follows.

C. This follows from B since continuity of  $h$  is equivalent to continuity of  $P \rightarrow \Sigma(P, T)$ .

D. Since the uniform asymptotic linearity of  $T$  and the boundedness of the eigenvalues of  $\Sigma(P, T)$  imply uniform consistency of  $T$ ,  $g$  is continuous. We will verify the continuity of  $P \rightarrow L_P(\psi(X_1, P))$  on  $\mathbf{P}$  with the total variation distance. Let  $\{P_m\}$  and  $P$  satisfy  $d_v(P_m, P) \rightarrow 0$ . It suffices to show that there exists a subsequence  $\{P_n\}$ , say, with  $L_{P_n}(\psi(X_1, P_n)) \rightarrow L_P(\psi(X_1, P))$ . Choose  $\{P_n\}$  such that  $nd_v(P_n, P) \rightarrow 0$ . Then contiguity and uniform asymptotic linearity imply that

$$n^{-1/2} \sum_{i=1}^n (\psi(X_i, P_n) + v(P_n) - \psi(X_i, P) - v(P)) = o_p(1),$$

and by lemma A.7.6 this yields

$$L_{P_n}(\psi(X_1, P_n) + v(P_n)) \rightarrow L_P(\psi(X_1, P) + v(P)).$$

Since  $d_v(P_n, P) \rightarrow 0$ , and  $g$  is continuous,  $v(P_n) \rightarrow v(P)$ , and hence  $L_{P_n}(\psi(X_1, P_n)) \rightarrow L_P(\psi(X_1, P))$ . The continuity of  $P \rightarrow L_P(\psi(X_1, P))$  follows. This continuity and that of  $P \rightarrow \Sigma(P, T)$  imply that the maps

$$P \rightarrow L_P \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi(X_i, P) \right)$$

converge uniformly to  $P \rightarrow N(0, \Sigma(P, T))$ . This can be verified by checking the Lindeberg condition; see (A.7.5). The result follows.  $\square$

The notions of uniform regularity and uniform Gaussian regularity we have introduced are suitable for parametric models but typically are somewhat too strong for the natural nonparametric and semiparametric models we are interested in. The following notion of local regularity at a point of a parametric model  $\mathbf{P} = \{P_\theta : \theta \in \Theta\}$  is a useful stepping stone to a suitable definition of regular estimates in nonparametric models. It is due to Hájek (1970), (1972) and implicit in the work of Le Cam.

**Definition 7.**  $T = \{T_n\}$  is *locally regular* at  $P_{\theta_0}$  if, whenever  $\sqrt{n}|\theta_n - \theta_0|$  stays bounded,

$$L_{\theta_n}(\sqrt{n}(T_n - v(P_{\theta_n}))) \rightarrow L_0,$$

where  $L_0$  does not depend on  $\{\theta_n\}$ . If  $L_0$  is Gaussian, we say that  $T_n$  is *locally Gaussian regular*, while *local asymptotic linearity* corresponds to uniformity of the expansion (1) for  $\theta$  depending on  $n$  with  $\sqrt{n}|\theta - \theta_0|$  bounded.

Necessary and sufficient conditions for an asymptotically linear estimator to be locally regular will be given in section 2.4.

**Convention.** We will frequently abbreviate *locally regular* to just *regular*, in keeping with most of the recent literature.

Of course, there are estimators which are not uniformly regular or even locally regular. The most well known of these is the ‘‘superefficient’’ estimator of Hodges; see, e.g., Lehmann (1983, pages 405, 407–408). Here is another example.

**Example 1. Stein's estimator of a normal mean.**

Suppose that  $X_1, \dots, X_n$  are i.i.d.  $N_k(\theta, I)$  in  $R^k$ ,  $k \geq 3$ . Consider the estimator

$$\hat{\theta}_n \equiv \left( 1 - \frac{k-2}{n|\bar{X}|^2} \right) \bar{X}$$

of  $\theta \in R^k$ . If  $\theta_0 = 0$  and  $\theta_n \equiv tn^{-1/2}$  with  $t \in R^k$ , then, since  $L_{\theta_n}(\sqrt{n}(\bar{X} - \theta_n)) = L_0(X_1) = N_k(0, I)$  under  $P_{\theta_n}$ ,

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta_n) &= \sqrt{n}(\bar{X} - \theta_n) \\ &\quad - \frac{k-2}{|\sqrt{n}(\bar{X} - \theta_n) + t|^2} \{ \sqrt{n}(\bar{X} - \theta_n) + t \} \\ &= {}_d X_1 - \frac{k-2}{|X_1 + t|^2} (X_1 + t) \end{aligned}$$

for every  $n \geq 1$ , the distribution of which is dependent on  $t$ . Hence  $\hat{\theta}_n$  is not locally regular at  $\theta_0 = 0$ .  $\square$

There also exist estimators which are (uniformly or locally) regular, but are *not* (uniformly or locally) Gaussian regular. Here are two examples.

**Example 2. Minimum Kolmogorov distance estimator of center of symmetry.**

Suppose that  $X_1, \dots, X_n$  are i.i.d. with distribution function  $F_\theta(x) = F_0(x - \theta)$  where  $F_0$  is continuous and symmetric about 0:  $1 - F_0(x) = F_0(-x)$  for all  $x$ . Thus  $X_1 - \theta, \dots, X_n - \theta$  are i.i.d.  $F_0$  and hence both

$$\frac{1}{n} \sum_{i=1}^n 1_{[X_i - \theta \leq x]} = F_n(x + \theta) \rightarrow_{\text{a.s.}} F_0(x)$$

and

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n 1_{[\theta - X_i \leq x]} &= 1 - F_n(\theta - x) \\ &\rightarrow_{\text{a.s.}} 1 - F_0(-x) = F_0(x), \end{aligned}$$

where

$$F_n(x) \equiv n^{-1} \sum_{i=1}^n 1_{[X_i \leq x]}$$

is the empirical df of the  $X$ 's. This suggests estimation of  $\theta$  as the value of  $t$  which minimizes some measure of distance between the empirical df of  $X_1 - t, \dots, X_n - t$  and the empirical df of  $t - X_1, \dots, t - X_n$ . The minimum Kolmogorov-Smirnov distance estimator  $\hat{\theta}_{KS}$  of  $\theta$  is any value of  $t$  which minimizes

$$D_n(t) \equiv \sup_x |F_n(x+t) - (1 - F_n(t-x-))|.$$

It was shown by Rao, Schuster, and Littell (1975) that

$$L_\theta(\sqrt{n}(\hat{\theta}_{KS} - \theta)) \rightarrow L(Z),$$

where  $L(Z)$  is *not* Gaussian, but is expressible in terms of a Brownian bridge process. Since the estimator is translation equivariant, the convergence is trivially uniform in  $\theta$ . If the Kolmogorov-Smirnov (supremum) distance is replaced by an  $L_2$ -distance, then the resulting estimator is uniformly Gaussian regular; see, e.g., Shorack and Wellner (1986, page 759, exercise 22.5.3).  $\square$

**Example 3. Bickel-Hodges estimate of center of symmetry.**

Suppose that  $X_1, \dots, X_n$  are i.i.d. with symmetric distribution as in example 2, and let  $\hat{\theta}_n \equiv \text{med} \{ \frac{1}{2}(X_{(i)} + X_{(n-i+1)}) : 1 \leq i \leq n \}$ . Then Bickel and Hodges (1967) show that

$$L_\theta(\sqrt{n}(\hat{\theta}_n - \theta)) \rightarrow L(Z),$$

where  $L(Z)$  is again not Gaussian, but is expressible in terms of a Brownian motion process. If the estimator is instead the Hodges-Lehmann estimator  $\tilde{\theta}_n = \text{med} \{ \frac{1}{2}(X_{(i)} + X_{(j)}) : 1 \leq i, j \leq n \}$ , then  $\tilde{\theta}_n$  is a uniformly Gaussian regular estimator of  $\theta$ ; see section 2.5.  $\square$

## 2.3 THE INFORMATION BOUND AND THE HÁJEK-LE CAM CONVOLUTION AND ASYMPTOTIC MINIMAX THEOREMS

Suppose  $v$  is a Euclidean parameter defined on a regular parametric model  $\mathbf{P} = \{P_\theta : \theta \in \Theta\}$ . We can identify  $v$  with the parametric function  $q: \Theta \rightarrow R^m$  defined by

$$q(\theta) = v(P_\theta).$$

Fix  $P = P_\theta$  and suppose  $q$  has a total differential matrix  $\dot{q}_{m \times k}$  at  $\theta$ . Define

(1)  $I^{-1}(P | v, P) = \dot{q}(\theta)I^{-1}(\theta)\dot{q}^T(\theta)$ , the *information bound* for  $v$ , and

(2)  $\tilde{I}(\cdot, P | v, P) = \dot{q}(\theta)I^{-1}(\theta)\dot{I}(\theta)$ , the *efficient influence function* for  $v$ .

As defined in (1) and (2), the information bound and influence function appear to depend on the parametrization  $\theta \rightarrow P_\theta$  of  $\mathbf{P}$ . However, as our notation indicates:

**Proposition 1.**  $I^{-1}(P | v, P)$  and  $\tilde{I}(\cdot, P | v, P)$  are invariant under smooth changes of parametrization.

Here is a formal calculation.

Suppose  $\gamma \rightarrow \theta(\gamma)$  is a one-to-one continuously differentiable mapping of an open subset  $\Gamma$  of  $R^k$  onto  $\Theta$  with nonsingular differential  $\dot{\theta}$ . We represent  $\mathbf{P} = \{Q_\gamma : \gamma \in \Gamma\}$ , where

$$Q_\gamma \equiv P_{\theta(\gamma)}.$$

Identify  $v$  by

$$v(\gamma) \equiv v(Q_\gamma) \equiv q(\theta(\gamma)).$$

Then, by the chain rule, the Fisher information matrix for  $\gamma$  is

$$\dot{\theta}(\gamma)I(\theta(\gamma))\dot{\theta}^T(\gamma),$$

while

$$\dot{v}(\gamma) = \dot{q}(\theta(\gamma))\dot{\theta}^T(\gamma).$$

Substituting back into (1) gives the same answer for  $\gamma \rightarrow Q_\gamma$  as for  $\theta \rightarrow P_\theta$ . A similar calculation works for  $\tilde{I}$ .  $\square$

**Theorem 1.** (Convolution Theorem) Suppose  $T$  is a uniformly regular estimate of  $v$  with corresponding limit law  $L_\theta$ . Then:

- A.  $L_\theta$  is representable as the convolution of a  $N(0, I^{-1}(P_\theta | v, P))$  distribution with that of another  $m$ -vector  $\Delta_\theta$ ;  $L_\theta = L(Z_\theta + \Delta_\theta)$ , where  $Z_\theta \sim N(0, I^{-1}(P_\theta | v, P))$  and  $\Delta_\theta$  are independent. More generally,

$$(3) \quad L_\theta \left( \begin{array}{c} \sqrt{n}(T_n - q(\theta)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i, P_\theta | v, P) \\ \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i, P_\theta | v, P) \end{array} \right) \rightarrow L_\theta \left( \begin{array}{c} \Delta_\theta \\ Z_\theta \end{array} \right).$$

- B. If the map  $\theta \rightarrow \dot{q}(\theta)$  is continuous, then (3) holds uniformly on compact subsets  $K$  of  $\Theta$ .
- C. Moreover, if  $\theta \rightarrow \dot{q}(\theta)$  is continuous, then  $L_\theta = N(0, I^{-1}(P_\theta | v, P))$  for all  $\theta$  if and only if  $T$  is uniformly asymptotically linear with efficient influence function  $\tilde{I} = \tilde{I}(\cdot, P_\theta | v, P)$ .

Since it is intuitively clear that an estimator  $T = \{T_n\}$  with limit law  $L_\theta = L(Z_\theta) = N(0, I^{-1}(P_\theta | v, P))$  and  $\Delta_\theta \equiv 0$  is "less spread out" than an estimator with nondegenerate  $\Delta_\theta$ , and since C implies that all uniformly Gaussian regular estimates  $T$  with  $\Sigma(P_\theta, T) = I^{-1}(P_\theta | v, P)$  are asymptotically equivalent, it is reasonable to use the result of theorem 1 to define efficiency.

**Definition 1.** If  $T = \{T_n\}$  is a uniformly Gaussian regular estimator of  $v(P_\theta) = q(\theta)$  with  $\Sigma(P_\theta, T) = I^{-1}(P_\theta | v, P)$ , we say that  $T$  is *uniformly efficient*. If this holds with "uniformly" replaced by "locally," then we say that  $T$  is *locally efficient* or just *efficient*.

**Proof of theorem 1.** Fix  $\theta$ . Recall that  $S_n(\theta) = n^{-1/2} \sum_{i=1}^n \dot{I}(X_i, \theta)$ . Regularity of  $\{T_n\}$  and regularity of the model  $P$  imply that

$$(U_n, V_n) \equiv (\sqrt{n}(T_n - q(\theta)), S_n(\theta))$$

are marginally convergent in law, hence marginally tight, and hence jointly tight, since marginal tightness implies joint tightness.

Fix a subsequence  $\{n'\}$ . By Prohorov's theorem A.7.5, there exists a further

subsequence  $\{n''\}$  such that  $L_\theta(U_{n''}, V_{n''})$  has a joint limit  $L(U, V)$  where  $V \sim N(0, I(\theta))$ . For convenience of notation, we now denote the subsequence  $\{n''\}$  by  $\{n\}$ . Let

$$W_n = L_n(\theta + \frac{t}{\sqrt{n}}) - L_n(\theta).$$

Then  $L_\theta(U_n, W_n)$  has limit  $L(U, t^T V - \frac{1}{2} t^T I(\theta) t)$  by the LAN property, proposition 2.1.2. Next, since the map  $\theta \rightarrow P_\theta$  is continuous at regular points  $\theta$ , the regularity of  $\{T_n\}$  implies by proposition 2.2.1.B that, for all  $t$ ,

$$L_{\theta+t/\sqrt{n}}(\sqrt{n}(T_n - q(\theta + \frac{t}{\sqrt{n}}))) \rightarrow L_\theta = L(U).$$

Hence, by differentiability of  $q$ ,

$$(a) \quad L_{\theta+t/\sqrt{n}}(U_n) \rightarrow L(U + \dot{q}(\theta)t).$$

By (a)

$$(b) \quad E_{\theta+t/\sqrt{n}}(\exp[ia^T U_n]) \rightarrow \exp[ia^T \dot{q}(\theta)t] E \exp[ia^T U].$$

On the other hand, by contiguity,

$$(c) \quad E_{\theta+t/\sqrt{n}}(\exp[ia^T U_n]) = E_\theta(\exp[ia^T U_n + W_n]) + o(1).$$

Finally, note that  $\{|\exp[ia^T U_n + W_n]|\} = \{\exp(W_n)\}$  is a uniformly integrable sequence of variables by proposition 2.1.2 and lemma A.7.2. Combining (b) and (c) we arrive, again by lemma A.7.2, at the identity, valid for all  $a \in R^m, t \in R^k$ ,

$$(d) \quad E \exp[ia^T U + t^T V - \frac{1}{2} t^T I(\theta) t] = \exp[ia^T \dot{q}(\theta)t] E \exp[ia^T U].$$

Both sides of (d), for fixed  $a$ , are functions of  $t$  analytic on  $C^k$ , where  $C$  is the complex plane. Therefore by analytic continuation, identity (d) holds for  $t^T = -i(a^T - b^T)\dot{q}(\theta)I^{-1}(\theta)$  as well. Then (d) becomes

$$(e) \quad E \exp[ia^T(U - \dot{q}(\theta)I^{-1}(\theta)V) + ib^T \dot{q}(\theta)I^{-1}(\theta)V] \\ = E \exp[ia^T U + \frac{1}{2} a^T \dot{q}(\theta)I^{-1}(\theta) \dot{q}^T(\theta) a] \\ \cdot \exp[-\frac{1}{2} b^T \dot{q}(\theta)I^{-1}(\theta) \dot{q}^T(\theta) b].$$

Note that formula (e) is the characteristic function of the limit law of  $(U_{n''} - \dot{q}(\theta)I^{-1}(\theta)V_{n''}, \dot{q}(\theta)I^{-1}(\theta)V_{n''})$  and is the same for every choice of initial subsequence  $\{n'\}$ . Consequently the full sequence, which is the left side of (3), has limit law with characteristic function (e). For  $b = 0$ , (e) reduces to

$$(f) \quad E \exp[ia^T(U - \dot{q}(\theta)I^{-1}(\theta)V)] \\ = E \exp[ia^T U + \frac{1}{2} a^T \dot{q}(\theta)I^{-1}(\theta) \dot{q}^T(\theta) a].$$

Together, (e) and (f) yield

$$(g) \quad E \exp[ia^T(U - \dot{q}(\theta)I^{-1}(\theta)V) + ib^T\dot{q}(\theta)I^{-1}(\theta)V] \\ = E \exp[ia^T(U - \dot{q}(\theta)I^{-1}(\theta)V)] \exp[-\frac{1}{2}b^T\dot{q}(\theta)I^{-1}(\theta)\dot{q}^T(\theta)b],$$

which is a product of a function of  $a$  with a function of  $b$ . The function of  $b$  is the characteristic function of a  $N(0, I^{-1}(P_\theta | v, P))$  distribution. From (g) with  $b = a$ , the first part of A follows. But (g) also yields (3).

To prove B and C, we copy the proof leading to (g), replacing  $\theta$  by a sequence  $\{\theta_n\}$  tending to  $\theta$ , e.g.,

$$(a') \quad L_{\theta_n + t/\sqrt{n}}(U_n) \rightarrow L(U + \dot{q}(\theta)t).$$

Note that for (a') the continuous differentiability of  $q$  is used. In this way, (g) holds, and hence (3), with  $\theta$  on the left side replaced by  $\theta_n$ , is proved.

If  $L_\theta = N(0, I^{-1}(P_\theta | v, P))$ , then part B of the theorem shows that  $\Delta_\theta = 0$  a.s. and hence that  $T$  is uniformly asymptotically linear with efficient influence function  $\tilde{I}$ . Conversely, if  $T$  is uniformly asymptotically linear with influence function  $\tilde{I}$ , then proposition 2.2.1.D implies uniform Gaussian regularity, since  $\theta \rightarrow I^{-1}(P_\theta | v, P)$  is continuous by (1), the continuity of  $\theta \rightarrow \dot{q}(\theta)$ , and the regularity of the model.  $\square$

**Information inequality.** If  $T$  is uniformly Gaussian regular, then

$$(4) \quad \Sigma(P_\theta, T) \geq I^{-1}(P_\theta | v, P)$$

in the order on nonnegative definite matrices. Equality holds if and only if  $T$  is uniformly efficient.

**Proof.** By the convolution theorem

$$\Sigma(P_\theta, T) = I^{-1}(P_\theta | v, P) + E(\Delta_\theta - E\Delta_\theta)(\Delta_\theta - E\Delta_\theta)^T.$$

If equality holds,  $\Delta_\theta$  is constant. Since the asymptotic mean of  $\sqrt{n}(T_n - v)$  is zero, we must have  $\Delta_\theta = 0$ .  $\square$

Note that, if  $q$  is the identity,

$$I^{-1}(P_\theta | v, P) = I^{-1}(\theta),$$

and (4) is the usual (asymptotic) information bound for  $\theta$ .

**Definition 2.** A function  $l: R^m \rightarrow R^+$  is called *bowl-shaped* if

$$l(x) = l(-x), \text{ and } \{x : l(x) \leq c\} \text{ is convex for every } c \geq 0.$$

Bowl-shaped functions  $l$  generate loss functions and risks via  $E_\theta l(\sqrt{n}(T_n - q(\theta)))$ .

**Asymptotic optimality theorem.** If  $T$  is uniformly regular and  $l$  is bowl-shaped, then

$$(5) \quad \liminf_{n \rightarrow \infty} E_\theta l(\sqrt{n}(T_n - q(\theta))) \geq El(Z_\theta),$$

where  $Z_\theta \sim N(0, I^{-1}(P_\theta | v, P))$ .

**Example 1. Quadratic loss.**

$$l(x) \equiv x^2.$$

The risk of  $T_n$  is  $n$  times its mean square error.  $\square$



**Example 2. Zero-one loss.**

$l(x) \equiv 1 - 1_C(x)$ , where  $C$  is a bounded symmetric convex set. The risk of  $T_n$  is the probability that the confidence region  $T_n + C/\sqrt{n}$  does not cover  $v(P_\theta) = q(\theta)$ .  $\square$

**Proof of the asymptotic optimality theorem.** Define

$$l_k(x) = 2^{-k} \sum_{i=1}^{k2^k} (1 - 1_{C_{ik}}(x)),$$

with  $C_{ik} \equiv \{y : l(y) \leq i2^{-k}\}$ . Note that  $l_k \uparrow l$  as  $k \rightarrow \infty$  and that  $l_k$  is bowl-shaped and continuous Lebesgue a.e. since the boundary of a convex set is a Lebesgue null set. Since  $T$  is uniformly regular, by the convolution theorem its limit law  $L_\theta$  is absolutely continuous and can be represented as that of  $Z_\theta + \Delta_\theta$  where  $\Delta_\theta$  is independent of  $Z_\theta$ . The boundedness and a.e. continuity of  $l_k$  then yield

$$\liminf_{n \rightarrow \infty} E_\theta l(\sqrt{n}(T_n - q(\theta))) \geq E l_k(Z_\theta + \Delta_\theta).$$

Apply Anderson's theorem (Anderson (1955)) and the independence of  $Z_\theta$  and  $\Delta_\theta$  to conclude

$$E l_k(Z_\theta + \Delta_\theta) \geq E l_k(Z_\theta).$$

The theorem follows by monotone convergence.  $\square$

**Note.** The Hájek–Le Cam convolution theorem, the information inequality, and the asymptotic optimality theorem hold if uniform (Gaussian) regularity and uniform efficiency are replaced by their local versions.

An extension of the asymptotic optimality theorem which applies to *all* (rather than just regular) estimates  $T$  is the

**Local asymptotic minimax theorem.** If  $\{T_n\}$  is any sequence of estimates, then

$$(6) \quad \lim_{M \rightarrow \infty} \liminf_{n \rightarrow \infty} \sup \{ E_{\theta'} l(\sqrt{n}(T_n - q(\theta'))) : \sqrt{n}|\theta' - \theta| \leq M \} \\ \geq E l(Z_\theta).$$

**Proof.** See Ibragimov and Has'minskii (1981, remark II.12.2).  $\square$

In this book we restrict ourselves to regular estimates and shall not refer to this result further.

## 2.4 NUISANCE PARAMETERS, ADAPTATION, AND SOME GEOMETRY

As we saw in chapter 1, parameters  $v$  are often defined implicitly through a parametrization  $\theta \rightarrow P_\theta$ , where  $\theta^T = (v^T, \eta^T)$ ,  $v \in N \subset R^m$ ,  $\eta \in H \subset R^{k-m}$ , where  $v$  is the parameter of interest and  $\eta$  is a *nuisance parameter*. If  $\theta_0 = (v_0, \eta_0) \in \Theta$ , let  $P_1(\eta_0) \equiv \{P_\theta : \eta = \eta_0, v \in N\}$ . This is the model when  $\eta = \eta_0$  is known. We want to assess the cost of not knowing  $\eta$  by com-

paring the information bounds and efficient influence functions for  $\nu$  at  $P_{\theta_0}$  in  $\mathbf{P}_1(\eta_0)$  and  $\mathbf{P}$ .

As before, we let  $\langle \cdot, \cdot \rangle_0$  be the inner product in  $L_2(P_{\theta_0})$ ,  $\|\cdot\|_0$  the norm, and write  $E_0$  for expectation under  $P_{\theta_0}$ .

Suppose the model is regular, and write  $\dot{\mathbf{i}}$  for the score function at  $\theta_0$  and  $\tilde{\mathbf{I}} = I^{-1}(\theta_0)\dot{\mathbf{i}}$  for the efficient influence function of the parameter  $\theta$  at  $P_{\theta_0}$  in  $\mathbf{P}$ . Decompose

$$\dot{\mathbf{i}} = \begin{pmatrix} \dot{\mathbf{i}}_1 \\ \dot{\mathbf{i}}_2 \end{pmatrix}, \quad \tilde{\mathbf{I}} = \begin{pmatrix} \tilde{\mathbf{I}}_1 \\ \tilde{\mathbf{I}}_2 \end{pmatrix},$$

with  $\tilde{\mathbf{I}}_1$ ,  $\dot{\mathbf{i}}_1$   $m$ -vectors,  $\tilde{\mathbf{I}}_2$ ,  $\dot{\mathbf{i}}_2$   $(k-m)$ -vectors. Write  $I(\theta_0)$  in block matrix form, suppressing dependence on  $\theta_0$ , as

$$I = [I_{ij}]_{i,j=1,2} = \begin{pmatrix} I_{11} & I_{12} \\ I_{21} & I_{22} \end{pmatrix},$$

with  $I_{11}$   $m \times m$ ,  $I_{12}$   $m \times (k-m)$ ,  $I_{21}$   $(k-m) \times m$ ,  $I_{22}$   $(k-m) \times (k-m)$ , and similarly decompose  $I^{-1}(\theta_0)$  into  $I^{ij}$ ,  $i, j = 1, 2$ . By well-known block matrix forms of matrix inverses we have

$$(1) \quad I^{-1}(\theta_0) = [I^{ij}]_{i,j=1,2} = \begin{pmatrix} I_{11}^{-1} & -I_{11}^{-1}I_{12}I_{22}^{-1} \\ -I_{22}^{-1}I_{21}I_{11}^{-1} & I_{22}^{-1} \end{pmatrix},$$

where

$$(2) \quad \begin{aligned} I_{11 \cdot 2} &\equiv I_{11} - I_{12}I_{22}^{-1}I_{21}, \\ I_{22 \cdot 1} &\equiv I_{22} - I_{21}I_{11}^{-1}I_{12}. \end{aligned}$$

By (2.3.1) and (2.3.2), the information bound for estimating  $\nu$  in  $\mathbf{P}$  is  $I^{11} = I_{11 \cdot 2}^{-1}$  and the efficient influence function for  $\nu$  in  $\mathbf{P}$  is

$$(3) \quad \begin{aligned} \tilde{\mathbf{I}}_1 &= I^{11}\dot{\mathbf{i}}_1 + I^{12}\dot{\mathbf{i}}_2 \\ &= I_{11 \cdot 2}^{-1}(\dot{\mathbf{i}}_1 - I_{12}I_{22}^{-1}\dot{\mathbf{i}}_2) \quad \text{by (1)} \\ &\equiv I_{11 \cdot 2}^{-1}\mathbf{I}_1^*. \end{aligned}$$

Since

$$\begin{aligned} I_{11 \cdot 2} &= E_0(\dot{\mathbf{i}}_1 - I_{12}I_{22}^{-1}\dot{\mathbf{i}}_2)(\dot{\mathbf{i}}_1 - I_{12}I_{22}^{-1}\dot{\mathbf{i}}_2)^T \\ &= E\mathbf{I}_1^*\mathbf{I}_1^{*T}, \end{aligned}$$

we see that (3) has the same form as  $\tilde{\mathbf{I}} = I^{-1}(\theta_0)\dot{\mathbf{i}}$  with  $\tilde{\mathbf{I}}$  replaced by  $\tilde{\mathbf{I}}_1$ ,  $I(\theta_0) = E_0(\dot{\mathbf{i}}\dot{\mathbf{i}}^T)$  replaced by  $I_{11 \cdot 2} = E_0(\mathbf{I}_1^*\mathbf{I}_1^{*T})$ , and  $\dot{\mathbf{i}}$  replaced by

$$(4) \quad \mathbf{I}_1^* \equiv \dot{\mathbf{i}}_1 - I_{12}I_{22}^{-1}\dot{\mathbf{i}}_2.$$

We therefore call  $\mathbf{I}_1^*$  the *efficient score function* for  $\nu$  in  $\mathbf{P}$ , and call  $I_{11 \cdot 2}$  the *information* for  $\nu$  in  $\mathbf{P}$ .

If, on the other hand,  $\eta = \eta_0$  is treated as known, the information bound (for  $v$  in  $P_1(\eta_0)$ ) is  $I_{11}^{-1}$  and the corresponding efficient influence curve (for  $v$  in  $P_1(\eta_0)$ ) is just

$$(5) \quad I_{11}^{-1} \dot{I}_1 .$$

From the block matrix formulas relating  $[I_{ij}]$  and  $[I^{ij}]$ , we can derive some important relations between these quantities. First note from (1) and (2) that

$$(6) \quad (I^{11})^{-1} = I_{11 \cdot 2} = I_{11} - I_{12} I_{22}^{-1} I_{21} ,$$

so not knowing  $\eta$  decreases the information for  $v$  by  $I_{12} I_{22}^{-1} I_{21}$ . Similarly,

$$I_{11}^{-1} = I^{11} - I^{12} (I^{22})^{-1} I^{21}$$

or

$$(7) \quad I^{11} = I_{11 \cdot 2}^{-1} = I_{11}^{-1} + I^{12} (I^{22})^{-1} I^{21} ,$$

so not knowing  $\eta$  increases the information bound (inverse information) by  $I^{12} (I^{22})^{-1} I^{21}$ . Moreover, from (6),

$$(8) \quad I_{11 \cdot 2} = I_{11} \quad \text{and} \quad I_{11 \cdot 2}^{-1} = I_{11}^{-1}$$

if and only if

$$(9) \quad I_{12} = 0 .$$

In this case it also follows from (3), (4), and (8) that

$$(10) \quad \tilde{I}_1 = I_{11}^{-1} \dot{I}_1 \quad \text{and} \quad I_1^* = \dot{I}_1 .$$

**Definition 1.**  $\{\hat{v}_n\}$  is an *adaptive estimate* of  $v$  in the presence of  $\eta$  if  $\hat{v}_n$  is regular on  $P$  and efficient for each of the models  $P_1(\eta)$ , for all  $\eta$ .

If an adaptive estimate exists, we can do as well not knowing  $\eta$  as knowing it. By (9) and (10), a necessary condition for the existence of adaptive estimates in regular parametric models is

$$(11) \quad I_{12}(\theta) = 0 \quad \text{for all } \theta .$$

Cox and Reid (1987) discuss reparametrization of  $\eta$  to achieve (11); see also Kass (1989, section 2.1.4, and pages 201, 202). In any case, adaptation is very much a feature of the parametrization, as the following examples show.

**Example 1. Normal location-scale.**

Suppose that  $P_\theta \equiv N(v, \eta)$ ,  $v \in R$ ,  $\eta > 0$ . Thus  $I_{12}(\theta) = 0$  for all  $\theta$ . In this, the normal location-scale model, we can estimate the mean equally well whether or not we know the variance.  $\square$

**Example 2. Reparametrization of normal location-scale.**

Now suppose that  $P_\theta = N(v, \eta - v^2)$ ,  $\eta > v^2$ . Then easy calculation shows that

$$I_{12}(\theta) = - \frac{v}{(\eta - v^2)^2}$$

by the classical formula

$$I(\theta) = - \left[ E_{\theta} \left( \frac{\partial^2 \mathbf{l}(\theta)}{\partial \theta_i \partial \theta_j} \right) \right]. \quad \square$$

We can think of  $I_1^*$  as the  $\dot{I}_1$  corresponding to the reparametrization  $(v, \eta) \rightarrow (v, \eta + I_{22}^{-1}(\theta_0)I_{21}(\theta_0)(v - v_0))$ . With this reparametrization, adaptation at  $\theta_0$  becomes possible since  $\dot{I}_2$  is unchanged and condition (3) is satisfied. If we can paste together these local reparametrizations and find  $(v, \eta) \rightarrow (v, \gamma(v, \eta))$  such that

$$\begin{aligned} \gamma(v, \eta) - \gamma(v_0, \eta_0) &= \eta - \eta_0 + I_{22}^{-1}(\theta_0)I_{21}(\theta_0)(v - v_0) \\ &\quad + o(|\theta - \theta_0|) \end{aligned}$$

for every  $\theta_0 = (v_0, \eta_0)$ , then under this reparametrization the necessary condition for adaptation holds. For instance in example 2 we can take  $\gamma(v, \eta) = \eta - v^2$ . These remarks have little practical significance since the initial parametrization is usually natural and the reparametrization is not.

### Some Geometry

The efficient influence function  $\tilde{I}_1$  and efficient score function  $I_1^*$  can be interpreted geometrically in the Hilbert space  $L_2(P_{\theta_0})$ ; see sections A.1 and A.2 for elementary Hilbert space theory. First suppose  $m = 1$ . Let  $[\dot{I}_2]$  be the linear span of the components of  $\dot{I}_2$  in  $L_2(P_{\theta_0})$ . Then by example A.2.1,  $I_{12}I_{22}^{-1}\dot{I}_2$  is the projection of  $\dot{I}_1$  on  $[\dot{I}_2]$ , and by (4) the efficient score function  $I_1^*$  is the projection of  $\dot{I}_1$  on the orthocomplement of  $[\dot{I}_2]$ .

We can also relate the efficient influence functions  $\tilde{I}_1$  and  $I_{11}^{-1}\dot{I}_1$  for  $v$  in  $\mathbf{P}$  and  $\mathbf{P}_1(\eta_0)$ . In particular,  $I_{11}^{-1}\dot{I}_1$  is the projection of  $\tilde{I}_1$  on  $[\dot{I}_1]$ . We need only check that  $\tilde{I}_1 - I_{11}^{-1}\dot{I}_1 = (I^{11} - I_{11}^{-1})\dot{I}_1 + I^{12}\dot{I}_2$  is orthogonal to  $\dot{I}_1$ , and this follows from  $I^{11}I_{11} + I^{12}I_{21} = 1$ .

If  $m > 1$  these relationships continue to hold if projection is interpreted componentwise. The following basic proposition can be viewed as providing the rationale for two different approaches to computing information bounds in semi-parametric models which will be presented in sections 3 and 4 of chapter 3.

#### Proposition 1.

- A. The efficient score function  $I_1^*(\cdot, P_{\theta_0} \mid v, \mathbf{P})$  is the projection of the score function  $\dot{I}_1$  on the orthocomplement of  $[\dot{I}_2]$  in  $L_2(P_{\theta_0})$ .
- B. The efficient influence function  $\tilde{I}(\cdot, P_{\theta_0} \mid v, \mathbf{P}_1(\eta_0))$  is the projection of the efficient influence function  $\tilde{I}_1(\cdot, P_{\theta_0} \mid v, \mathbf{P})$  on  $[\dot{I}_1]$  in  $L_2(P_{\theta_0})$ .

See figures 1 and 3.

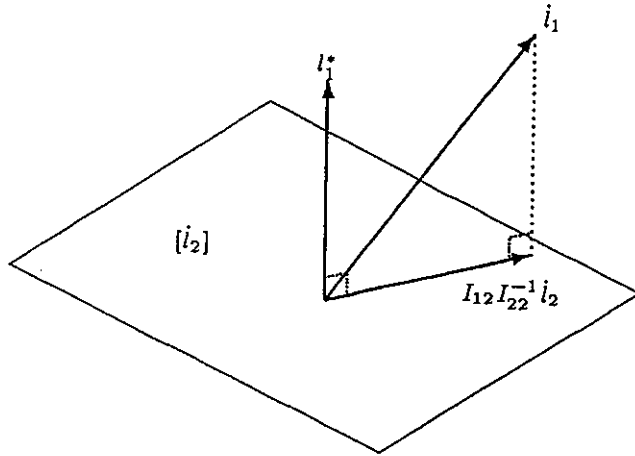


FIGURE 1. Projection of score functions.

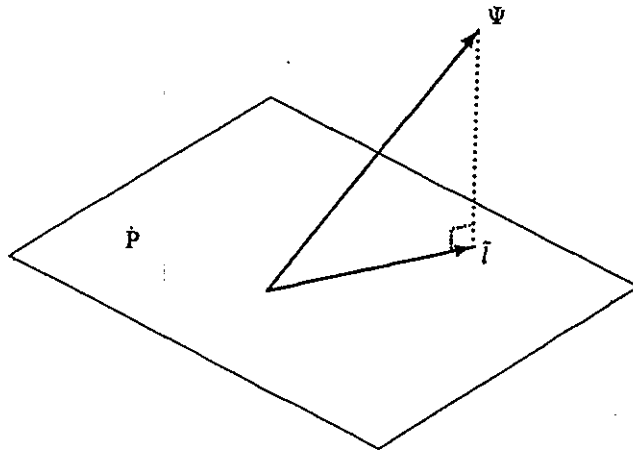


FIGURE 2. Projection of influence functions.

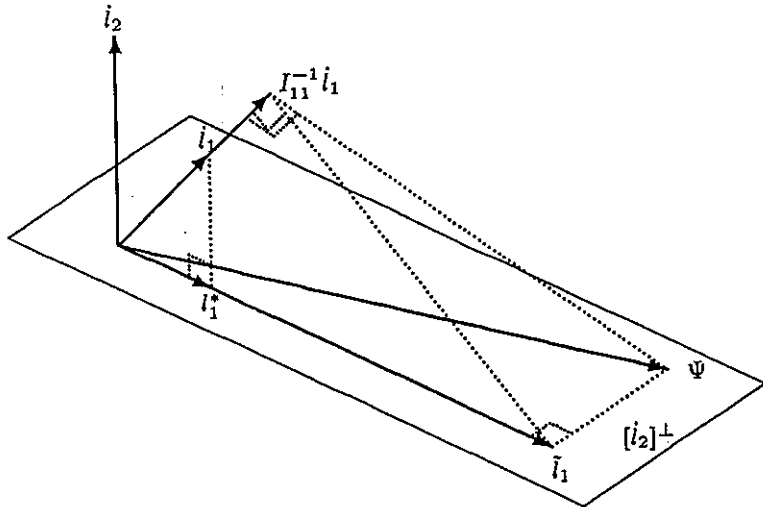


FIGURE 3. Score and influence function projections.

Here is another relationship between the influence and score functions of  $P_1(\eta_0)$  and its companion  $P_2(v_0) \equiv \{P_{(v_0, \eta)} : \eta \in H\}$ . We use the subscript 2 for score and influence function in the companion model. The efficient influence function  $\tilde{I}_1$  can be written as

$$(12) \quad \tilde{I}_1 = I_{11}^{-1} \dot{I}_1 - I_{11}^{-1} I_{12} \tilde{I}_2 .$$

This relationship was implicit in section 4 of Begun, Hall, Huang, and Wellner (1983); it will appear again in the context of semiparametric models (with  $v$  infinite-dimensional and  $\eta$  finite-dimensional) in section 5.4. Note that (12) provides an immediate proof, by orthogonality of  $\tilde{I}_2$  to  $[\dot{I}_1]$ , of the formula

$$(13) \quad I_{11,2}^{-1} = I_{11}^{-1} + I_{11}^{-1} I_{12} I_{22,1}^{-1} I_{21} I_{11}^{-1} ,$$

which is another way of writing (7).

**Proof of (12).** From (1),

$$\begin{aligned} \tilde{I}_1 + I_{11}^{-1} I_{12} \tilde{I}_2 &= I^{11} \dot{I}_1 + I^{12} \dot{I}_2 + I_{11}^{-1} I_{12} (I^{21} \dot{I}_1 + I^{22} \dot{I}_2) \\ &= I_{11}^{-1} \left\{ (I_{11} I^{11} + I_{12} I^{21}) \dot{I}_1 + (I_{11} I^{12} + I_{12} I^{22}) \dot{I}_2 \right\} \\ &= I_{11}^{-1} \dot{I}_1. \quad \square \end{aligned}$$

Table 1 summarizes the efficient score functions, efficient influence functions, information, and inverse information for the two models  $P$  and  $P_1(\eta_0)$ .

Name	Notation	Model	
		$P$	$P_1(\eta_0)$
Efficient score	$I_1^*(\cdot, P   v, \cdot)$	$I_1^* = \dot{I}_1 - I_{12} I_{22}^{-1} \dot{I}_2$	$\dot{I}_1$
Information	$I(P   v, \cdot)$	$E I_1^* I_1^{*T} = I_{11} - I_{12} I_{22}^{-1} I_{21}$ $\equiv I_{11,2}$	$I_{11}$
Efficient influence function	$\tilde{I}_1(\cdot, P   v, \cdot)$	$\tilde{I}_1 = I^{11} \dot{I}_1 + I^{12} \dot{I}_2$ $= I_{11,2}^{-1} I_1^*$ $= I_{11}^{-1} \dot{I}_1 - I_{11}^{-1} I_{12} \tilde{I}_2$	$I_{11}^{-1} \dot{I}_1$
Information bound	$I^{-1}(P   v, \cdot)$	$I^{11} = I_{11,2}^{-1}$ $= I_{11}^{-1} + I_{11}^{-1} I_{12} I_{22,1}^{-1} I_{21} I_{11}^{-1}$	$I_{11}^{-1}$

Table 1

We now use several examples to illustrate the relationships between score functions, efficient score functions, and efficient influence functions given in proposition 1 and table 1. In our first two examples, the efficient score functions are intuitively plausible.

**Example 3. The bivariate normal distribution.**

Suppose that  $X = (X_1, X_2) \sim N_2(\theta, \Sigma)$ , where  $\theta = (v, \eta) \in R^2$ , and  $\Sigma$  is the covariance matrix with ones on the diagonal and correlation  $\rho$ ; we will suppose here that  $\rho$  is known. The joint density is

$$p(x, \theta) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)}\{z_1^2 - 2\rho z_1 z_2 + z_2^2\}\right\},$$

where  $z_1 \equiv x_1 - \nu$ ,  $z_2 \equiv x_2 - \eta$ , and the model is  $\mathbf{P} = \{P_\theta : \theta \in R^2\}$ . Hence, the scores for  $\nu = \theta_1$  and  $\eta = \theta_2$  are

$$\begin{aligned} \dot{\mathbf{i}}_1(x) &= \frac{1}{1-\rho^2}\{(x_1 - \nu) - \rho(x_2 - \eta)\}, \\ \dot{\mathbf{i}}_2(x) &= \frac{1}{1-\rho^2}\{(x_2 - \eta) - \rho(x_1 - \nu)\}. \end{aligned}$$

It is straightforward to calculate

$$I(\theta) = \frac{1}{1-\rho^2} \begin{pmatrix} 1 & -\rho \\ -\rho & 1 \end{pmatrix}.$$

Hence the efficient score for  $\nu = \theta_1$  in the presence of the nuisance parameter  $\eta = \theta_2$  is

$$\mathbf{I}_1^*(x) = x_1 - \nu,$$

and the information for  $\nu$  is  $I(P_\theta | \nu, \mathbf{P}) = 1$ . Of course the information bound is achieved by the sample mean,  $\hat{\nu} = \bar{X}_1 \equiv n^{-1} \sum_{i=1}^n X_1^{(i)}$ ; here  $X^{(1)}, \dots, X^{(n)}$  denote the observations which are i.i.d. as  $X$ . Note that indeed  $\mathbf{I}_1^* \perp \dot{\mathbf{I}}_2$ .

Now consider the submodel  $\mathbf{P}_1 = \mathbf{P}_1(\eta_0)$  in which  $\eta$  is known; without loss of generality suppose that  $\eta_0 = 0$ . For this submodel  $\mathbf{P}_1$ , the score function  $\dot{\mathbf{i}}_1$  for  $\nu$  is also the efficient score function  $\mathbf{I}_1^* \equiv \mathbf{I}_1^*(\cdot, P_\theta | \nu, \mathbf{P}_1)$ , so that

$$I(P_\theta | \nu, \mathbf{P}_1) = \frac{1}{1-\rho^2}, \quad I^{-1}(P_\theta | \nu, \mathbf{P}_1) = 1-\rho^2.$$

The efficient influence function for  $\nu$  in the submodel  $\mathbf{P}_1$  is

$$\tilde{\mathbf{I}}_1(x, P_\theta | \nu, \mathbf{P}_1) = x_1 - \nu - \rho x_2.$$

Thus knowledge of  $\eta$  reduces the information bound for estimation of  $\nu$  from 1 to  $1-\rho^2$ . An estimator achieving this bound (assuming known covariance matrix  $\Sigma$ ) is  $\hat{\nu}^0 = n^{-1} \sum_{i=1}^n (X_1^{(i)} - \rho X_2^{(i)})$ .  $\square$

**Example 4. The multinomial distribution.**

Suppose that  $X = (X_1, \dots, X_{k+1}) \sim \text{Mult}_{k+1}(1, (p_1, \dots, p_{k+1}))$ , and let  $\theta = (p_1, \dots, p_k)$  so that  $p_{k+1} = 1 - \sum_{i=1}^k p_i = 1 - \sum_{i=1}^k \theta_i$ . The density function is

$$p(x, \theta) = \left\{ \prod_{i=1}^k \theta_i^{x_i} \right\} (1 - \sum_{i=1}^k \theta_i)^{x_{k+1}}$$

for  $x_j \in \{0, 1\}$ ,  $j=1, \dots, k+1$ . Hence the scores for  $\theta$  are easily calculated to be

$$\dot{\mathbf{i}}_j(x) = \frac{x_j}{\theta_j} - \frac{x_{k+1}}{p_{k+1}}, \quad j=1, \dots, k,$$

and the information matrix for  $\theta$  is

$$I(P_\theta | \theta, \mathbf{P}) \equiv I(\theta) = \text{diag}\left(\frac{1}{\theta}\right) + \left(\frac{1}{p_{k+1}}\right) \underline{1} \underline{1}^T,$$

where  $\underline{1}$  is a  $k$ -vector of 1's. Therefore the information bound is given by

$$I^{-1}(\theta) = \text{diag}(\theta) - \theta\theta^T.$$

Of course, the usual estimator  $\hat{\theta}$  given by the first  $k$  coordinates of  $\hat{p} = n^{-1} \sum_{i=1}^n X^{(i)} = \bar{X}$  (where  $X^{(1)}, \dots, X^{(n)}$  are i.i.d. as  $X$ ) achieves this bound.

Now consider estimation of  $v = (\theta_1, \dots, \theta_m)$  with  $m < k$ . The efficient scores for  $v$ , with  $(\theta_{m+1}, \dots, \theta_k) \equiv \eta$  as nuisance parameters are

$$I_j^*(x) = \frac{x_j}{\theta_j} - \frac{\sum_{l=m+1}^{k+1} x_l}{\sum_{l=m+1}^{k+1} p_l}, \quad j=1, \dots, m,$$

which can be easily checked via orthogonality. This is the same as the efficient score for  $v$  in the multinomial model for observation of

$$Y = (X_1, \dots, X_m, \sum_{j=m+1}^{k+1} X_j) \sim \text{Mult}_{m+1}(1, (v, 1 - \sum_{j=1}^m v_j)),$$

and corresponds with intuition. Consequently the information bound for estimation of  $v$  is just the upper left corner of  $I^{-1}(\theta)$ :

$$I^{-1}(v) = I^{-1}(P_\theta | v, \mathbf{P}) = \text{diag}(v) - vv^T.$$

Again the bound is achieved by the natural estimator  $\hat{v}$  given by the first  $m$  coordinates of  $\hat{p} = \bar{X}$  (which is exactly equal to the first  $m$  coordinates of  $\bar{Y}$  where  $Y^{(1)}, \dots, Y^{(n)}$  are defined in terms of  $X^{(1)}, \dots, X^{(n)}$  exactly as  $Y$  was defined in terms of  $X$  above).

On the other hand, consider the submodel  $\mathbf{P}_1 = \mathbf{P}_1(\eta_0)$  in which  $\eta = (\theta_{m+1}, \dots, \theta_k) = \eta_0$  is known. Then the efficient scores  $I_j^*(\cdot, P_\theta | v, \mathbf{P}_1)$  equal the original scores  $I_j$ ,  $j=1, \dots, m$ , and the  $(m \times m)$ -information matrix is given by the upper left corner of  $I(\theta)$ :

$$I(P_\theta | v, \mathbf{P}_1) = \text{diag}\left(\frac{1}{v}\right) + \frac{1}{p_{k+1}} \underline{1} \underline{1}^T.$$

Therefore the information bound is given by

$$I^{-1}(P_\theta | v, \mathbf{P}) = \text{diag}(v) - \frac{1}{c} vv^T,$$

where  $c \equiv 1 - \sum_{j=m+1}^k \theta_j = \sum_{i=1}^m v_i + p_{k+1}$ . This bound is achieved by the estimator

$$\hat{v}_j^0 = \frac{c \bar{X}_j}{\bar{X}_{k+1} + \sum_{l=1}^m \bar{X}_l}, \quad j=1, \dots, m,$$



where  $\bar{X}_j \equiv n^{-1} \sum_{i=1}^n X_j^{(i)}$ ,  $j = 1, \dots, m, k + 1$ . This estimator is the maximum likelihood estimator with respect to both the unconditional and conditional (given  $X_{m+1}, \dots, X_k$ ) likelihoods for the submodel  $P_1$ .  $\square$

**Example 5. Gaussian linear regression model.**

Let  $X = (Z, Y)$  where  $Y$  is scalar,  $Z$  and  $\theta$  are  $k$ -vectors,  $Y = \theta^T Z + e$ , and  $e \sim N(0, 1)$  independent of  $Z$ . Here  $Z \sim H$ , and we assume that  $E(ZZ^T)$  is nonsingular. This is the model specified by (1.3.6), (1.3.7) with  $\sigma^2 = 1$  for simplicity. If  $\theta \equiv (\theta_1, \dots, \theta_k)^T$ , identify  $v$  with  $\theta_1$ ,  $\eta$  with  $(\theta_2, \dots, \theta_k)$ . Then, if  $Z \equiv (Z_1, \dots, Z_k)^T$ ,

$$\dot{1} = Z e,$$

and hence

$$\dot{1}_1 = Z_1 e, \quad \dot{1}_2 = Z_{[2]} e,$$

where  $Z_{[2]} \equiv (Z_2, \dots, Z_k)^T$ . Therefore

$$(14) \quad \begin{aligned} I &= E(ZZ^T e^2) = E(ZZ^T), \\ \tilde{I} &= I^{-1} Z e, \end{aligned}$$

and

$$(15) \quad \dot{1}_1^* = (Z_1 - I_{12} I_{22}^{-1} Z_{[2]}) e,$$

where  $I_{12} = E(Z_1 Z_{[2]}^T)$ ,  $I_{22} = E(Z_{[2]} Z_{[2]}^T)$ , and

$$(16) \quad \tilde{I}_1 = (Z_1 - I_{12} I_{22}^{-1} Z_{[2]}) e / E(Z_1 - I_{12} I_{22}^{-1} Z_{[2]})^2.$$

If we observe  $Y = (Y_1, \dots, Y_n)^T$ ,  $Z_{k \times n} = (Z^{(1)}, \dots, Z^{(n)})$  where  $(Z^{(j)}, Y_j)$  is the  $j$ th observation, then the least squares (maximum likelihood) estimate of  $\theta$  is

$$\hat{\theta} = [ZZ^T]^{-1} ZY.$$

Since  $n^{-1} ZZ^T = n^{-1} \sum_{i=1}^n Z^{(i)} [Z^{(i)}]^T = E(ZZ^T) + O_p(n^{-1/2})$ , we see that

$$\hat{\theta} = \theta + n^{-1} \sum_{i=1}^n I^{-1} Z^{(i)} e_i + o_p(n^{-1/2})$$

indeed has the (efficient) influence function  $\tilde{I}$ . If we replace  $I_{12}$ ,  $I_{22}$  by the corresponding blocks  $\hat{I}_{12}$ ,  $\hat{I}_{22}$  in  $n^{-1} ZZ^T$ , we obtain the least squares estimate of  $v$  as

$$(17) \quad \begin{aligned} \hat{v} &\equiv \sum_{i=1}^n (Z_1^{(i)} - \hat{I}_{12} \hat{I}_{22}^{-1} Z_{[2]}^{(i)}) Y_i / \sum_{j=1}^n (Z_1^{(j)} - \hat{I}_{12} \hat{I}_{22}^{-1} Z_{[2]}^{(j)})^2 \\ &= v + n^{-1} \sum_{i=1}^n \tilde{I}_1(Z^{(i)}, Y_i) + o_p(n^{-1/2}). \end{aligned}$$

Note that the influence function merely replaces the coefficients of the regression of  $Z_1^{(i)}$  on  $Z_{[2]}^{(i)}$  based on the  $n$  observations by the corresponding population quantities, or, equivalently, the empirical measure orthogonality condition

$$(18) \quad n^{-1} \sum_{i=1}^n (Z_1^{(i)} - \hat{I}_{12} \hat{I}_{22}^{-1} Z_{[2]}^{(i)}) Z_{[2]}^{(i)} = 0$$

is replaced by the population orthogonality condition

$$(19) \quad E(\mathbf{1}_1^* \dot{\mathbf{i}}_2) = E[(Z_1 - I_{12} I_{22}^{-1} Z_{[2]}) Z_{[2]}] = 0,$$

which corresponds to proposition 1.A. Similarly if  $\eta$  is assumed known, the least squares estimate of  $\nu$  is

$$\hat{\nu}^* \equiv \sum_{i=1}^n Z_1^{(i)} Y_i / \sum_{j=1}^n [Z_1^{(j)}]^2$$

with influence function  $[I_{11}]^{-1} \dot{\mathbf{i}}_1$ . Also from (17) and (18)

$$(20) \quad \hat{\nu}^* - \hat{\nu} = \sum_{i=1}^n \left\{ \frac{Z_1^{(i)}}{|Z_1^{(i)}|^2} - \frac{Z_1^{(i)} - \hat{I}_{12} \hat{I}_{22}^{-1} Z_{[2]}^{(i)}}{|Z_1^{(i)} - \hat{I}_{12} \hat{I}_{22}^{-1} Z_{[2]}^{(i)}|^2} \right\} e_i$$

where  $|Z_1^{(i)}|$  is the Euclidean norm of  $(Z_1^{(i)}, \dots, Z_{[2]}^{(i)})^T$ . Now

$$(21) \quad \sum_{i=1}^n \left\{ \frac{Z_1^{(i)}}{|Z_1^{(i)}|^2} - \frac{Z_1^{(i)} - \hat{I}_{12} \hat{I}_{22}^{-1} Z_{[2]}^{(i)}}{|Z_1^{(i)} - \hat{I}_{12} \hat{I}_{22}^{-1} Z_{[2]}^{(i)}|^2} \right\} Z_1^{(i)} = 0$$

by (18). Again this orthogonality relation is the sample version of the population relation

$$(22) \quad E\left(\left(\frac{\dot{\mathbf{i}}_1}{\|\dot{\mathbf{i}}_1\|_0^2} - \tilde{\mathbf{I}}_1\right) \dot{\mathbf{i}}_1\right) = E\left(\frac{Z_1}{\|Z_1\|_0^2} - \frac{Z_1 - I_{12} I_{22}^{-1} Z_{[2]}}{\|Z_1 - I_{12} I_{22}^{-1} Z_{[2]}\|_0^2}\right) Z_1 = 0,$$

which corresponds to proposition 1.B. □

Here is another basic example illustrating proposition 1.

**Example 6. The bivariate normal distribution, continued.**

Suppose that  $X = (X_1, X_2) \sim N_2(\nu, \Sigma)$  where  $\nu \in R^2$  and

$$\Sigma = \begin{pmatrix} \eta_1^2 & \rho \eta_1 \eta_2 \\ \rho \eta_1 \eta_2 & \eta_2^2 \end{pmatrix};$$

here  $\theta = (\nu_1, \nu_2, \eta_1^2, \eta_2^2, \rho) \in R^2 \times R^{+2} \times (-1, 1)$ . The joint density is, with  $z_i \equiv z_i(\theta) = (x_i - \nu_i)/\eta_i$ ,  $i = 1, 2$ ,

$$p(x, \theta) = \frac{1}{2\pi\eta_1\eta_2\sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)}\{z_1^2 - 2\rho z_1 z_2 + z_2^2\}\right\},$$

and the model is  $\mathbf{P} = \{P_\theta : \theta \in R^2 \times R^{+2} \times (-1, 1)\}$ , the family of (nondegenerate) bivariate normal distributions with all five parameters unknown. Hence, the scores for  $\theta$  are

$$\dot{\mathbf{i}}_1(x) = \frac{1}{(1-\rho^2)\eta_1} \{z_1 - \rho z_2\},$$

$$\dot{\mathbf{i}}_2(x) = \frac{1}{(1-\rho^2)\eta_2} \{z_2 - \rho z_1\},$$

$$\dot{i}_3(x) = -\frac{1}{2\eta_1^2} \left\{ 1 - \frac{1}{1-\rho^2} \{z_1^2 - \rho z_1 z_2\} \right\},$$

$$\dot{i}_4(x) = -\frac{1}{2\eta_2^2} \left\{ 1 - \frac{1}{1-\rho^2} \{z_2^2 - \rho z_1 z_2\} \right\},$$

$$\dot{i}_5(x) = \frac{1}{(1-\rho^2)^2} \{ \rho(1-\rho^2) - \rho \{z_1^2 + z_2^2\} + (1+\rho^2)z_1 z_2 \}.$$

It is straightforward to calculate that  $[\dot{i}_1, \dot{i}_2] \perp [\dot{i}_3, \dot{i}_4, \dot{i}_5]$ , and hence the information matrix  $I(P_\theta | \theta, \mathbf{P})$  is block diagonal: the upper left  $2 \times 2$  block is

$$I(v) = \frac{1}{1-\rho^2} \begin{pmatrix} \frac{1}{\eta_1^2} & -\frac{\rho}{\eta_1 \eta_2} \\ -\frac{\rho}{\eta_1 \eta_2} & \frac{1}{\eta_2^2} \end{pmatrix},$$

and the lower right  $3 \times 3$  block for  $(\eta_1^2, \eta_2^2, \rho)$  is

$$I(\eta_1^2, \eta_2^2, \rho) = \frac{1}{1-\rho^2} \begin{pmatrix} \frac{2-\rho^2}{4\eta_1^4} & \frac{-\rho^2}{4\eta_1^2 \eta_2^2} & \frac{-\rho}{2\eta_1^2} \\ \frac{-\rho^2}{4\eta_1^2 \eta_2^2} & \frac{2-\rho^2}{4\eta_2^4} & \frac{-\rho}{2\eta_2^2} \\ \frac{-\rho}{2\eta_1^2} & \frac{-\rho}{2\eta_2^2} & \frac{1+\rho^2}{1-\rho^2} \end{pmatrix}.$$

Thus, knowledge of  $(\eta_1^2, \eta_2^2, \rho)$  does not affect how well we can estimate  $(v_1, v_2)$ , and vice versa. Furthermore,  $I^{-1}(P_\theta | \theta, \mathbf{P})$  is also block diagonal with

$$I^{-1}(P_\theta | v, \mathbf{P}) = \begin{pmatrix} \eta_1^2 & \rho \eta_1 \eta_2 \\ \rho \eta_1 \eta_2 & \eta_2^2 \end{pmatrix}$$

and

$$I^{-1}(P_\theta | \eta_1^2, \eta_2^2, \rho, \mathbf{P}) = \begin{pmatrix} 2\eta_1^4 & 2\rho^2 \eta_1^2 \eta_2^2 & \rho(1-\rho^2)\eta_1^2 \\ 2\rho^2 \eta_1^2 \eta_2^2 & 2\eta_2^4 & \rho(1-\rho^2)\eta_2^2 \\ \rho(1-\rho^2)\eta_1^2 & \rho(1-\rho^2)\eta_2^2 & (1-\rho^2)^2 \end{pmatrix}.$$

These information bounds are achieved by the usual maximum likelihood estimators  $(\bar{X}, \hat{\Sigma})$ , where

$$\hat{\eta}_j^2 = n^{-1} \sum_{i=1}^n (X_j^{(i)} - \bar{X}_j)^2, \quad j = 1, 2,$$

$$\hat{\rho} = n^{-1} \sum_{i=1}^n (X_1^{(i)} - \bar{X}_1)(X_2^{(i)} - \bar{X}_2) / (\hat{\eta}_1 \hat{\eta}_2),$$

and  $X^{(1)}, \dots, X^{(n)}$  are i.i.d. copies of  $X$ . Note that the efficient score function for  $v_2$  in the model  $\mathbf{P}$  is

$$I_2^*(x, P_\theta | v_2, P) = \dot{i}_2(x) - I_{21} I_{11}^{-1} \dot{i}_1(x) = \frac{z_2}{\eta_2},$$

and the efficient influence function for  $v_2$  in the model  $P$  is

$$\tilde{I}_2(x, P_\theta | v_2, P) = \eta_2 z_2 = x_2 - v_2.$$

Now consider the submodel  $P_1 \equiv P_1(v_2) \subset P$  in which the mean  $v_2$  of  $X_2$  is known to be zero:

$$P_1 = \{P_\theta \in P : v_2 = 0\}.$$

For this submodel  $P_1$ , the results from example 3 generalize straightforwardly: the efficient score function  $I_1^*$  for  $v_1$  is

$$I_1^*(x, P_\theta | v_1, P_1) = \dot{i}_1(x, P_\theta | v_1, P_1) = \frac{1}{(1 - \rho^2)\eta_1} (z_1 - \rho z_2),$$

so that

$$I(P_\theta | v_1, P_1) = \frac{1}{(1 - \rho^2)\eta_1^2}, \quad I^{-1}(P_\theta | v_1, P_1) = (1 - \rho^2)\eta_1^2,$$

and the efficient influence function (for  $v_1$  in the submodel  $P_1$ ) is

$$\tilde{I}_1(x, P_\theta | v_1, P_1) = \eta_1(z_1 - \rho z_2) = x_1 - v_1 - \rho \frac{\eta_1}{\eta_2} x_2.$$

Thus knowledge of  $v_2$  reduces the information bound for estimation of  $v_1$  from  $\eta_1^2$  to  $(1 - \rho^2)\eta_1^2$ . If  $\eta_1^2, \eta_2^2, \rho$  are also known, an estimator achieving this bound is

$$\hat{v}_1^0 \equiv \frac{1}{n} \sum_{i=1}^n (X_1^{(i)} - \rho \frac{\eta_1}{\eta_2} X_2^{(i)}).$$

If  $\eta_1^2, \eta_2^2, \rho$  are unknown, then replacing them by their natural estimators also yields an efficient estimate in the model  $P_1$ .

Similarly, if  $\eta_1^2$  and  $\eta_2^2$  are unknown, the information lower bound for estimation of  $\rho$  is  $(1 - \rho^2)^2$ . But if  $\eta_1^2 = \eta_2^2 = 1$  are both known, then the information lower bound for  $\rho$  is  $(1 - \rho^2)^2 / (1 + \rho^2)$ . □

Proposition 1 can be put in a broader context.

**Proposition 2.** Let  $m = 1$  and suppose that  $T_n$  is an asymptotically linear estimator of  $v$  with influence function  $\psi$ . Then:

A.  $T_n$  is Gaussian regular if and only if

$$(23) \quad \psi - \tilde{I}_1 \perp \dot{P} = [\dot{i}_1, \dot{i}_2],$$

or, equivalently, if and only if both

$$(24) \quad \langle \psi, \dot{i}_1 \rangle_0 = 1$$

and

$$(25) \quad \psi \perp [\dot{i}_2].$$

B. If  $T_n$  is regular, then  $\psi \in \dot{P} = [\dot{I}_1, \dot{I}_2]$  if and only if  $\psi = \tilde{I}_1$ .

See figures 2 and 3 on page 31.

We note in passing that (24) and (25) are asymptotic versions of the equations leading to the Cramér-Rao information bound. Consider the problem of minimizing  $\Sigma(P_{\theta_0}, T) = E_0 \psi^2$  subject to (24) and (25). For simplicity take  $k = 2$ . If we write

$$\psi = c\dot{I}_1 + d\dot{I}_2 + \Delta,$$

where  $\Delta \perp [\dot{I}_1, \dot{I}_2]$ , then (25) holds if and only if

$$\psi = c(\dot{I}_1 - I_{12}I_{22}^{-1}\dot{I}_2) + \Delta = c\dot{I}_1^* + \Delta,$$

while (24) forces

$$c = \|\dot{I}_1 - I_{12}I_{22}^{-1}\dot{I}_2\|_0^{-2}.$$

Finally

$$\|\psi\|_0^2 = \|\dot{I}_1^*\|_0^{-2} + \|\Delta\|_0^2.$$

Therefore, the minimizing  $\Delta = 0$ , and, as expected, the minimizing  $\psi$  is the efficient influence function. This argument makes clear the characterizing features of the efficient influence function implied in proposition 1.B:

- (i)  $\tilde{I}_1$  and all other influence functions are orthogonal to  $[\dot{I}_2]$ .
- (ii)  $\tilde{I}_1$  is the unique influence function belonging to  $[\dot{I}_1, \dot{I}_2]$ .
- (iii)  $\tilde{I}_1$  can be obtained by projecting any influence function  $\psi$  corresponding to a regular estimate for  $v$  into  $[\dot{I}_1, \dot{I}_2]$ .

Here is a slight generalization of proposition 2 to a general function  $v(P_\theta) = q(\theta)$ .

**Proposition 3.** Suppose that  $T_n$  is an asymptotically linear estimator at  $\theta_0$  of  $v(P_\theta) = q(\theta)$  with influence function  $\psi$  where  $q : \Theta \rightarrow R^m$ . Then:

A.  $T_n$  is (Gaussian) regular at  $\theta_0$  if and only if  $q(\theta)$  is differentiable at  $\theta_0$  with derivative  $\dot{q}(\theta_0)$  and, with  $\tilde{I} \equiv \tilde{I}(\cdot, P_{\theta_0} | v, P)$ ,

$$(26) \quad \psi - \tilde{I} \perp \dot{P} = [\dot{I}_1, \dot{I}_2],$$

where (26) is equivalent to

$$(27) \quad E_0 \psi \dot{I}^T = \dot{q}(\theta_0).$$

B. If  $T_n$  is regular, then  $\psi \in \dot{P}^m$  if and only if

$$(28) \quad \psi = \tilde{I} = \dot{q}(\theta_0)I^{-1}(\theta_0)\dot{I}(\theta_0).$$

**Proof.** By asymptotic linearity of  $T_n$  and proposition 2.1.2,

$$(a) \quad \mathbf{L}_{\theta_0} \begin{pmatrix} \sqrt{n}(T_n - q(\theta_0)) \\ L_n(\theta_0 + t_n / \sqrt{n}) - L_n(\theta_0) \end{pmatrix} \rightarrow N \left( \begin{pmatrix} 0 \\ -\Sigma_{22} / 2 \end{pmatrix}, \Sigma \right)$$

where

$$(b) \quad \Sigma = [\Sigma_{ij}], \quad \Sigma_{11} = E_0 \Psi \Psi^T, \quad \Sigma_{12} = E_0 \Psi \dot{\mathbf{I}}^T t, \\ \Sigma_{22} = t^T I(\theta_0) t, \quad t_n \rightarrow t.$$

Consequently, by Le Cam's third lemma (lemma A.9.3)

$$(c) \quad \mathbf{L}_{\theta_0 + t_n / \sqrt{n}}(\sqrt{n}(T_n - q(\theta_0))) \rightarrow N(\Sigma_{12}, \Sigma_{11}).$$

Assume now that  $T_n$  is regular. Then

$$(d) \quad \mathbf{L}_{\theta_0 + t_n / \sqrt{n}}(\sqrt{n}(T_n - q(\theta_0 + \frac{t_n}{\sqrt{n}}))) \rightarrow N(0, \Sigma_{11})$$

and from (c) and (d) we conclude

$$(e) \quad \sqrt{n}(q(\theta_0 + \frac{t_n}{\sqrt{n}}) - q(\theta_0)) \rightarrow \Sigma_{12} = E_0 \Psi \dot{\mathbf{I}}^T t.$$

But this implies that  $q$  is differentiable at  $\theta_0$  with derivative  $\dot{q}(\theta_0)$  satisfying (27) and hence (26).

On the other hand, if  $q$  is differentiable and (27) holds, then (e) is valid, which together with (c) implies (d) and hence Gaussian regularity. The proof of A is complete.

As for B, note that A implies that  $\dot{q}$  and hence  $\tilde{\mathbf{I}}$  are well defined and that (26) holds. Since  $\tilde{\mathbf{I}} \in \tilde{\mathbf{P}}^m$ , (26) yields  $\psi \in \tilde{\mathbf{P}}^m$  if and only if  $\psi - \tilde{\mathbf{I}} = 0$ .  $\square$

Choosing  $q(\theta) = q(\nu, \eta) = \nu$  in proposition 3 immediately yields a generalization of proposition 2 to  $m > 1$ . Now (27) becomes

$$(29) \quad E_0 \Psi \dot{\mathbf{I}}_1^T = J_{m \times m},$$

$$(30) \quad E_0 \Psi \dot{\mathbf{I}}_2^T = 0,$$

where  $J$  is the identity. In particular, if  $m = k$  we obtain that the influence function of any linear and Gaussian regular estimate of  $\theta$  has

$$(31) \quad E_0 \Psi \dot{\mathbf{I}}^T = J_{k \times k}.$$

For elementary versions of propositions 2 and 3 under different hypotheses, see Hall and Mathiason (1990, section 3.1).

## 2.5 CONSTRUCTION OF $\sqrt{N}$ -CONSISTENT AND EFFICIENT ESTIMATES

There are many ways of constructing estimates in smooth parametric models which under appropriate regularity conditions are efficient. The most popular construction is Fisher's method of maximum likelihood. Although this method

can fail spectacularly—a famous example is in Kiefer and Wolfowitz (1956, page 905)—there are closely related  $M$ -estimate methods which work under minimal conditions. We shall discuss these further below and in considerable detail in chapter 7. Bayes estimates corresponding to smooth prior distributions and bounded bowl-shaped loss functions also work quite generally. The best recent result is due to Ibragimov and Has'minskii (1981, theorem III.3.1, page 185). Other methods include minimum Hellinger distance estimates (see, e.g., Beran (1977b)), and a variety of approaches suitable in particular models such as minimum  $\chi^2$  in discrete data models, and  $L$ - and  $R$ -estimates for location and scale; see Huber (1981, chapter 3), for example.

### $\sqrt{n}$ -Consistent Preliminary Estimators

We begin by studying the construction of regular, but not in general efficient, minimum distance estimates. These procedures were introduced by Wolfowitz (1957). Using these  $\sqrt{n}$ -consistent starting points we then show how to construct efficient estimates. There are, of course, many methods of constructing  $\sqrt{n}$ -consistent estimates. These include the direct efficient constructions as well as other portmanteau methods, such as the method of moments. We specialize to the method of minimum distance because it can be applied to general (not just parametric) models and is simple to describe and analyze, though not necessarily to implement.

Let  $\mathbf{P} \subset \mathbf{M}$  be the set of all probability measures on  $\mathbf{X}$ , and let  $\theta : \mathbf{P} \rightarrow R^k$  be a parameter. A natural way of constructing estimates of  $\theta$  is to find a "smooth" extension  $\bar{\theta}$  of  $\theta$  to  $\mathbf{M}$  and let

$$(1) \quad T_n = \bar{\theta}(I_n),$$

where  $I_n$  is the empirical distribution of  $X_1, \dots, X_n$  i.i.d.  $P_0$  given by (A.6.6). One way of obtaining a smooth extension  $\bar{\theta}$  of  $\theta$  is by "minimum distance" as follows. Suppose that:

(D1)  $\rho$  is a metric compatible with  $I_n$  in the sense of (A.6.7).

(D2) There exists a map  $\Pi : \mathbf{M} \rightarrow \mathbf{P}$  with  $\rho(\Pi(Q), Q) = \inf \{ \rho(P, Q) : P \in \mathbf{P} \}$  for every  $Q \in \mathbf{M}$ .

(D3)  $\theta$  is  $\rho$ -continuous on  $\mathbf{P}$ . That is,  $\rho(P_m, P) \rightarrow 0$  with  $P_m, P \in \mathbf{P}$  implies  $|\theta(P_m) - \theta(P)| \rightarrow 0$ .

Then the  $\rho$ -extension  $\bar{\theta}$  of  $\theta$  is given by  $\bar{\theta}(Q) = \theta(\Pi(Q))$ , and (1) becomes

$$(2) \quad T_n = \theta(\Pi(I_n)).$$

**Lemma 1.** If (D1)–(D3) hold, then  $T_n$  given by (2) is consistent.

**Proof.** By definition of  $\Pi$  and by (D1)

$$(a) \quad \begin{aligned} \rho(\Pi(I_n), P_0) &\leq \rho(\Pi(I_n), I_n) + \rho(I_n, P_0) \\ &\leq 2\rho(I_n, P_0) \rightarrow_p 0 \quad \text{as } n \rightarrow \infty. \end{aligned}$$

The lemma follows from (D3). □

**Note.** Here and in the future we ignore measurability questions involving  $T_n$  which are unimportant in practice.

**Definition 1.**  $\theta$  is  $\rho$ -Lipschitz on  $\mathbf{P}$  if for every  $P_0 \in \mathbf{P}$  there exists  $c(P_0) < \infty$ ,  $\varepsilon(P_0) > 0$ , such that

$$(3) \quad |\theta(P) - \theta(P_0)| \leq c(P_0)\rho(P, P_0) \quad \text{if } \rho(P, P_0) \leq \varepsilon(P_0).$$

Of course,  $\rho$ -Lipschitz implies  $\rho$ -continuity.

**Lemma 2.** If (D1)–(D2) hold and  $\theta$  is  $\rho$ -Lipschitz on  $\mathbf{P}$ , then  $T_n$  given by (2) is  $\sqrt{n}$ -consistent.

**Proof.** By (2), (3) and then (a) of lemma 1,

$$|T_n - \theta(P_0)| \leq c(P_0)\rho(\Pi(IP_n), P_0) \leq 2c(P_0)\rho(IP_n, P_0). \quad \square$$

We use the approach of lemma 2 to prove the following important theorem, concerning the existence of uniformly  $\sqrt{n}$ -consistent estimates, due to Le Cam (1956). For a generalization to non-Euclidean sample spaces, see Le Cam (1986, section 17.6).

**Theorem 1.** If  $\mathbf{P} = \{P_\theta : \theta \in \Theta\}$  is a regular parametric model on a Euclidean space  $\mathbf{X}$  and  $\theta$  is identifiable, then there exist uniformly  $\sqrt{n}$ -consistent estimates of  $\theta$ .

**Proof.** It follows from proposition A.5.3.C that, uniformly for  $\theta \in K$  compact and uniformly for all Borel sets  $A$ ,

$$(a) \quad P_{\theta+h}(A) - P_\theta(A) = \int_A \dot{p}^T(x, \theta)h \, d\mu(x) + o(|h|).$$

Hence for any class  $\mathbf{A}$  of Borel sets, it follows that

$$(b) \quad \liminf |h|^{-1} \sup \left\{ |P_{\theta+h}(A) - P_\theta(A)| : A \in \mathbf{A} \right\} \\ \geq \inf_{|e|=1} \sup \left\{ \left| \int_A \dot{p}^T(x, \theta)e \, d\mu(x) \right| : A \in \mathbf{A} \right\}.$$

The right side of (b) vanishes if and only if there exists  $e$  such that

$$(c) \quad \int_A \dot{p}^T(x, \theta)e \, d\mu(x) = 0 \quad \text{for all } A \in \mathbf{A}.$$

If  $\mathbf{A}$  is the class of all shifted quadrants, (c) implies

$$(d) \quad \frac{\dot{p}^T(x, \theta)}{p(x, \theta)}e = 0 \quad \text{a.s. } P_\theta,$$

which cannot hold if  $I(\theta)$  is nonsingular. Therefore, for a regular parametrization,

$$(e) \quad d_K(P_\theta, P_{\theta_0}) \geq c^{-1}(\theta_0)|\theta - \theta_0| \quad \text{if } |\theta - \theta_0| \leq \varepsilon(\theta_0).$$

Moreover, the uniformity in (a) and  $L_1$  continuity of  $\theta \rightarrow \dot{p}(\theta)$  imply that we can bound  $c(\theta)$  and  $1/\varepsilon(\theta)$  on compacts.



Since  $\Theta$  is open, there exist compacts  $K_j$  with  $\Theta = \cup_{j=1}^{\infty} K_j$ ,  $K_{j+1} \supset K_j$ ,  $j \geq 1$ . Define  $T_{nj}$  to minimize  $d_K(P_{\theta}, P_n)$  for  $\theta \in K_j$ .  $d_K$  continuity of  $\theta \rightarrow P_{\theta}$  on  $K_j$  guarantees the existence of  $T_{nj}$ . Since  $\theta$  is identifiable, the map  $P_{\theta} \rightarrow \theta$  is  $d_K$  continuous on  $K_j$  and (e) implies (3) with  $\rho = d_K$ . As in lemma 2, it follows from the compatibility of  $d_K$  with  $P_n$  (see section A.6) that  $T_{nj}$  is uniformly  $\sqrt{n}$ -consistent on  $K_j$ . Then let  $T_n = T_{nj}$  with  $d_K(P_{T_n}, P_n) \leq n^{-1/4}$  and  $j$  minimal. It is easy to see that  $P_{\theta}(T_n = T_{nj_0}) \rightarrow 1$ , where  $K_{j_0}$  is the first  $K_j$  such that  $\theta \in K_j$ . Uniform  $\sqrt{n}$ -consistency of  $T_n$  follows.  $\square$

*Asymptotically Efficient Estimators*

The classical method of estimation in regular parametric models is maximum likelihood. If  $\mathbf{P} = \{P_{\theta} : \theta \in \Theta\}$  is a (regular) parametric model, a *maximum likelihood estimate*  $\hat{\theta}_n$  of  $\theta$  satisfies

$$L_n(\hat{\theta}_n) = \max\{L_n(\theta) : \theta \in \Theta\},$$

where  $L_n(\theta)$  is the log-likelihood for  $\theta$  as defined in (2.1.10). Of course, as noted at the beginning of this section,  $\hat{\theta}_n$  may not exist (as in the example of Kiefer and Wolfowitz (1956)), or it may exist, but be inconsistent (see, e.g., Kraft and Le Cam (1956), Ferguson (1982), or Le Cam (1990)). However, the ‘‘usual’’ Cramér-type smoothness assumptions, typically involving boundedness of third derivatives of  $l$ , imply that if  $\hat{\theta}_n$  is well defined it satisfies

$$(4) \quad n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \hat{\theta}_n) \equiv S_n(\hat{\theta}_n) = 0,$$

and, if it is consistent, it is also asymptotically efficient; see, e.g., Cramér (1946, section 33.3), Lehmann (1983, sections 6.2 and 6.3), or Ibragimov and Has’minskii (1981, section III.3). (The conditions of the latter authors for asymptotic efficiency of the maximum likelihood estimate are the weakest.)

Even if the maximum likelihood estimate  $\hat{\theta}_n$  does not exist, we can define a one-step Newton-Raphson approximate ‘‘solution’’ of (4) by

$$(5) \quad \hat{\theta}_n^{\text{approx}} = \tilde{\theta}_n + \left[ -\frac{1}{n} \sum_{i=1}^n \ddot{l}(X_i, \tilde{\theta}_n) \right]^{-1} \frac{1}{n} \sum_{i=1}^n \dot{l}(X_i, \tilde{\theta}_n)$$

(assuming existence of the Hessian matrix  $\ddot{l}$ ). This is the basis of the construction of an efficient estimator given below which, by avoiding use of  $\ddot{l}$ , always works for regular parametric models.

In fact, it follows from the work of Le Cam (1960), (1969), (1970), that regularity of the model together with existence of a  $\sqrt{n}$ -consistent preliminary estimator is enough to guarantee the existence of an efficient estimator. When the sample space is Euclidean, existence of uniformly  $\sqrt{n}$ -consistent estimators is guaranteed by theorem 1. Here is an admittedly artificial construction that is motivated by and refines the one-step approximate solution (5).

- (i) Construct  $\tilde{\theta}_n$  uniformly  $\sqrt{n}$ -consistent as in theorem 1.

(ii) Form a grid of cubes with sides of length  $c n^{-1/2}$  over  $R^k$ , and, given  $\tilde{\theta}_n$ , define  $\theta_n^*$  to be the midpoint of the cube into which  $\tilde{\theta}_n$  has fallen (with some consistent rule for the boundaries of cubes); then  $\theta_n^*$  is also uniformly  $\sqrt{n}$ -consistent. This discretization was introduced by Le Cam (1956, page 144).

(iii) As in section 2.3, let

$$\tilde{I}(\cdot, \theta) = \tilde{I}(\cdot, P_\theta | \theta, \mathbf{P}) = I^{-1}(\theta) \dot{I}(\cdot, \theta),$$

and define

$$(6) \quad \hat{\theta}_n = \theta_n^* + n^{-1} \sum_{i=1}^n \tilde{I}(X_i, \theta_n^*).$$

**Theorem 2.** If  $\mathbf{P}$  is a regular parametric model and if there exists a uniformly (respectively locally)  $\sqrt{n}$ -consistent estimator  $\tilde{\theta}_n$  of  $\theta$ , then the estimator  $\hat{\theta}_n$  given in (6) is a uniformly (respectively locally) efficient estimator of  $\theta$ .

**Proof.** Suppose that  $\theta_n \rightarrow \theta$  for  $\theta_n, \theta \in \Theta$ . In view of theorem 2.3.1 and (2.1.14) of proposition 2.1.2, it suffices to show that, for all  $\varepsilon > 0$ ,

$$(a) \quad P_{\theta_n}(|\sqrt{n}(\hat{\theta}_n - \theta_n) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i, \theta_n)| \geq \varepsilon) \rightarrow 0$$

as  $n \rightarrow \infty$ . Now it follows from the definition (6) of  $\hat{\theta}_n$  that for any  $M > 0$  the left side of (a) is bounded by

$$(b) \quad P_{\theta_n}(A_n) + P_{\theta_n}(|\sqrt{n}(\theta_n^* - \theta_n)| + \frac{1}{\sqrt{n}} \sum_{i=1}^n \{|\tilde{I}(X_i, \theta_n^*) - \tilde{I}(X_i, \theta_n)|\} \geq \varepsilon, A_n^c)$$

$$(c) \quad \leq P_{\theta_n}(A_n) + \sum_{\theta'_n} P_{\theta'_n}(|\sqrt{n}(\theta'_n - \theta_n)| + \frac{1}{\sqrt{n}} \sum_{i=1}^n \{|\tilde{I}(X_i, \theta'_n) - \tilde{I}(X_i, \theta_n)|\} \geq \varepsilon),$$

where

$$A_n \equiv [|\sqrt{n}|\theta_n^* - \theta_n| > M],$$

and where the sum in (c) is over all  $\theta'_n$  in the grid which are at distance at most  $Mn^{-1/2}$  from  $\theta_n$ . Now choose  $M$  so large that the limsup on  $n$  of the first term is arbitrarily small; this is possible by the uniform  $\sqrt{n}$ -consistency of  $\tilde{\theta}_n$  and the choice of  $\theta_n^*$ . Since the number of summands in the second term is bounded, it suffices to show that for all  $\theta'_n$  satisfying  $|\theta'_n - \theta_n| \leq Mn^{-1/2}$  we have

$$(d) \quad \sqrt{n}(\theta'_n - \theta_n) + I^{-1}(\theta'_n)S_n(\theta'_n) - I^{-1}(\theta_n)S_n(\theta_n) = o_{P_{\theta_n}}(1).$$

But this follows from (2.1.15) of proposition 2.1.2 and the continuity of  $\theta \rightarrow I(\theta)$ .

For the local version of the theorem, the argument is the same with  $\theta_n = \theta + t_n/\sqrt{n}$  with  $|t_n|$  bounded.  $\square$

The construction (6) uses  $I(\theta_n^*)$  as an estimator of  $I(\theta)$ . The theorem remains valid if  $I(\theta)$  is estimated instead by

$$\hat{I}_n \equiv \hat{I}_n(\theta_n^*) \equiv n^{-1} \sum_{i=1}^n \dot{\ell}^T(X_i, \theta_n^*),$$

as is easily proved by use of (d) of the proof of lemma A.9.5 and Chung's uniform law of large numbers, theorem A.7.3.

Another (also artificial) construction of an efficient estimator will be given in section 7.8 as a corollary of a different construction using sample splitting in place of discretization.

# 3 | Information Bounds for Euclidean Parameters in Infinite-Dimensional Models

## 3.1 INTRODUCTION AND OVERVIEW

Our point of view is that of Stein (1956) as developed by Koshevnik and Levit (1976), Levit (1975), (1978), Pfanzagl and Wefelmeyer (1982), and Begun, Hall, Huang, and Wellner (1983). We consider a general  $\mathbf{P} \subset \mathbf{M}$ , a parameter  $v: \mathbf{P} \rightarrow R^m$  which we wish to estimate, and a fixed  $P_0 \in \mathbf{P}$ . Recall that for a fixed sample space  $(\mathbf{X}, \mathcal{A})$ ,  $\mathbf{M}$  is the set of all probability measures dominated by the  $\sigma$ -finite measure  $\mu$ .

**Definition 1.** We call any subset  $\mathbf{Q}$  of  $\mathbf{P}$  which has a regular Euclidean parametrization a *regular parametric submodel* of  $\mathbf{P}$ .  $P_0 \in \mathbf{P}$  is regular if it belongs to a regular parametric submodel of  $\mathbf{P}$ .

For example,  $\{N(\mu, \sigma^2) : \mu \in R, \sigma^2 > 0\}$  is a regular parametric submodel of  $\{\text{all absolutely continuous symmetric distributions on } R\}$ .

We say that an estimate  $T$  is *locally regular on  $\mathbf{P}$*  if it is locally regular on all regular parametric submodels  $\mathbf{Q}$  of  $\mathbf{P}$ . *Local* and *uniform Gaussian regularity* and *linearity* are defined similarly.

**Convention.** We will frequently abbreviate *locally regular* to just *regular*, in keeping with most of the recent literature.

Suppose  $m = 1$ . If  $T$  is locally Gaussian regular as an estimate of  $v$  on  $\mathbf{P}$  with variance  $\Sigma(P_0, T)$ , and  $\mathbf{Q}$  is a regular parametric submodel which contains  $P_0$ , then

$$(1) \quad \Sigma(P_0, T) \geq I^{-1}(P_0 | v, \mathbf{Q})$$

if  $v$  considered as a function of the parameter indexing  $\mathbf{Q}$  is smooth enough so that  $I^{-1}(P_0 | v, \mathbf{Q})$  is defined. It is then natural to define

$$(2) \quad I^{-1}(P_0 | v, \mathbf{P}) \equiv \sup \{ I^{-1}(P_0 | v, \mathbf{Q}) : \text{regular parametric } \mathbf{Q} \subset \mathbf{P} \}$$

as a measure of the best we can achieve in estimating  $v$  over  $\mathbf{P}$ .

Definition (2) creates no contradictions if  $\mathbf{P}$  is itself a regular parametric model for which  $I^{-1}(P_0 | v, \mathbf{P})$  is defined in the sense of (2.3.1). Indeed, by theorems 2.5.1 and 2.5.2, if  $\mathbf{X}$  is Euclidean and  $v$  is identifiable, we can then construct an efficient estimate  $T^*$  of  $v$  on  $\mathbf{P}$ . Such an estimate is Gaussian

regular on any regular parametric submodel  $\mathbf{Q}$  of  $\mathbf{P}$ . Therefore,

$$I^{-1}(P_0 | v, \mathbf{P}) = \Sigma(P_0, T^*) \geq I^{-1}(P_0 | v, \mathbf{Q}).$$

Since  $\mathbf{P}$  is a regular parametric submodel of itself, definitions (2.3.1) and (2) are consistent.

Our main goal can now be stated as exhibiting or showing how to construct estimates  $T$  of  $v$  which

- (i) are locally Gaussian regular on  $\mathbf{P}$ ;
- (ii) have  $\Sigma(P_0, T) = I^{-1}(P_0 | v, \mathbf{P})$ .

A further goal is to construct locally Gaussian regular estimates  $T$  which satisfy (ii) for all  $P_0 \in \mathbf{P}$ .

In fact, there may *not* exist any locally regular  $T$  which achieve the bounds. Ritov and Bickel (1990) give an example in which  $I^{-1}(P_0 | v, \mathbf{P})$  is well defined and positive, but no  $\sqrt{n}$ -consistent, much less locally regular, estimate of  $v$  exists.

In this chapter we shall give conditions on  $\mathbf{P}$  and  $v$  which enable us, in principle, to find regular parametric submodels  $\mathbf{Q}$  for which  $I^{-1}(P_0 | v, \mathbf{P})$  is assumed or at least a sequence  $\{\mathbf{Q}_j\}$  of such models such that  $I^{-1}(P_0 | v, \mathbf{Q}_j) \rightarrow I^{-1}(P_0 | v, \mathbf{P})$ . Note that such a least favorable  $\mathbf{Q}$  will typically not be unique and will depend on  $P_0, v, \mathbf{P}$ . We are then able to define  $I^{-1}$  for  $m > 1$  as well. In chapter 5 we extend these notions to infinite-dimensional  $v$ . In chapters 4 and 6 we discuss the construction of different types of semiparametric models and apply the methods of chapters 3 and 5 to a large number of examples.

Exhibition of  $\mathbf{Q}$  or  $\{\mathbf{Q}_j\}$  still leaves the major problem of constructing or exhibiting  $T$  which satisfy (i) and (ii). In particular models, locally regular  $T$  have often been proposed on other grounds and can be checked against (ii). Examples appear in this chapter and chapters 4–6. Various general methods of estimation are discussed and applied in chapter 7.

To calculate least favorable  $\mathbf{Q}$  we pick up the theme of section 2.4. We characterize  $\mathbf{Q}$  and associated efficient score and influence functions by projection on tangent spaces, appropriately defined closed linear subspaces of  $L_2(\mu)$  or  $L_2(P_0)$  associated with  $P_0$  and  $\mathbf{P}$ . Tangent spaces are introduced in section 2. Their calculation is discussed there and for special cases throughout chapters 3–6. The extension of the ideas of section 2.4 is the topic of sections 3 and 4.

The machinery we develop in theorems 3.3.1 and 3.4.1 is sufficient but not necessary to achieve our goals of obtaining  $T$  satisfying (i) and (ii). These methods of finding  $\mathbf{Q}$  can and shall be used heuristically without necessarily checking regularity conditions. The  $\mathbf{Q}$  so obtained if it is a regular parametric submodel can always be used in providing a bound for the asymptotic variance of any locally regular estimate. This point of view is developed in the final example of section 4.

We conclude this section with a cautionary example showing that  $I^{-1}$  can be infinite, implying that there are parameters not estimable at rate  $n^{-1/2}$ .

There is an enormous literature on the subject of best rates of convergence and how to achieve them in various contexts. The geometry behind the spectrum of such rates is explored by Donoho and Liu (1987), (1991a,b), who also put the  $n^{-1/2}$ -rate domain in context.

**Example 1. Density estimation.**

Let  $\mathbf{P}$  consist of all distributions on  $R$  with continuous densities. Let

$$v(P) = p(0),$$

where  $p$  is the density of  $P$ . Then

$$(3) \quad I^{-1}(P \mid v, \mathbf{P}) = \infty$$

for all  $P$ . Simple nonparametric restrictions on  $p$  such as bounds on higher derivatives, analyticity do not help. This is in line with various results in density estimation showing only rates of convergence slower than  $\sqrt{n}$  are possible in such models. Surprisingly, however, there does exist a nonparametric model  $\mathbf{P}$  (Ibragimov and Has'minskii (1982)) for which  $I^{-1}$  is finite. We establish (3) for  $p(0) > 0$ . Define a parametric submodel, valid for  $|\eta| < 1$  by

$$(4) \quad p(x, \eta) = p(x)(1 + \eta h(x)),$$

where

$$(5) \quad \sup_x |h(x)| \leq 1, \quad \int h(x)p(x) dx = 0.$$

Then

$$v(P_\eta) = p(0)(1 + \eta h(0)).$$

It is easy to check regularity of  $\mathbf{Q} = \{P_\eta\}$  by, e.g., proposition 2.1.1 and calculate

$$I^{-1}(P_0 \mid v, \{P_\eta\}) = p^2(0)h^2(0) / \int h^2(x)p(x) dx.$$

It is also easy to exhibit  $h_j$  satisfying (5) and  $h_j(0) = 1$  with  $\int h_j^2(x)p(x) dx \rightarrow 0$  as  $j \rightarrow \infty$ , so that the sup in (2) is  $\infty$ . Bounds on derivatives of any order on  $p(\cdot, \eta)$  only translate into smoothness conditions on  $h$  for small  $|\eta|$  and thus (3) holds quite generally. This method can be extended to yield optimal rates of convergence; see, e.g., Ibragimov and Has'minskii (1981, section IV.5, pp. 237–240). □

**Notation.** Quantities such as  $I(P \mid v, \mathbf{P})$ ,  $\tilde{I}(\cdot, P \mid v, \mathbf{P})$  occur frequently in the sequel. We will often suppress one or more of their arguments, the state of nature,  $P$ , the parameter  $v$ , or the model  $\mathbf{P}$ , when convenient.

### 3.2 TANGENT SPACES

In this section we will study  $\mathbf{P}$  as a subset  $\mathbf{S}$  of the Hilbert space  $L_2(\mu)$  and as a subset of  $L_2(P_0)$ , via the correspondences  $P \longleftrightarrow s$  and  $P \longleftrightarrow r$  of chapter 2. Our aim is to identify regular parametric submodels and score functions, as

geometrical objects. We need some definitions whose geometrical origins are evident.

Let  $\mathbf{H}$  be a Hilbert space with norm  $\|\cdot\|$  and inner product  $\langle \cdot, \cdot \rangle$ , and  $\mathbf{V}$  a subset of  $\mathbf{H}$ . We establish several conventions:

**Convention 1.** If  $W$  is a subset or vector of elements of  $\mathbf{H}$ , let  $\text{lin}(W)$ ,  $\overline{W}$ , and  $[W]$  denote the linear span of  $W$ , the closure of  $W$ , and the closed linear span of  $W$  respectively.

**Convention 2.** If we write

$$a_n = b_n + O(\varepsilon_n),$$

where  $a_n, b_n \in \mathbf{H}$ ,  $\varepsilon_n \geq 0$  we mean  $\|a_n - b_n\| = O(\varepsilon_n)$ . The same convention applies to  $o$  notation.

**Definition 1.**  $\mathbf{V}$  is a ( $k$ -dimensional) *surface* in  $\mathbf{H}$  if it can be represented as the image of the open unit sphere in  $R^k$  under a continuously Fréchet differentiable map which is of rank  $k$ . That is, we can write

$$\mathbf{V} = \{v(\eta) : |\eta| < 1\},$$

where

- (i)  $v(\eta + \Delta) = v(\eta) + \Delta^T \dot{v}(\eta) + o(|\Delta|)$ ,
- (ii)  $\dot{v} = (\dot{v}_1, \dots, \dot{v}_k)^T \in \mathbf{H}^k$ ,
- (iii)  $\dim[\dot{v}] = k$ .

Any surface or curve ( $k = 1$ ) has many representations (parametrizations); however if  $\eta \rightarrow v(\eta)$  and  $\gamma \rightarrow g(\gamma)$  are two representations of  $\mathbf{V}$ ,  $g(\gamma_0) = v(\eta_0)$ , and  $v^{-1}(v(\eta_0)) = \eta_0$ , then  $\dot{v}(\eta_0) = M \dot{g}(\gamma_0)$  where  $M_{k \times k}$  is nonsingular. This follows by an application of the chain rule and inverse function theorem (e.g. Dieudonné (1960, theorem 10.2.5)). Therefore  $[\dot{v}(\eta_0)] = [\dot{g}(\gamma_0)]$  is independent of the parametrization. We call it the *tangent space* of  $\mathbf{V}$  at  $v(\eta_0)$  and write it  $\dot{\mathbf{V}}$  (or  $\dot{\mathbf{V}}(v_0)$  when we need to identify  $v_0 = v(\eta_0)$ ). For convenience we will identify a surface or curve with a representation  $v(\cdot)$  and for a curve call  $\dot{v}(\eta)$  its tangent at  $v(\eta)$ .

If we take  $\mathbf{H} = L_2(\mu)$  and  $\mathbf{V} = \mathbf{S}$  our discussion and example A.2.1 lead to:

**Proposition 1.**

- A.  $\mathbf{P} = \{P_\theta : \theta \in \Theta \subset R^k\}$  is a regular parametric model and  $\Theta$  is a surface in  $R^k$  if and only if  $\mathbf{S}$  is a  $k$ -dimensional surface in  $L_2(\mu)$ .
- B. In this case we can represent the projection operator onto  $\dot{\mathbf{S}}$  by

$$(1) \quad \Pi(h | \dot{\mathbf{S}}) = 4 \langle h, \dot{s} \rangle^T I^{-1} \dot{s},$$

where we define  $\dot{s} = (\dot{s}_1, \dots, \dot{s}_k)^T$  and hence  $\langle h, \dot{s} \rangle = (\langle h, \dot{s}_1 \rangle, \dots, \langle h, \dot{s}_k \rangle)^T$ .

Now consider the image of  $\mathbf{S}$  under the mapping

$$(2) \quad s \rightarrow r = 2\left(\frac{s}{s_0} - 1\right)1_{[s_0 > 0]},$$

which maps  $s_0$  into zero. This is a subset of  $L_2(P_0)$ . The mapping (2) may not be one-to-one, and the location of the image obviously depends on  $P_0$ . Nevertheless because all information bound calculations are local we do not lose anything by identifying this image with  $\mathbf{P}$ . By a familiar abuse of notation, we even call this set  $\mathbf{P}$ . Then

$$(3) \quad \dot{\mathbf{P}} = \left[ \frac{2\dot{\mathbf{S}}}{s_0} \right] = [\dot{i}_1, \dots, \dot{i}_k]$$

by (2.1.4). Moreover

$$\Pi_0(h | \dot{\mathbf{P}}) = \langle h, \dot{\mathbf{i}} \rangle_0 I^{-1} \dot{\mathbf{i}}$$

where  $\dot{\mathbf{i}} = (\dot{i}_1, \dots, \dot{i}_k)^T$ ,  $\langle h, \dot{\mathbf{i}} \rangle_0$  is now the vector of inner products and the subscript 0 refers to operations in the Hilbert space  $L_2(P_0)$ . Note that for any  $\dot{\mathbf{S}}$  and corresponding  $\dot{\mathbf{P}}$ , the projection operators in the two Hilbert spaces are related by the identities

$$(4a) \quad \Pi_0(h | \dot{\mathbf{P}}) = s_0^{-1} \Pi(h s_0 | \dot{\mathbf{S}}) \quad \text{for } h \in L_2(P_0),$$

$$(4b) \quad \Pi(t | \dot{\mathbf{S}}) = s_0 \Pi_0\left(\frac{t}{s_0} | \dot{\mathbf{P}}\right) \quad \text{for } t \in L_2(\mu).$$

Note that (4a) and (4b) are equivalent to the isomorphism of  $\dot{\mathbf{P}}$  and  $\dot{\mathbf{S}}$ .

We return to the general nonparametric case.

**Definition 2.** If  $v_0 \in \mathbf{V} \subset \mathbf{H}$ , let the *tangent set* at  $v_0$  be the union of all the (one-dimensional) tangent spaces of curves  $\mathbf{C} \subset \mathbf{V}$  passing through  $v_0$ , and denote it  $\dot{\mathbf{V}}^0$ . We call the closed linear span  $[\dot{\mathbf{V}}^0]$  of the tangent set  $\dot{\mathbf{V}}^0$  the *tangent space* of  $\mathbf{V}$ , and denote it by  $\dot{\mathbf{V}}$ .

If  $\mathbf{V}$  is a  $k$ -dimensional surface,  $\dot{\mathbf{V}}^0 = \dot{\mathbf{V}}$  and  $\dot{\mathbf{V}}$  agrees with our previous definition. In general,  $\dot{\mathbf{V}}^0$  is a union of one-dimensional spaces. For example, let

$$\begin{aligned} \mathbf{P} = \{ & N(\mu, 1) : \mu \in R \} \\ & \cup \{ (1 - \varepsilon)N(0, 1) + \varepsilon U(-\frac{1}{2}, \frac{1}{2}) : |\varepsilon| < \phi(\frac{1}{2}) / (1 - \phi(\frac{1}{2})) \} \end{aligned}$$

where  $U(-\frac{1}{2}, \frac{1}{2})$  is the uniform distribution on  $(-\frac{1}{2}, \frac{1}{2})$ . At  $P_0 = N(0, 1)$ ,  $\dot{\mathbf{P}}^0$  is the union of the linear space generated by the identity and that generated by  $x \rightarrow 1_{[-1/2, 1/2]}(x) / \phi(x) - 1$ .

Note that, in general,  $\dot{\mathbf{V}}^0$  is linear if any two curves contained in  $\mathbf{V}$  passing through  $v_0$  are both contained in some surface contained in  $\mathbf{V}$  (at least in a neighborhood of  $v_0$ ). The following proposition will be very useful in later chapters in calculating tangent spaces  $\dot{\mathbf{S}}$  and  $\dot{\mathbf{P}}$ .

**Proposition 2.** If  $\overline{\dot{\mathbf{V}}^0}$  is linear, then  $\overline{\dot{\mathbf{V}}^0} = \dot{\mathbf{V}}$ .

**Proof.** The inclusion  $\overline{\dot{\mathbf{V}}^0} \subset \dot{\mathbf{V}}$  is trivial since  $\dot{\mathbf{V}}^0 \subset \text{lin}(\dot{\mathbf{V}}^0)$ , and hence



$\overline{\dot{V}^0} \subset \overline{\text{lin}(\dot{V}^0)} \equiv [\dot{V}^0] \equiv \dot{V}$ . To prove the reverse inclusion, note that if  $\overline{\dot{V}^0}$  is linear, then  $\overline{\dot{V}^0} \supset \text{lin}(\dot{V}^0)$ . Hence  $\overline{\dot{V}^0} \supset \overline{\text{lin}(\dot{V}^0)} = \dot{V}$ .  $\square$

In all examples of interest to us, the closures  $\overline{\dot{P}^0}$  and  $\overline{\dot{S}^0}$  of  $\dot{P}^0$  and  $\dot{S}^0$  are linear spaces; thus proposition 2 will be of frequent use in the subsequent developments. An important feature of general  $\dot{S}$  and  $\dot{P}$  is given by:

**Proposition 3.** If  $t \in \dot{S}$

$$\langle t, s_0 \rangle = \int t s_0 d\mu = 0,$$

and, if  $h \in \dot{P}$ ,

$$\langle h, 1 \rangle_0 = \int h dP_0 = 0.$$

For  $t \in \dot{S}^0$ ,  $h \in \dot{P}^0$ , these are just restatements of (2.1.8)–(2.1.9). Note that (4a) and (4b) continue to hold for general  $\dot{S}$  and  $\dot{P}$ . We continue to deal with both of these isomorphic spaces because each has its advantages;  $\dot{S}$  in examples such as 3.3.2,  $\dot{P}$  in most of the examples we deal with.

### Calculation of Tangent Spaces

If  $h$  is tangent at  $P_0$  and  $h \longleftrightarrow \{P_\eta : |\eta| < 1\}$ , then  $h$  is the derivative in  $P_0$ -measure of  $\log(dP_\eta/d\mu)$  at  $\eta = 0$ . Therefore (by, e.g., Loève (1977, section 9.1, page 153)),

$$h = \lim_j \left\{ \log(dP_{\eta_j}/d\mu) - \log(dP_0/d\mu) / \eta_j \right\} \text{ a.s. } P_0$$

for some sequence  $\eta_j \rightarrow 0$ . This suggests that we approach the calculation of  $\dot{P}$  heuristically by

- (i) Calculating  $(\partial/\partial\eta)\log(dP_\eta/d\mu) |_{\eta=0}$  for smooth families  $\{P_\eta\}$ .
- (ii) Intersecting the set of all functions obtained in (i) with  $L_2(P_0)$  and forming the closed linear span of the result.

Although this closed linear space may be larger than  $\dot{P}$  and, conceivably, since  $h$  need only be a derivative in measure, could be smaller, it coincides with  $\dot{P}$  in most of the substantive examples we study. The class of  $\{P_\eta\}$  which need be considered can often be reduced. For instance, suppose  $\mathbf{P} = \{P_G : G \in \mathbf{G}\}$  and  $\mathbf{G}$  is a convex subset of a linear space. Then the closed linear span of

$$\left\{ \frac{\partial}{\partial\eta} \log p(\cdot, (1-\eta)G_0 + \eta G_1) |_{\eta=0} : G_1 \in \mathbf{G} \right\}$$

usually agrees with  $\dot{P}$  at  $P_{G_0}$ .

Here is another useful principle. If  $\mathbf{P} = \{P_{(\alpha,\beta)} : \alpha \in A, \beta \in B\}$ ,  $P_0 = P_{(\alpha_0,\beta_0)}$ , and  $\dot{P}_1$  is the tangent space at  $P_0$  to the set  $\mathbf{P}_1 = \{P_{(\alpha,\beta_0)} : \alpha \in A\}$  and  $\dot{P}_2$  is the tangent space at  $P_0$  to  $\mathbf{P}_2 = \{P_{(\alpha_0,\beta)} : \beta \in B\}$ , then  $\dot{P}$  typically agrees with  $\dot{P}_1 + \dot{P}_2$ . This principle is discussed further below in the context of

the symmetric location problem. All of these heuristics are used at least implicitly in our subsequent calculations of tangent spaces.

Here are the first important examples in which  $\dot{S} = \overline{\dot{S}^0}$  and  $\dot{P} = \overline{\dot{P}^0}$  can be calculated explicitly. We calculate either  $\dot{S}$  and  $\Pi(\cdot | \dot{S})$  and then pass to  $\dot{P}$  and  $\Pi_0(\cdot | \dot{P})$  via (3) and (4) or calculate in the opposite direction, depending on which option is most convenient.

**Example 1. All probabilities dominated by  $\mu$ ,  $P = M_\mu$ .**

In this case, by proposition 3,

$$(5) \quad \dot{P} \subset \{h \in L_2(P_0) : \int h dP_0 = 0\} \equiv L_2^0(P_0).$$

In fact equality holds in (5).

The claimed identity follows if we exhibit a dense subset of  $\{h \in L_2(P_0) : \int h dP_0 = 0\}$  which is contained in  $\dot{P}^0$ . In particular, consider  $h$  such that  $h$  is bounded and  $\int h dP_0 = 0$ . To see that this set is dense, approximate general  $h$  with  $\int h dP_0 = 0$  by

$$h_M - \int h_M dP_0,$$

where  $h_M = h 1_{\{|h| \leq M\}}$ . For such  $h$  define

$$(6) \quad p(\eta) = \exp[\eta h - b(\eta)] p_0, \quad \eta \in R,$$

where

$$b(\eta) = \log \int \exp[\eta h] p_0 d\mu.$$

Now  $\eta \rightarrow p(\eta) \in \mathbf{P}$  is an exponential family, and, for  $h \neq 0$ , regular with

$$(7) \quad \dot{s}(0) = \frac{1}{2} h s_0,$$

since  $\dot{b}(0) = \int h dP_0 = 0$ . Consequently equality holds in (5) in view of (2.1.4). Equivalently

$$(8) \quad \dot{S} = \{t \in L_2(\mu) : t = t 1_{\{s_0 > 0\}}, \langle t, s_0 \rangle = 0\}.$$

For future reference we note that, by (A.2.16),

$$(9) \quad \begin{aligned} \Pi(t | \dot{S}) &= t 1_{\{s_0 > 0\}} - \langle t, s_0 \rangle s_0 \quad \text{and} \\ \Pi_0(h | \dot{P}) &= h - \int h dP_0. \end{aligned}$$

It is not necessary to restrict attention to bounded  $h$  for this example (and the next). Given  $h \in L_2(P_0)$ ,  $\int h dP_0 = 0$ , let  $\psi: R \rightarrow (0, \infty)$  be a bounded continuously differentiable function with bounded derivative  $\psi'$  and with  $\psi(0) = \psi'(0) = 1$  and  $\psi'/\psi$  bounded; for example  $\psi(x) = 2(1 + e^{-2x})^{-1}$ . Define

$$p^*(\eta) = p_0 \psi(\eta h) / \int \psi(\eta h) dP_0.$$

Now,  $p^*(0) = p_0$ , and since

$$\frac{\partial}{\partial \eta} \int \psi(\eta h) dP_0 |_{\eta=0} = \int h dP_0 = 0,$$

we have

$$\frac{\partial}{\partial \eta} \log p^*(\eta) |_{\eta=0} = h.$$

It is easy to see, using proposition 2.1.1, that for  $h \neq 0$  and  $\varepsilon$  sufficiently small  $\{P_\eta : dP_\eta/d\mu = p^*(\eta), |\eta| < \varepsilon\}$  is a regular parametric model with score function  $h$ , and we again arrive at equality in (5). Note that this second construction in fact shows that  $\dot{P}^0 = \dot{P}$ . □

**Example 2. All probabilities symmetric about  $\theta_0$  fixed.**

If  $s(\eta) = p^{1/2}(\eta)$  given by (6) is symmetric about  $\theta_0$  for all  $\eta$ , so is the  $L_2$ -limit of  $(s(\eta) - s(0))/\eta$  since the  $L_2$ -limit is an a.e. ( $\mu$ ) limit for some sequence  $\eta_j \rightarrow 0$ . The second construction of example 1 now yields

$$\begin{aligned} \dot{P} &= \{h \in L_2(P_0) : \int h dP_0 = 0, h(x) = h(2\theta_0 - x) \text{ a.s. } P_0\} \\ &= \dot{P}^0. \end{aligned}$$

Here we have

$$(10) \quad \Pi_0(h | \dot{P})(x) = \frac{1}{2}(h(x) + h(2\theta_0 - x)) - \int h dP_0$$

and

$$\Pi(t | \dot{S}) = \frac{1}{2}(t(x) + t(2\theta_0 - x)) 1_{[s_0 > 0]} - \langle t, s_0 \rangle s_0.$$

It is easy to verify (10), using proposition A.3.1, upon noting that

$$\{h : h \text{ symmetric about } \theta_0\} = \{h : h \text{ measurable } \mathcal{B}_0\},$$

where  $\mathcal{B}_0$  is the  $\sigma$ -field induced by the function  $x \rightarrow |x - \theta_0|$ . □

**Example 3. Constraint defined models.**

This is an example where it is convenient to first calculate  $\dot{S}$ . Suppose that

$$(11) \quad P = \{P \ll \mu : \gamma_i(P) = 0, i = 1, \dots, r\}.$$

Suppose  $\gamma_i$  is pathwise Fréchet differentiable with derivative  $\dot{\gamma}_i$  when viewed as a map from  $S$  to  $R$ . That is, abusing notation in a familiar way, as  $s \in S$  tends to  $s_0$ ,

$$(12) \quad \gamma_i(s) = \gamma_i(s_0) + \langle \dot{\gamma}_i(s_0), s - s_0 \rangle + o(\|s - s_0\|)$$

with  $\dot{\gamma}_i(s_0) \in L_2(\mu)$ . An example of such a constraint for  $P$  absolutely continuous Lebesgue and supported on a fixed compact is knowledge that the variance of  $X$  is  $\sigma_0^2$ , so that

$$\gamma(P) = \text{Var}_P(X) - \sigma_0^2.$$

A similar, but more complicated constraint is studied in section 3.3.

Note that

$$(13) \quad \dot{S} \subset \{t \in L_2^0(\mu) : t = t 1_{[s_0 > 0]}, \langle t, \dot{\gamma}_i(s_0) \rangle = 0, i = 1, \dots, r\},$$

since for any curve  $s(\eta) = s_0 + \eta t + o(\eta) \in S$ ,

$$0 = \gamma_i(s(\eta)) - \gamma_i(s_0) = \eta \langle \dot{\gamma}_i(s_0), t \rangle + o(\eta).$$

Equality in (13) holds under further conditions. Here is the case  $r = 1$ .

Suppose  $s \rightarrow \dot{\gamma}(s)$  is continuous at  $s_0$  and  $\dot{\gamma}(s_0) \notin [s_0]$ . There exists  $t_1 \in L_2(\mu)$  such that

$$(14) \quad \langle t_1, s_0 \rangle = 0, \quad \langle t_1, \dot{\gamma}(s_0) \rangle = 1.$$

Given  $t$  such that

$$\langle t, s_0 \rangle = 0, \quad \langle t, \dot{\gamma}(s_0) \rangle = 0,$$

let

$$p(\eta, \varepsilon) = p_0 \Psi \left( \varepsilon \frac{t_1}{s_0} + \eta \frac{t}{s_0} \right) / b(\eta, \varepsilon)$$

where  $\Psi$  is as in example 1,

$$b(\eta, \varepsilon) = \int p_0 \Psi \left( \varepsilon \frac{t_1}{s_0} + \eta \frac{t}{s_0} \right) d\mu.$$

But, if  $s(\eta, \varepsilon) = p^{1/2}(\eta, \varepsilon)$  we have, by  $\langle t_1, s_0 \rangle = \langle t, s_0 \rangle = 0$ ,

$$s(\eta, \varepsilon) = s_0 + \varepsilon t_1 + \eta t + o(\varepsilon) + o(\eta).$$

Then, by (12)

$$\begin{aligned} \gamma(s(\eta, \varepsilon)) &= \langle \dot{\gamma}(s_0), \varepsilon t_1 + \eta t \rangle + o(\varepsilon) + o(\eta) \\ &= \varepsilon + o(\varepsilon) + o(\eta). \end{aligned}$$

Consequently  $\varepsilon \rightarrow \gamma(s(\eta, \varepsilon))$  has a root  $\varepsilon(\eta)$  with

$$\varepsilon(\eta) = o(\eta).$$

Then for  $|\eta|$  small,  $\eta \rightarrow s(\eta, \varepsilon(\eta))$  is a curve in  $S$  with the required tangent  $t$ , and equality in (13) and  $\dot{S} = \dot{S}^0$  follow.

Here, by example A.2.1 and formula (A.2.11), if  $t 1_{[s_0 > 0]} = t$ ,

$$\begin{aligned} (15) \quad \Pi(t | \dot{S}) &= t - \Pi(t | [s_0, \dot{\gamma}(s_0)]) \\ &= t - \langle t, s_0 \rangle_{s_0} - \Pi(t | [\dot{\gamma}(s_0) - \langle \dot{\gamma}(s_0), s_0 \rangle_{s_0}]) \\ &= t - \langle t, s_0 \rangle_{s_0} - c(t)(\dot{\gamma}(s_0) - \langle \dot{\gamma}(s_0), s_0 \rangle_{s_0}), \end{aligned}$$

where

$$c(t) = \frac{\langle t, \dot{\gamma}(s_0) - \langle \dot{\gamma}(s_0), s_0 \rangle s_0 \rangle}{\|\dot{\gamma}(s_0) - \langle \dot{\gamma}(s_0), s_0 \rangle s_0\|^2}$$

It follows that  $\dot{\mathbf{P}} = \overline{\dot{\mathbf{P}}^0} = \{h \in L_2(P_0) : \int h dP_0 = 0, \int h \dot{\lambda} dP_0 = 0\}$  and

$$(16) \quad \Pi_0(h | \dot{\mathbf{P}}) = h - \int h dP_0 - \frac{\text{Cov}_0(h, \dot{\lambda})}{\text{Var}_0(\dot{\lambda})}(\dot{\lambda} - \int \dot{\lambda} dP_0),$$

where  $\dot{\lambda} \equiv \dot{\gamma}(s_0)/2s_0$ .

This argument can be extended to  $r > 1$  if the maps  $s \rightarrow \dot{\gamma}_i(s), i = 1, \dots, r$ , are continuous at  $s_0$  and  $\dot{\gamma}_i(s_0), i = 1, \dots, r$ , and  $s_0$  are linearly independent.  $\square$

Characterizing tangent spaces for semiparametric models is much harder. If  $\mathbf{P} = \{P_{(\theta, G)} : \theta \in \Theta, G \in \mathbf{G}\}$ , let

$$\mathbf{P}_1(G_0) = \{P \in \mathbf{P} : \theta \in \Theta, G = G_0\},$$

$$\mathbf{P}_2(\theta_0) = \{P \in \mathbf{P} : G \in \mathbf{G}, \theta = \theta_0\}$$

be the two sections of  $\mathbf{P}$  at  $(\theta_0, G_0) \in \Theta \times \mathbf{G}$ .

Define  $\dot{\mathbf{P}}_1, \dot{\mathbf{P}}_2$  to be the tangent spaces corresponding to  $\mathbf{P}_1, \mathbf{P}_2$ . By definition,  $\dot{\mathbf{P}}_1, \dot{\mathbf{P}}_2$  are both contained in  $\dot{\mathbf{P}}$ . Therefore

$$(17) \quad \dot{\mathbf{P}} \supset \dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2.$$

Formally we expect the inclusion to be an equality. To see this, suppose

$$p(\eta) = p(\theta(\eta), G(\eta))$$

represents a curve through  $P_0$ . Differentiating pointwise formally we get

$$\begin{aligned} \frac{1}{p_0} \frac{\partial}{\partial \eta} p(\eta) |_{\eta=0} &= \frac{1}{p_0} \frac{\partial}{\partial \eta} p(\theta_0, G(\eta)) |_{\eta=0} \\ &+ \frac{\partial \theta(\eta)}{\partial \eta} \frac{1}{p_0} \frac{\partial}{\partial \theta} p(\theta, G_0) |_{\theta=\theta_0}. \end{aligned}$$

The right side formally belongs to  $\dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2$ , and so we expect that

$$(18) \quad \dot{\mathbf{P}} = \dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2.$$

In fact it is possible that  $\dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2$  is not closed if both  $\dot{\mathbf{P}}_1$  and  $\dot{\mathbf{P}}_2$  are infinite-dimensional; see section A.4. We can and do verify (18) and  $\dot{\mathbf{P}} = \dot{\mathbf{P}}^0$  for the symmetric location model.

**Example 4. The symmetric location model.**

We make  $P \longleftrightarrow (\theta, G)$  where  $\theta$  is the center of symmetry of  $P$  and  $G$  is the distribution of  $X - \theta$ , symmetric about zero. We restrict  $G$  so that its density  $G' = g$  is absolutely continuous with derivative  $g'$  and the Fisher information for location  $I_g$  is finite:

$$I_g = \int \frac{[g']^2}{g}(x) dx.$$

If  $s(\eta)$  is a curve passing through  $s_0$ , we write

$$s(\eta) = s_\eta(\cdot - \theta(\eta)).$$

Without loss of generality take  $\theta(0) = 0$  and let  $t = \dot{s}(0)$  so that

$$(19) \quad s_\eta(\cdot - \theta(\eta)) - s_0(\cdot) = \eta t + o(\eta).$$

By corollary A.5.1

$$\dot{S}_1 = [s'_0] = \dot{S}_1^0,$$

where  $s'_0 = -g'g^{-1/2}/2$ , and by example 2 and the second construction of example 1,

$$\begin{aligned} \dot{S}_2 &= \{t \in L_2(\mu) : t \text{ symmetric about } 0, t1_{[s_0 > 0]} = t, \langle t, s_0 \rangle = 0\} \\ &= \dot{S}_2^0 \supset \{t \in \dot{S}_2 : t \text{ absolutely continuous with derivative } t' \in L_2(\mu)\}. \end{aligned}$$

Of course,  $\mu$  is Lebesgue measure here. Now (18) is equivalent to  $\dot{S} \subset \dot{S}_1 + \dot{S}_2$ , which holds if,

$$(20) \quad t \in \dot{S}_1 + \dot{S}_2$$

for all  $t \in \dot{S}^0$ . Note that  $\dot{S}_1 + \dot{S}_2$  is closed since  $\dot{S}_1$  is one-dimensional.

**Proof of (20).** Now

$$(a) \quad \theta(\eta) \rightarrow 0$$

by weak convergence of  $P_{(\theta(\eta), G_\eta)}$  to  $P_{(0, G_0)}$ , and

$$(b) \quad s_\eta = s_0(\cdot + \theta(\eta)) + \eta t + o(\eta),$$

in the  $L_2$  sense. To see (b), note that

$$\begin{aligned} \|s_\eta - s_0(\cdot + \theta(\eta)) - \eta t\| & \\ &= \|s_\eta(\cdot - \theta(\eta)) - s_0 - \eta t(\cdot - \theta(\eta))\| \\ &\leq \|s_\eta(\cdot - \theta(\eta)) - s_0 - \eta t\| + \|\eta t\| \|\theta(\eta)\| \\ &= o(\eta) \end{aligned}$$

by (19), (a), and the  $L_2$ -continuity theorem A.1.1 applied to  $t$ . Also

$$(c) \quad \dot{S}_2^\perp = [s_0] + L_1 + L_2,$$

where

$$L_1 = \{h \in L_2(\mu) : h(x) = -h(-x) \text{ a.e. } \mu\},$$

$$L_2 = \{h \in L_2(\mu) : h1_{[s_0 > 0]} = 0\}.$$

Therefore, since

$$\Pi(h | L_1)(x) = \frac{1}{2}(h(x) - h(-x))$$

by (A.2.5) (note that  $h(x) = \Pi(h | L_1) = \frac{1}{2}(h(x) + h(-x)) \perp L_1$ ),

$$\begin{aligned} & \Pi\left(\frac{s_\eta(\cdot) - s_0(\cdot + \theta(\eta))}{\eta} | L_1\right) \\ &= \frac{s_\eta(\cdot) - s_\eta(-\cdot) - s_0(\cdot + \theta(\eta)) + s_0(-\cdot + \theta(\eta))}{2\eta} \\ (d) \quad &= -\frac{s_0(\cdot + \theta(\eta)) - s_0(\cdot - \theta(\eta))}{2\eta} \end{aligned}$$

by symmetry of  $s_\eta$  and  $s_0$ . By (b) and (d)

$$(e) \quad D_\eta \equiv \left\| \frac{s_0(\cdot + \theta(\eta)) - s_0(\cdot - \theta(\eta))}{2\eta} \right\| = O(\|t\|) = O(1).$$

But by corollary A.5.1 again

$$(f) \quad D_\eta = \frac{1}{2} \left| \frac{\theta(\eta)}{\eta} \right| I^{1/2}(G_0)(1 + o(1)),$$

and comparison of (e) and (f) yields  $|\theta(\eta)/\eta| = O(1)$ . Since  $t \perp s_0$  and  $t = t1_{[s_0 > 0]}$ , (c), (b), (d), (f), and example A.3.1 yield

$$\begin{aligned} \Pi(t | \dot{S}_2^\perp) &= \Pi(t | L_1) \\ &= -\frac{s_0(\cdot + \theta(\eta)) - s_0(\cdot - \theta(\eta))}{2\eta} + o(1) \\ &= -\frac{\theta(\eta)}{\eta} s'_0(\cdot) + o(1), \end{aligned}$$

and hence

$$\Pi(t | \dot{S}_2^\perp) \in \dot{S}_1.$$

Thus

$$t = \Pi(t | \dot{S}_2^\perp) + \Pi(t | \dot{S}_2) \in \dot{S}_1 + \dot{S}_2. \quad \square$$

### 3.3 INFORMATION BOUNDS VIA DERIVATIVES OF FUNCTIONS: THE NONPARAMETRIC APPROACH

In this section we show how to obtain information bounds and efficient influence functions for parameters which are explicitly defined as functions on  $\mathbf{P}$ . In the next section we consider parameters which are defined implicitly through a parametrization. The approaches are equivalent, but each has its own advantages depending on the example considered.

Suppose that  $\mathbf{P}$  is a regular parametric, semiparametric, or nonparametric model, and that  $v : \mathbf{P} \rightarrow R^m$  is a Euclidean parameter. Let  $v$  also denote the corresponding map from  $\mathbf{S}$  to  $R^m$  given by  $v(s(P)) = v(P)$ .

**Definition 1.** Let  $m = 1$ .  $v$  is *pathwise differentiable* on  $\mathbf{S}$  at  $s_0$  if there exists a bounded linear functional  $\dot{v}(s_0) : \dot{\mathbf{S}} \rightarrow R$  such that

$$(1) \quad v(s(\eta)) = v(s_0) + \eta \dot{v}(s_0)(t) + o(|\eta|)$$

for any curve  $s(\cdot)$  in  $\mathbf{S}$  which passes through  $s_0 = s(0)$  and has  $\dot{s}(0) = t$ . We define the bounded linear functional  $\dot{v}(P_0) : \dot{\mathbf{P}} \rightarrow R$  by

$$(2) \quad \dot{v}(P_0)(h) = \dot{v}(s_0)(\frac{1}{2} h s_0).$$

Then (1) holds if and only if

$$(3) \quad v(P_\eta) = v(P_0) + \eta \dot{v}(P_0)(h) + o(|\eta|),$$

where  $\{P_\eta\}$  is the curve in  $\mathbf{P}$  corresponding to  $\{s(\eta)\}$  in  $\mathbf{S}$  and  $h = 2t/s_0$ ; see (3.2.3). We call (3) *pathwise differentiability* of  $v$  on  $\mathbf{P}$  at  $P_0$ . For convenience in what follows we ignore the dependence of  $\dot{v}(P_0)$  on  $P_0$  and write  $\dot{v}(h)$  for  $\dot{v}(P_0)(h)$ , and likewise  $\dot{v}(t)$  for  $\dot{v}(s_0)(t)$ .

The functional  $\dot{v}$  is by (3) uniquely defined on  $\dot{\mathbf{P}}^0$  and hence on  $\dot{\mathbf{P}}$ . Often the parameter  $v$  is described most naturally as the restriction to  $\mathbf{P}$  of a parameter  $v_e$  defined on a larger model  $\mathbf{M}_0 \supset \mathbf{P}$ . Suppose  $v_e$  is pathwise differentiable on  $\mathbf{M}_0$  with derivative  $\dot{v}_e : \dot{\mathbf{M}}_0 \rightarrow R$ . Necessarily

$$(4) \quad \dot{v}_e = \dot{v} \quad \text{on } \dot{\mathbf{P}}.$$

By the Riesz representation theorem (see example A.1.8) there exists a unique element of  $\dot{\mathbf{P}}$  which, abusing notation we also call  $\dot{v}$ , such that

$$(5) \quad \dot{v}(h) = \langle \dot{v}, h \rangle_0$$

for all  $h \in \dot{\mathbf{P}}$ . Strictly speaking the element is  $\dot{v}^T(1)$  where  $\dot{v}^T : R \rightarrow \dot{\mathbf{P}}$  is the adjoint of  $\dot{v} : \dot{\mathbf{P}} \rightarrow R$ . If  $\dot{v}_e \in \dot{\mathbf{M}}_0$  is similarly defined, (4) implies that

$$(6) \quad \dot{v} = \Pi_0(\dot{v}_e | \dot{\mathbf{P}}).$$

This follows since for all  $h \in \dot{\mathbf{P}}$ ,

$$\begin{aligned} \langle \dot{v} - \Pi_0(\dot{v}_e | \dot{\mathbf{P}}), h \rangle_0 &= \langle \dot{v} - \dot{v}_e, h \rangle_0 \\ &= \dot{v}(h) - \dot{v}_e(h) = 0 \quad \text{by (4)}. \end{aligned}$$

Viewed as an element of  $L_2(P_0)$ ,  $\dot{v}_e$  is called a gradient of  $v$  by Koshevnik, Levit, and Pfanzagl, among others, and  $\dot{v}$  is the canonical gradient. Figure 1 gives the geometry.

Note that  $\dot{v}_e \in \dot{\mathbf{M}}$  implies

$$(7) \quad \langle \dot{v}_e, 1 \rangle_0 = 0.$$

Of course this is true for  $\dot{v}$  as well. As we shall see, the information bounds and efficient influence function are determined by the canonical gradient, but computation of  $\dot{v}$  is usually done via (6). We illustrate this with:

**Example 1. The symmetric location model.**

We use the assumptions and notation of example 3.2.4. The parameter  $v(P)$  we want to consider is the center of symmetry. There are many representations



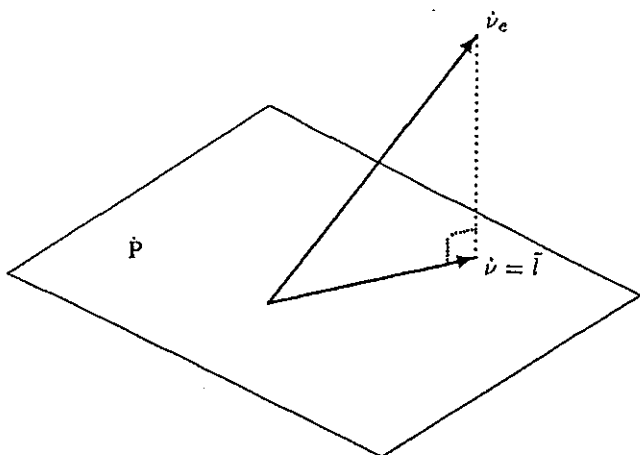


FIGURE 1. Projection of pathwise derivatives.

of  $v$  which can be thought of as parameters on larger models, for instance, the median of  $P$  defined as  $(F^{-1}(\frac{1}{2}+) + F^{-1}(\frac{1}{2}))/2$  where  $F$  is the distribution function corresponding to  $P$  and  $F^{-1}(t) \equiv \inf\{x : F(x) \geq t\}$ . A particularly convenient representation is given by

$$(8) \quad v(P) = \int_0^1 F^{-1}(t) w(t) dt,$$

where  $w$  is symmetric about  $1/2$ , is infinitely differentiable, vanishes off  $[1/4, 3/4]$ , say, and  $\int_0^1 w(s) ds = 1$ . Formula (8) evidently defines a parameter  $v_e$  on  $M_\mu$  with  $\mu$  Lebesgue measure which agrees with the center of symmetry on  $P$ .

We claim that  $v_e$  is pathwise differentiable on  $M_\mu$ , and at  $P_0$  has gradient

$$(9) \quad \dot{v}_e(x) = - \int_0^1 (1_{[u \geq F_0(x)]} - u) w(u) dF_0^{-1}(u).$$

To see this, write

$$v_e(P) = \int_{-\infty}^{\infty} x w(F(x)) dF(x).$$

Then

$$\begin{aligned} v_e(P_\eta) - v_e(P_0) &= \int_{-\infty}^{\infty} x w(F_0(x)) (p_\eta(x) - p_0(x)) dx \\ &\quad + \int_{-\infty}^{\infty} x (w(F_\eta(x)) - w(F_0(x))) dF_0(x) \\ &\quad + \int_{-\infty}^{\infty} x (w(F_\eta(x)) - w(F_0(x))) (p_\eta(x) - p_0(x)) dx. \end{aligned}$$

If  $s_\eta = s_0 + \eta s_0 h/2 + o(\eta)$  for  $h \in L_2(P_0)$ , then

$$\int |p_\eta - p_0 - \eta h p_0|(x) dx = o(\eta)$$

and

$$\sup_x |F_\eta(x) - F_0(x) - \eta \int_{-\infty}^x h(y) p_0(y) dy| = o(\eta).$$

Since  $w(F_0(x))$  vanishes outside  $\{x : 1/4 \leq F_0(x) \leq 3/4\}$  and  $\sup_x |F_\eta(x) - F_0(x)| \rightarrow 0$  it is easy to see that

$$\begin{aligned} v_e(P_\eta) - v_e(P_0) &= \eta \left\{ \int_{-\infty}^{\infty} x \left[ w(F_0(x)) h(x) \right. \right. \\ &\quad \left. \left. + w'(F_0(x)) \left( \int_{-\infty}^x h(y) p_0(y) dy \right) \right] p_0(x) dx \right\} \\ &\quad + o(\eta^2) \\ &= -\eta \int_{-\infty}^{\infty} \left( \int_x^{\infty} w(F_0(y)) dy \right) h(x) p_0(x) dx + O(\eta^2). \end{aligned}$$

We deduce (recalling our convention given in (7)) that

$$\dot{v}_e(x) = - \int_x^{\infty} w(F_0(y)) dy + E_0 \left( \int_x^{\infty} w(F_0(y)) dy \right)$$

and (9) follows. If  $P_0 \in \mathbf{P}$ , say  $F_0(x) = G_0(x - \theta)$ , then it is easy to see (using  $\int (1 - 2u) w(u) du = 0$ ) that

$$(10) \quad \dot{v}_e(\theta + x) = -\dot{v}_e(\theta - x).$$

For simplicity and without loss of generality suppose  $\theta = 0$ . Then from example 3.2.4 we know that

$$\dot{\mathbf{P}} = \begin{bmatrix} g'_0 \\ g_0 \end{bmatrix} + \dot{\mathbf{P}}_2,$$

where

$$\dot{\mathbf{P}}_2 = \{h \in L_2(G_0) : h \text{ symmetric about } 0, \int h dG_0 = 0\}.$$

Let  $\psi_0 \equiv -g'_0/g_0$ . By (10),  $\dot{v}_e \perp \dot{\mathbf{P}}_2$ , and hence

$$(11) \quad \dot{v} = \Pi_0(\dot{v}_e | [\psi_0]) = \frac{\langle \dot{v}_e, \psi_0 \rangle_0}{\|\psi_0\|_0^2} \psi_0.$$

A simple calculation yields  $\langle \dot{v}_e, \psi_0 \rangle_0 = 1$ , in agreement, as we shall see, with (3.4.5).  $\square$

Now suppose  $\mathbf{v} = (v_1, \dots, v_m)^T$  is  $m$ -dimensional. If each  $v_i$  is pathwise differentiable with derivative  $\dot{v}_i = \dot{v}_i(P_0)$ , we call  $\mathbf{v}$  *pathwise differentiable on  $\mathbf{P}$  at  $P_0$*  with derivative  $\dot{\mathbf{v}} = (\dot{v}_1, \dots, \dot{v}_m)^T$ . The remarks we have just made extend obviously to  $m > 1$ .

For a regular  $k$ -dimensional parametric model  $\mathbf{Q} = \{P_\theta : \theta \in \Theta\}$ , pathwise differentiability of  $\mathbf{v}(P_\theta) = q(\theta) : R^k \rightarrow R^m$  at  $P_0 = P_{\theta_0}$  holds if  $q$  is differentiable at  $\theta_0$  with total derivative matrix  $\dot{q}_{m \times k}$ . To see this, note that  $\dot{\mathbf{Q}} = [\dot{\mathbf{i}}(\theta_0)] \equiv [\dot{\mathbf{I}}]$ , so all  $h \in \dot{\mathbf{Q}}$  are of the form  $h = a^T \dot{\mathbf{I}}$  for some  $a \in R^k$ .

But with  $p_0(\eta) = p(\theta_0 + \eta a)$ ,  $a \in R^k$ ,  $p_0(\eta)$  has tangent  $a^T \dot{\mathbf{i}}$ , and since

$$(12) \quad \langle I^{-1}(\theta_0) \dot{\mathbf{i}}, \dot{\mathbf{i}}^T \rangle_0 = J \equiv \text{identity matrix},$$

it follows that

$$\begin{aligned} v(P_0(\eta)) &= q(\theta_0 + \eta a) \\ &= q(\theta_0) + \eta \dot{q}(\theta_0) a + o(|\eta|) \\ (13) \quad &= v(P_0) + \eta \langle \dot{q}(\theta_0) I^{-1}(\theta_0) \dot{\mathbf{i}}, \dot{\mathbf{i}}^T a \rangle_0 + o(|\eta|). \end{aligned}$$

Comparison of (3), (5), and (13) yields

$$(14) \quad \dot{v} = \dot{q}(\theta_0) I^{-1}(\theta_0) \dot{\mathbf{i}}(\theta_0) = \tilde{\mathbf{I}}(\cdot, P_0 | v, \mathbf{Q}),$$

the efficient influence function for  $v$  in  $\mathbf{Q}$  as defined in (2.3.2), and hence

$$(15) \quad \langle \dot{v}, \dot{v}^T \rangle_0 = \dot{q}(\theta_0) I^{-1}(\theta_0) \dot{q}^T(\theta_0) = I^{-1}(P_0 | v, \mathbf{Q}),$$

the information bound for  $v$  in  $\mathbf{Q}$ . (Note that  $\dot{v}^T$  here is a row vector of elements of  $L_2(P_0)$  rather than the adjoint of the mapping  $\dot{v}: \dot{\mathbf{P}} \rightarrow R^m$ .)

**Theorem 1.** Let  $v$  be pathwise differentiable on  $\mathbf{P}$  at  $P_0$  with  $m = 1$ . For any regular parametric submodel  $\mathbf{Q}$  for which  $I^{-1}(P_0 | v, \mathbf{Q})$  is defined, the efficient influence function  $\tilde{\mathbf{I}}(\cdot, P_0 | v, \mathbf{Q})$  satisfies

$$(16) \quad \tilde{\mathbf{I}}(\cdot, P_0 | v, \mathbf{Q}) = \Pi_0(\dot{v}(P_0) | \dot{\mathbf{Q}}),$$

and consequently

$$(17) \quad I^{-1}(P_0 | v, \mathbf{Q}) = \|\Pi_0(\dot{v}(P_0) | \dot{\mathbf{Q}})\|_0^2 \leq \|\dot{v}(P_0)\|_0^2$$

with equality if and only if

$$(18) \quad \dot{v}(P_0) \in \dot{\mathbf{Q}}.$$

**Proof.** Fix  $P_0 \in \mathbf{P}$ . For any regular  $k$ -dimensional parametric submodel  $\mathbf{Q} = \{Q_\eta: \eta \in H\}$  with  $Q_0 = P_0 \in \mathbf{P}$  for which  $I^{-1}(P_0 | v, \mathbf{Q})$  is defined,  $\eta \rightarrow v(s(\eta))$  is differentiable at  $\eta = 0$ , and by the argument leading to (14), the efficient influence function for  $v$  in  $\mathbf{Q}$  is

$$(a) \quad \tilde{\mathbf{I}}(\cdot, P_0 | v, \mathbf{Q}) = \Pi_0(\dot{v}(P_0) | \dot{\mathbf{Q}}).$$

Equation (16) has been proved, and (17) follows from it by Pythagoras since  $\dot{\mathbf{Q}} \subset \dot{\mathbf{P}} \subset L_2(P_0)$ .  $\square$

If  $\mathbf{Q}$  is a regular parametric submodel of  $\mathbf{P}$  satisfying (18) and hence giving equality in (17), we call  $\mathbf{Q}$  *least favorable*. It follows immediately from theorem 1 that under the condition  $\dot{v}(P_0) \in \dot{\mathbf{P}}^0$  there exists a least favorable  $\mathbf{Q}$  which may be chosen one-dimensional.

These results extend to  $m \geq 1$  as follows. Let

$$(19) \quad \begin{aligned} \tilde{\mathbf{I}}_j &= \dot{v}_j(P_0), \quad j = 1, \dots, m, \\ \tilde{\mathbf{I}} &= (\tilde{\mathbf{I}}_1, \dots, \tilde{\mathbf{I}}_m)^T. \end{aligned}$$

**Corollary 1.** Suppose that:

- (i)  $v$  is pathwise differentiable on  $\mathbf{P}$  at  $P_0$ .
- (ii)  $\mathbf{Q}$  is any regular parametric submodel for which  $I^{-1}(P_0 | v, \mathbf{Q})$  is defined.

Then the efficient influence function  $\tilde{\mathbf{I}}(\cdot, P_0 | v, \mathbf{Q})$  for  $v$  in  $\mathbf{Q}$  satisfies (16) and, consequently,

$$(20) \quad I^{-1}(P_0 | v, \mathbf{Q}) = \langle \Pi_0(\dot{v}(P_0) | \dot{\mathbf{Q}}), \Pi_0(\dot{v}^T(P_0) | \dot{\mathbf{Q}}) \rangle_0 \\ \leq \langle \tilde{\mathbf{I}}, \tilde{\mathbf{I}}^T \rangle_0$$

in the order  $A \leq B$  if  $B - A$  is nonnegative definite. Equality holds if and only if

$$(21) \quad \tilde{\mathbf{I}}_j \in \dot{\mathbf{Q}}, \quad j = 1, \dots, m,$$

and, in that case, the efficient influence function  $\tilde{\mathbf{I}}(\cdot, P_0 | v, \mathbf{Q})$  equals  $\tilde{\mathbf{I}}$ . Furthermore, if

$$(iii) \quad [\tilde{\mathbf{I}}_j : j = 1, \dots, m] \subset \dot{\mathbf{P}}^0,$$

then for each  $a \in R^m$  there exists a one-dimensional regular parametric submodel  $\mathbf{Q}$  such that

$$(22) \quad a^T I^{-1}(P_0 | v, \mathbf{Q}) a = a^T E(\tilde{\mathbf{I}} \tilde{\mathbf{I}}^T) a.$$

**Proof.** Apply theorem 1 to  $v(P) \equiv \sum_{i=1}^m a_i v_i(P)$  and note that  $\Pi_0(\sum_{i=1}^m a_i \dot{v}_i(P_0) | \dot{\mathbf{Q}}) = \sum_{i=1}^m a_i \Pi_0(\dot{v}_i(P_0) | \dot{\mathbf{Q}})$ . □

A least favorable model in the sense of corollary 1 will exist under our assumptions if any  $m$  curves intersecting at  $s_0$  are all contained in an  $m$ -dimensional surface passing through  $s_0$  at least in a neighborhood of  $s_0$ . If a least favorable submodel exists, there is an  $m$ -dimensional least favorable submodel.

As we have seen in theorem 1, condition (iii) yields, for  $m = 1$ , a least favorable parametric model  $\mathbf{Q}$  in which the information and efficient influence function for  $v$  can naturally be identified with the corresponding quantities for  $\mathbf{P}$ . We extend this idea for general  $m$  to situations in which no least favorable  $\mathbf{Q}$  exists, but a least favorable sequence  $\{\mathbf{Q}_j\}$  does.

Suppose that  $v$  is pathwise differentiable as in theorem 1 and corollary 1, and

$$(iii) \quad [\dot{v}_j(P_0) : j = 1, \dots, m] \subset \overline{\dot{\mathbf{P}}^0}.$$

Under these conditions we define efficient influence functions and the information matrix for  $v$ .

**Definition 2.** We call  $\tilde{\mathbf{I}} = \tilde{\mathbf{I}}(\cdot, P_0 | v, \mathbf{P})$  defined by (19) the *efficient influence function* for  $v$ :

$$(23) \quad \tilde{\mathbf{I}} = \tilde{\mathbf{I}}(\cdot, P_0 | v, \mathbf{P}) = \dot{v}(P_0).$$

The information  $I(P_0 | v, P)$  for  $v$  in  $P$  is defined as the inverse of

$$(24) \quad I^{-1}(P_0 | v, P) = \langle \tilde{I}, \tilde{I}^T \rangle_0 = E_0(\tilde{\Pi}^T).$$

**Definition 3.** If  $\hat{v}_n$  is a locally Gaussian regular (on  $P$ ) estimate of  $v$  with  $\Sigma(P_0, \hat{v}_n) = I^{-1}(P_0 | v, P)$  we say  $\hat{v}_n$  is *efficient at  $P_0$* . If  $\hat{v}_n$  is efficient at all regular  $P_0$ , then we say  $\hat{v}_n$  is *efficient*.

If an efficient  $\hat{v}_n$  exists, its influence function must be  $\tilde{I}$ , and conversely any locally Gaussian regular and linear estimate with influence function  $\tilde{I}$  is efficient. See also proposition 1 below.

**Corollary 2.** Suppose that conditions (i) and  $\overline{\text{(iii)}}$  hold for  $v: P \rightarrow R^m$ . Let  $q: R^m \rightarrow R^d$  and assume that  $q$  is continuously differentiable with derivative  $\dot{q}$  (at  $v(P_0)$ ). Then  $q(v(P))$  also satisfies (i) and  $\overline{\text{(iii)}}$  and has influence function

$$(25) \quad \tilde{I}(\cdot, P_0 | q(v), P) = \dot{q}(v(P_0))\tilde{I}(\cdot, P_0 | v, P).$$

**Proof.** Pathwise differentiability of  $q(v(s))$  with derivative  $\dot{q}(v(s_0))\dot{v}(s_0)$  follows from the chain rule for Fréchet derivatives, and the conclusion follows from theorem 1.  $\square$

Note that corollary 2 is just a generalization of (2.3.2). Also, note that  $\overline{\text{(iii)}}$  holds for *all*  $v$  satisfying (i) if  $\overline{P^0} = \dot{P}$ , or, equivalently (see proposition 3.2.2), if  $\overline{P^0}$  is a linear space.

The following theorem justifies definitions 2 and 3.

**Theorem 2.** Suppose that  $v: P \rightarrow R^m$  is pathwise differentiable on  $P$  at  $P_0$ ,  $\overline{\text{(iii)}}$  holds, and  $T_n$  is a locally regular estimate of  $v$  with corresponding limit law  $L_{P_0}$ . Then:

- A.  $L_{P_0}$  can be represented as the convolution of a  $N(0, I^{-1}(P_0 | v, P))$  distribution with another distribution on  $R^m$ . More generally,

$$(26) \quad L_{P_0} \left( \begin{array}{c} \sqrt{n}(T_n - v(P_0)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i, P_0 | v, P) \\ \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i, P_0 | v, P) \end{array} \right) \rightarrow L \left( \begin{array}{c} \Delta_0 \\ Z_0 \end{array} \right)$$

where  $Z_0 \sim N(0, I^{-1}(P_0 | v, P))$  and  $\Delta_0$  is independent of  $Z_0$ . More generally still, if  $h \in (\overline{P^0})'$  then

$$(27) \quad L_{P_0} \left( \begin{array}{c} \sqrt{n}(T_n - v(P_0)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i, P_0 | v, P) \\ \frac{1}{\sqrt{n}} \sum_{i=1}^n h(X_i) \end{array} \right) \rightarrow L \left( \begin{array}{c} \Delta_0 \\ W_0 \end{array} \right)$$

where  $W_0 \sim N(0, Ehh^T)$  and  $\Delta_0$  is independent of  $W_0$ .

B.  $\Delta_0 = 0$  if and only if  $T_n$  is asymptotically linear with influence function  $\tilde{I}$ , and then  $T_n$  is efficient.

The extension (27) of (26) was noticed by Pfanzagl (1989).

**Proof.** Note that if  $T_n$  is locally regular on  $\mathbf{P}$ , it is locally regular on any one-dimensional submodel  $\mathbf{Q}$ . Suppose  $\dot{I}$  is the score function for such a model  $\{P_\gamma : |\gamma| < 1\}$ ; that is,

$$(a) \quad s(P_\gamma) = s(P_0) + \frac{1}{2} \gamma \dot{I} s_0 + o(|\gamma|).$$

Let  $v(P_\gamma) \equiv q(\gamma)$ . By (i),

$$(b) \quad q(\gamma) = q(0) + \gamma E(\tilde{I}\dot{I}) + o(\gamma),$$

where  $\tilde{I}$  is defined by (19). Here and in the sequel, expectations are under  $P \equiv P_0$ .

Let

$$(c) \quad U_n = \sqrt{n} (T_n - v(P)),$$

$$(d) \quad V_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n \dot{I}(X_i).$$

Then, by theorem 2.3.1,

$$\mathbf{L}(U_n, V_n) \rightarrow \mathbf{L}(U, V)$$

and using (a) and (b), identity (d) in the proof of theorem 2.3.1 becomes, for all  $a \in R^m, t \in R$ ,

$$(e) \quad E \exp[ia^T U + tV - \frac{1}{2} t^2 E(\dot{I}^2)] = \exp[ia^T E(\tilde{I}\dot{I})t] E \exp[ia^T U].$$

By (iii), for each  $b \in R^m, c \in R^l$ , there exists a sequence  $\{\dot{I}_{bcj}\}$  of score functions of regular one-dimensional models such that

$$(f) \quad \|b^T \dot{I}^* + c^T h - \dot{I}_{bcj}\|_0 \rightarrow 0 \quad \text{as } j \rightarrow \infty,$$

where

$$(g) \quad \dot{I}^* \equiv I(P | v, \mathbf{P}) \tilde{I}.$$

Let

$$V_{nj} \equiv \frac{1}{\sqrt{n}} \sum_{i=1}^n \dot{I}_{bcj}(X_i)$$

and let  $\mathbf{L}(U, V_j)$  denote the limit of  $\mathbf{L}(U_n, V_{nj})$ . Now let

$$V_n = n^{-1/2} \sum_{i=1}^n \dot{I}^*(X_i), \quad W_n = n^{-1/2} \sum_{i=1}^n h(X_i).$$

Suppose, for a subsequence which we identify with  $\{n\}$ ,  $\mathbf{L}(U_n, V_n, W_n) \rightarrow \mathbf{L}(U, V, W)$ . Such a subsequence always exists since  $\{U_n\}, \{V_n\}, \{W_n\}$  are tight. Since

$$\sup_n E (b^T V_n + c^T W_n - V_{nj})^2 = \|b^T \Gamma^* + c^T h - \dot{\mathbf{i}}_{bcj}\|_0^2 \rightarrow 0$$

as  $j \rightarrow \infty$ , we must have

$$L(U, b^T V + c^T W) = \lim_{j \rightarrow \infty} L(U, V_j).$$

Take  $\dot{\mathbf{i}} = \dot{\mathbf{i}}_{bcj}$ ,  $t = 1$ ,  $V = V_j$  in (e). Since

$$\begin{aligned} E e^{V_j} &= \exp\left[\frac{1}{2} E(\dot{\mathbf{i}}_{bcj}^2)\right] \rightarrow \exp\left[\frac{1}{2} E(b^T \Gamma^* + c^T h)^2\right] \\ &= E e^{b^T V + c^T W}, \end{aligned}$$

we can pass to the limit as  $j \rightarrow \infty$  to get

$$\begin{aligned} \text{(h)} \quad E \exp\left\{ia^T U + b^T V + c^T W - \frac{1}{2} b^T I b - b^T E \Gamma^* h^T c - \frac{1}{2} c^T E h h^T c\right\} \\ = \exp[ia^T b + ia^T E \tilde{\mathbf{l}} h^T c] E \exp[ia^T U], \end{aligned}$$

where  $I \equiv I(P_0 | \mathbf{v}, \mathbf{P})$ . Use analytic continuation in (h) and take  $b = -iI^{-1}a$ ,  $c = id$  to get, after some algebra,

$$\begin{aligned} \text{(i)} \quad E \exp[ia^T(U - I^{-1}V) + id^T W] \\ = E \exp[ia^T U + \frac{1}{2} a^T I^{-1} a] \exp[-\frac{1}{2} d^T E h h^T d]. \end{aligned}$$

Since

$$\begin{aligned} L(U - I^{-1}V, W) \\ = \lim_n L(\sqrt{n}(T_n - \mathbf{v}(P_0)) - n^{-1/2} \sum_{i=1}^n \tilde{\mathbf{l}}(X_i), n^{-1/2} \sum_{i=1}^n h(X_i)), \end{aligned}$$

the various parts of the theorem follow from (i) in the same manner as theorem 2.3.1 followed from (e) of that proof.  $\square$

Using theorem 2 we strengthen corollary 1 along the lines of proposition 2.4.3.

**Proposition 1.** Suppose that  $T_n$  is an asymptotically linear estimate of  $\mathbf{v}: \mathbf{P} \rightarrow R^m$  at  $P_0$ , with influence function  $\psi$ . Then:

A.  $T_n$  is regular at  $P_0$  if and only if  $\mathbf{v}$  is pathwise differentiable with derivative  $\dot{\mathbf{v}}$  and

$$(28) \quad \psi - \dot{\mathbf{v}} = \psi - \tilde{\mathbf{l}} \perp \dot{\mathbf{P}}.$$

B. Let  $T_n$  be regular at  $P_0$  and suppose that  $\overline{\text{(iii)}}$  holds. Then  $T_n$  is efficient at  $P_0$  if and only if  $\psi \in \dot{\mathbf{P}}^m$ , and then  $\psi = \tilde{\mathbf{l}}$ .

The geometrical content of this proposition is apparent from figure 2.

**Remark 1.** Suppose  $T_n = \mathbf{v}_e(IP_n)$  for an extension  $\mathbf{v}_e: \mathbf{M}_0 \rightarrow R$  where

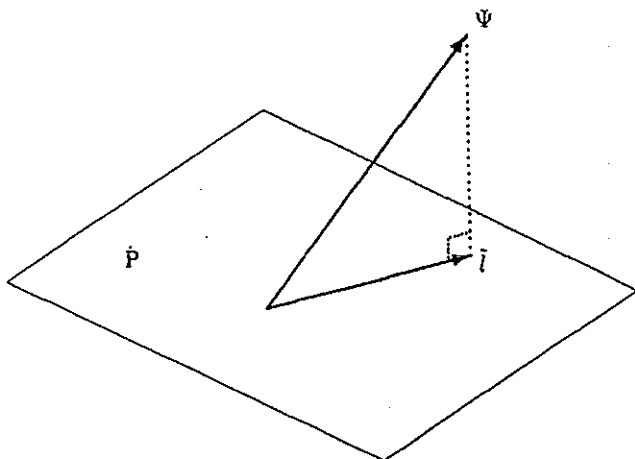


FIGURE 2. Projection of influence functions.

$M_0 \supset P \cup \{\text{discrete distributions}\}$ . If  $v_e$  is pathwise differentiable and Gâteaux differentiable on  $M_0$ , as in (2.2.3), a formal computation suggests that  $T_n$  is asymptotically linear with influence function  $\psi$  such that

$$\dot{v}_e = \psi.$$

By proposition 1, this implies that  $T_n$  is regular.

**Proof of proposition 1.** The idea for the proof of A is the same as that for the proof of proposition 2.4.3.A. Let  $\{P_\eta\}$  be an  $l$ -dimensional regular parametric submodel with score function  $h \in (\dot{P}^0)'$  at  $\eta = 0$ . Let  $P_n = P_{\eta_n}$ ,  $\eta_n = t_n/\sqrt{n}$ ,  $t_n \rightarrow t$ . By the asymptotic linearity of  $\{T_n\}$  at  $P_0$  and proposition 2.1.2,

$$(a) \quad L_{P_0} \begin{pmatrix} \sqrt{n}(T_n - v(P_0)) \\ L_n(\eta_n) - L_n(0) \end{pmatrix} \rightarrow N \left( \begin{pmatrix} 0 \\ -\Sigma_{22}/2 \end{pmatrix}, \Sigma \right),$$

with

$$(b) \quad \Sigma = [\Sigma_{ij}], \quad \Sigma_{11} = E\psi\psi^T, \quad \Sigma_{12} = E\psi h^T t, \quad \Sigma_{22} = t^T E h h^T t.$$

Consequently, by Le Cam's third lemma (lemma A.9.3)

$$(c) \quad L_{P_n}(\sqrt{n}(T_n - v(P_0))) \rightarrow N(\Sigma_{12}, \Sigma_{11}).$$

Assume now that  $T_n$  is regular. Then

$$(d) \quad L_{P_n}(\sqrt{n}(T_n - v(P_n))) \rightarrow N(0, \Sigma_{11}),$$

and from (c) and (d) we conclude that

$$(e) \quad \sqrt{n}(v(P_n) - v(P_0)) \rightarrow \Sigma_{12} = E\psi h^T t$$

or

$$v(P_{\eta_n}) = v(P_0) + \langle \psi, h^T \rangle_0 \eta_n + o(|\eta_n|)$$



for any sequence  $\eta_n$  as above. Now take  $l = 1$ . Given any sequence of reals,  $\tilde{\eta}_m \rightarrow 0$ , if we take  $n_m = \lceil \tilde{\eta}_m^{-2} \rceil$ , the integer part of  $\tilde{\eta}_m^{-2}$ , then  $\tilde{\eta}_m = c_m / \sqrt{n_m}$  where  $c_m \rightarrow 1$ . This implies that  $v$  is pathwise differentiable with derivative  $\dot{v}$  satisfying

$$\dot{v}(h) = \langle \psi, h \rangle_0$$

for all  $h \in \dot{P}^0$ . Consequently,  $\dot{v}$  can be defined on all of  $\dot{P}$  by

$$(f) \quad \dot{v}(h) = \langle \psi, h \rangle_0 = \langle \dot{v}, h \rangle_0 \quad \text{with } \dot{v} \in \dot{P}.$$

But (f) implies (28).

On the other hand, if  $v$  is pathwise differentiable and (28) holds, then (e) is valid, which together with (c) implies (d). This completes the proof of A.

For the proof of B, note that A implies that (i) holds. First suppose that  $T_n$  is efficient at  $P_0$ . Then definition 3 and theorem 2.A imply that  $\Delta_0 = 0$ , and hence  $T_n$  has influence function  $\psi = \tilde{I}$  by theorem 2.B.

Now suppose that  $\psi \in \dot{P}^m$ . Then  $\psi - \tilde{I} \in \dot{P}^m$ . Together with (28) this yields  $\psi = \tilde{I}$ , and hence efficiency of  $T_n$  at  $P_0$ .  $\square$

If we examine the proof of this convolution theorem 2, we see that we can weaken the definition of  $\dot{P}^0$ ,  $\dot{S}^0$ , and the corresponding  $\dot{P}$ ,  $\dot{S}$ , and still obtain the conclusions of the theorem. Specifically, let

$$\dot{P}_w^0 = \{h \in L_2(P_0) : r(\gamma) = \gamma h + o(\gamma)\}$$

for some mapping  $\gamma \rightarrow P_\gamma$ ,  $|\gamma| < 1$ , and  $r(\cdot)$  as defined in (3.2.2). That is, we replace the requirement of continuous Fréchet differentiability of  $\gamma \rightarrow r(\gamma)$  by Fréchet differentiability at  $\gamma = 0$ . Evidently  $\dot{P}_w^0 \supset \dot{P}^0$ , and if we define  $\dot{P}_w = [\dot{P}_w^0]$ , we have  $\dot{P}_w \supset \dot{P}$ . If we replace  $\dot{P}$  by  $\dot{P}_w$  in (19), (23), and so forth, theorem 2 continues to hold.

Here are three examples in which we exhibit efficient estimates.

**Example 2. Estimation of the mean,  $P$  unconstrained.**

Suppose that  $\mu$  concentrates on an interval  $[-M, M]$ . We want to estimate

$$v(P) = \int x dP(x).$$

Identify  $v$  on  $S$  as

$$v(s) = \int x s^2(x) d\mu(x).$$

Fix  $s_0$ . Then  $v$  is pathwise differentiable, and using (7) in the form  $\dot{v}(s_0) \perp s_0$  we obtain

$$(29) \quad \dot{v}(s_0)(x) = 2s_0(x)(x - E_0 X) \quad \text{a.e. } \mu.$$

To see (29), check

$$(30) \quad \begin{aligned} v(s) - v(s_0) - \langle \dot{v}, s - s_0 \rangle &= \int (x - E_0 X)(s - s_0)^2(x) d\mu(x) \\ &= O(\|s - s_0\|^2). \end{aligned}$$

Hence by (3.2.8), (23), and (29) it follows that

$$\tilde{I}(x, P_0) = x - E_0 X.$$

Not surprisingly, the sample mean is efficient:

$$\Sigma(P_0, \bar{X}) = \int \tilde{I}^2(x) dP_0(x).$$

More generally, suppose that  $\mathbf{P}$  satisfies

$$(31) \quad \sup_{P \in \mathbf{P}} E_P X^2 < \infty.$$

Then  $v(P) \equiv \int x dP(x) = E_P X$  is pathwise differentiable on  $\mathbf{P}$  with  $\dot{v}(s_0)(x) = 2s_0(x)(x - E_0 X)$  or  $\dot{v}(P_0) = x - E_0 X = \tilde{I}(x, P_0)$  by proposition A.5.2. The mean continues to be efficient and locally regular in view of Le Cam's third lemma A.9.3.  $\square$

**Example 3. Estimation of the mean,  $\mathbf{P}$  constrained.**

Suppose that  $\mathbf{X} = [-M, M]$ ,  $v$  is as in example 1, but now as in example 3.2.3, we constrain  $\mathbf{P}$  by fixing the coefficient of variation at  $\sqrt{c_0}$ . This is equivalent to

$$\gamma(P) = \int x^2 dP(x) - (1 + c_0)v^2(P) = 0.$$

The mean  $v$  is, of course, pathwise differentiable on the smaller tangent space of this model, and  $\dot{v}(P_0)$  is given by (6), where  $\dot{v}_e(P_0)(x) = x - E_0 X$  by example 1. Check that

$$\dot{\gamma}(P)(x) = (x^2 - \int x^2 dP(x)) - 2(1 + c_0)v(P)(x - v(P)),$$

where  $P \rightarrow \dot{\gamma}(P)$  is continuous at  $P_0$ . We can apply example 3.2.3, which shows that  $\dot{\mathbf{P}} = \{h \in L_2^0(P) : h \perp \dot{\gamma}(P)\}$ . Thus (23), (3.2.16), and algebra yield

$$\begin{aligned} (32) \quad \tilde{I}(x, P_0) &= \Pi_0(\dot{v}_e \mid \dot{\mathbf{P}}) \\ &= \dot{v}_e(P_0) - \frac{\langle \dot{v}_e(P_0), \dot{\gamma}(P_0) \rangle_0}{\|\dot{\gamma}(P_0)\|_0^2} \dot{\gamma}(P_0) \\ &= [1 + 2(1 + c_0)(E_0 X)a(P_0)](x - E_0 X) \\ &\quad - a(P_0)(x^2 - E_0 X^2) \end{aligned}$$

where

$$(33) \quad a(P) \equiv \frac{E(\dot{v}(P)\dot{\gamma}(P))}{E(\dot{\gamma}^2(P))} = \frac{\text{Cov}[X - EX, X^2 - 2(1 + c_0)(EX)X]}{\text{Var}[X^2 - 2(1 + c_0)(EX)X]}$$

Consider the estimate

$$(34) \quad \hat{v}_n = \bar{X} + a(IP_n)\gamma(IP_n),$$

where  $IP_n$  is the empirical df. It is straightforward to show that  $\hat{v}_n$  is uniformly Gaussian regular and efficient. For this model,  $\hat{v}_n$  improves  $\bar{X}$  :

$$\begin{aligned} \Sigma(P_0, \hat{v}) &= \Sigma(P_0, \bar{X}) - a^2(P_0)Var_0[X^2 - 2(1+c_0)(E_0 X)X] \\ &= I^{-1}(P_0 | v, P). \end{aligned}$$

Estimates such as (34) were introduced by Levit (1975). Alternative procedures are discussed in Haberman (1984) and Sheehy (1987), (1988).

Note that  $\hat{v}_n$  can be identified with  $\psi(IP_n)$  where

$$(35) \quad \psi(P) \equiv \int x dP + a(P)\gamma(P).$$

We can think of  $\psi(P)$  as being a rather unobvious extension of the mean from our constrained  $\mathbf{P}$  to  $\mathbf{M}$ . It is the "right" extension in the sense that  $\dot{\psi} \in \dot{\mathbf{P}}$ .  $\square$

We can generalize example 1 considerably.

**Example 4.**  $d_K$ -differentiable parameters on  $\mathbf{M} = \{\text{all } P \text{ on } \mathbf{X}\}$ .

Suppose that  $\mathbf{X} \subset R$  and that  $v$  is a parameter on  $\mathbf{M}$  satisfying the following regularity conditions:

- (i) For all  $P_0 \in \mathbf{M}$ ,  $v$  is continuously Fréchet differentiable at  $P_0$  with respect to  $d_K$  given by (A.6.8).
- (ii) For all  $P_0 \in \mathbf{P}$ , the derivative  $\dot{v}$  has the representation

$$\dot{v}(P_0)(P) = \int \psi(x, P_0) dP(x),$$

where  $\psi$  is continuous, bounded in  $x$ , continuous in  $P_0$  with respect to the total variation metric  $d_v$ , and

$$\int \psi(x, P_0) dP_0(x) = 0.$$

If  $\mathbf{X} = [-M, M]$  the mean has these properties as well as the functional  $\psi(P)$  of (35).

Let  $IP_n$  be the empirical distribution. Then  $v(IP_n)$  is an efficient estimate of  $v$ . To see this, note that by (i) and (ii)

$$(36) \quad v(IP_n) = v(P) + n^{-1} \sum_{i=1}^n \psi(X_i, P) + o_p(n^{-1/2}).$$

Since  $\psi$  is continuous, bounded, and continuous in  $P$ , the conditions of proposition 2.2.1.D hold, and  $v(IP_n)$  is regular. By (3.2.9),  $\psi(\cdot, P_0) \in \dot{\mathbf{M}}$ , and efficiency follows. Note that we do not really need (i) and (ii) for  $v(IP_n)$  to be efficient here, but merely  $v(IP_n)$  asymptotically linear.  $\square$

### 3.4 INFORMATION BOUND CALCULATIONS VIA SCORES: THE SEMIPARAMETRIC APPROACH

In semiparametric models parameters are defined implicitly through the parametrization rather than directly as functions on  $\mathbf{P}$ . Part A of proposition 2.4.1 generalizes to such models and gives a way of computing information and influence functions. It is convenient to present this approach in terms of  $\mathbf{P}$  rather than  $\mathbf{S}$ . Let

$$(1) \quad \mathbf{P} = \{P_{(v,G)} : v \in N, G \in \mathbf{G}\},$$

where  $N$  is an open subset of  $R^m$  and  $\mathbf{G}$  is general. Fix  $P_0 = P_{(v_0, G_0)}$ . In agreement with chapter 2, let  $\dot{\mathbf{I}}_1$  denote the vectors of partial derivatives of  $\log p(\cdot, v, G)$  with respect to  $v$  evaluated at  $P_0$ , and this is just the score function for the parametric model

$$\mathbf{P}_1 = \{P_{(v, G_0)} : v \in N\}.$$

We also let

$$\mathbf{P}_2 = \{P_{(v_0, G)} : G \in \mathbf{G}\}$$

and follow a similar convention for parametric submodels  $\mathbf{Q}$ .

**Theorem 1.** Suppose that:

- (i)  $\mathbf{P}_1$  is regular.
- (ii)  $v$  is a 1-dimensional parameter.

Let

$$(2) \quad \mathbf{I}_1^* = \dot{\mathbf{I}}_1 - \Pi_0(\dot{\mathbf{I}}_1 \mid \dot{\mathbf{P}}_2).$$

Then:

A. If  $\mathbf{Q} = \{P_{(v, G_\gamma)} : v \in N, \gamma \in \Gamma\}$  is a regular parametric submodel of  $\mathbf{P}$  with  $P_0 \in \mathbf{Q}$ , then

$$(3) \quad I(P_0 \mid v, \mathbf{Q}) \geq \|\mathbf{I}_1^*\|_0^2 = E_0(\mathbf{I}_1^{*2})$$

with equality if and only if  $\mathbf{I}_1^* \in \dot{\mathbf{Q}}$ .

B. If  $\dot{\mathbf{P}} = \dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2$ ,  $\mathbf{I}_1^* \neq 0$  and the assumptions of theorem 3.3.1 apply to  $\mathbf{P}$  and  $v$ , then the efficient influence function for  $v$ , as defined in section 3, is given by

$$(4) \quad \tilde{\mathbf{I}}_1 = \|\mathbf{I}_1^*\|_0^{-2} \mathbf{I}_1^* \equiv I^{-1}(P_0 \mid v, \mathbf{P}) \mathbf{I}_1^*.$$

In the same way as (2.4.4) led to the definition of efficient score function in regular parametric models, this theorem leads to

**Definition 1.** We call  $\mathbf{I}_1^*$  the *efficient score function* for  $v$  in  $\mathbf{P}$ , and write it  $I^*(\cdot, P_0 \mid v, \mathbf{P})$ .

The geometry of the relationship between  $\dot{\mathbf{I}}_1$  and  $\mathbf{I}_1^*$  is given in figure 3.

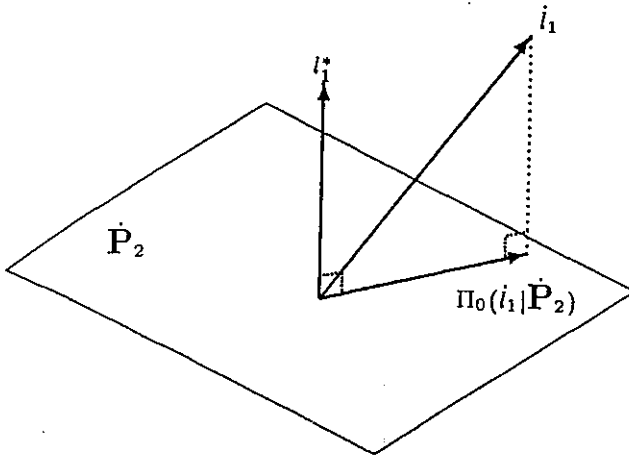


FIGURE 3. Projection of score functions.

**Proof of theorem 1.**

A. By proposition 2.4.1.A for any regular parametric submodel  $\mathbf{Q}$  as above,

$$\begin{aligned} I(P_0 \mid v, \mathbf{Q}) &= \|\dot{\mathbf{i}}_1 - \Pi_0(\dot{\mathbf{i}}_1 \mid \dot{\mathbf{Q}}_2)\|_0^2 \\ &= \|\dot{\mathbf{i}}_1 - \Pi_0(\dot{\mathbf{i}}_1 \mid \dot{\mathbf{P}}_2) + \Pi_0(\dot{\mathbf{i}}_1 \mid \dot{\mathbf{P}}_2) - \Pi_0(\dot{\mathbf{i}}_1 \mid \dot{\mathbf{Q}}_2)\|_0^2 \\ &= \|\dot{\mathbf{i}}_1^*\|_0^2 + \|\Pi_0(\dot{\mathbf{i}}_1 \mid \dot{\mathbf{Q}}_2) - \Pi_0(\dot{\mathbf{i}}_1 \mid \dot{\mathbf{P}}_2)\|_0^2 \end{aligned}$$

since  $\dot{\mathbf{i}}_1^* \perp \dot{\mathbf{P}}_2$  and, since  $\dot{\mathbf{Q}}_2 \subset \dot{\mathbf{P}}_2$  and both  $\Pi_0(\dot{\mathbf{i}}_1 \mid \dot{\mathbf{P}}_2)$  and  $\Pi_0(\dot{\mathbf{i}}_1 \mid \dot{\mathbf{Q}}_2) \in \dot{\mathbf{P}}_2$ . The result (3) follows.

Part B of the theorem follows from theorem 3.3.1 and the following proposition, which isolates the essential features of any differentiable function  $v(P)$  which identifies  $v$  on  $\mathbf{P}$  of (1).  $\square$

**Proposition 1.** Suppose that  $\mathbf{P}$  is as in (1), that  $v$  is identified by  $v : \mathbf{P} \rightarrow R$  so  $v(P_{(v,G)}) = v$  for  $P_{(v,G)} \in \mathbf{P}$ , that  $v(P)$  is pathwise differentiable on  $\mathbf{P}$  with derivative  $\dot{v}(P_0)$ , and that  $\mathbf{P}_1$  is regular. Then

$$(5) \quad \langle \dot{v}(P_0), \dot{\mathbf{i}}_1 \rangle_0 = 1,$$

and

$$(6) \quad \dot{v}(P_0) \perp \dot{\mathbf{P}}_2.$$

If also  $\dot{\mathbf{P}} = \dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2$ , then

$$(7) \quad \dot{v}(P_0) = \tilde{I}(\cdot, P_0 \mid v, \mathbf{P}) = \tilde{I}_1$$

where  $\tilde{I}_1 \equiv I^{-1}(P_0 \mid v, \mathbf{P})\dot{\mathbf{i}}_1^*$  as in (4) with efficient score function  $\dot{\mathbf{i}}_1^*$  given in (2).

**Proof.** First consider  $v$  restricted to the model  $\mathbf{P}_1$ . Then for  $P_{(v,G_0)} \equiv P_v \in \mathbf{P}_1$ , it follows on the one hand that

$$(a) \quad v(P_{(v,G_0)}) - v(P_0) = v - v_0 \equiv t$$

by definition of  $v(P)$ , while, on the other hand, pathwise differentiability of  $v$  yields

$$(b) \quad v(P_{(v, G_0)}) - v(P_0) = t \langle \dot{v}(P_0), \dot{I}_1 \rangle_0 + o(|t|).$$

Equality of (a) and (b) for all  $t$  yields

$$1 = \langle \dot{v}(P_0), \dot{I}_1 \rangle_0,$$

proving (5). To prove (6), note that  $h \in \dot{P}_2^0$  implies that there exists a one-dimensional regular parametric submodel  $Q_2 = \{P_{(v_0, G_\gamma)} : |\gamma| < 1\}$  having  $h$  as a tangent. Consideration of  $v$  restricted to  $Q_2$  yields, if  $P_\gamma \equiv P_{(v_0, G_\gamma)}$ ,

$$0 = v(P_\gamma) - v(P_0) = \gamma \langle \dot{v}(P_0), h \rangle_0 + o(\gamma)$$

and hence

$$0 = \langle \dot{v}(P_0), h \rangle_0 \quad \text{for } h \in \dot{P}_2^0,$$

which implies (6).

To prove (7), note that, since  $\dot{v}(P_0) \in \dot{P}$ ,

$$\dot{v}(P_0) = a I_1^* + bh,$$

where  $h \in \dot{P}_2$ . Here we have used the assumptions that  $\dot{P} = \dot{P}_1 + \dot{P}_2$  and  $P_1$  is regular. Then (6) implies that

$$\dot{v}(P_0) = a I_1^*,$$

and (5) then yields

$$1 = a \langle I_1^*, \dot{I}_1 \rangle_0 = a \|I_1^*\|_0^2.$$

Hence (7) holds. □

Suppose  $v$  is  $m$ -dimensional,  $m > 1$ , with vector efficient score function

$$(8) \quad I_1^* = (I_{11}^*, \dots, I_{1m}^*)^T,$$

where  $I_{11}^*, \dots, I_{1m}^*$  are the score functions of the components of  $v$  given by (2).

**Corollary 1.** Suppose  $P_1$  is regular.

A. If  $Q$  is a regular parametric submodel as in theorem 1, for which  $I^{-1}(P_0 | v, Q)$  is defined, then

$$(9) \quad I(P_0 | v, Q) \geq E_0(I_1^* I_1^{*T}),$$

with equality if and only if  $[I_1^*]$  is contained in  $\dot{Q}$ .

B. If  $\dot{P} = \dot{P}_1 + \dot{P}_2$  and the conditions of theorem 3.3.2 hold, then the efficient influence function is given by

$$(10) \quad \tilde{I}_1 = I^{-1}(P_0 | v, P) I_1^*$$

and

$$(11) \quad I(P_0 | v, P) = E_0(I_1^* I_1^{*T}).$$

**Proof.** Apply proposition 2.4.1.A to  $Q = \{P_{(v, G_\gamma)} : v \in N, \gamma \in \Gamma\}$  to deduce that for all  $a \in R^m$ ,

$$\begin{aligned} a^T I(P_0 | v, Q) a &= \| a^T \dot{I}_1 - \Pi_0(a^T \dot{I}_1 | \dot{Q}_2) \|_0^2 \\ &= \| a^T I_1^* \|_0^2 \\ &\quad + \| a^T (\Pi_0(\dot{I}_1 | \dot{Q}_2) - \Pi_0(\dot{I}_1 | \dot{P}_2)) \|_0^2 \\ &\geq a^T E_0(I_1^* I_1^{*T}) a, \end{aligned}$$

and A follows.

To prove B we generalize proposition 1. □

**Proposition 2.** Suppose  $P_1$  is regular,  $v: P \rightarrow R^m$  is identified as in proposition 1, and that  $v$  is pathwise differentiable on  $P$  with derivative  $\dot{v}(P_0)_{m \times 1}$  at  $P_0$  and  $E_0 \dot{v}(P_0) = 0$ . Then

$$(12) \quad E_0(\dot{v}(P_0) \dot{I}_1^T) = J_{m \times m} \equiv \text{the } m \times m \text{ identity}$$

and

$$(13) \quad [\dot{v}(P_0)] \perp \dot{P}_2.$$

If also  $\dot{P} = \dot{P}_1 + \dot{P}_2$  holds, then

$$(14) \quad \dot{v}(P_0) = \tilde{I}(P_0 | v, P) = \tilde{I}_1,$$

where  $\tilde{I}_1$  is given by (3.3.16) and

$$(15) \quad \tilde{I}_1 = I^{-1}(P_0 | v, P) I_1^*,$$

where  $I^{-1}(P_0 | v, P)$  is given in (3.3.17) and

$$(16) \quad I(P_0 | v, P) = E_0(I_1^* I_1^{*T}).$$

**Proof.** Claims (12) and (13) are proved as for proposition 1. Claims (14)–(16) follow by writing  $\dot{v}(P_0) = A I_1^* + h$ , where  $A_{m \times m}$  is a matrix and  $h$  is an  $m$ -vector whose components belong to  $\dot{P}_2$ . Then, for any  $b \in R^m$ ,

$$(a) \quad \dot{v}^T(P_0) b = I_1^{*T} A^T b + h^T b.$$

As before, we get from (13) that  $h^T b = 0$ . Then, by taking the inner product of both sides of (a) with  $b^T h$  in  $L_2(P_0)$ ,

$$[E_0(\dot{I}_1 \dot{v}^T(P_0))] b = E_0(\dot{I}_1 I_1^{*T}) A^T b$$

or from (12) and (13)

$$b = E_0(I_1^* I_1^{*T}) A^T b \quad \text{for every } b \in R^m.$$

Hence,  $A$  and  $E_0(I_1^* I_1^{*T})$  are nonsingular and, by (12),

$$A = [E_0(I_1^* I_1^{*T})]^{-1}.$$

So,  $\tilde{I}_1 = \Pi_0(\dot{v}(P_0) | \dot{P}) = [E_0(I_1^* I_1^{*T})]^{-1} I_1^*$ , and the proposition follows. □

Now suppose we have Euclidean nuisance parameters present. That is,  $\mathbf{P} = \{P_{(\theta, G)} : \theta = (v, \eta), v \in N, \eta \in H, G \in \mathbf{G}\}$ ,  $N \subset R^m$ ,  $H \subset R^{k-m}$ , with  $N, H$  open. We follow the usual convention and define

$$(17) \quad \begin{aligned} \mathbf{P}_{12} &= \{P_{(\theta, G_0)} : \theta \in \Theta\}, \\ \mathbf{P}_{23} &= \{P_{(v_0, \eta, G)} : \eta \in H, G \in \mathbf{G}\}, \\ \mathbf{P}_1 &= \{P_{(v, \eta_0, G_0)} : v \in N\}, \end{aligned}$$

and so forth, and follow the same convention for parametric submodels  $\mathbf{Q}$ . Suppose the parametric model  $\mathbf{P}_{12}$  is regular. Let  $\dot{\mathbf{I}}_1, \dot{\mathbf{I}}_2$  be the score functions for  $\mathbf{P}_1, \mathbf{P}_2$ , and define

$$(18) \quad \dot{\mathbf{I}}_1^* \equiv \dot{\mathbf{I}}_1 - \Pi_0(\dot{\mathbf{I}}_1 | \dot{\mathbf{P}}_2 + \dot{\mathbf{P}}_3).$$

Note that, since  $\dot{\mathbf{P}}_2$  is finite-dimensional,  $\dot{\mathbf{P}}_2 + \dot{\mathbf{P}}_3$  is closed, and if  $\mathbf{Q} = \{P_{(\theta, G_\gamma)} : \theta \in \Theta, \gamma \in \Gamma\}$  is a regular parametric submodel then

$$\dot{\mathbf{Q}}_{23} = \dot{\mathbf{Q}}_2 + \dot{\mathbf{Q}}_3 \subset \dot{\mathbf{P}}_2 + \dot{\mathbf{P}}_3.$$

We can deduce, as in corollary 1, that

$$I(P_0 | v, \mathbf{Q}) \geq E_0(\dot{\mathbf{I}}_1^* \dot{\mathbf{I}}_1^{*T}).$$

If further

$$(19) \quad \dot{\mathbf{P}} = \dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2 + \dot{\mathbf{P}}_3$$

and the conditions of theorem 3.3.2 hold, then the efficient influence function is given by (10) and (11) with  $\dot{\mathbf{I}}_1^*$  defined in (18). We can use the same argument as in proposition 2.

Note, by applying (A.2.11), that (18) can be calculated by first projecting on  $\dot{\mathbf{P}}_2$  and then on the orthocomplement of  $\dot{\mathbf{P}}_2$  in  $\dot{\mathbf{P}}_2 + \dot{\mathbf{P}}_3$ , yielding

$$(20) \quad \begin{aligned} \dot{\mathbf{I}}_1^* &= \dot{\mathbf{I}}_1 - \Pi_0(\dot{\mathbf{I}}_1 | \dot{\mathbf{P}}_2) - \Pi_0(\dot{\mathbf{I}}_1 - \Pi_0(\dot{\mathbf{I}}_1 | \dot{\mathbf{P}}_2) | (\dot{\mathbf{P}}_2 + \dot{\mathbf{P}}_3) \cap \dot{\mathbf{P}}_2^\perp) \\ &= \dot{\mathbf{I}}_1 - \Pi_0(\dot{\mathbf{I}}_1 | \dot{\mathbf{P}}_2) - \Pi_0(\dot{\mathbf{I}}_1 | (\dot{\mathbf{P}}_2 + \dot{\mathbf{P}}_3) \cap \dot{\mathbf{P}}_2^\perp), \end{aligned}$$

since  $\dot{\mathbf{P}}_2 \perp \dot{\mathbf{P}}_2^\perp$ . Alternatively, we can reverse the roles of  $\dot{\mathbf{P}}_2, \dot{\mathbf{P}}_3$ . In fact, it is often convenient to proceed as follows. Begin by computing the efficient score function of  $\theta = (v, \eta)$ , which we shall call  $\dot{\mathbf{I}}^*$ . So

$$(21) \quad \dot{\mathbf{I}}^* = \dot{\mathbf{I}} - \Pi_0(\dot{\mathbf{I}} | \dot{\mathbf{P}}_3) \equiv \dot{\mathbf{I}}^*(\cdot, P_0 | \theta, \mathbf{P}).$$

If we write  $\dot{\mathbf{I}}^* = (\dot{\mathbf{I}}_1^*, \dots, \dot{\mathbf{I}}_k^*)^T$ , it is *not* true, in general, that, say,

$$\dot{\mathbf{I}}^*(\cdot, P_0 | v_1, \mathbf{P}) = \dot{\mathbf{I}}_1^*.$$

However, by reversing the roles of  $\dot{\mathbf{P}}_2$  and  $\dot{\mathbf{P}}_3$  in (20) and taking  $m = 1$  we obtain

$$(22) \quad \dot{\mathbf{I}}^*(\cdot, P_0 | v_1, \mathbf{P}) = \dot{\mathbf{I}}_1^* - \Pi_0(\dot{\mathbf{I}}_1^* | [\dot{\mathbf{I}}_i^* : i \neq 1]).$$



The definition 2.4.1 of full adaptation extends naturally to semiparametric models.

**Definition 2.** If

$$(23) \quad I(P_0 | \nu, \mathbf{P}) = I(P_0 | \nu, \mathbf{P}_{12})$$

and there exists a locally Gaussian regular (on  $\mathbf{P}$ ) estimate  $\hat{\nu}_n$  which is efficient at  $P_0$ , we say that  $\hat{\nu}_n$  is *adaptive at  $P_0$* . If  $\hat{\nu}_n$  is efficient for all  $P_0$  such that  $\mathbf{P}_{12}$  is regular we say  $\hat{\nu}_n$  is *adaptive*.

**Proposition 3.** Under the conditions of corollary 1 the adaptation condition (23) is equivalent to

$$(24) \quad [\dot{\mathbf{I}}_1 - \Pi_0(\dot{\mathbf{I}}_1 | \dot{\mathbf{P}}_2)] \perp \dot{\mathbf{P}}_3.$$

**Proof.** By (11) and (20),

$$I(P_0 | \nu, \mathbf{P}) = I(P_0 | \nu, \mathbf{P}_{12}) - E_0(WW^T),$$

where

$$W = \Pi_0(\dot{\mathbf{I}}_1 - \Pi_0(\dot{\mathbf{I}}_1 | \dot{\mathbf{P}}_2) | (\dot{\mathbf{P}}_2 + \dot{\mathbf{P}}_3) \cap \dot{\mathbf{P}}_2^\perp).$$

$W = 0$  if and only if (24) holds. □

A useful sufficient condition for (24) is given by:

**Corollary 2.** If the conditions of corollary 1 hold and  $\Pi_0(\dot{\mathbf{I}}_1 | \dot{\mathbf{P}}_2) = \Pi_0(\dot{\mathbf{I}}_1 | \dot{\mathbf{P}}_3)$ , then (23) and (24) hold.

**Corollary 3.** Let  $\mathbf{F}$  be the set of all real-valued functions  $q$  of  $\theta$  such that

(i)  $q$  is a function of  $\nu$  only;

(ii)  $q$  has a derivative  $\dot{q}_{1 \times k}$ .

If  $\mathbf{P}_{12}$  is regular, then (24) is equivalent to

$$(25) \quad (\dot{\mathbf{I}}_1^T, \dot{\mathbf{I}}_2^T) I^{-1}(\theta) \dot{q}^T \perp \dot{\mathbf{P}}_3$$

for all  $q \in \mathbf{F}$ .

**Proof.** The left-hand side of (25) is, by (2.3.2),  $\tilde{\mathbf{I}}(\cdot | q, \mathbf{P}_{12})$ , which is just a nonsingular transformation of the left side of (24), and hence the conditions are equivalent. □

### Examples

Now we illustrate the methods developed in this section by considering two important examples. Many more examples will be given in chapter 4.

**Example 1. The symmetric location model.**

By example 3.2.4 the conditions of theorem 1 are satisfied here. If  $I(G) < \infty$ , then, with the notation of example 3.3.1,

$$(26) \quad \dot{\mathbf{I}}_1(x) = \psi_0(x - \theta_0) = -\frac{g_0'}{g_0}(x - \theta_0)$$

is the efficient score function in  $\mathbf{P}_1$ . There is no parametric nuisance parameter  $\eta$  and

$$(27) \quad \dot{P}_3 = \{a(\cdot - \theta_0) : a \text{ symmetric, } E_0 a(X - \theta_0) = 0\}.$$

Since  $X - \theta_0$  is symmetric about 0,  $\psi_0$  is odd and

$$E_0 \psi_0(X - \theta_0) a(X - \theta_0) = 0 \quad \text{for all } a(\cdot - \theta_0) \in \dot{P}_3.$$

Thus,  $\dot{I}_1 \perp \dot{P}_3$ ,  $\tilde{I}_1$  agrees with (3.3.11), and full adaptation to shape is possible in this model. Examples of locally Gaussian regular adaptive estimates have been constructed by Van Eeden (1970), Beran (1974), (1978), Sacks (1975), and Stone (1975) among others. We give a construction in chapter 7.  $\square$

Further applications of this method (or at least the point of view) of Begun, et al. (1983) are given in chapter 4. For a systematic development of *score operators* as introduced by Begun, et al. (1983), see sections 5.4 and 5.5. Huang (1982), (1984), Choi (1989), and Choi and Hall (1988) have developed connections between score operators and least favorable submodels as defined in section 3.3. We will not pursue their approach here.

In the following example, and many other examples of interest in chapters 4 and 6, it is difficult to determine  $\dot{P}$  exactly. We now expand the argument in section 3.1 to the effect that such lack of knowledge is not necessarily troublesome for our theory. Suppose that  $T$  is a subspace of  $L_2^0(P_0)$  which is a candidate for  $\dot{P}$ . Typically we can easily prove either

$$(28) \quad T \subset \dot{P}$$

or

$$(29) \quad T \supset \dot{P}.$$

Usually the latter inclusion is the most difficult one to establish. For example, if the model involves group structure, mixing, or missing data, then it is easy to verify that  $T \equiv R(\dot{I}) \subset \dot{P}$  where  $\dot{I}$  is the score operator for the model. (See example 2 below and sections 5.4 and 5.5 for careful definitions and a full development.) But the reverse inclusion  $T \supset \dot{P}$  is much more difficult; recall the difficulty in showing  $\dot{P} = \dot{P}_1 + \dot{P}_2$  in the symmetric location model, example 3.2.4. On the other hand, in the case of constraint-defined models such as in examples 3.2.3, 3.3.3, and section 6.2, (29) is easy to prove for a natural choice of  $T$ , but (28) is more difficult and holds only under further hypotheses; recall example 3.2.3.

We argue here that if we are in case (28), have established an information bound based on  $T$ , and can produce a regular estimator which *achieves* the bound based on  $T$ , then proving the reverse inclusion (29) and hence  $\dot{P} = T$  is largely irrelevant. Here is a simple argument showing why this is true. Note that if  $T \subset \dot{P}$ , then for  $h \in L_2^0(P_0)$

$$(30) \quad \|\Pi_0(h | T)\|_0 \leq \|\Pi_0(h | \dot{P})\|_0.$$

Thus the conjectural or putative information bound  $\|\Pi_0(\dot{v} | T)\|_0^2$  based on the subspace  $T$  is less than or equal to the honest, true information bound  $I^{-1}(\dot{v} | P_0, P)$  defined by (3.3.24) which we may not be able to calculate

because of lack of knowledge of  $\dot{\mathbf{P}}$ . Now suppose that  $\hat{v}_n$  is a Gaussian regular estimator of  $v$  with asymptotic variance  $\Sigma(P_0, \hat{v}_n)$ . Then by the convolution theorem and (30)

$$(31) \quad \Sigma(P_0, \hat{v}_n) \geq I^{-1}(v | P_0, \mathbf{P}) = \|\Pi_0(\dot{v} | \dot{\mathbf{P}})\|_0^2 \geq \|\Pi_0(\dot{v} | \mathbf{T})\|_0^2.$$

But if  $\hat{v}_n$  achieves the bound provided by  $\mathbf{T}$ ; i.e., if

$$\Sigma(P_0, \hat{v}_n) = \|\Pi_0(\dot{v} | \mathbf{T})\|_0^2,$$

then equality must hold throughout (31), and both the bound is sharp and the estimate  $\hat{v}_n$  is *asymptotically efficient*. Alternatively, if  $\hat{v}_n$  is asymptotically linear with influence function  $\psi \in \mathbf{T}$ , then  $\psi \in \dot{\mathbf{P}}$  by (28), and by proposition 3.3.1.B,  $\hat{v}_n$  is again asymptotically efficient. These arguments justify our attitude toward exact determination of the tangent space  $\dot{\mathbf{P}}$  in example 2 below and in the examples in chapters 4 and 6 for which the inclusion (28) holds. Of course, the argument should be completed by producing a regular estimator which achieves the bound, as is done in the following example, but not universally for the examples considered in chapters 4 and 6.

**Example 2. Cox model without censoring.**

We consider a simple version of example 1.3.7 with  $C = \infty$  (no censoring) and  $m = k = 1$ . We observe  $(Z, T)$  with  $T$  having conditional hazard function given  $Z = z \in R$  specified by

$$(32) \quad \lambda(t | z) = r(vz)\lambda(t),$$

where  $v \in R$ ,  $\lambda = g/\bar{G}$ ,  $r(z) = \exp(z)$ ,  $g$  is the density of  $G$  on  $[0, \infty)$ , and  $\bar{G} = 1 - G$ . The distribution  $H$  of  $Z$  is assumed known with density  $h$  with respect to some fixed measure  $m$  while, with  $\mu \equiv$  Lebesgue measure,

$$\mathbf{G} = \{G \ll \mu\}.$$

Then, if we write  $r = r(vz)$ ,

$$P[T \geq t | Z = z] = \bar{G}^r(t) = [\bar{G}(t)]^r$$

and the density with respect to  $\mu \times m$  is

$$(33) \quad f(z, t, v, G) = r g(t) \bar{G}^{r-1}(t) h(z).$$

Calculation of scores for  $v$  and "for  $g$ " is straightforward: differentiation yields

$$(34) \quad \dot{h}_1(z, t) = z[1 - \Lambda_r(t)],$$

where

$$\Lambda_r(t) \equiv r(vz)\Lambda(t) \equiv r(vz) \int_0^t \frac{dG}{1-G}.$$

Similarly, letting  $\{g_\eta\}$  be a regular parametric family through  $g \equiv g_0$  with

$$a \in \{a \in L_2(G) : \int a dG = 0\} \equiv \mathbf{H}_0$$

as score for  $\eta$ , the score (operator) for  $g$  is

$$(35) \quad \dot{I}_2 a(z, t) = a(t) + (r-1) \frac{\int_t^\infty a(s) dG(s)}{\bar{G}(t)}.$$

For example,

$$g_\eta = \psi(\eta a) g / \int \psi(\eta a) dG$$

will do, if  $\psi : R \rightarrow R^+$  is bounded and continuously differentiable with bounded derivative  $\psi'$  satisfying  $\psi(0) = \psi'(0) = 1$  and with  $\psi'/\psi$  bounded; see also example 3.2.1.

Formula (35) defines a linear transformation  $\dot{I}_2$  from functions of  $t$  to functions of  $(z, t)$  which is well defined on  $H_0$ , but is, in general, unbounded, and hence its range is not necessarily in  $L_2(P)$ . (We will show that  $\dot{I}_2$  is a bounded operator if  $v \geq 0$ ,  $Z \geq 0$ , and  $Er(vZ) < \infty$ , and that this can always be arranged by a reparametrization if (46) holds ( $Z$  bounded).) Under the same assumption we can show, by arguing as in example 3.1.1, that  $\text{Range}(\dot{I}_2) \equiv R(\dot{I}_2) \subset \dot{P}_2$ . Equality of these two sets has not yet been established.

It will be very helpful to rewrite the scores in (34) and (35) in terms of the operators  $R$  and  $L$  introduced in section A.1. Let  $H \equiv L_2(G)$ , and let  $H_0 \equiv \{a \in H : \int a dG = 0\}$  as above. Then  $L : H \rightarrow H_0$  and  $R : H_0 \rightarrow H_0$  are defined by

$$(36) \quad \begin{aligned} R a(t) &\equiv a(t) - \frac{\int_t^\infty a dG}{1 - G(t)} \\ &= -E\{a(Y) - a(t) \mid Y > t\} \equiv -e_a(t) \\ &= -\text{mean residual life of } a(Y) \text{ at } t, \end{aligned}$$

where  $Y \sim G$ , and

$$(37) \quad L a(t) \equiv a(t) - \int_0^t a \frac{dG}{1-G} = a(t) - \int_0^t a d\Lambda,$$

where  $\Lambda(t) \equiv \int_0^t (1-G)^{-1} dG$  is the cumulative hazard function corresponding to  $G$ . The operators  $R$  and  $L$  arise as the logarithmic derivatives of the maps  $Hg \equiv g/\bar{G} \equiv \lambda$  and  $D\lambda \equiv \lambda \exp(-\int_0^\cdot \lambda(s) ds) \equiv g$ , respectively. In view of (32) and (33), it is not at all surprising that the scores (34) and (35) can be expressed in terms of  $R$  and  $L$  and the operators  $R_r$  and  $L_r$  corresponding to the conditional (given  $Z = z$ )  $df(1-G)^r$  with hazard rate  $\lambda_r \equiv r\lambda$ :

$$(38) \quad \begin{aligned} R_r a(t) &= a(t) - \frac{\int_t^\infty a dF(\cdot | z)}{1 - F(t | z)} \\ &= -E\{a(Y) - a(t) \mid Z = z, Y > t\}, \end{aligned}$$

and

$$(39) \quad L_r a(t) \equiv a(t) - \int_0^t a \, d\Lambda(\cdot | z) = a(t) - r(vz) \int_0^t a \, d\Lambda.$$

Comparison of (34) with (39) yields

$$(40) \quad \dot{I}_1(z, t) = z(L_r 1)(t).$$

Furthermore

$$(41) \quad \dot{I}_2 a(z, t) = (L_r R a)(z, t)$$

follows from (35), (36), and (39) since  $L \circ Ra = a$  implies

$$a(t) = Ra(t) - \int_0^t Ra \, d\Lambda$$

or, by definition of  $R$ ,

$$(42) \quad \frac{\int_0^t a \, dG}{1 - G(t)} = a(t) - Ra(t) = - \int_0^t Ra \, d\Lambda.$$

Hence, the right side of (35) equals

$$\begin{aligned} Ra(t) + r \frac{\int_0^t a \, dG}{1 - G(t)} &= Ra(t) - r \int_0^t Ra \, d\Lambda \quad \text{by (42)} \\ &= Ra(t) - \int_0^t Ra \, d\Lambda_r \\ &= L_r Ra(t) \quad \text{by (39)}. \end{aligned}$$

Based on these score calculations, we can take several different approaches to establishing information bounds for estimation of  $v$ . Here is a brief sketch of three different approaches:

**Method 1.** (Via a candidate efficient estimator). Consider a particular estimator  $\hat{v}_n$  of  $v$ , for example the Cox (1972) partial likelihood estimator. If we show that the given estimator  $\hat{v}_n$  is (locally) regular and asymptotically linear with an influence function  $\psi$  in the tangent space  $\dot{P}$  of the model, then in view of proposition 3.3.1,  $\psi = \tilde{I}_1$  and  $\hat{v}_n$  is efficient.

**Method 2.** (Via orthogonality calculations). This approach proceeds without a candidate efficient estimator. Instead, we try to compute  $I_1^*$ , and hence  $\tilde{I}_1$ , via orthogonality considerations: we try to find  $a^* \in H_0$  so that

$$(43) \quad \dot{I}_1 - \dot{I}_2 a^* \perp \dot{I}_2 a \quad \text{for all } a \in H_0.$$

If  $\dot{P}_2 = \overline{R(\dot{I}_2) \cap L_2^0(P)}$ , this gives  $\dot{I}_1 - \dot{I}_2 a^* \perp \dot{P}_2$ , and  $I_1^* = \dot{I}_1 - \dot{I}_2 a^*$  as in (2). If not, we at least obtain a lower bound.

**Method 3.** (Via inversion of  $\dot{I}_2 \dot{I}_2$ ). This approach involves calculation of the projection  $\Pi_0(\dot{I}_1 | \dot{P}_2)$  by inversion of  $\dot{I}_2^T \dot{I}_2$  and application of theorem A.2.2.

This was the approach taken for the Cox model (with censoring) by Begun, et al. (1983). Inversion of  $\dot{\mathbf{I}}_2^T \dot{\mathbf{I}}_2$  is complicated however, and gives much more than is needed for calculation of a bound for estimation of  $\nu$ ; in fact, this is essentially related to calculation of bounds for estimation of  $G$ , a theme which will be pursued in section 5.4. Again the sharpness of the bound depends on the validity of  $\dot{\mathbf{P}}_2 = \mathbf{R}(\dot{\mathbf{I}}_2) \cap L_2^0(P)$ .

We now illustrate methods 1 and 2 in the present case of the Cox model without censoring.

**Method 1. (Via estimation).**

The well-known partial likelihood estimator  $\hat{\nu}_n$  of Cox (1972) is our natural candidate for an efficient estimator. If  $EZ^2 r^2(\nu Z)$  is bounded uniformly in a neighborhood of  $\nu_0$ ,  $\hat{\nu}_n$  is locally linear and asymptotically normal as can be seen by generalizing the proof of Tsiatis (1981); see example 7.4.4. The influence function of  $\hat{\nu}_n$  is  $\mathbf{I}_1^* \|\mathbf{I}_1^*\|_0^{-2}$ , where

$$(44) \quad \mathbf{I}_1^*(z, t) = \dot{\mathbf{I}}_1(z, t) - \left( \frac{S_1}{S_0}(t) - r(\nu z) \int_0^t \frac{S_1}{S_0} d\Lambda \right),$$

and

$$(45) \quad S_i(t) = E\{Z^i r(\nu Z) 1_{[T \geq t]}\}, \quad i = 0, 1.$$

Under the assumption

$$(46) \quad |Z| \leq K < \infty \text{ a.s.},$$

we shall show that  $\mathbf{I}_1^* = \dot{\mathbf{I}}_1 - \dot{\mathbf{I}}_2 a^*$  for  $a^* \in \mathbf{H}_0$ , and hence  $\mathbf{I}_1^* \in \dot{\mathbf{P}}$ . In view of proposition 1, this yields the efficiency of Cox's estimator.

Note that (46) implies  $|S_1| \leq K S_0$ , so  $S_1/S_0$  is bounded, and hence  $L(S_1/S_0) \in \mathbf{H}_0$ . Since  $R \circ L a = a$  for  $a \in \mathbf{H} \equiv L_2(G)$  we can write  $S_1/S_0 = R \circ L(S_1/S_0)$ . But then, in view of the definition (39) of  $L_r$ ,  $\mathbf{I}_1^*$  in (44) can be written as

$$\begin{aligned} \mathbf{I}_1^*(z, t) &= \dot{\mathbf{I}}_1(z, t) - L_r \frac{S_1}{S_0}(z, t) \\ &= \dot{\mathbf{I}}_1(z, t) - L_r R(L \frac{S_1}{S_0})(z, t) \\ &= \dot{\mathbf{I}}_1(z, t) - \dot{\mathbf{I}}_2 a^*(z, t) \quad \text{by (41)} \end{aligned}$$

with

$$a^*(t) \equiv L\left(\frac{S_1}{S_0}\right) \in \mathbf{H}_0.$$

This establishes the claim:  $\mathbf{I}_1^* \in \dot{\mathbf{P}}$  and hence  $\hat{\nu}_n$  is efficient.

**Method 2. (Via orthogonality).**

We want to find an  $a^* \in \mathbf{H}_0$  so that (43) holds. To do this, we first express

the scores (40) and (41) in terms of a counting process martingale. Now  $N_r$  defined by

$$N_r(t) \equiv 1_{[T \leq t]}, \quad t \geq 0,$$

is a counting process. Conditional on  $Z = z$ , it has a compensator

$$\begin{aligned} A_r(t) &\equiv \int_0^t 1_{[T \geq s]} d\Lambda_r(s) \\ &= r \int_0^t 1_{[T \geq s]} d\Lambda(s); \end{aligned}$$

thus

$$M_r(t) \equiv N_r(t) - A_r(t)$$

is a martingale with respect to the  $\sigma$ -fields  $\mathcal{F}_t \equiv \sigma\{Z, 1_{[T \leq s]} : s \leq t\}$ . Suppose  $b$  is a fixed function. Then, by direct calculation

$$\begin{aligned} \int_0^\infty b(s) dM_r(s) &= b(T) - \int_0^\infty b(s) 1_{[T \geq s]} d\Lambda_r(s) \\ &= b(T) - \int_0^T b(s) d\Lambda_r(s) \\ (47) \qquad \qquad \qquad &= L_r b(T) \quad \text{by (39)}. \end{aligned}$$

Comparison of (47) with (40) and (41) yields

$$(48) \quad \dot{I}_1(Z, T) = Z \int_0^\infty dM_r(s)$$

and

$$(49) \quad \dot{I}_2 a(Z, T) = \int_0^\infty Ra(s) dM_r(s).$$

To calculate the efficient score function  $I_1^*$  for  $v$ , we want to find a function  $a^*$  with  $\int a^* dG = 0$  so that

$$I_1^* \equiv \dot{I}_1 - \dot{I}_2 a^* \perp \dot{I}_2 a \quad \text{in } L_2(P)$$

for all functions  $a \in \mathbf{H}_0$ ; i.e.,

$$(50) \quad E\{[\dot{I}_1 - \dot{I}_2 a^*] \dot{I}_2 a\} = 0 \quad \text{for all } a \in \mathbf{H}_0.$$

This is just as in Begun, et al. (1983), except that here we are working in  $L_2(P)$  rather than  $L_2(\mu)$  and have replaced  $A$  by  $\dot{I}_2$ ,  $\beta^*$  by  $a^*$ , and  $\beta$  by  $a$ .

Now we use a bit of martingale theory. By conditioning on  $Z$ , the expectation in (50) is easily calculated as the expectation of the predictable covariation process of the martingale transforms in (48) and (49); see, e.g., theorem B.3.1, page 891 of Shorack and Wellner (1986). Thus the left side of (50) equals

$$\begin{aligned} &EE\{[\dot{I}_1 - \dot{I}_2 a^*] \dot{I}_2 a \mid Z\} \\ &= EE\left\{\int (Z - Ra^*) Ra 1_{[T \geq s]} r(vZ) d\Lambda(s) \mid Z\right\} \end{aligned}$$

$$\begin{aligned}
 &= \int \{ E[Z r(vZ) 1_{[T \geq s]}] - E[r(vZ) 1_{[T \geq s]}] Ra^*(s) \} Ra(s) d\Lambda(s) \\
 (51) \quad &= \int \{ S_1(s) - S_0(s) Ra^*(s) \} Ra(s) d\Lambda(s)
 \end{aligned}$$

where  $S_1$  and  $S_0$  are as defined in (45). From (51) it is easy to make the right choice of  $a^*$ : set

$$(52) \quad a^* \equiv L\left(\frac{S_1}{S_0}\right).$$

Since  $R \circ L = \text{identity}$  by proposition A.1.8, it follows that

$$Ra^* = \frac{S_1}{S_0},$$

and hence the integrand of (51) is zero identically, and (50) holds. Thus the efficient score function for estimation of  $v$  is

$$\begin{aligned}
 I_1^*(Z, T) &= \dot{I}_1(Z, T) - \dot{I}_2 a^*(Z, T) \\
 &= \int_0^\infty \left[ Z - \frac{S_1(t)}{S_0(t)} \right] dM_r(t) \\
 (53) \quad &= \int_0^\infty [Z - E(Z | T = t)] dM_r(t),
 \end{aligned}$$

since

$$(54) \quad \frac{S_1}{S_0}(t) = E(Z | T = t)$$

by straightforward calculations. Hence the information for  $v$  is, by an easy martingale calculation

$$\begin{aligned}
 I(v) &= E[I_1^*(Z, T)^2] \\
 &= E E \left\{ \int_0^\infty [Z - E(Z | T = t)]^2 1_{[T \geq t]} r(vZ) d\Lambda(t) \mid Z \right\} \\
 &= E \int_0^\infty [Z - E(Z | T = t)]^2 r(vZ) \bar{G}(t)^{r(vZ)-1} dG(t) \\
 (55) \quad &= E[Z - E(Z | T)]^2 = E \text{Var}(Z | T).
 \end{aligned}$$

This argument can be extended to the censored case. For more on the Cox model, see examples 4.7.1, 5.5.2, and 7.4.4. For information calculations in a family of models containing the Cox model, see Sasieni (1992).  $\square$



# 4 | Euclidean Parameters: Further Examples

## 4.1 INTRODUCTION: MODELS

The examples of the last chapter illustrate that both the form of tangent spaces and the ease with which the projections associated with information bounds for Euclidean parameters can be calculated varies widely. Our aim in this chapter is to provide a classification of most of the interesting semiparametric models according to common features of their tangent spaces. Not surprisingly, this classification arises by considering the different methods one can employ to construct parametric models.

Semiparametric models are typically constructed by eliminating some of the restrictions specifying a parametric model. For instance, by eliminating the requirement that errors are normal in the normal measurement model, we obtain the symmetric location model. Occasionally, they arise by parametrically restricting some aspects of a nonparametric model. An example of this type is the model considered in example 3.3.2 in which the set of all distributions on a bounded set in  $R$  is restricted by specifying the coefficient of variation.

Therefore, to construct semiparametric models, we first need to think about the usual ways of constructing parametric models.

### *Parametric Models*

One major way is to think of the observation  $X$  as having been obtained from an  $X$ -valued variable  $\varepsilon$  with known distribution  $G_0$  by means of one of a parametric group of transformations  $A = \{a(\cdot, \theta) : \theta = (v, \eta), v \in N, \eta \in H\}$  of  $X$ . Thus

$$X = a(\varepsilon, \theta),$$

where both  $a$  and  $\varepsilon$  take values in  $X$ . These are the *parametric group models* considered, for example, by Lehmann (1983, section 1.3, pages 19–26). If  $A$  is the location or location and scale group on  $R$ , or more generally a subgroup of the affine group on  $R^d$ , we obtain classical models which are used as building blocks in the construction of more complex parametric and, as we shall see, semiparametric models. Here are several particular cases of parametric group models:

**Normal measurement model:**  $\mathbf{X} = R$ ,  $G_0 = N(0,1)$ ,

$$a(\varepsilon, \theta) = v + \eta\varepsilon, \quad \eta > 0.$$

A generalization is the

**Multivariate Gaussian model:**  $\mathbf{X} = R^d$ ,  $\theta = (v, \Sigma)$  with  $v$   $d \times 1$ ,  $\Sigma$  positive definite,  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_d)^T$  with the  $\varepsilon_j$  i.i.d.  $N(0,1)$ ,

$$a(\varepsilon, \theta) = v + S\varepsilon,$$

where  $S$  is the unique lower triangular root of  $\Sigma$  with positive diagonal entries. The existence of such a root is given, for instance, in Graybill (1983, theorem 8.6.2).

**Exponential lifetime model:**  $\mathbf{X} = R^+$ ,  $G_0 = \text{Exponential}(1)$ , the standard exponential distribution, and

$$a(\varepsilon, v) = v\varepsilon, \quad v > 0.$$

The second major class of basic models is that of the full rank  $k$ -parameter exponential families discussed in Lehmann (1983, section 1.4, pages 26–36). These families are determined by a dominating measure  $\mu$  on  $\mathbf{X}$  and a statistic  $T(\mathbf{X})$  which is  $k \times 1$  and such that  $\mu\{\theta^T T(\mathbf{X}) \neq c\} > 0$  for all  $c$ . Then  $P_\theta$  belongs to the family generated by  $T$  and  $\mu$  if and only if

$$\frac{dP_\theta(x)}{d\mu} = \exp[\theta^T T(x) - b(\theta)]$$

with

$$\Theta \equiv \left\{ \theta \in R^k : \int \exp(\theta^T T(x)) d\mu(x) < \infty \right\}$$

and

$$b(\theta) = \log \left( \int \exp(\theta^T T(x)) d\mu(x) \right).$$

The common building block examples include the multivariate Gaussian model, the two-parameter gamma family, the multinomial family including the binomial and the Poisson family. The three basic group models we have cited are also exponential families (but of course this is not true in general).

More complex models can be constructed from basic models in a variety of ways which we list and describe below. These methods of construction can be alternated and iterated as we please.

### 1. Restricting the parameter space

If  $\mathbf{Q} = \{Q_\gamma : \gamma \in \Gamma\}$  is the basic model and  $\theta \rightarrow \gamma(\theta)$  is a map from  $\Theta$  into  $\Gamma$  with  $\gamma(\Theta) \subset \Gamma$ , we obtain a submodel

$$\mathbf{P} = \{Q_{\gamma(\theta)} : \theta \in \Theta\}.$$

If, for instance,  $\mathbf{X}$  has  $d$  distinct points which we identify with the  $d$  standard basis vectors in  $R^d$  and  $\mathbf{Q}$  is the *Multinomial*(1,  $\gamma_1, \dots, \gamma_d$ ) family then  $\mathbf{Q}$  corresponds to the set of *all* probability distributions on  $\mathbf{X}$ . Any model on  $\mathbf{X}$  can be obtained as a submodel of this  $\mathbf{Q}$ .

## 2. Linking to covariates

If  $X = (Z, Y)$  with  $Z$  a vector of variables (covariates, dummy variables, etc.) with known distribution  $H_0$  on  $Z$  and  $Q = \{Q_\gamma : \gamma \in \Gamma\}$  is a given model for  $Y$ , we obtain *generalized regression models* by linking  $\gamma$  to  $Z$  via a regression function  $r: Z \times \Theta \rightarrow \Gamma$  so that the conditional distribution of  $Y$  given  $Z = z$  under  $P_\theta$  is  $Q_{r(z, \theta)}$ , or

$$P_\theta(Z \in A, Y \in B) = \int_A Q_{r(z, \theta)}(Y \in B) dH_0(z).$$

Examples include the

**Normal linear regression model:**  $Q = \{N(\mu, \sigma^2) : \gamma = (\mu, \sigma) \in \Gamma\}$ ,  $\Gamma = R \times R^+$ . Then,  $\theta = (v^T, \sigma)^T$ ,  $v \in R^m$ ,  $\sigma > 0$ , and

$$r(z, \theta) = (z^T v, \sigma)$$

gives the normal linear regression model. If  $z^T v$  is replaced in  $r(z, \theta)$  by a general parametric family we obtain nonlinear regression. If  $\sigma$  is also made to depend on  $z$  we model heteroscedasticity.

**Logistic regression:** In this model

$$Q = \{\text{Bernoulli}(\gamma) \text{ distributions, } 0 < \gamma < 1\}$$

and

$$r(z, \theta) = \exp(\theta^T z) / (1 + \exp(\theta^T z)).$$

Both of these models are special cases of so-called GLIM models; see McCullagh and Nelder (1983).

## 3. Mixing

If  $Q$  is as above, we can generate models by viewing  $\gamma$  as random  $\gamma \sim G_\theta \in G$  where  $G \equiv \{G_\theta : \theta \in \Theta\}$  is itself a parametric model. Thus,

$$P_\theta = \int Q_\gamma dG_\theta(\gamma).$$

Here is a well-known example.

**Negative binomial:** Suppose that  $Q = \{\text{Poisson}(\gamma) : \gamma > 0\}$  so that  $Q_\gamma(\{k\}) = \gamma^k \exp(-\gamma) / k!$ . Also let

$$G = \{\text{Gamma}(p, \lambda) : \lambda > 0, p > 0\},$$

where  $\text{Gamma}(p, \lambda)$  has density

$$g_{(p, \lambda)}(\gamma) = \frac{\lambda^p \gamma^{p-1}}{\Gamma(p)} \exp(-\lambda \gamma).$$

Then

$$P_{(p, \lambda)}(\{k\}) = \frac{\Gamma(k+p)}{\Gamma(p)\Gamma(k+1)} \left(\frac{\lambda}{\lambda+1}\right)^p \left(\frac{1}{\lambda+1}\right)^k, \quad k = 0, 1, 2, \dots,$$

is the negative binomial family.

We may combine linking with mixing. An example is:

**ANOVA Model II:** Take  $Z = (\delta_1, \dots, \delta_p)$ , with  $\delta_j$  independent 0 or 1. We

get the one-way ANOVA model by linking the mean  $v$  of the normal measurement model to  $Z$  via  $v(z, \gamma) = \gamma^T z$ . If we mix so that  $\gamma_1, \dots, \gamma_p$  are independent  $N(0, \theta)$ , we obtain model II.

Other models which can be generated in this way include Gaussian errors in variables and, in general, situations in which the initial unknown parameter  $\gamma$  can be thought of at least in part as a hidden random factor affecting the observation.

#### 4. State space transformation

We can generate new models through families of transformations  $\mathbf{A} = \{a(\cdot, v) : v \in N\}$  on  $\mathbf{X}$ . For such an  $\mathbf{A}$ , we can define

$$\mathbf{P} = \{Q_\gamma a^{-1}(\cdot, v) : \gamma \in \Gamma, v \in N\}.$$

Group models are generated in this way from  $\mathbf{Q} = \{G_0\}$ . We are led to this formulation when we suppose there is a "correct" but unknown scale  $a^{-1}(\cdot, v)$  on which the observations are distributed according to  $\mathbf{Q}$ , usually an interpretable group model. An example is the

**Box-Cox model:** Suppose  $\mathbf{Q} =$  Normal linear regression model and let

$$a(z, y, v) = (z, h_v(y)),$$

$$h_v(y) = \begin{cases} v^{-1}(|y|^v \operatorname{sgn}(y) - 1), & v \neq 0 \\ \log |y|, & v = 0. \end{cases}$$

That is, if  $v$  were known, we could by applying  $h_v^{-1}$  to  $Y$  get an ordinary normal linear model.

In some important situations, the observation  $X$  may be modeled as a reduction  $a(Y)$  of a complete but unobservable variable  $Y \sim Q \in \mathbf{Q}$ . An important example is *censoring*. In particular, we consider the:

**Exponential proportional hazards model with right censoring:**

Suppose that

$$Y = (Z, C, T), \quad (C, Z) \sim H_0$$

and that

$$T = \exp(\theta^T Z) \varepsilon, \quad \varepsilon \sim \text{Exponential}(1).$$

Then an observation  $X$  in this model is

$$X = a(Z, T, C) \equiv (Z, T \wedge C, 1_{[T \leq C]}).$$

That is, we begin by modeling lifetimes  $T$  of patients, items, etc., parametrically according to an exponential lifetime model with expected lifetime linked to covariates  $Z$  via  $\exp(\theta^T Z)$ . An ideal observation  $Y = (Z, T, C)$  consists of the independent variable  $Z$ ,  $T$ , and the potential censoring time  $C$  assumed independent of  $T$ . Mixture models can be viewed in this way too. Here the ideal observation is  $(\gamma, X)$ , where  $\gamma$  is given above. State space transformations figure prominently in the formulation of the EM algorithm.

#### 5. Biased sampling

Given a parametric model  $\mathbf{Q} = \{Q_\gamma : \gamma \in \Gamma\}$  on  $\mathbf{X}$ , known weight functions  $w_i(\cdot) \geq 0$ ,  $1 \leq i \leq s$ , and probabilities  $\lambda_1, \dots, \lambda_s$ ,  $\sum_{i=1}^s \lambda_i = 1$ , we

generate the *biased sampling model*

$$\mathbf{P} = \{P_\gamma : \gamma \in \Gamma\}$$

on  $\{1, \dots, s\} \times \mathbf{X}$  with observation  $(I, X)$  via

$$P_\gamma(I=i, X \in A) = \frac{\lambda_i \int_A w_i(x) dQ_\gamma(x)}{W_i(\gamma)},$$

where

$$W_i(\gamma) = \int w_i(x) dQ_\gamma(x).$$

That is, given  $I = i$ , we sample  $X$  not with its proper probability  $dQ_\gamma(x)$ , but with probability proportional to  $w_i(x)dQ_\gamma(x)$ . Models of this kind arise quite naturally in econometric and epidemiological contexts and, more generally, when parametric models are introduced into survey sampling. They are discussed further in section 4. The following important special case occurs when  $s = 1$  and  $w(x) = 1_A(x)$ .

**Truncation:** Given a parametric model  $\mathbf{Q} = \{Q_\gamma : \gamma \in \Gamma\}$  on  $\mathbf{X}$  and a measurable subset  $A$  of  $\mathbf{X}$ , we obtain a truncation model on  $\mathbf{X}$  via  $P_\gamma(X \in B) = Q_\gamma(X \in B \mid X \in A)$ . Truncated regression models arise quite naturally; see section 4.6.

### *Semiparametric Models*

We can construct semiparametric models by the same methods if we introduce nonparametric = "large" aspects to one or more of the operations 1–5. Thus,

In section 2 we consider semiparametric group models where  $G_0 \in \mathbf{G}$ , a large set of distributions, for instance all symmetric distributions or all  $k$ -variate distributions with independent components.

In section 3 we consider semiparametric regression models. These are obtained in two ways:

- Linking covariates to a semiparametric group model.
- Linking covariates to a parametric group model, with a semiparametric regression function.

For instance, if  $Z = (Z_1, Z_2)$  and we put

$$r(z_1, z_2, \theta, \nu) = \theta^T z_1 + \nu(z_2)$$

where  $\nu$  is arbitrary, we obtain the Engle, Granger, Rice, and Weiss (1986) model discussed further in section 3. If we take

$$r(z, \theta, \nu) = \nu(\theta^T z),$$

we obtain a projection pursuit model as in Huber (1985).

Semiparametric extensions of biased sampling models obtained by taking the starting  $Q$  as semi- or nonparametric are discussed in section 4. If in mixture models we permit mixing on only part of  $\gamma$  and the family of mixing distributions to be non- or semiparametric we obtain semiparametric mixture models. For instance, let  $\gamma = (\theta, \eta)$  and

$$P_{(\theta, G)} = \int Q_{(\theta, \eta)} dG(\eta),$$

where  $G$  ranges over all distributions of  $\eta$ . A special case is the Reiersøl errors in variables model. Semiparametric mixture models are discussed in section 5.

We also obtain semiparametric models by considering a set of state-space transformations  $A$  which is partly nonparametric. For instance, if in the exponential proportional hazards model we take  $T^* = a(T)$  with  $A = \{\text{all monotone transformations}\}$ , and then

$$X = (Z, T^* \wedge C, 1_{[T^* \leq C]}),$$

we obtain the Cox proportional hazards model discussed in section 3.4. Such models are discussed further in sections 6 and 7.

Many semiparametric models have been developed in the biostatistics and econometrics literature. General surveys in econometrics discussing many of the models and methods we shall study as well as many others are in Robinson (1988b), Newey (1990), and Stoker (1991).

## 4.2 SEMIPARAMETRIC GROUP MODELS

Let  $A$  be a group of measurable transformations of  $X$  onto itself with identity denoted by  $e$ , and  $G$  a probability distribution on  $X$ . Write  $a^{-1}$  for the inverse of the transformation  $a$ ,  $a_1 \circ a_2$  for the composition of the transformations  $a_1$  and  $a_2$ , and  $ax$  or  $a(x)$  for the image of  $x$  under  $a$ . Let  $G_a$  denote the distribution of  $a(X)$  when  $X$  has distribution  $G$ . Following Lehmann (1983, p. 19), we give:

**Definition 1.** The group model induced by  $(A, G)$  is

$$P = \{G_a : G_a = Ga^{-1}, a \in A\}.$$

We call  $P$  a *parametric group model* if

$$A = \{a(\cdot, \theta) : \theta \in \Theta\}, \quad \Theta \subset R^k,$$

and the maps  $\theta \rightarrow a(x, \theta)$  are smooth. Then we represent

$$P = \{P_\theta : \theta \in \Theta\}, \quad P_\theta = Ga^{-1}(\cdot, \theta).$$

### Conventions

- We write  $a_\theta$  for  $a(\cdot, \theta)$  and  $a$  for  $a_\theta$  when there is no chance of confusion.
- Define  $\theta_e$  as any solution of  $a_{\theta_e} = e$ . Note that  $P_{\theta_e} = G$ , even though  $\theta_e$  need not be uniquely defined.

If  $\mu$  is a measure on  $X$  such that  $\mu a^{-1}$  is equivalent to  $\mu$  for all  $a$  and  $\mu$  dominates  $G$ , write

$$(1) \quad g = \frac{dG}{d\mu}, \quad j(\cdot, \theta) = \frac{d\mu a^{-1}}{d\mu}.$$

Then the density of  $P_\theta$  with respect to  $\mu$  is given by (see also (b) of the proof of lemma 1 below),

$$(2) \quad p(x, \theta) = g(a^{-1}x)j(x, \theta).$$

Semiparametric extensions of parametric group models are obtained by simply letting  $G$  range over  $\mathbf{G}$ , a large subset of  $\mathbf{M}$ . We will call these *semiparametric group models*.

We have seen essentially two examples of parametric group models in section 1, the multivariate Gaussian model and the exponential lifetime model. The Gaussian linear model of example 2.4.5 is of this type also. We discuss it in detail next.

**Example 1. The Gaussian linear model.**

Let  $X = (Z^T, Y)^T$ ,  $Z \in R^m$ ,  $Y$  real,  $\Theta = R^{m+1} \times R^+$ . Write  $\theta = (v_1, \dots, v_m, \Delta, \sigma)^T$ , and define the group action by

$$a(z, y, \theta) = (z^T, \Delta + \sum_{j=1}^m v_j z_j + \sigma y)^T.$$

Suppose that  $Z \sim H_0$  is known and that  $Z$  is independent of  $Y \sim N(0, 1)$  under  $G$ . Then  $G$  is known, and under  $P_\theta$ ,  $Z \sim H_0$ , and given  $Z$

$$Y \sim N(\Delta + \sum_{j=1}^m v_j Z_j, \sigma^2).$$

This defines  $P_\theta$ , and hence  $P$ . □

Here are two semiparametric group models.

**Example 2. The general linear model.**

Under  $G$ , we take  $Y$  to have density  $q$  with respect to Lebesgue measure  $\lambda$  and  $Z$  to have a density  $h$  with respect to  $\mu_0$ , and let  $\mathbf{G}$  be the set of all product densities  $h q$ . The group action  $a$  is specified as in example 1, so that

$$a^{-1}(z, y) = (z, \sigma^{-1}(y - \Delta - \sum_{j=1}^m v_j z_j))^T, \quad j(\cdot, \theta) = \sigma^{-1}.$$

The model is then specified by densities with respect to  $\mu = \mu_0 \times \lambda$  by

$$p(z, y, \theta, g) = \sigma^{-1} h(z) q(\sigma^{-1}(y - \Delta - \sum_{j=1}^m v_j z_j)).$$

Note that the full parametrization  $(\theta, g)$  is not identifiable but  $v_1, \dots, v_m$  are bona fide parameters if  $h$  does not concentrate on a hyperplane. □

**Example 3. The elliptic model.**

Here is a generalization of the multivariate Gaussian model studied by many authors, see Cambanis, Huang, and Simons (1981) for a review. We can write

this model as

$$X = v + S\varepsilon,$$

where  $\varepsilon \sim G \in \mathbf{G} = \{\text{all Lebesgue absolutely continuous probability distributions on } R^d \text{ invariant under orthogonal transformations}\}$ ,  $S \equiv S(\Sigma)$  and  $a$  are identified as in the multivariate Gaussian model of section 1 with  $v \in R^d$ ,  $\Sigma$  positive definite. Then this is a semiparametric group model. If  $g$  is the density of  $G$ , it follows from the definition of  $\mathbf{G}$  that

$$g(x) = \tilde{g}(|x|^2)$$

for  $\tilde{g}: R^+ \rightarrow R^+$ . Hence, the density of  $P_{(\theta, G)}$  is given by

$$p(x, \theta, G) = \tilde{g}(|S^{-1}(x - \Delta)|^2) |\det(S^{-1})|.$$

where “det” denotes determinant. Note that here  $e = a(0, J)$  where  $J$  is the  $d \times d$  identity. Again the full parametrization is not identifiable, but  $(\Delta, \Sigma/\text{trace}(\Sigma))$  is a proper parameter. □

Calculation of information bounds and estimation in group models is greatly simplified by two observations discussed in the two subsections which follow.

- (a) For fixed  $G$  the efficient influence functions at  $\theta$  are linear transforms of the efficient influence function at  $\theta_e$  whose coefficients don't depend on  $G$ ; see lemmas 1 and 2.
- (b) The tangent space  $\dot{P}_2(\theta)$  obtained by fixing  $\theta$  is determined by  $\dot{P}_2(\theta_e)$ ; see lemma 3.

### Two Formulae for Efficient Score Functions for Parametric Group Models

Suppose that the parametrization  $\theta \rightarrow a_\theta$  of the group model has the following properties:  $\Theta$  is open  $\subset R^k$  and:

- (i)(a) Given any pair of points  $\theta_0, \theta_1$ , let  $\theta_{01}$  be such that  $a_{\theta_{01}} = a_{\theta_0}^{-1} \circ a_{\theta_1}$ . There exist open neighborhoods  $V_0, V_1$ , and  $V_{01}$  of  $\theta_0, \theta_1, \theta_{01}$ , respectively such that,  $\theta' \in V_0, \theta \in V_1$ , implies  $a_{\theta'}^{-1} \circ a_\theta = a_{\theta^*}$  for some  $\theta^* \in V_{01}$ .
- (b) The mapping  $\theta \rightarrow a_\theta$  is one-to-one on  $V_{01}$ . Let  $i$  denote its (local) inverse on the image of  $V_{01}$ .
- (c) Define a mapping from  $V_0 \times V_1 \rightarrow V_{01}$  by

$$\lambda(\theta', \theta) = i(a_{\theta'}^{-1} \circ a_\theta).$$

Suppose that  $\lambda$  has a derivative at  $(\theta_0, \theta_1)$ . It follows from these conditions that if  $\theta_0 \in V_1$  then the mapping  $\omega: V_1 \rightarrow V_{01}$  given by  $\omega(\theta) = \lambda(\theta_0, \theta)$  has a derivative  $\dot{\omega}(\theta_0)$  at  $\theta_0$  which we identify with a  $k \times k$  matrix. Assume that  $\dot{\omega}(\theta)$  is continuous in  $\theta$ .



- (ii) Suppose that the model generated by  $G$  and  $A$  is dominated so that its densities take the form (2) and that  $l(x, \theta)$  is differentiable at any  $\theta_e$  and its derivative  $\dot{l}(x)$  does not depend on  $\theta_e$ .

The existence of a smooth parametrization compatible with the group operation as in (i) makes the group  $A$  a Lie group and  $\theta \rightarrow a_\theta$  a compatible parametrization; see, for example, the definition given by Pontryagin (1966, page 284).

**Convention:** In this section if  $v: \Theta \rightarrow E$ , where  $E$  is Euclidean, we may treat the derivative  $\dot{v}(\theta_0)$  as an operator from  $R^k$  to  $E$ . If  $d\theta \in R^k$ ,  $\dot{v}(\theta_0)[d\theta]$  will denote the image of  $d\theta$  under  $\dot{v}(\theta_0)$ . For example,  $\dot{l}(\theta_0)[d\theta] = \dot{l}^T(\theta_0) d\theta$  if  $\dot{l}$ ,  $d\theta$  are represented as vectors. This convention applies whenever  $\dot{v}$  has an argument  $[\cdot]$ .

The following lemma essentially appears in Pfanzagl and Wefelmeyer (1982, page 269 of section 18.3), in a different form. It reveals the central role of  $\omega(\theta)$  in group model calculations.

**Lemma 1.** Assume (i) and (ii) hold. Then the score function  $\dot{l}(x, \theta)$  satisfies

$$(5) \quad \dot{l}(x, \theta) = \dot{\omega}^T(\theta)\dot{l}(a^{-1}x),$$

where  $a \equiv a_\theta$ . The information matrix  $I(\theta)$  satisfies

$$(6) \quad I(\theta) = \dot{\omega}^T(\theta)I\dot{\omega}(\theta),$$

where  $I = I(\theta_e)$ . So  $I(\theta)$  is nonsingular if and only if  $I$  and  $\dot{\omega}(\theta)$  are. In that case the model is regular and the efficient influence function  $\tilde{l}(\theta)$  is given by

$$(7) \quad \tilde{l}(x, \theta) = \dot{\omega}^{-1}(\theta)\tilde{l}(a^{-1}x),$$

where  $\tilde{l}(x) = I^{-1}\dot{l}(x)$ .

**Proof.** Fix an arbitrary  $\theta_0 \in \Theta$ . Let  $a_0 \equiv a_{\theta_0}$  and let  $V_0$  be as in (i). We claim that for all  $\theta \in V_0$ ,

$$(a) \quad p(x, \theta) = p(a_0^{-1}x, \omega(\theta))j(x, \theta_0).$$

To see this, first note that

$$P_{\omega(\theta)} = G(a_0^{-1} \circ a)^{-1} = Ga^{-1}(a_0).$$

Therefore, if  $X \sim P_{\omega(\theta)}$ , then  $a_0(X) \sim P_{\theta_0}$ . Now (a) follows from

$$\frac{dP_\theta}{d\mu}(x) = \frac{dP_{\omega(\theta)}a_0^{-1}}{d\mu}(x)$$

and the formula (see for example Halmos (1950, theorem 39.D, page 164))

$$(b) \quad \frac{dQT^{-1}}{d\mu}(x) = \frac{dQ}{d\mu}(T^{-1}(x))\frac{d\mu T^{-1}}{d\mu}(x),$$

which is valid for any measurability preserving point transformation  $T$ , and  $Q$  and  $\mu T^{-1}$  dominated by  $\mu$ .

Apply the chain rule to (a) to get (5) at  $\theta = \theta_0$ . Regularity of the model follows from proposition 2.1.1 since, by (6),  $\theta \rightarrow I(\theta)$  is continuous.  $\square$

Another representation can be obtained if we assume

- (iii)  $\mathbf{X} = R^d$ ,  $\mu$  is Lebesgue measure, and  $\mathbf{A}$  is a subgroup of the affine group on  $R^d$ . Then we can write  $a_\theta = (\Delta(\theta), S(\theta))$  with  $\Delta$  a  $d$ -vector,  $S$  a nonsingular  $(d \times d)$ -matrix, and

$$a_\theta x = \Delta(\theta) + S(\theta)x.$$

- (iv) The map  $\theta \rightarrow (\Delta(\theta), S(\theta))$  has continuous first partial derivatives, or, equivalently, has a  $\theta$  continuous derivative  $(\dot{\Delta}(\theta), \dot{S}(\theta))$ .
- (v) The density  $g$  of  $G$  is continuously differentiable with derivative  $g'$ , a  $d$ -vector. Let  $\psi \equiv -g'/g$ .

Note that examples 1 through 3 satisfy (iii) and (iv).

**Lemma 2.** If (iii)–(v) hold for the model generated by  $G$  and  $\mathbf{A}$  and if

$$(8) \quad \int (1 + |z|^2) |\psi(z)|^2 dG(z) < \infty,$$

then

$$(9) \quad \begin{aligned} \dot{\mathbf{i}}(x, \theta) [d\theta] &= \psi^T(\varepsilon) S^{-1}(\theta) \{ \dot{S}(\theta) [d\theta] \varepsilon + \dot{\Delta}(\theta) [d\theta] \} \\ &\quad - \text{trace}(S^{-1}(\theta) \dot{S}(\theta) [d\theta]) \\ &= - \text{trace} [ (J - \varepsilon \psi^T(\varepsilon)) S^{-1}(\theta) \dot{S}(\theta) [d\theta] \\ &\quad - S^{-1}(\theta) \dot{\Delta}(\theta) [d\theta] \psi^T(\varepsilon) ], \end{aligned}$$

where  $\varepsilon = S^{-1}(\theta)(x - \Delta(\theta))$  and  $J$  is the  $d \times d$  identity matrix. If  $I(\theta)$  is nonsingular, the model is regular.

**Proof.** From (iii) and (2)

$$p(x, \theta) = g(S^{-1}(\theta)(x - \Delta(\theta)))j(\theta),$$

where  $j(\theta) = |\det(S^{-1}(\theta))|$ .

Applying (iv) and (v), we see that  $\dot{\mathbf{i}}(x, \theta)$  is well defined a.e.  $\mu$  and

$$(a) \quad \begin{aligned} \dot{\mathbf{i}}(x, \theta) [d\theta] &= -\psi^T(\varepsilon) \{ \dot{S}^{-1}(\theta) [d\theta] (x - \Delta(\theta)) - S^{-1}(\theta) \dot{\Delta}(\theta) [d\theta] \} \\ &\quad + \frac{\dot{j}}{j}(\theta) [d\theta] \\ &= \psi^T(\varepsilon) S^{-1}(\theta) \{ \dot{S}(\theta) [d\theta] \varepsilon + \dot{\Delta}(\theta) [d\theta] \} + \frac{\dot{j}}{j}(\theta) [d\theta] \end{aligned}$$

since  $\dot{S}^{-1}(\theta) [d\theta] = -S^{-1}(\theta) \dot{S}(\theta) [d\theta] S^{-1}(\theta)$ .

Now, by (iv) and (a), there exists a continuous function  $c(\theta)$  such that

$$(b) \quad \|\dot{\mathbf{i}}(x, \theta)\|^2 \leq c(\theta) |\psi|^2(\varepsilon) (1 + |\varepsilon|^2) + 2 \frac{\|\dot{j}(\theta)\|^2}{j^2(\theta)},$$

where  $\| \cdot \|$  is any Euclidean norm. If we apply theorem A.7.2 using (b) and (8), we see that the map  $\theta \rightarrow I(\theta)$  is continuous. The argument of (2.1.8) and proposition 2.1.1 now yield

$$(c) \quad E_{\theta} \dot{I}(\theta) = 0$$

and, if  $I(\theta)$  is nonsingular, that the model is regular. If we apply (c) to (a) and use  $\int g'(z) dz = 0$ , we see that

$$\begin{aligned} \frac{j}{j}(\theta)[d\theta] &= \int g'^T(\epsilon)[S^{-1}(\theta)\dot{S}(\theta)[d\theta]\epsilon] d\epsilon \\ &= \text{trace} \{ S^{-1}(\theta)\dot{S}(\theta)[d\theta] \int \epsilon g'^T(\epsilon) d\epsilon \}. \end{aligned}$$

By (8) and an argument as in the proof of lemma I.2.4.a and b of Hájek and Šidák (1967),

$$(d) \quad \int (Jg(\epsilon) + \epsilon g'^T(\epsilon)) d\epsilon = 0.$$

Hence

$$\int \epsilon g'^T(\epsilon) d\epsilon = -J,$$

and (9) follows. □

### *Tangent Spaces in Semiparametric Group Models*

Suppose  $\theta = (v, \eta)$ ,  $v \in N \subset R^m$ ,  $\eta \in H \subset R^{k-m}$ . Designate  $v$ ,  $\eta$ , and  $g$  as variables 1, 2, and 3 respectively. Write  $\dot{P}_3$  for  $\dot{P}_3(\theta_e)$ .

**Lemma 3.** If  $P$  is a semiparametric group model dominated by  $\mu$  and if  $\mu a^{-1}$  is equivalent to  $\mu$ ,  $\dot{P}_3$  determines  $\dot{P}_3(\theta)$  for all  $\theta$ . In fact,

$$(10) \quad \dot{P}_3(\theta) = \{ h(a^{-1} \cdot) : h \in \dot{P}_3 \}.$$

Therefore,

$$(11) \quad \Pi_0(h | \dot{P}_3(\theta)) = \Pi_0(h \circ a | \dot{P}_3) \circ a^{-1}.$$

**Proof.** Note that since  $\mu$  and  $\mu a^{-1}$  are equivalent,

$$\begin{aligned} 2 \left( \frac{s(\cdot, \theta, g_{\eta})}{s(\cdot, \theta, g_0)} - 1 \right) &= 2 \left( \left( \frac{g_{\eta}}{g_0} \right)^{1/2} (a^{-1} \cdot) - 1 \right) \\ &= \eta h(a^{-1} \cdot) + o(\eta) \end{aligned}$$

in  $L_2(P_{\theta})$  if and only if

$$2 \left( \frac{s(\cdot, \theta_e, g_{\eta})}{s(\cdot, \theta_e, g_0)} - 1 \right) = \eta h(\cdot) + o(\eta)$$

in  $L_2(G)$ . □

The set  $G$  is, in the examples we consider, so large that projection on  $\dot{P}_3$  is easy to determine. One interesting class arises by postulating another group  $T$  of transformations on  $X$  and taking

$$G = \{G : GT^{-1} = G \text{ for all } T \in T\}.$$

For instance, in example 3,  $G$  is the set of all  $G$  invariant under orthogonal transformations. The projection on the right side in (11) can be calculated in these cases; see proposition A.3.2. Calculations of this type are carried out in section 6.2.

In fact, however, the explicit form of the projection is not needed in situations where adaptation is possible, which is the case in all the examples of this section. Group models yield the principal examples of situations where adaptation is possible. We proceed to derive a set of sufficient conditions for the possibility of adaptation. We then verify these conditions in a number of important examples. All of these examples are of the type studied in lemma 2. Now (9) reveals that what we essentially need for adaptive estimation are a preliminary estimate of  $\theta$  and an estimate of  $\psi = -g'/g$  which is consistent in an appropriate sense. Construction of the resulting adaptive estimates is discussed in section 7.8.

### Adaptation

Fix  $G$ . Assume  $P_{12}$  is regular. Suppose first that  $k = m$ . Then proposition 3.4.2 yields adaptation for  $v$  at  $\theta$  if and only if

$$(12) \quad [\dot{I}_1(\cdot, \theta)] \perp \dot{P}_3(\theta).$$

By lemma 1, (12) is equivalent to

$$[\dot{I}_1(a^{-1} \cdot)] \perp \dot{P}_3(\theta)$$

and hence, by lemma 3, to

$$\int \dot{I}_1(a^{-1}x)h(a^{-1}x) dP_\theta(x) = 0$$

or, equivalently,

$$(13) \quad \int \dot{I}_1(x)h(x) dG(x) = 0$$

for all  $h \in \dot{P}_3$ .

If  $m < k$ , corollary 3.4.3 and lemmas 1 and 3 similarly yield as a necessary condition for adaptation at  $\theta$  that, with  $\dot{I}^T = (\dot{I}_1^T, \dot{I}_2^T)$  and  $I$  the information matrix for  $P_{12}$ ,

$$(14) \quad \int \dot{I}^T(x)h(x) dG(x)I^{-1}(\dot{\omega}^{-1}(\theta))^T \dot{q}^T(\theta) = 0$$

whenever  $q(\theta)$  is a function of  $v$  only with derivative  $\dot{q}(\theta)_{1 \times k}$ . We can rewrite (14) in operator form in various ways. In particular we shall use

$$(15) \quad \dot{q}(\theta)\dot{\omega}^{-1}(\theta)[I^{-1} \int \dot{I}(x)h(x) dG(x)] = 0.$$

**Convention:** In several of the following examples  $X = (Z, Y)$ . We shall write  $h(Y)$ ,  $f(Z)$ , etc., to indicate the domain of functions  $h$ ,  $f$ . Unsubscripted expectations and probabilities are calculated under  $G$ .

**Example 2. The general linear model, continued.**

Fix  $G$ . It may be shown as in Hájek and Šidák (1967, pages 210–213) that  $P_{12}$  is regular if  $q$  is absolutely continuous,

$$\int (1 + u^2) \frac{[q']^2}{q}(u) du < \infty, \quad E|Z|^2 < \infty,$$

and  $E(Z - E(Z))(Z - E(Z))^T$  is nonsingular. Under these assumptions,

$$(16) \quad \dot{l}(x, \theta, g) = \sigma^{-1}(z^T \psi(u), \psi(u), u\psi(u) - 1)^T,$$

where  $u = \sigma^{-1}(y - \Delta - \sum_{j=1}^m v_j z_j)$  and  $\psi = -(q'/q)$ . By arguing as in example 3.2.1,

$$(17) \quad \dot{P}_3 = \{a(Y) + b(Z) : Ea(Y) = Eb(Z) = 0\}.$$

From (16), (17) and lemma 3, the score functions for both location  $\Delta$  and scale  $\sigma$  are in  $\dot{P}_3$  as one would expect since both parameters are unidentifiable. For the regression parameter  $v$ , by lemma 3,

$$\Pi_0(\dot{l}_1(X, \theta) | \dot{P}_3(\theta)) = \sigma^{-1} \Pi_0(\psi(Y)Z | \dot{P}_3) \cdot a^{-1}.$$

But  $Y$  and  $Z$  are independent when  $G = G_{\theta}$ . Therefore,

$$\{a(Y) : Ea(Y) = 0\} \perp \{b(Z) : Eb(Z) = 0\},$$

and by proposition A.2.3.C

$$\begin{aligned} \Pi_0(\psi(Y)Z | \dot{P}_3) &= \psi(Y)E(Z | Y) + ZE(\psi(Y) | Z) \\ &= \psi(Y)EZ \end{aligned}$$

since  $E\psi(Y) = 0$ . Hence,

$$\Pi_0(\dot{l}_1(X, \theta) | \dot{P}_3(\theta))(x) = \sigma^{-1} EZ \psi(u)$$

which is proportional to the score function for  $\Delta$ ; see example 2.4.5 for the Gaussian case. We conclude that  $\Pi_0(\dot{l}_1 | \dot{P}_3) = \Pi_0(\dot{l}_1 | \dot{P}_2)$  and hence that adaptation for  $v$  is possible in this example. The efficient score function is

$$I_1^*(x, v) = \frac{z - EZ}{\sigma} \psi(u), \quad u \equiv \sigma^{-1}(y - \Delta - v^T z),$$

and hence

$$\tilde{I}_1(x, v) = \frac{\sigma}{I_q} [E(Z - EZ)(Z - EZ)^T]^{-1} (z - EZ) \psi(u),$$

where

$$I_q \equiv E \psi^2(Y) = \int \frac{[q']^2}{q}(u) du.$$

Estimates  $T_n$  of  $\nu$  which are asymptotically normal with

$$\Sigma(P_0, T) = \frac{\sigma^2}{I_q} [E(Z - E(Z))(Z - E(Z))^T]^{-1}$$

for all  $P_0$  such that  $\Sigma(P_0, T)$  is defined, have been constructed by Bickel (1982), Koul and Susarla (1983), and under stronger conditions by Dionne (1981). Hsieh and Manski (1987) study issues arising in the practical implementation of these adaptive estimates. Local regularity of such estimates which is necessary for the definition of efficiency will be discussed in chapter 7. Bounds for estimation of the distribution  $Q$ , and hence of functionals of it such as the mean or median will be discussed in section 6.2. If  $G$  is restricted so that  $q$  is symmetric, then  $\Delta$  is also adaptively estimable much as in example 3.4.1. Construction of adaptive estimates in this and the next example are given in section 7.8.  $\square$

**Example 3. The elliptic model, continued.**

We show that adaptive estimation of the parameters  $(\Delta, \Sigma / \text{trace}(\Sigma))$  is possible. Note that any homogeneous function  $f$  of  $\Sigma$  which is of degree 0,  $f(\lambda\Sigma) = f(\Sigma)$  for all  $\lambda > 0$ , is a function of  $\Sigma / \text{trace}(\Sigma)$  so that our claim applies to any function of the mean, variances, and covariances which is homogeneous of degree zero in the latter.

We need to verify (15) for appropriate  $q(\theta)$ . Note first that

$$a_{\theta_0}^{-1} \circ a_{\theta} : x \rightarrow S_0^{-1} Sx + S_0^{-1} (\Delta - \Delta_0)$$

and

$$\omega(\theta) = (S_0^{-1} (\Delta - \Delta_0), S_0^{-1} \Sigma(S_0^{-1})^T).$$

Therefore,

$$\dot{\omega}(\theta_0)[d\theta] = (S_0^{-1}[d\Delta], S_0^{-1}[d\Sigma](S_0^{-1})^T)$$

and

$$(18) \quad \dot{\omega}^{-1}(\theta_0)[d\theta] = (S_0[d\Delta], S_0[d\Sigma](S_0)^T).$$

Next write

$$p(x, \theta, g) = \tilde{g}[(x - \Delta)^T \Sigma^{-1}(x - \Delta)][\det(\Sigma)]^{-1/2}$$

and note that, under  $G$ :

- (i)  $|X|$  has density  $c(d)r^{d-1}\tilde{g}(r^2)$  where  $c(1) = 2$  and, for  $d > 1$ ,  $c(d)$  is the surface area of the unit sphere in  $R^d$ .
- (ii)  $|X|$  is independent of  $U = X/|X|$  which is uniformly distributed on the surface of the unit sphere in  $R^d$ . Therefore, under  $G$ ,  $(|X_1|, \dots, |X_d|)$  and  $(\text{sgn}(X_1), \dots, \text{sgn}(X_d))$  are independent,  $\text{sgn}(X_1), \dots, \text{sgn}(X_d)$  are i.i.d., and  $\text{sgn}(X_1) = \pm 1$  with probability  $\frac{1}{2}$ .

If  $\Delta = (\Delta_1, \dots, \Delta_k)$ ,  $\Sigma = [\sigma_{ij}]$ , we use these symmetries and the approach of lemma 2 to calculate the components  $\dot{l}_i = \partial l / \partial \Delta_i$ ,  $\dot{l}_{ij} = \partial l / \partial \sigma_{ij}$  at  $(0, J)$  as follows:

$$\dot{\mathbf{I}}[d\Delta, d\Sigma] = \psi(|x|^2)(2x^T d\Delta + x^T d\Sigma x) - \frac{1}{2}\text{trace}(d\Sigma),$$

where  $\psi = -\tilde{g}'/\tilde{g}$ . Therefore,

$$(19) \quad \dot{\mathbf{I}}_i(x) = 2\psi(|x|^2)x_i$$

and

$$(20) \quad \dot{\mathbf{I}}_{ij}(x) = (2 - \delta_{ij})\left(x_i x_j \psi(|x|^2) - \frac{\delta_{ij}}{2}\right).$$

Properties (i) and (ii) yield further that  $I$  is diagonal with diagonal entries

$$(21) \quad \begin{aligned} E \dot{\mathbf{I}}_i^2(X) &\equiv a_0, \\ E \dot{\mathbf{I}}_{ij}^2(X) &= b_0 \delta_{ij} + b_1(1 - \delta_{ij}), \end{aligned}$$

where

$$(22) \quad \begin{aligned} a_0 &= 4E\left(|X_1|^2 \psi^2(|X|^2)\right) = 4E(U_1^2)I_0 = \frac{4}{d}I_0, \\ b_0 &= \text{Var}\left(|X_1|^2 \psi(|X|^2)\right) = E\left(|X_1|^4 \psi^2(|X|^2)\right) - \frac{1}{4} \\ &= E(U_1^4)I_1 - \frac{1}{4} \\ &= \frac{3}{d(d+2)}I_1 - \frac{1}{4} > 0, \\ b_1 &= 4E\left(X_1^2 X_2^2 \psi^2(|X|^2)\right) = 4E\left(U_1^2 U_2^2\right)I_1 = \frac{4}{d(d+2)}I_1, \end{aligned}$$

and

$$(23) \quad \begin{aligned} I_0 &= c(d) \int_0^\infty r^{d+1} \psi^2(r^2) \tilde{g}(r^2) dr, \\ I_1 &= c(d) \int_0^\infty r^{d+3} \psi^2(r^2) \tilde{g}(r^2) dr. \end{aligned}$$

The final identities in (22) are obtained from the moments of the *Dirichlet*  $(1/2, \dots, 1/2)$  distribution of  $(U_1^2, \dots, U_k^2)$ ; see, e.g., Wilks (1962, formula (7.7.7), page 179).

Next, calculate

$$\dot{\mathbf{P}}_3 = \{h \in L_2(G) : h \text{ a function of } |x|^2, \int h(x) dG(x) = 0\}.$$

For  $h \in \dot{\mathbf{P}}_3$ , we obtain from (19), (20), using symmetry again, that

$$(24) \quad \begin{aligned} \int \dot{\mathbf{I}}_i(x) h(x) dG(x) &= 0, \\ \int \dot{\mathbf{I}}_{ij}(x) h(x) dG(x) &= c_o \delta_{ij} \end{aligned}$$

for some  $c_o$ . From the structure of  $I$  and (24), the components of  $I^{-1} \int \dot{\mathbf{I}}(x) h(x) dG(x)$  have the same structure as (24). Now apply (18) to get that

$$(25) \quad \dot{\omega}^{-1}(\theta) \left[ I^{-1} \int \dot{h}(x) h(x) dG(x) \right] = (0, b_0^{-1} c_0 \Sigma).$$

Finally, we note that to verify (15) it is enough to consider

$$(a) \quad q(\theta) = \Delta$$

and

$$(b) \quad q(\theta) = f(\Sigma)$$

where  $f: R^{d(d+1)/2} \rightarrow R$  is homogeneous of order zero. We are done since, by (25), condition (15) evidently holds in case (a), while in case (b),  $q(\theta) = f(\Sigma) = f(\lambda\Sigma)$  for all  $\lambda$  implies

$$(26) \quad 0 = \frac{\partial}{\partial \lambda} f(\lambda\Sigma)|_{\lambda=1} = \sum_{i,j} \sigma_{ij} \frac{\partial f}{\partial \sigma_{ij}}(\Sigma)$$

(which is Euler's equation for homogeneous functions), and then (25) and (26) yield (15).

Adaptive estimates have been constructed in this example by Bickel (1982) under the minimal conditions suggested by lemma 2 and (23), that is,  $\tilde{g}$  absolutely continuous and

$$\int_0^\infty r^{d+1} (1+r^2) \psi^2(r^2) \tilde{g}(r^2) dr < \infty.$$

However, adaptive estimates with reasonable small sample behavior still have to be constructed.  $\square$

There is at least one natural alternative to the elliptic model as a generalization of the multivariate Gaussian model. We retain the affine group to generate  $P_1$ , but let  $G$  range over all distributions  $G$  such that, under  $G$ ,  $X_1, \dots, X_d$  are independent, identically distributed, and symmetric about 0. Thus,

$$X = \Delta + S \varepsilon,$$

where  $S = [s_{ij}]$  is nonsingular  $d \times d$ ,  $\varepsilon_1, \dots, \varepsilon_d$  are independent and identically and symmetrically distributed about 0.

To address identifiability issues we choose  $\tilde{\Delta}$  and  $\tilde{S}$  such that the components of  $\tilde{S}^{-1}(X - \tilde{\Delta})$  are independent and identically distributed and symmetric about 0. Since

$$(27) \quad \tilde{S}^{-1}(X - \tilde{\Delta}) = \tilde{S}^{-1}(\Delta - \tilde{\Delta}) + \tilde{S}^{-1}S\varepsilon,$$

$\tilde{\Delta} = \Delta$  and  $\Delta$  is identifiable. Furthermore, (27) implies that for some constant  $c$  and permutation matrix  $\Pi$ ,

$$(28) \quad \tilde{S}^{-1}S = c\Pi.$$

Since  $\Pi\Pi^T = J$ , the identity, we have  $S\Pi^T = c\tilde{S}$  and  $SS^T = c^2\tilde{S}\tilde{S}^T$ . It follows that, e.g.,  $SS^T / \text{trace}(SS^T)$  is identifiable.

If we drop the assumption of symmetry, (27) still yields (28) in view of the Darmois-Skitovich theorem (Kagan, Linnik, and Rao (1973, theorem 3.1.1, page 89)), provided  $\varepsilon_1$  is not normal. Consequently,  $SS^T / \text{trace}(SS^T)$  is still



identifiable. However,  $\Delta$  is not. (27) implies the existence of some constant  $\gamma$  such that

$$(29) \quad \tilde{S}^{-1} (\Delta - \tilde{\Delta}) = \gamma 1,$$

where  $1 = (1, \dots, 1)^T$ . By (28) we have  $S1 = c\tilde{S}1$ , which together with (29) yields the identifiability of

$$\Delta - 1^T \Delta (1^T S 1)^{-1} S 1 = \tilde{\Delta} - 1^T \tilde{\Delta} (1^T \tilde{S} 1)^{-1} \tilde{S} 1.$$

Note that for  $S = J$ , this is the natural parameter

$$\Delta - (d^{-1} \sum_{i=1}^d \Delta_i) 1.$$

Furthermore, the necessary condition for adaptive estimation is satisfied, for identifiable parameters. These claims which may be verified by applying (15) are left as an exercise. Adaptive estimates or even good  $\sqrt{n}$ -consistent estimates remain to be constructed in these situations.  $\square$

### Two-Sample Models

Suppose  $X = (Y, Z)$ ,

$$Y = \alpha + \sigma e_1, \quad Z = v_1 + v_2(\alpha + \sigma e_2)$$

where  $v_2, \sigma > 0$ ,  $e_1$  and  $e_2$  are independent, each with fixed (Lebesgue) density  $q$ . This is a parametric group model—the *two-sample location and scale model*. If we vary  $q$  freely, we obtain a semiparametric group model in which  $\alpha, \sigma$  are no longer identifiable. However,  $v_1$  and  $v_2$  are identifiable, and Stein (1956) showed that the orthogonality condition (3.4.24) for fully adaptive estimation of  $v_1$  and  $v_2$  holds. This example was greatly generalized by Pfanzagl and Wefelmeyer (1982). Here is a generalization of their result.

Let  $\mathbf{X} = \mathbf{X}_0 \times \mathbf{X}_0$  and write  $X = (Y, Z)$ . Let  $\mathbf{V}$  be a group of transformations on  $\mathbf{X}_0$ . A general two-sample parametric group model is obtained by taking  $G$  on  $\mathbf{X}$  making  $Y$  and  $Z$  independent and identically distributed and choosing  $\mathbf{A} = \mathbf{V} \times \mathbf{V}$  operating in the usual way on  $\mathbf{X}$ . Suppose  $\mathbf{V}$  is parametrized by  $\eta \rightarrow v(\cdot, \eta), \eta \in H \subset R^m$ . Parametrize  $\mathbf{A}$  by

$$(30) \quad (v, \eta) \rightarrow (v(\cdot, \eta), v(\cdot, v) \circ v(\cdot, \eta)).$$

Let  $\mu$  be a  $\sigma$ -finite measure on  $\mathbf{X}_0$ , and let  $\mathbf{G}$  consist of all product densities  $g(y, z) = q(y)q(z)$  on  $\mathbf{X}$ . We want to study estimation of  $v$  defined by (30) in the semiparametric group model  $\mathbf{P} = \{P_{(v, \eta, G)}\}$  generated by  $\mathbf{A}$  and  $\mathbf{G}$ . Assume that condition (i) applies to  $\mathbf{V}$  and its parametrization with  $\lambda$  having a nonsingular derivative. Suppose, further, that the model generated by  $q$  and  $\mathbf{A}$  satisfies condition (ii) of this section. Then

**Theorem 1.** The necessary condition (3.4.24) for adaptive estimation of  $v$  holds under these conditions.

Pfanzagl and Wefelmeyer (1982, sections 18.3, 18.5) established this result when  $\mathbf{X}_0$  is Euclidean and  $\mathbf{A}$  is a subgroup of the affine group.

**Proof.** Fix  $(v_0, \eta_0)$ . To simplify notation write  $v(v)$  for  $v_v$ . In view of lemma 1 we can assume without loss of generality that  $v(v_0) = e$ . Since  $\dot{\omega}$  is nonsingular, this can be achieved by reparametrizing  $V \times V = A$  locally by  $(v', \eta')$ , where  $(v', \eta') = (i(v^{-1}(v_0) \circ v(v)), \eta)$ . For convenience we continue to refer to the parameter as  $(v, \eta)$ .

Since  $\lambda$  defined in (i)(c) has a nonsingular derivative, we can further reparametrize  $V \times V$  locally by  $(v', \eta')$ , where

- (a)  $\eta = i(v(v_0 - v') \circ v(\eta'))$ ,  
 (b)  $v = i(v(v_0 + v') \circ v^{-1}(v_0 - v'))$ .

Note that the derivative of this mapping is nonsingular, and hence  $(v', \eta')$  is well defined. The tangent spaces  $\dot{P}'_2, \dot{P}'_3$  for the new parametrization coincide with  $\dot{P}_2$  and  $\dot{P}_3$ , and  $\dot{P}'_{12}$  agrees with  $\dot{P}_{12}$ . However,  $\dot{P}'_1 \neq \dot{P}_1$ , and in fact, as we shall see,  $\dot{P}'_1$  is orthogonal to  $\dot{P}_2$ ; i.e.,  $v'$  can be adaptively estimated in  $P_{12}$ . We also will show that the score function for  $v'$  in  $P_1$  is orthogonal to  $\dot{P}_3$ . Therefore the condition (3.4.24) for adaptation holds for the efficient score function, and hence the efficient influence function  $\tilde{I}(\cdot, P_0 | v', P_{12})$  is also orthogonal to  $\dot{P}_3$ . Since  $v$  is a smooth function of  $v'$  when the domain of  $(v', \eta')$  is restricted to a sufficiently small open neighborhood of  $(0, \eta_0)$ , we can apply corollary 3.3.2 to conclude that  $\tilde{I}(\cdot, P_0 | v, P_{12})$  is a linear transform of  $\tilde{I}(\cdot, P_0 | v', P_{12})$ . Hence (3.4.24) will hold. Here are the details of this argument.

With the new  $(v', \eta')$  parametrization, the density of  $X$  is

$$(c) \quad q(v^{-1}(\eta') \circ v^{-1}(v_0 - v')y)q(v^{-1}(\eta') \circ v^{-1}(v_0 + v')z) \\ j(y, \omega_R(v - v', \eta'))j(z, \omega_R(v + v', \eta')),$$

where  $\omega_R(\eta, \eta^*) \equiv i(v(\eta) \circ v(\eta^*))$ . By our assumptions  $\omega_R$  has a derivative; see Pontryagin (1966 pages 96; 284).

Now it is easily seen from the symmetries of (c) that the derivatives of the log likelihood for the reparametrized model are of the form

$$(d) \quad \tilde{I}_1(y, z, v'_0, \eta'_0) = \Psi_1(y) - \Psi_1(z) \\ (e) \quad \tilde{I}_2(y, z, v'_0, \eta'_0) = \Psi_1(y) + \Psi_2(z).$$

Under  $(v'_0, \eta'_0) = (0, \eta_0)$ ,  $Y$  and  $Z$  are i.i.d. Hence  $\dot{P}'_1$  and  $\dot{P}'_2 = \dot{P}_2$  are orthogonal. On the other hand, by formally varying  $q$  in (c), and for rigor arguing as in example 3.2.1,

$$(f) \quad \dot{P}_3(\theta_0) = \{b(Y) + b(Z) : \int b q_0 d\mu = 0\}.$$

Hence  $\dot{P}'_1$  is orthogonal to  $\dot{P}_3$  as well. The theorem follows.  $\square$

#### Example 4. Two-sample location and scale.

Stein's (1956) example falls under theorem 1. To see this, we parametrize the affine group  $V$  on  $R$  by  $\eta = (\alpha, \sigma) \longleftrightarrow [y \rightarrow \alpha + \sigma y]$ . Then, if  $\eta = (\alpha, \sigma)$ ,  $v = (v_1, v_2)$  we can see that  $Y$  and  $(Z - v_1)/v_2$  have identical density  $\sigma^{-1}q(\sigma^{-1}(\cdot - \alpha))$ , and the example follows. Fully adaptive estimates

have been constructed for a “regular” subset of this model by Weiss and Wolfowitz (1970). Lemma 2 suggests that constructions exist which provide estimates efficient at all  $q$  which are absolutely continuous and have

$$\int_{-\infty}^{\infty} (1 + t^2) \frac{[q']^2}{q}(t) dt < \infty.$$

Park (1990) constructs adaptive estimates under the minimal conditions in this model and the two-sample rotation model.

We get a natural generalization of the two-sample model by weakening the assumption that  $\epsilon_1, \epsilon_2$  are i.i.d. under  $G \in \mathbf{G}$  to the assumption that  $\epsilon_1$  and  $\epsilon_2$  are exchangeable:

$$G = \{ \text{all distributions } G \ll \mu \text{ such that } GT^{-1} = G \}$$

where  $T(y, z) = (z, y)$  and  $\mu$  is exchangeable. We then observe  $X \sim P_{(\theta, G)} \in \mathbf{P}$  with  $\theta = (v, \eta)$  and  $X = (Y, Z)$ , where  $(v_{\eta}^{-1}Y, v_{\eta}^{-1} \circ v_v^{-1}Z)$  has distribution  $G \in \mathbf{G}$ . Here  $v_{\eta}, v_v \in \mathbf{V}$ . Then  $X$  has a density with respect to  $\mu$  given by

$$g(v_{\eta}^{-1}y, v_{\eta}^{-1} \circ v_v^{-1}z) j(y, \eta) j(z, \omega_R(v, \eta)),$$

where  $\omega_R(v, \eta) = i(v_v \circ v_{\eta})$ .

An example of this type of paired observations model is given by

$$Y = U + \alpha + \sigma e_1, \quad Z = v_1 + v_2(U + \alpha + \sigma e_2),$$

where  $e_1, e_2, U$  are unobserved random variables with  $e_1, e_2$  i.i.d. and independent of  $U$ .

**Proposition 1.** Suppose that conditions (i) and (ii) are satisfied by  $\mathbf{A} = \mathbf{V} \times \mathbf{V}$  and  $\mathbf{G}$  with  $\lambda$  having nonsingular derivative. Then the necessary condition (3.4.24) for adaptive estimation of  $v$  holds in this model.

**Proof.** We proceed as in the proof of theorem 1, but replace (c)–(f) by

$$(c') \quad g(v^{-1}(\eta') \circ v^{-1}(v_0 - v')y, v^{-1}(\eta') \circ v^{-1}(v_0 + v')z) \\ \cdot j(y, \omega_R(v_0 - v', \eta')) j(z, \omega_R(v_0 + v', \eta')),$$

$$(d') \quad \dot{I}'_1(y, z, v'_0, \eta'_0) = \psi_1(y, z), \quad \text{where } \psi_1(y, z) = -\psi_1(z, y),$$

$$(e') \quad \dot{I}'_2(y, z, v'_0, \eta'_0) = \psi_2(y, z), \quad \text{where } \psi_2(y, z) = \psi_2(z, y),$$

$$(f') \quad \dot{P}'_3(\theta_0) \subset \{ a(Y, Z) \in L_2(P_0) : a(y, z) = a(z, y) \}.$$

In view of example A.3.2, the same argument as in the proof of theorem 1 completes the proof of the proposition.  $\square$

**Example 5. Two-sample scale mixture model.**

Suppose that  $U \sim G$  and that, given  $U = \eta$ ,  $\eta Y$  and  $\eta v Z$  are i.i.d. each with density  $q_0$ . Then

$$(a) \quad p(y, z, v, G) = v \int \eta^2 q_0(\eta y) q_0(v \eta z) dG(\eta),$$

and, at  $v_0 = 1$ ,

$$\dot{I}_1(Y, Z) = 1 + E\left[UZ \frac{q'_0}{q_0}(UZ) \mid Y, Z\right].$$

By example A.3.2 with  $\mathbf{H}_0$  the right-hand side of ( $f'$ ) in the proof of proposition 1 above,

$$\dot{I} - \Pi_0(\dot{I}_1 \mid \mathbf{H}_0) = \dot{I}_1(Y, Z) - \frac{1}{2} [\dot{I}_1(Y, Z) + \dot{I}_1(Z, Y)].$$

It is easily checked that

$$I^*(Y, Z) = \dot{I}_1 - \Pi_0(\dot{I}_1 \mid \mathbf{H}_0) = \frac{1}{2} E\left( UZ \frac{q'_0}{q_0}(UZ) - UY \frac{q'_0}{q_0}(UY) \mid Y, Z \right)$$

is attained by the least favorable submodel  $G_\gamma = G(\cdot/\gamma)$ :

$$\begin{aligned} \text{(b)} \quad p(y, z; v, G_\gamma) &= v \int \eta^2 q_0(\eta y) q_0(v\eta z) dG\left(\frac{\eta}{\gamma}\right) \\ &= v \gamma^2 \int \eta^2 q_0(\gamma\eta y) q_0(v\eta \gamma z) dG(\eta). \end{aligned}$$

Alternatively, estimation of  $v$  can be based on the maximal invariant  $R \equiv Z/Y = (UZ)/(UY)$ . Now

$$P_v(R \leq r) = P_v(vUZ \leq vUYr) = \int Q_0(vry) q_0(y) dy,$$

so  $R$  has the density

$$h(r) = v \int q_0(vry) q_0(y) y dy$$

and hence, formally, at  $v_0 = 1$ ,

$$\dot{I}_R(r) = 1 + \frac{r \int q'_0(ry) q_0(y) y^2 dy}{\int q_0(ry) q_0(y) y dy}.$$

The distribution of  $R$  is independent of  $G$ , and hence the score function for  $v$  based on  $R$  is the same for any assumed parametric submodel  $P_{(v, G)}$ . In particular, we can take  $P_{(v, G)}$  to be the least favorable submodel for  $v$  since for this model  $\partial l(\cdot, v, G_\gamma)/\partial v = I^*$ . Hence an application of proposition A.5.5 to this model yields

$$\dot{I}_R(r) = E[I^*(Y, Z) \mid R = r].$$

The loss of information in using the efficient influence function based on  $R$  is governed by  $I^* - E[I^* \mid R]$ , which in general doesn't vanish.

Finally, note that  $I^*$  can be rewritten as

$$I^*(y, z) = - \frac{1}{2p(y, z)} \left[ y \frac{\partial}{\partial y} p(y, z) - z \frac{\partial}{\partial z} p(y, z) \right],$$

where  $p(y, z) \equiv p(y, z, v, G)$ . Hence  $I^*$  can be estimated directly, without first estimating  $G$ , and efficient estimators can easily be constructed without prior knowledge of  $q_0$ . See Van der Vaart (1988b) for a far-reaching generalization.

The special case of this example with  $q_0$  an exponential distribution is studied in more detail in example 4.5.4, but we do not explore adaptation to  $q_0$  further.  $\square$

### 4.3 REGRESSION MODELS

Regression models arise when we observe  $X = (Z, Y)$  with  $Z \in \mathbf{Z}$ , a vector of covariates and  $Y \in \mathbf{Y}$ . Suppose  $Z \sim H \in \mathbf{H}$  and let  $\mathbf{Q} = \{Q_\gamma : \gamma \in \Gamma\}$  be a model on  $\mathbf{Y}$ . Finally, let  $\mathbf{R}$  be a family of (regression) functions. If  $r \in \mathbf{R}$ , then  $r : \mathbf{Z} \rightarrow \Gamma$ .

The general regression model  $\mathbf{P}$  is now defined by,  $P \in \mathbf{P}$  if

$$P(Z \in A, Y \in B) = \int_A Q_{r(z)}(B) dH(z).$$

Here,  $\mathbf{R}$  is parametrized wholly or partly by a Euclidean parameter  $\theta$ . In this generality every model is a regression model. However, if we specialize we arrive at models which can be analyzed simply.

We consider two cases:

- I.  $\mathbf{Q}$  is a semiparametric group model,  $\mathbf{R}$  parametric. The prototype here is the nonlinear regression model with arbitrary error distribution.
- II.  $\mathbf{Q}$  is parametric,  $\mathbf{R}$  semiparametric. The prototype here is the projection pursuit model. We begin with:

#### Regression Model I

Suppose  $\mathbf{Q}$  is a semiparametric group model satisfying assumptions (i) and (ii) of lemma 4.2.1,

$$\mathbf{Q} = \{Ga^{-1}(\cdot, \gamma) : \gamma \in \Gamma \subset R^s, G \in \mathbf{G}\},$$

where  $a(\cdot, \gamma)$  is a group of transformations. Assume that:

- (i)  $\varepsilon$  and  $Z$  are independent,  $\varepsilon \sim G \in \mathbf{G}$ ,  $Z \sim H \in \mathbf{H}$ , and  $Y = a(\varepsilon, r(Z, \theta))$  so that  $L(Y | Z) = Ga^{-1}(\cdot, r(Z, \theta))$ .
- (ii)  $\mathbf{R}$  is purely parametric  $r = r(\cdot, \theta)$ ,  $\theta \in \Theta \subset R^k$ . Further,  $\theta \rightarrow r(\cdot, \theta)$  is differentiable with derivative  $\dot{r}$  and  $|\dot{\omega}(r(\cdot, \theta)) \dot{r}(\cdot, \theta)| \in L_2(H)$ , where  $\dot{\omega}$  is as in assumption (i) of lemma 4.2.1.

Then  $X = (Z, Y)$  has density

$$(1) \quad p(x, \theta, g, h) = g(a^{-1}(y, r(z, \theta)))j(y, r(z, \theta))h(z)$$

with respect to a product measure  $\mu \times m$ ,  $\mathbf{Q} \ll \mu$ ,  $\mathbf{H} \ll m$ . If  $P_{(\theta, G, H)}$  has density given by (1), then regression model I is

$$\mathbf{P} = \{P_{(\theta, G, H)} : \theta \in \Theta, G \in \mathbf{G}, H \in \mathbf{H}\}.$$

By lemma 4.2.3 the tangent space at  $Q = Q_{(\gamma, G)}$  is

$$\dot{\mathbf{Q}}_2(\gamma, G) = \{ha^{-1}(\cdot, \gamma) : h \in \dot{\mathbf{G}}\}.$$

Suppose that

- (iii)  $\dot{G} = \{h \in L_2(G) : h \text{ is } \mathcal{B}_0 \text{ measurable, } Eh = 0\}$  for some sub- $\sigma$ -field  $\mathcal{B}_0$  of the Borel field on  $Y$ .

Here are two examples in which we assume (i) and (ii) are satisfied.

**Example 1. Nonlinear regression.**

If  $Q$  is the semiparametric location model and  $r(Z, \theta)$  is general, as above, we obtain the nonlinear regression model with arbitrary errors,

$$Y = r(Z, \theta) + \varepsilon.$$

If  $G = \{\text{all distributions on } Y\}$ , of course

$$\dot{G} = \{h \in L_2(G) : Eh = 0\},$$

and (iii) holds. If  $G = \{\text{all distributions symmetric about } 0\}$ , again  $\dot{G}$  is of the form given in (iii) with  $\mathcal{B}_0$  the sub- $\sigma$ -field of symmetric sets.  $\square$

**Example 2. Heteroscedasticity.**

It is often useful to suppose that not only the mean  $\mu$  of a dependent variable  $Y$  is a function  $\mu(Z, \theta)$  of covariates  $Z$  but that its scale  $\sigma$  can be similarly represented,  $\sigma(Z, \theta) = \exp[r_2(Z, \theta)]$ , say. The appropriate parametric regression model can be written

$$Y = \mu(Z, \theta) + \sigma(Z, \theta)\varepsilon,$$

where  $\varepsilon$  and  $Z$  are independent and only  $\theta$  is unknown. The corresponding semiparametric model starts with  $Q$ , the semiparametric location and scale model, and has  $r(Z, \theta) \equiv (\mu(Z, \theta), \sigma(Z, \theta))$ . We can again naturally consider the two possibilities for  $G$  from example 1.  $\square$

For type I models the essential formulae take a simple form.

**Proposition 1.** If the regression model  $P$  satisfies assumptions (i)–(iii), then:

- A. The score function for  $\theta$  in  $P_1$  is
- (2)  $\dot{l}_1(X, \theta, G) = \dot{r}^T(Z, \theta) \dot{\omega}^T(r(Z, \theta)) \dot{l}(\varepsilon)$ ,  
where  $\dot{l}(\varepsilon)$  is the score function for  $\gamma$  in the model  $Q$  at  $G$  and  $\gamma_e$ , as in (ii) of lemma 4.2.1.
- B. The information  $I(P \mid \theta, P_1)$  for  $\theta$  is
- (3)  $I(P \mid \theta, P_1) = E[\dot{r}^T(Z, \theta) \dot{\omega}^T(r(Z, \theta)) I \dot{\omega}(r(Z, \theta)) \dot{r}(Z, \theta)]$ ,  
where  $I = E \dot{l}(\varepsilon) \dot{l}^T(\varepsilon)$ .
- C. If  $I(P \mid \theta, P_1)$  is nonsingular and the map
- $$\theta \rightarrow \text{trace}\{E[\dot{r}^T(Z, \theta) \dot{\omega}^T(r(Z, \theta)) \dot{\omega}(r(Z, \theta)) \dot{r}(Z, \theta)]\}$$
- is continuous, then  $P_1$  is regular.
- D. The efficient score function for  $\theta$  in  $P$  is given by
- (4)  $I^*(X) = \dot{r}^T(Z, \theta) \dot{\omega}^T(r(Z, \theta)) \dot{l}(\varepsilon) - E[\dot{r}^T(Z, \theta) \dot{\omega}^T(r(Z, \theta))] E(\dot{l}(\varepsilon) \mid \mathcal{B}_0)$ .

Note that  $I^*$  does not depend on the structure of  $H$ .

**Proof.** Let  $\theta, G, H$  be labelled as variables 1, 2, 3, respectively, and let  $P_0 \longleftrightarrow (\theta, G, H)$ . Then A follows from lemma 4.2.1 and (1), B follows from (2), and C is immediate.

D. Assumption (iii) and lemma 4.2.3 yield that (in random variable notation),  $\dot{P}_2 = \{a(\varepsilon) : a \text{ is } \mathcal{B}_0\text{-measurable, } E a(\varepsilon) = 0\}$ . It is easy to see that  $\dot{P}_{23} = \dot{P}_2 + \dot{H}$ . By (i),  $\dot{P}_2 \perp \dot{H}$ . Now  $I^*$  is the projection in  $L_2(P_0)$  of  $\dot{I}_1(\cdot, \theta, G)$  on the orthocomplement of  $\dot{P}_{23}$ . Writing

$$\dot{I}_1(X, \theta, G) = \dot{I}_1(Z, a(\varepsilon, r(Z, \theta)), \theta, G),$$

we obtain from proposition A.3.1 and (2) that the projection of  $\dot{I}_1$  on  $\dot{P}_2$  is

$$\begin{aligned} E(\dot{I}_1(Z, a(\varepsilon, r(Z, \theta)), \theta, G) \mid \mathcal{B}_0)(\varepsilon) \\ = E[\dot{r}^T(Z, \theta)\dot{\omega}^T(r(Z, \theta))]E[\dot{I}(\varepsilon) \mid \mathcal{B}_0]. \end{aligned}$$

On the other hand, by (i),

$$E(\dot{I}_1(X, \theta, G) \mid Z) = 0.$$

Equation (4) follows from proposition A.2.3.C. □

**Example 1. Nonlinear regression, continued.**

If  $G$  ranges over all distributions, by proposition 1

$$(5) \quad I^*(X) = (\dot{r}^T(Z, \theta) - E(\dot{r}^T(Z, \theta)))\dot{I}(\varepsilon).$$

In particular, if  $\theta = (v^T, \eta)^T$ ,

$$r(Z, \theta) = \eta + r_1(Z, v),$$

and  $r_1(\cdot, v)$  is not constant, we get

$$(6) \quad I^*(X) = (-[\dot{r}_1(Z, v) - E(\dot{r}_1(Z, v))]^T \frac{g'}{g}(\varepsilon), 0)^T.$$

The first component is just the score function for  $v$  and agrees, as in the special case of the linear regression model, with the score function for  $v$  in the parametric nonlinear regression model,  $P_1$ . Therefore, we can adaptively estimate  $v$  in the presence of  $\eta$  just as in the linear model. This was noted by Manski (1984). If  $\varepsilon$  is assumed symmetric about 0, we can also adaptively estimate  $\mu$  in  $r(Z, \theta) \equiv \mu + r_1(Z, \eta)$  and  $\theta^T = (\mu, \eta^T)$ . □

**Example 2. Heteroscedasticity, continued.**

If  $\theta = (v_1^T, v_2^T, \eta_1, \eta_2)^T$  and

$$\mu(Z, \theta) = \eta_1 + r_1(Z, v_1),$$

$$\sigma(Z, \theta) = \exp\{\eta_2 + r_2(Z, v_2)\},$$

then

$$\dot{I}(X, \theta, G) = -\dot{r}^T(Z, \theta) \exp\left[-(\eta_2 + r_2(Z, v_2))\right] \left(\frac{g'}{g}(\varepsilon), \varepsilon \frac{g'}{g}(\varepsilon) + 1\right)^T,$$

where

$$\dot{r}(Z, \theta) = \begin{pmatrix} \dot{r}_1(Z, v_1) & 0 & 1 & 0 \\ 0 & \sigma \dot{r}_2(Z, v_2) & 0 & \sigma \end{pmatrix}.$$

Suppose first that  $\mathbf{G} = \{\text{all distributions symmetric about } 0\}$ . Then

$$E\left(\frac{g'}{g}(\varepsilon) \mid \mathcal{B}_0\right) = 0 \quad \text{and} \quad E\left(\varepsilon \frac{g'}{g}(\varepsilon) \mid \mathcal{B}_0\right) = \varepsilon \frac{g'}{g}(\varepsilon),$$

and, after some algebra, we obtain, writing  $\mathbf{I}^* = (\mathbf{I}_1^{*T}, \mathbf{I}_2^{*T}, \mathbf{I}_3^*, \mathbf{I}_4^*)^T$ ,

$$\mathbf{I}_1^* = -\frac{1}{\sigma(Z)} \dot{r}_1^T(Z, v_1) \frac{g'}{g}(\varepsilon),$$

$$\mathbf{I}_2^* = -[\dot{r}_2(Z, v_2) - E \dot{r}_2(Z, v_2)]^T (\varepsilon \frac{g'}{g}(\varepsilon) + 1),$$

$$\mathbf{I}_3^* = -\frac{1}{\sigma(Z)} \frac{g'}{g}(\varepsilon),$$

$$\mathbf{I}_4^* = 0.$$

Since  $(\mathbf{I}_1^*, \mathbf{I}_3^*) \perp \mathbf{I}_2^*$  and  $\mathbf{I}_1^*, \mathbf{I}_3^*$  agree with  $\dot{\mathbf{I}}_1, \dot{\mathbf{I}}_3$  we see that  $(v_1, \eta_1)$  can be estimated as well not knowing  $(v_2, \eta_2)$ ,  $G$ ,  $H$  as knowing them, i.e., that  $(v_1, \eta_1)$  should be estimable adaptively. Similarly, since  $\Pi_0(\dot{\mathbf{I}}_2 \mid \dot{\mathbf{P}}_2) = \Pi_0(\dot{\mathbf{I}}_2 \mid [\dot{\mathbf{I}}_4])$  by the formulas for  $\dot{r}$  and  $\mathbf{I}_2^*$ ,  $v_2$  should be estimable adaptively (in the presence of  $\eta_2$ ).

If  $\mathbf{G}$  ranges over all distributions,  $\mathbf{I}_2^*$  and  $\mathbf{I}_4^*$  are unchanged, but

$$\mathbf{I}_1^* = -\left(\frac{\dot{r}_1(Z, v_1)}{\sigma(Z)} - E\left(\frac{\dot{r}_1(Z, v_1)}{\sigma(Z)}\right)\right)^T \frac{g'}{g}(\varepsilon),$$

$$(7) \quad \mathbf{I}_3^* = -\left(\frac{1}{\sigma(Z)} - E\left(\frac{1}{\sigma(Z)}\right)\right) \frac{g'}{g}(\varepsilon).$$

Since  $\mathbf{I}_1^*$  and  $\mathbf{I}_2^*$  may now be correlated, estimation of  $v_1$  is not as precise as if  $v_2$  is known. However, it should still be possible to estimate  $(v_1, v_2)$  adaptively in the presence of  $(\eta_1, \eta_2)$ .  $\square$

### Regression Model II

Suppose  $\mathbf{Q}$  is a regular parametric model, and that  $\Gamma \subset R^s$ .

- (iv) Assume that  $\mathbf{R} = \{r(\cdot, \theta, t) : \theta \in \Theta \subset R^k, t \in \mathbf{T}\}$  where  $\mathbf{T}$  is a nice subset of a function space. Necessarily  $r(z, \theta, t) \in R^s$ . Suppose,  $Z \sim H \in \mathbf{H}$  and  $\mathbf{R} \subset [L_2(H)]^s$  for all  $H \in \mathbf{H}$ .
- (v)  $\theta \rightarrow r(\cdot, \theta, t)$  is continuously Fréchet differentiable as a map to  $[L_2(H)]^s$  with derivative  $\dot{r}_1$ .
- (vi) Since  $\mathbf{R}$  is a subset of a Hilbert space, we can define its tangent space  $\dot{\mathbf{R}}$  at  $r_0 = r(\cdot, \theta_0, t_0) \in \mathbf{R}$  and the tangent space  $\dot{\mathbf{R}}_2$  of  $\mathbf{R}_2 = \{r(\cdot, \theta_0, t) : t \in \mathbf{T}\}$  in the sense of section 3.2. By (v),  $\dot{\mathbf{R}}_1 = [\dot{r}_1(\cdot, \theta_0, t_0)]$ . Suppose  $\dot{\mathbf{R}}_2$ , the tangent space of  $\mathbf{R}$  with  $\theta$  fixed, is of the form



$\dot{R}_2 = \{a(Z) \in [L_2(H)]^r : a \text{ is } \mathcal{B}_0\text{-measurable}\}$ , for some sub- $\sigma$ -field  $\mathcal{B}_0$  of the Borel field on  $Z$ .

Finally suppose that, as we would expect,

(vii)  $\dot{R} = \dot{R}_1 + \dot{R}_2$ .

Here are some examples.

**Example 3. Projection pursuit regression.**

Friedman and Stuetzle (1981) introduced a class of regression models intermediate between nonparametric and fully parametric linear regression. The simplest example of the model is

(8)  $Y = t(\theta^T Z) + \varepsilon$ ,

where  $\varepsilon \sim N(0, \sigma^2)$  and is independent of  $Z$ ,  $t$  is differentiable, and  $\sigma$  is assumed known. If  $t$  is known, this is an ordinary nonlinear regression model. However, if  $t$  ranges over (say) all differentiable  $L_2$  functions of  $\theta^T Z$ , this becomes a regression model of type II with  $Q$  the Gaussian location model,  $R = \{t(\theta^T \cdot)\}$ . So,  $\mathcal{B}_0$  is generated by  $\theta^T Z$ . Note that  $\theta$  is not identifiable, but  $[\theta]$  is. That is, the direction of  $\theta$  is identifiable, but not its length nor its orientation. One might identify the line  $[\theta]$  by a nonzero element of it, for example,

$$\frac{\theta}{|\theta|} \text{sign}(\theta^T e_j),$$

where  $e_1, \dots, e_k$  are the standard unit vectors and  $j$  is minimal with  $\text{sign}(\theta^T e_j) \neq 0$ . □

**Example 4. Periodic function regression.**

McDonald (1983) has made use of a model similar to (8) in which  $\theta$  itself is estimable, namely

(9)  $Y = t(\theta^{-1} Z) + \varepsilon$ ,

where  $Z$  is real,  $t$  is an unknown function of period one, and  $\varepsilon \sim N(0, \sigma^2)$  independent of  $Z$  as before. This is a generalization of the classical parametric problem of estimating the period  $\theta$  of a sinusoid:

$$Y = A \sin(2\pi Z / \theta) + \varepsilon.$$

Here  $Z$  is usually time and typically the  $Z_i$ 's are preselected points. However, the asymptotic properties of efficient estimates in the preselected and random time models can be shown to be the same. This model is discussed, for instance, in Ibragimov and Has'minskii (1982, section II.6, pages 141–142, and section III.5, pages 209–211). The semiparametric model (9) is a type II regression model with  $\mathcal{B}_0 = \{A : A = A + \theta\}$  and the  $\mathcal{B}_0$ -measurable functions are the measurable functions of period  $\theta$ . □

**Example 5. "Partial spline" regression models.**

The following class of models was introduced by Engle et al. (1986) to deal with situations where the regression of the dependent variable on some of the covariates is arbitrary while on others it may be assumed linear. That is,

$X = (Z, W, Y)$  with  $Z, W$  vector and  $Y$  a scalar,  $\theta = (v^T, \mu)^T$ , and

$$(10) \quad Y = \mu + t(Z) + v^T W + \varepsilon,$$

where  $\varepsilon$  is independent of  $(Z, W)$ ,  $\varepsilon \sim N(0, \sigma^2)$  with  $\sigma$  known. We want to permit  $t$  to vary as freely as possible. Rewrite (10) as

$$(11) \quad Y = \tilde{\mu} + \tilde{t}(Z) + v^T(W - E(W|Z)) + \varepsilon,$$

where

$$\tilde{\mu} \equiv E t(Z) + v^T E W + \mu,$$

$$\tilde{t}(Z) \equiv t(Z) - E t(Z) + v^T E(W|Z) - v^T E W.$$

Then  $\tilde{\mu}$  and  $\tilde{t}$  are identifiable via

$$\tilde{\mu} = E(Y)$$

and

$$\tilde{t}(Z) = E(Y - \tilde{\mu} | Z),$$

while  $v$  is identifiable provided the distribution of  $W - E(W|Z)$  does not concentrate on an  $(s - 1)$  dimensional hyperplane. Suppose the distribution of  $(Z^T, W^T)^T$  and  $\sigma$  are known and  $t$  varies freely. This is again a regression model of type II with  $Q$  the Gaussian location model and with

$$r(z, w, v, t) = t(z) + v^T w.$$

If  $T$  is large enough,  $\mathcal{B}_0$ , given in (vi), is generated by  $Z$ . □

Here are the essential formulae for regression models of type II. Write  $P = \{P_{(\theta, t, H)} : \theta \in \Theta, t \in T, H \in H\}$  for such a model. Then  $P_{(\theta, t, H)}$  has density

$$(12) \quad p(x, \theta, t, H) = q(y, r(z, \theta, t)) h(z),$$

where  $q(\cdot, \gamma)$  is the density of  $Q_\gamma$  and  $h$  of  $H$  with respect to a suitable dominating measure.

**Proposition 2.** For the regression model II satisfying assumptions (iv)–(vii) of this section:

A. The score function for  $\theta$  in  $P_1$  is given by

$$\dot{1}_1(X, \theta, t, H) = \dot{r}_1^T(Z, \theta, t) \dot{1}(Y, r(Z, \theta, t)),$$

where  $\dot{1}(\cdot, \gamma)$  is the score function for  $\gamma$  in  $Q$ .

B. The efficient score function of  $\theta$  in  $P$  is

$$(13) \quad 1_1^*(X) = \tilde{r}^T(Z, \theta, t) \dot{1}(Y, r(Z, \theta, t)),$$

where

$$(14) \quad \tilde{r}^T(Z, \theta, t) = \dot{r}_1^T(Z, \theta, t) - E(\dot{r}_1^T(Z, \theta, t) \dot{1}^T(Y, r(Z, \theta, t)) | \mathcal{B}_0) I^{-1}(Z)$$

and

$$I(Z) = E(\dot{1}^T(Y, r(Z, \theta, t)) | \mathcal{B}_0).$$

C. If  $Q$  is a parametric location model based on  $G_0$ ,  $\varepsilon = Y - r(Z, \theta, t)$ , then

$$(15) \quad I_1^*(X) = - \left[ \dot{r}_1^T(Z, \theta, t) - E(\dot{r}_1^T(Z, \theta, t) | \mathcal{B}_0) \right] \frac{g_0'}{g_0}(\varepsilon).$$

**Proof.** A is immediate.

B. Note first that from (12) by an application of the chain rule for Fréchet derivatives,

$$\dot{P}_2 = \{ a^T(Z) \dot{i}(Y, r(Z, \theta, t)) : a \in [L_2(H_0)]^s, a \text{ is } \mathcal{B}_0\text{-measurable} \}.$$

Apply proposition A.3.5 with  $b \equiv \dot{i}(Y, r(Z, \theta, t))$  and  $H_0 = \dot{P}_2$  to get the formula for the projection on  $\dot{P}_2$ . Then, since  $E[\dot{i}(Y, r(Z, \theta, t)) | Z] = 0$ ,  $I^*$  in (13) is orthogonal to  $\dot{P}_3$ . C follows by specializing (14).  $\square$

**Remark 1.** Regularity of  $P_1$  holds if the map  $\theta \rightarrow \dot{r}_1^T(\cdot, \theta)$  is  $L_2$ -continuous and the matrix

$$I(\theta, t, H) = E[\dot{r}_1^T(Z, \theta, t) \dot{i}^T(Y, r(Z, \theta, t)) \dot{r}_1^T(Z, \theta, t)]$$

is nonsingular.

**Example 3. Projection pursuit regression, continued.**

By (15) of proposition 2

$$(16) \quad I_1^*(X) = [Z - E(Z | \theta^T Z)] t'(\theta^T Z) \frac{\varepsilon}{\sigma^2}.$$

The information matrix from (3.4.16) for  $\theta$  in  $\mathbf{P}$  is singular since  $\theta^T I^* = 0$ . This is to be expected since  $\theta$  is estimable only up to a constant factor. However a sufficient condition for identifiability of homogeneous functions of  $\theta$  at  $(\theta_0, t_0, H_0)$  is that  $H_0$  have a density function  $h_0$  with respect to Lebesgue measure which is positive in a neighborhood where  $t'_0(\theta_0^T z) \neq 0$ . This is not a situation where adaptation to  $t$  is possible.

A natural comparison can be made to the situation where  $t(x) = t_0(\eta x)$ , where  $t_0$  is known. The efficient score function for  $v$  in this model is easily seen to be

$$\left( Z - \frac{E\{Z[t'(\theta^T Z)]^2 \theta^T Z\}}{E[(\theta^T Z)t'(\theta^T Z)]^2} \theta^T Z \right) t'(\theta^T Z) \frac{\varepsilon}{\sigma^2}.$$

The loss in information is governed by the difference between the linear regression (without intercept) of  $Z t'(\theta^T Z)$  on  $\theta^T Z t'(\theta^T Z)$  and the conditional expectation of  $Z$  given  $\theta^T Z$ . There is no loss if and only if  $E(Z | \theta^T Z)$  is a vector multiple of  $\theta^T Z$ , for instance, if  $Z$  has a multivariate normal distribution with mean zero.

It is in fact surprising that estimation of  $[\theta]$  at rate  $n^{-1/2}$  should be possible since, if  $\mathbf{T}$  is effectively nonparametric, for instance, consists of all thrice continuously differentiable  $t$ , then  $t'$  can only be estimated at a rate slower than  $n^{-1/2}$ , for instance  $n^{-2/5}$  in our example. Nevertheless, it follows from a theorem of Brillinger (1982) that if  $H$  is spherically symmetric Gaussian and  $\hat{\theta}$  is the least squares estimate of  $\theta$  for the linear regression,  $Y = \theta^T Z + \varepsilon$ , then  $[\hat{\theta}]$

is a  $\sqrt{n}$ -consistent estimate of  $[\theta]$ . Results for this model have been obtained by Ruud (1986) and Duan and Li (1989), and for general  $H$  by Stoker (1986), (1991), Härdle and Stoker (1989), and others.

In projection pursuit regression we may elaborate by adding functions of other functions of  $Z$ , for example,

$$Y = t_1(\theta_1^T Z) + t_2(\theta_2^T Z) + \varepsilon.$$

Parameters  $[\theta_j]$  are in general identifiable provided that  $\theta_1$  is not a multiple of  $\theta_2$ , that is  $[\theta_1] \neq [\theta_2]$ . Algorithms for fitting the  $t_j$  and  $[\theta_j]$  are known, but the asymptotic theory of these estimates is unexplored.  $\square$

**Example 4. Periodic function regression, continued.**

If  $H$  has density  $h$  with respect to some translation invariant  $\sigma$ -finite measure  $\mu$ , then for any  $v \in L_1(H)$ ,

$$E(v(Z) | \mathcal{B}_0) = \frac{\sum_{j=-\infty}^{\infty} v(Z + j\theta)h(Z + j\theta)}{\sum_{j=-\infty}^{\infty} h(Z + j\theta)}.$$

Indeed, if  $B \in \mathcal{B}_0$ , then

$$\begin{aligned} & E\left(1_B(Z) \frac{\sum_{j=-\infty}^{\infty} v(Z + j\theta)h(Z + j\theta)}{\sum_{j=-\infty}^{\infty} h(Z + j\theta)}\right) \\ &= \int_{[0,\theta)} 1_B(z) \frac{\sum_{j=-\infty}^{\infty} v(z + j\theta)h(z + j\theta)}{\sum_{j=-\infty}^{\infty} h(z + j\theta)} \sum_{j=-\infty}^{\infty} h(z + j\theta) d\mu(z + j\theta) \\ &= \int 1_B(z) v(z) h(z) d\mu(z) = E 1_B(Z) v(Z). \end{aligned}$$

Since  $\dot{r}_1(Z, \theta, t) = \theta^{-2} t'(\theta^{-1} Z)$ , where  $t'(\theta^{-1} Z)$  is periodic of period  $\theta$ , it follows that

$$E(\dot{r}_1(Z, \theta, t) | \mathcal{B}_0) = \theta^{-2} t'(\theta^{-1} Z) E(Z | \mathcal{B}_0),$$

and hence, using the above projection formula to calculate  $E(Z | \mathcal{B}_0)$ ,

$$\begin{aligned} (17) \quad \Gamma^*(X, \theta, t) &= \theta^{-2} [Z - E(Z | \mathcal{B}_0)] t'(\theta^{-1} Z) \frac{\varepsilon}{\sigma^2} \\ &= \left\{ \frac{\sum_{j=-\infty}^{\infty} j h(Z + j\theta)}{\sum_{j=-\infty}^{\infty} h(Z + j\theta)} \right\} \theta^{-1} t'(\theta^{-1} Z) \frac{\varepsilon}{\sigma^2}. \end{aligned} \quad \square$$

**Example 5. Partial spline regression models, continued.**

Here

$$(18) \quad \Gamma^*(X, v, t) = [W - E(W | Z)] \frac{\varepsilon}{\sigma^2}.$$

There is a loss in information as compared to the linear regression model

$$Y = \mu + \eta^T Z + v^T W + \varepsilon,$$

which is governed by the difference between the linear regression of  $W$  on  $Z$  and  $E(W|Z)$ .

The parameter  $\mu$  of model (10) is unidentifiable. However, if we restrict  $t$  by  $E t(Z) = 0$ , the situation changes. The score function of  $\mu$  is just  $\varepsilon/\sigma^2$ , which is orthogonal to  $\dot{P}_2$ , so that adaptive estimation of  $\mu$  is possible. An adaptive estimate of  $\mu$  for this situation has been given by Schick (1986). Estimates of  $v$  which are efficient have been constructed by Chen (1985) if  $W$  and  $Z$  are independent, and for general  $W, Z$  by Chen (1988) and Robinson (1988a). See also Ritov and Bickel (1990) for the impossibility of attaining the information bound here for arbitrary  $W, Z$ .  $\square$

**Example 6. Logistic partial spline model.**

Suppose that

$$Q = \left\{ \text{Bernoulli} \left( \frac{e^\gamma}{1+e^\gamma} \right), \gamma \in R \right\}$$

and

$$r(z, w, \theta, t) = t(z) + \theta^T w$$

is a partial spline version of the logistic regression model. Suppose that  $t$  ranges over a class of functions of  $Z$  so large that  $\mathcal{B}_0$  given in condition (vi) is generated by  $Z$ . Then

$$\dot{I}(Y, r(Z, W, v, t)) = Y - p(Z, W)$$

and

$$\tilde{r}^T(Z, W, v, t) = W - E(W(Y - p(Z, W))^2 | Z) I^{-1}(Z),$$

where

$$I(Z) = E(p(Z, W)(1 - p(Z, W)) | Z)$$

and

$$p(Z, W) = \frac{\exp(t(Z) + \theta^T w)}{1 + \exp(t(Z) + \theta^T w)}.$$

Therefore, from (13) of proposition 2,

$$(19) \quad I^*(X, \theta, t) = \left( W - E(W p(Z, W)(1 - p(Z, W)) | Z) I^{-1}(Z) \right) (Y - p(Z, W)).$$

If  $W$  and  $Z$  are not functionally related,  $I^*$  given by (19) does not, in general, vanish so that we expect that, in this type of exponential family partial spline models, it should also be possible to estimate the parameters  $v$  at rate  $n^{-1/2}$ . For comparison, note that the efficient score function for  $v$  in the parametric model  $P_0$  with  $t(z) = \eta^T z$  is

$$(20) \quad I^*(X, P_0 | v, P_0) = (W - AB^{-1}ZC^{-1})(Y - p(Z, W)),$$

where

$$\begin{aligned}
 A &\equiv E(Wp(Z,W)(1 - p(Z,W))Z^T), \\
 B &\equiv EZZ^T, \\
 C &\equiv \{E(p(Z,W)(1 - p(Z,W))Z^T)\}B^{-1}Z,
 \end{aligned}$$

and the loss of information for  $v$  in  $P$  relative to  $P_0$  is determined by the difference between (19) and (20). □

*Extensions*

In the first three examples of regression model II,  $Q$  is the Gaussian location model. It is natural to combine types I and II and permit the distribution  $G$  of  $\epsilon$  to be arbitrary. This adds to  $\dot{P}_{23}$  of the model II the space  $\{b(\epsilon) : Eb(\epsilon) = 0\}$ .  $I^*$  given by (15) is easily seen to be orthogonal also to the last one of these spaces. Therefore, in examples 3–5 and, more generally, in situations in which  $Q$  is the semiparametric location (group) model and assumptions (iv)–(vii) hold, adaptation to lack of knowledge of the distribution of  $\epsilon$  is possible.

This result does not apply to estimation of  $\mu$  in example 5, since we need to add the side condition  $E t(Z) = 0$ . It may be shown that adaptation to lack of knowledge of the distribution  $G$  of  $\epsilon$  is still possible at least if  $G$  is taken symmetric about 0, so that  $\mu$  is identifiable. However, lack of knowledge of the distribution of  $Z$  may be shown to lead to an information bound of  $Var[\tilde{t}(Z)] + \sigma^2/I_g$ , where  $I_g \equiv \int [g']^2/g$ , rather than the  $\sigma^2/I_g$  achievable in principle if  $H$  is known. This more conservative bound is achievable, at least if  $G$  is normal, by ignoring the  $Z_i$ —that is, using the estimate  $\bar{Y}$ .

Suppose that  $Q$  is a semiparametric location and scale model,

$$Q = \{G(\frac{\cdot - \mu}{\sigma}) : G \in G, \mu \in R, \sigma > 0\}$$

and  $r(Z, \theta) = (\mu(Z, v, t), \sigma(Z, \eta, t))^T$ , but otherwise we are in the situation of proposition 2. Then  $\dot{P}_2$  has added to it a linear space of functions of the form  $(\epsilon(g'/g)(\epsilon) + 1)b(Z)$ . If  $G$  is restricted to symmetric distributions, then  $I^*$  given by (15) is orthogonal to any such function and adaptation to total lack of knowledge of the functional form of  $\sigma(Z)$  should be possible. A result of this type with  $\mu$  linear in  $Z$ ,  $G$ , Gaussian, and smoothness conditions on  $\sigma$  has been proved by Carroll (1982). We conjecture that neither the linearity nor the smoothness conditions, nor the requirement of Gaussianity are essential. This appears to contradict an example of Fuller and Rao (1978), who show that, in the Gaussian case, when there are several observations for each observed value of  $Z$  the estimate obtained by naively substituting the natural estimates of  $\sigma^2(Z)$  in the weighted least squares estimate of  $v$  may not be efficient. All that this result indicates is that smoothing is required when estimating  $\sigma^2(Z)$ . We conjecture that the use of nonparametric regression estimates of  $\sigma^2(\cdot)$  such as those of Stone (1977) will lead to adaptive estimation of  $v$  without any condition on  $\sigma$  other than finiteness and nonsingularity of the information matrix of  $v$  for this

model. On the other hand, symmetry is crucial for adaptation. If  $\sigma(Z)$  is "arbitrary" and  $G$  is not symmetric, the efficient score function for  $v$  involves the projection of  $I^*$  given by (15) on the space of all functions of the form  $\{\varepsilon(g'/g)(\varepsilon) + 1\}b(Z)$ . An uninformative explicit form can be given in this case also.

#### 4.4 BIASED SAMPLING MODELS

As we indicated in section 1, a biased sampling model requires three ingredients:  $\mathbf{Q} = \{Q_{(\theta,G)} : \theta \in \Theta, G \in \mathbf{G}\}$ , a semiparametric or parametric model on  $\mathbf{X}$ , a set of "stratum" weight functions  $w_i : \mathbf{X} \rightarrow R^+, i = 1, \dots, s$ , and selection probabilities  $\lambda_i, i = 1, \dots, s$ , (with  $\sum_{i=1}^s \lambda_i = 1$ ). The  $w_i$  and  $\lambda_i$  are assumed known. In biased sampling we first select the stratum  $i$  with probability  $\lambda_i$  and then observe  $x \in \mathbf{X}$  not with probability  $dQ_{(\theta,G)}(x)$ , but rather with the weighted probability

$$\frac{w_i(x)}{W_i(\theta,G)} dQ_{(\theta,G)}(x),$$

where

$$(1) \quad W_i(\theta,G) = \int w_i(x) dQ_{(\theta,G)}(x).$$

Formally, the *biased sampling model*  $\mathbf{P} = \{P_{(\theta,G)} : \theta \in \Theta, G \in \mathbf{G}\}$  is defined on  $\{1, \dots, s\} \times \mathbf{X}$ . Suppose that  $\mathbf{Q} \ll \mu$  and write  $q(\cdot, \theta, G)$  for  $dQ/d\mu$ . Then, an observation from  $P$ , written  $(I, X)$ , has density with respect to counting measure  $\times \mu$  given by

$$(2) \quad p(i,x,\theta,G) = \lambda_i \frac{w_i(x)}{W_i(\theta,G)} q(x,\theta,G).$$

The  $W_i(\theta,G)$  may, in some cases, be known. An early view of this class of models in econometrics is given by Heckman (1976). They have also arisen earlier in the sample survey literature; see Cochran (1977), for example.

#### Examples

We shall see that the tangent spaces of  $\mathbf{Q}$  and  $\mathbf{P}$  can be simply related and projections for  $\mathbf{P}$  can be calculated in many interesting cases. Before doing so, we identify a number of important examples from the econometric, survival analysis, and related literatures, which fit into this framework.

##### Example 1. Case control studies.

Prentice and Pyke (1979) describe a biased sampling model arising in case control studies. Patients for whom a variety of covariates  $Z$  are available are classified into one of  $s$  disease categories. If  $Y$  denotes the disease category, it is customary to assume a parametric model for the conditional distribution of  $Y$  given  $Z = z$ ,  $P(Y = y | Z = z) = f(z,y,\theta)$ . The most common choice is the logistic regression model

$$(3) \quad f(z, y, \theta) = \exp(\mu_y + v_y^T z) / \sum_j \exp(\mu_j + v_j^T z),$$

$$\theta = (v^T, \mu^T)^T, \quad \mu = (\mu_1, \dots, \mu_s)^T \in R^s, \quad v = (v_1, \dots, v_s) \in R^{p \times s},$$

where, as usual, we take  $\mu_s = 0$  and  $v_s = 0$  for identifiability. If we suppose  $Z_{p \times 1}$  is distributed according to  $H \in \mathbf{H}$ , with corresponding density  $h$  with respect to  $m$ , we arrive at the semiparametric random sampling model  $\mathbf{Q}$  with density

$$(4) \quad q(z, y, \theta, H) = f(z, y, \theta)h(z).$$

In a case control study we sample from the population of disease categories with specified probabilities and then observe the resulting covariates. That is, we let  $I = Y$ , select category  $i$  with probability  $\lambda_i$  and then  $Z$  according to the conditional distribution of  $Z$  given  $Y = i$ . This is a biased sampling model with  $x = (z, y)$ . Here

$$(5) \quad p(i, z, y, \theta, H) = \lambda_i \frac{w_i(z, y)}{W_i(\theta, H)} f(z, y, \theta)h(z),$$

where

$$w_i(z, y) = 1_{[y=i]}, \quad W_i(\theta, H) = \int f(z, i, \theta)h(z) dm(z).$$

If  $f(\cdot, \theta)$  is given by (3),  $\mu_1, \dots, \mu_{s-1}$  which are identifiable in the random sampling model  $\mathbf{Q}$ , are no longer identifiable in  $\mathbf{P}$  since we can always rewrite (5) to have  $\mu_1 = \dots = \mu_{s-1} = 0$  by replacing  $h(z)$  with  $\tilde{h}(z) = h(z) / \sum_j \exp\{\mu_j + v_j^T z\}$ , and  $h$  is only specified up to a scale factor. However, comparative odds ratios,

$$\begin{aligned} \left\{ \frac{f(z, i, \theta)}{f(z, j, \theta)} \right\} \left\{ \frac{f(z', i, \theta)}{f(z', j, \theta)} \right\}^{-1} &= \exp\{(v_i - v_j)^T(z - z')\} \\ &= \left\{ \frac{p(i, z, y, \theta, H)}{p(j, z, y, \theta, H)} \right\} \left\{ \frac{p(i, z', y, \theta, H)}{p(j, z', y, \theta, H)} \right\}^{-1} \end{aligned}$$

are identifiable, and hence so are differences  $v_i - v_j$ .

In this model it is often reasonable to suppose  $W_i(\theta, H)$ , the proportion of individuals in category  $i$ , to be known from previous studies. This knowledge leads to a lower information bound and facilitates the construction of  $\sqrt{n}$ -consistent estimates.  $\square$

### Example 2. Choice-based sampling.

The model of example 1 has been developed and pursued in the econometric literature, by Manski and Lerman (1977) and Cosslett (1981). The categories of the variable  $Y$  now correspond to different economic or other choices,  $Z$  is a vector of covariates as before, and  $P(Y = j | Z = z) = f(z, j, \theta)$ ,  $1 \leq j \leq m$ . However, strata are now found in the space of categories which may contain more than one member. That is, if the range of  $Y$  is  $\{1, \dots, m\}$ , subsets  $S_1, \dots, S_s$ , not



necessarily disjoint, are formed such that  $\cup_{i=1}^s S_i = \{1, \dots, m\}$ . Sampling then proceeds by first choosing  $S_i$  with specified probability  $\lambda_i$  and then choosing  $(Y, Z)$  with probability

$$p(i, z, y, \theta, H) = \lambda_i \frac{1_{[y \in S_i]} f(z, y, \theta) h(z)}{W_i(\theta, H)}$$

where

$$W_i(\theta, H) = \sum_{y \in S_i} \int f(z, y, \theta) dH(z).$$

This model is of the form (2) with

$$w_i(z, y) = 1_{[y \in S_i]}.$$

□

### Example 3. Truncated regression and extensions.

Regression models where the dependent variable is truncated arise naturally in a variety of fields. An application to astronomy and theoretical contributions appear in Bhattacharya, Chernioff, and Yang (1983). A variety of other applications and an extension of the model are given in Jewell (1985).

The model specifies an ideal unobservable  $X' = (Z', Y')$  which follows a semiparametric linear regression model as in example 4.2.2,

$$(6) \quad Y' = \eta + v^T Z' + \varepsilon,$$

$Z' \sim H \in \mathbf{H}$ , with  $Z'$  independent of  $\varepsilon \sim G \in \mathbf{G}$ , the set of all absolutely continuous distributions. Truncation at  $y_0$  means that we only observe  $X'$  if  $Y' \leq y_0$ . In the astronomy example,  $Z'$  is the log velocity of a distant celestial object, while  $Y'$  is the negative log of the luminosity. Objects of low luminosity are not observed. In terms of (2),  $\mathbf{Q}$  is the model (6),  $s = 1$ , and  $w(z, y) = 1_{[y \leq y_0]}$ . Hence  $P_{(\theta, G, H)}$  has a density with respect to Lebesgue measure  $\times m$ , where  $m$  is a measure dominating  $\mathbf{H}$ , given by

$$(7) \quad p(z, y, \theta, G, H) = 1_{[y \leq y_0]} \frac{g(y - \eta - v^T z) h(z)}{W(\theta, G, H)},$$

where  $\theta = (v^T, \eta)^T$ ,  $g$  is the density of  $G$  and  $h$  that of  $H$ .  $W$  in this case is naturally unknown. Here  $\eta$  is unidentifiable. Jewell (1985) considers the more general model in which we sample from the  $i$ th of  $s$  strata  $S_i = (y_{i-1}, y_i]$  with probability  $\lambda_i$ , where  $y_0 = -\infty$ ,  $y_s = \infty$ . This is just model (2) with  $w_i(z, y) = w_i(y) = 1_{[y_{i-1} < y \leq y_i]}$ , and  $\mathbf{Q}$  the semiparametric linear regression model. In the examples he considers,  $W_i$  is often known. He gives a number of interesting examples from biostatistics and econometrics where these models seem appropriate. The theory of estimation in this model is developed in Bickel and Ritov (1991). □

### Example 4. Vardi's model and stratified sampling.

A final biased sampling model, which we shall explore further in chapter 6, is due to Vardi (1985). Here,  $\mathbf{Q}$  is taken as the set of all distributions on  $\mathbf{X}$  and interest centers on estimation of real- or vector-valued functions of  $\mathcal{Q}$  or  $\mathcal{Q}$  itself

when we observe  $(I, X) \sim P$  with

$$P(I = i, X \in A) = \lambda_i \frac{\int_A w_i(x) dQ(x)}{W_i(Q)}.$$

An important special case is *stratified sampling*; then  $w_i(x) = 1_{A_i}(x)$ ,  $i = 1, \dots, s$ , where  $A_1, \dots, A_s$  is a partition of the sample space  $\mathbf{X}$  ( $A_i \cap A_j = \emptyset$  for  $i \neq j$  and  $A_1 \cup \dots \cup A_s = \mathbf{X}$ ). In this case it is clear that  $Q$  is *not* identifiable if  $W_i(Q) = Q(A_i)$ ,  $i = 1, \dots, s$ , are unknown. If, however, the  $W_i(Q) = Q(A_i)$  are known, (as is essentially assumed in classical sampling theory), then  $Q$  is identifiable and functionals of  $Q$  such as the distribution function are regularly and efficiently estimable.  $\square$

We will return to these examples later in this section after developing formulae which simplify our calculation of efficient score functions.

### Score Functions and Tangent Spaces

We turn to estimation of Euclidean parameters in these models. Formally, from (2), if  $\dot{\mathbf{l}}(\cdot | \mathbf{P}) \equiv \dot{\mathbf{l}}(\cdot | \theta, \mathbf{P})$  is the derivative with respect to  $\theta$  of the log-likelihood for  $\mathbf{P}$  and  $\dot{\mathbf{l}}(\cdot | \mathbf{Q}) \equiv \dot{\mathbf{l}}(\cdot | \theta, \mathbf{Q})$  is defined similarly and we assume that  $W_i(\theta, G)$  is unknown, then

$$\begin{aligned} (8) \quad \dot{\mathbf{l}}(i, x | \mathbf{P}) &= \frac{\dot{p}}{p}(i, x, \theta, G) = \frac{\dot{q}}{q}(x, \theta, G) - \frac{\dot{W}_i}{W_i}(\theta, G) \\ &= \frac{\dot{q}}{q}(x, \theta, G) - \frac{\int (\dot{q}/q)(x, \theta, G) w_i(x) q(x, \theta, G) d\mu(x)}{W_i(\theta, G)} \\ &= \dot{\mathbf{l}}(x | \mathbf{Q}) - E_P(\dot{\mathbf{l}}(X | \mathbf{Q}) | I = i). \end{aligned}$$

Proceeding formally in the same way we expect that if  $v \in \dot{\mathbf{Q}}_2$ , then

$$(9) \quad v - E_P(v(X) | I) \in \dot{\mathbf{P}}_2$$

and, conversely, we expect every element of  $\dot{\mathbf{P}}_2$  to be of the form (9) (if the  $W_i$  are unknown).

For (8) and (9) to make sense in general we need to know that  $v \in \dot{\mathbf{Q}}$  implies that  $v \in L_2(P)$ . We make the blanket assumption

$$L_2(P) = L_2(Q)$$

in the sense that  $P$  and  $Q$  dominate each other and  $\int a^2 dP < \infty$  if and only if  $\int a^2 dQ < \infty$ . This assumption holds if

$$(10) \quad \delta < \sum_{i=1}^s \frac{\lambda_i}{W_i} w_i(x) < \delta^{-1}$$

for all  $x$  and some  $\delta > 0$ . All our models save the basic truncated regression

model satisfy this condition. Truncated regression is most easily treated directly; see section 6.6.

It is important to note that in (8) and (9) all expectations are carried out under the model  $\mathbf{P}$  and that the distribution of  $\dot{\mathbf{I}}(X | \mathbf{Q})$  is not the same as it would be under  $\mathbf{Q}$ ; in particular,  $\dot{\mathbf{I}}(X | \mathbf{Q})$  may not have mean zero under  $P$ . Further, objects such as  $\dot{\mathbf{Q}}$  are viewed simply as subspaces of  $L_2(P)$ .

It will be convenient to use the scores for the model  $\mathbf{Q}$  centered under  $P \in \mathbf{P}$ . We define

$$(11) \quad \dot{\mathbf{I}}^c(x) \equiv \frac{\dot{q}}{q}(x, \theta, G) - E_P \left( \frac{\dot{q}}{q}(X, \theta, G) \right)$$

and in general use the superscript  $c$  to denote centering by  $E_P$ . So

$$\dot{\mathbf{Q}}_2^c \equiv \{v - E_P(v(X)) : v \in \dot{\mathbf{Q}}_2\}.$$

Note that from (8)

$$(12) \quad \dot{\mathbf{I}}(i, x | \theta, \mathbf{P}) = \dot{\mathbf{I}}^c(x) - E_P(\dot{\mathbf{I}}^c | I = i),$$

and, more generally, from (9), if  $v \in \dot{\mathbf{Q}}_2^c$ , we expect

$$v - E_P(v | I = i) \in \dot{\mathbf{P}}_2.$$

We assume the correctness of these formal calculations. In theorem 3 we argue that a strengthening of (10) partially justifies them.

**Theorem 1.** Suppose that:

- (i) (12) holds.
- (ii)  $\dot{\mathbf{P}}_2 = \{v - E_P(v(X) | I) : v \in \dot{\mathbf{Q}}_2^c\}$ .

Let  $\mathbf{V} \equiv \dot{\mathbf{Q}}_2^c + \mathbf{L}$ , where  $\mathbf{L} \equiv \{a(I) : a \in L_2^0(P)\}$ . Then  $\mathbf{V} = \dot{\mathbf{P}}_2 + \mathbf{L}$ , and the efficient score function for  $\theta$  in  $\mathbf{P}$  is given by

$$(13) \quad \dot{\mathbf{I}}^*(I, X | \theta, \mathbf{P}) = \dot{\mathbf{I}}^c(X) - \Pi_0(\dot{\mathbf{I}}^c(X) | \mathbf{V}),$$

where  $\dot{\mathbf{I}}^c$  is given in (11).

**Proof.** By hypothesis we can represent

$$\mathbf{V} = \mathbf{L} + \{v(X) - E_P(v(X) | I) : v \in \dot{\mathbf{Q}}_2^c\} = \mathbf{L} + \dot{\mathbf{P}}_2.$$

But  $\mathbf{L} \perp \dot{\mathbf{P}}_2$ , so that

$$(a) \quad \begin{aligned} \Pi_0(\dot{\mathbf{I}}^c(X) | \mathbf{V}) &= \Pi_0(\dot{\mathbf{I}}^c(X) | \mathbf{L}) + \Pi_0(\dot{\mathbf{I}}^c(X) | \dot{\mathbf{P}}_2) \\ &= E_P(\dot{\mathbf{I}}^c(X) | I) + \Pi_0(\dot{\mathbf{I}}^c(X) | \dot{\mathbf{P}}_2) \end{aligned}$$

by (A.2.9) and (A.3.1). Now  $\mathbf{L} \perp \dot{\mathbf{P}}_2$ , (i) and (12) imply that  $\Pi_0(\dot{\mathbf{I}}(I, X | \theta, \mathbf{P}) | \dot{\mathbf{P}}_2) = \Pi_0(\dot{\mathbf{I}}^c(X) | \dot{\mathbf{P}}_2)$ . Hence it follows from (a) that

$$\begin{aligned} \dot{\mathbf{I}}^*(I, X | \theta, \mathbf{P}) &= \dot{\mathbf{I}}(I, X | \theta, \mathbf{P}) - \Pi_0(\dot{\mathbf{I}}(I, X | \theta, \mathbf{P}) | \dot{\mathbf{P}}_2) \\ &= \dot{\mathbf{I}}^c(X) - E_P(\dot{\mathbf{I}}^c | I) - \{\Pi_0(\dot{\mathbf{I}}^c | \mathbf{V}) - E_P(\dot{\mathbf{I}}^c | I)\} \end{aligned}$$

$$= \dot{I}^c(X) - \Pi_0(\dot{I}^c | V).$$

Thus (13) holds. □

**Corollary 1.** Under the assumptions of proposition 1,  $I^*$  is unchanged if we assume the  $\lambda_i$  unknown.

**Proof.** If the  $\lambda_i$  are assumed unknown and we add them as a third parameter then

$$\dot{P}_{23} = \overline{\dot{Q}_2^c} + L = V.$$

The result follows from theorem 1. □

We apply theorem 1 to compute efficient score functions and information bounds in our examples. We introduce the notation

$$ACE(h | U, V) = \Pi_0(h | L_U + L_V),$$

where  $L_U \equiv \{a(U) : a(U) \in L_2^0(P)\}$ , and

$$ACE(h | U, V) = ACE_1(h) + ACE_2(h),$$

where  $ACE_1$  is a function of  $U$ ,  $ACE_2$  a function of  $V$ . The notation can and will be extended to more than two variables. Further,  $ACE(h | U) = E(h | U) - E(h)$ . ACE is a useful mnemonic for the Breiman-Friedman algorithm, which is one of the ways of computing  $\Pi_0(h | L_U + L_V)$ . For further details see appendix A.4.

**Examples 1 and 2 continued.**

In this case the assumptions of theorem 1 are easily verified. Note  $\dot{I}^c = (\dot{f}/f)(X, \theta) - E(\dot{f}/f)(X, \theta)$  and  $\dot{Q}_2^c = \{a(Z) \in L_2^0(P)\}$ . Therefore,  $V$  of theorem 1 is just  $L_Z + L_I$ . We conclude that

$$(14) \quad I^*(I, X | \theta, P) = \frac{\dot{f}}{f}(X, \theta) - E_P \frac{\dot{f}}{f}(X, \theta) - ACE\left(\frac{\dot{f}}{f}(X, \theta) | Z, I\right)$$

by the definition of ACE and  $ACE(\text{constant} | Z, I) = 0$ .

We specialize to the logistic regression model where  $Y = I$ . If  $j \leftrightarrow \mu_{j-s}$  for  $j = s+1, \dots, 2s-1$ , we get as expected

$$I_j^*(I, X | \theta, P) = 0.$$

This follows from (14) since

$$\frac{\partial}{\partial \mu_j} \log f(Z, Y, \theta) = \varepsilon_j - f(Z, j, \theta), \quad \text{where } \varepsilon_j \equiv 1_{[I=j]},$$

and  $ACE(a(I) + b(Z) | I, Z) = a(I) + b(Z)$  for all  $a(I) + b(Z) \in L_2^0(P)$ . On the other hand, if  $I_j^*$  is the component of  $I^*(\cdot | \theta, P)$  corresponding to  $v_j$ ,  $j = 1, \dots, s-1$ ,

$$(15) \quad \begin{aligned} I_j^*(I, X) &= Z(\varepsilon_j - f(Z, j, \theta)) - E_P Z(\varepsilon_j - f(Z, j, \theta)) \\ &\quad - ACE(Z(\varepsilon_j - f(Z, j, \theta)) | Z, I) \\ &= Z \varepsilon_j - E_P Z \varepsilon_j - ACE(Z \varepsilon_j | Z, I). \end{aligned}$$

For simplicity, we consider the binomial case,  $s = 2$ . By (15) the matrix  $I(\theta | v_1, P)$  is nonsingular unless there exist  $\lambda_{p \times 1} \neq 0$ ,  $a(Z)$ ,  $b(\epsilon_1)$  such that

$$(16) \quad \lambda^T Z \epsilon_1 = a(Z) + b(\epsilon_1).$$

This cannot hold unless the distribution of  $Z$  is singular (concentrated on a  $(p - 1)$ -dimensional hyperplane). To check this, set  $\epsilon_1 = 0$  and  $1$  in (16) and get

$$\lambda^T Z = b(1) - b(0).$$

Since  $Y = I$  is finite we can exhibit  $ACE(Z_i \epsilon_j | Z, I)$  explicitly by theorem A.4.5. The formula is uninformative.  $\square$

**Example 3. Truncated regression and extensions, continued.**

In the Jewell model it is easy to see (compare example 4.2.2) that, if  $\epsilon = Y - \eta - v^T Z$ ,

$$\dot{Q}_2^c = \{a(\epsilon) + b(Z) : a(\epsilon), b(Z) \in L_2(P), Ea(\epsilon) = Eb(Z) = 0\},$$

while, if  $\theta = (v^T, \eta)^T$ ,

$$\frac{\dot{q}}{q}(X, \theta, G, H) = (Z^T \psi(\epsilon), \psi(\epsilon))^T,$$

where  $\psi \equiv -g'/g$ . If  $L_2(P) = L_2(Q)$ , we can apply proposition 1 to get

$$(17) \quad I^*(I, X | \theta, P) = (Z^T \psi(\epsilon) - ACE(Z^T \psi(\epsilon) | \epsilon, Z, I), 0)^T.$$

It is easy to see directly that (17) also holds for truncated regression but simplifies since  $I \equiv 1$ ,  $w_1(Z, Y) = 1_{[\epsilon \leq y_0 - \eta - v^T Z]}$ , and  $\epsilon$  and  $Z$  are independent under  $Q$ . An explicit computation for this case is given in section 6.4. In the general Jewell model it appears that explicit formulae for (17) are available only if  $Z$  is also finite. In that case we may again apply theorem A.4.5.  $\square$

*The Effect of Knowledge of the  $W_i$*

Knowledge of the  $W_i$  does make a difference. Suppose that  $W_i = W_{i0}$ ,  $i = 1, \dots, s$ , are known. We want to calculate the efficient score function or influence function for  $\theta$  in the submodel

$$(18) \quad P_0 \equiv \{P_{(\theta, G)} \in P : W_i(\theta, G) = W_{i0}, i = 1, \dots, s\}.$$

To see the form of  $\dot{P}_0$ , suppose that  $q_\eta \equiv q(\cdot, \theta_\eta, G_\eta)$  is a path in  $Q$  with tangent  $a \in \dot{Q}$  and  $P_{(\theta_\eta, G_\eta)} \in P_0$ . Then, we find formally that

$$(19) \quad \int a(x) w_i(x) q(x, \theta_0, G_0) d\mu(x) = 0, \quad \text{for } i = 1, \dots, s.$$

Let

$$(20) \quad \bar{w}_i(x) \equiv \frac{w_i(x)}{\sum_{j=1}^s \lambda_j w_j(x) W_{j0}^{-1}} = P(I = i | X = x).$$

Then (19) can be rewritten as

$$(21) \quad E_P(a(X)\bar{w}_i(X)) = 0 \quad \text{for } i = 1, \dots, s,$$

or,

$$(22) \quad E_P(a(X) | I) = 0.$$

This suggests that

$$(23) \quad \dot{P}_0 = \{a \in \dot{Q} : \langle a, \bar{w}_i \rangle_0 = 0, \quad i = 1, \dots, s\}.$$

Note that if  $a$  is a member of the right side in (23) then (22) yields  $E_P a(X) = 0$ , so that  $\dot{Q}$  may be replaced in (23) by  $\dot{Q}^c$ .

**Theorem 2.** Let  $\dot{I} = (\dot{q}/q)(\cdot, \theta, G)$ , and assume that conditions (i) and (ii) of theorem 1 hold. Suppose that:

- (i)  $\dot{Q}^c = [\dot{I}^c] + \dot{Q}_2^c$ .
- (ii)  $\dot{P}_0$  is given in (23).
- (iii)  $\dim[\Pi_0(\bar{w}_i | \dot{Q}_2)] = \dim[\bar{w}_i]$ .

Let

$$(24) \quad \Gamma^*(X | \theta, P_0) = \dot{I}(X) - \Pi_0(\dot{I}(X) | \dot{Q}_2^c) - \sum_{j=1}^s c_j \Pi_0(\bar{w}_j | \dot{Q}_2^c),$$

where  $c_j \in R^p$ ,  $j = 1, \dots, s$ , satisfy

$$(25) \quad \sum_{j=1}^s c_j E_P \{ \bar{w}_i \Pi_0(\bar{w}_j | \dot{Q}_2^c) \} = E_P \left( \bar{w}_i \left\{ \dot{I} - \Pi_0(\dot{I} | \dot{Q}_2^c) \right\} \right)$$

for  $i = 1, \dots, s$ , and let

$$I(\theta, P_0) = E_P(\Gamma^* \Gamma^{*T}).$$

Then, the efficient influence function for  $\theta$  is given by

$$\tilde{I}(X | \theta, P_0) = I(\theta, P_0)^{-1} \Gamma^*(X | \theta, P_0).$$

**Proof.** That (25) has a solution follows from the dimensionality assumption (iii). (The solution is unique if the  $\bar{w}_i$ ,  $i = 1, \dots, s$ , are linearly independent or if  $\dim[\Pi_0(\bar{w}_i | \dot{Q}_2^c)] = s - 1$ .) From (i), (ii), (24), and (25) it follows that  $\Gamma^* \in \dot{P}_0$ . If  $\{P_{(\theta, G_0)}\}$  is a parametric submodel of  $P_0$  we obtain in the usual way the score function for  $\theta$ ,  $\dot{I} + v$  where  $v \in \dot{Q}_2$  and

$$(a) \quad \langle \dot{I} + v, \bar{w}_i \rangle_0 = 0, \quad i = 1, \dots, s.$$

Since, by (ii), any  $v \in \dot{Q}_2$  such that (a) holds is possible we obtain the efficient score function by minimizing  $\|\dot{I} + v\|_0^2$  subject to (a). Write

$$\begin{aligned} \dot{I} + v &= \dot{I} - \Pi_0(\dot{I} | \dot{Q}_2^c) + w \\ &= \dot{I} - \Pi_0(\dot{I} | \dot{Q}_2^c) - \sum_{j=1}^s c_j \Pi_0(\bar{w}_j | \dot{Q}_2^c) + z, \end{aligned}$$

where  $w$  ranges over  $\dot{Q}_2^c$ ,  $c_j$  over  $R^p$ , and  $z$  over

$$[\Pi_0(\bar{w}_i | \dot{Q}_2^c), 1 \leq j \leq s]^\perp \cap \dot{Q}_2^c.$$

Then (a) becomes for  $i = 1, \dots, s$ ,

$$(b) \quad \langle \dot{I} - \Pi_0(\dot{I} | \dot{Q}_2^c), \bar{w}_i \rangle_0 = \sum_{j=1}^s c_j \langle \Pi_0(\bar{w}_j | \dot{Q}_2^c), \bar{w}_i \rangle_0,$$

which is just (25). So

$$\dot{I} + v = I^* + z$$

where  $z \perp I^*$ . The result follows. □

**Example 1 continued: Known stratum weights.**

Suppose that the  $W_i(\theta, H) = W_{i0}$ ,  $i = 1, \dots, s$ , are known. In this example,

$$\bar{w}_i(X) = \frac{1_{[Y=i]}}{\sum_j 1_{[Y=j]} \lambda_j / W_j} = \frac{W_i}{\lambda_i} 1_{[Y=i]}.$$

Since  $\dot{Q}_2^c = \{a(Z) \in L_2^0(P)\}$ ,

$$[\Pi_0(\bar{w}_i | \dot{Q}_2^c)] = [p(i | Z) - p(i)],$$

where

$$(26) \quad p(i | Z) = \frac{(\lambda_i / W_i) f(Z, i, \theta)}{\sum_j (\lambda_j / W_j) f(Z, j, \theta)}, \quad p(i) = E p(i | Z) = W_i.$$

We now specialize to the logistic regression model with  $s = 2$  and  $W_1 = 1 - W_2$  known. We assume that  $\mu_1 = 0$ . This does not restrict our argument since we can always replace  $Z$  by  $(1, Z)^T$  and treat  $\mu$  as  $v_1$ . Let  $\tilde{p}(Z) = p(1 | Z)$  and assume that  $Var[\tilde{p}(Z)] > 0$ . Then

$$\begin{aligned} \dot{I}_1 - \Pi_0(\dot{I}_1 | \dot{Q}_2^c) &= Z(\varepsilon_1 - f(Z, 1, \theta)) - E(Z(\varepsilon_1 - f(Z, 1, \theta)) | Z) \\ &= Z(\varepsilon_1 - \tilde{p}(Z)). \end{aligned}$$

Hence, by theorem 2,

$$\begin{aligned} I_1^* &= Z(\varepsilon_1 - \tilde{p}(Z)) - c_1 \{p(1 | Z) - p(1)\} - c_2 \{p(2 | Z) - p(2)\} \\ &= Z(\varepsilon_1 - \tilde{p}(Z)) - c_1^* \tilde{p}(Z) - c_2^*, \end{aligned}$$

where  $c_1^*$  and  $c_2^*$  are to be chosen so that  $I_1^* \perp [\bar{w}_1, \bar{w}_2] = [\varepsilon_1, \varepsilon_2] = [1, \varepsilon_1]$  since  $\varepsilon_1 + \varepsilon_2 = 1$ . This yields (dropping the  $P$  subscripts)

$$(27) \quad I_1^* = Z(\varepsilon_1 - \tilde{p}(Z)) - \frac{E[Z\tilde{p}(Z)(1 - \tilde{p}(Z))]}{Var[\tilde{p}(Z)]} \{\tilde{p}(Z) - E\tilde{p}(Z)\}.$$

The two terms on the right of (27) are orthogonal, so that

$$(28) \quad \begin{aligned} I(\theta, P) &= E\{ZZ^T \tilde{p}(Z)(1 - \tilde{p}(Z))\} \\ &\quad + \frac{E[Z\tilde{p}(Z)(1 - \tilde{p}(Z))]E[Z^T \tilde{p}(Z)(1 - \tilde{p}(Z))]}{Var[\tilde{p}(Z)]}. \end{aligned}$$

The difference between  $I_1^*$  given by (27) and that given by (15) can in principle be used to measure the information gained by knowledge of the  $W_i$ . We can also make two qualitative remarks.

**Remark 1.** The intercept term, which was unidentifiable in the unconstrained model  $\mathbf{P}$  is identifiable in the constrained model  $\mathbf{P}_0$ . This follows from (28) under the assumption that  $Z$  is not concentrated on a hyperplane.

**Remark 2.** If  $v$  and  $Z$  are both real,  $P \in \mathbf{P}_0$ ,  $v(P) = 0$ , and  $EZ \neq 0$ , then  $P_0 \in \mathbf{P}_0$  is not regular. Formally,  $I(P | v, \mathbf{P}_0) = \infty$ , since then  $\tilde{p}(Z)$  is constant. In this case there is no sequence of distributions in  $\mathbf{P}_0$  contiguous to  $P$ . Thus, independence of  $Z$  and  $Y$  is particularly easy to detect in the constrained model. See also Jagers, Oden, and Trulsson (1985).  $\square$

Our last example, in addition to having its own interest, shows how we can finesse the calculation of the  $c_j$  of theorem 2.

**Example 4. Vardi's model and stratified sampling, continued.**

Consider estimation of  $\mu \equiv E_Q X$  in the case that all  $W_i = W_{i0} = P(A_i)$  are known. (Recall that  $P$ , and hence  $\mu$ , is not identifiable if the  $W_i$  are unknown.) For simplicity, suppose  $P(|X| \leq M) = 1$  for all  $P \in \mathbf{P}$ . Save for the fact that a random rather than fixed number of observations is taken in each stratum this is just the classical stratified sampling model. The usual estimator of  $\mu$  for this model is  $\hat{\mu}_n \equiv \sum_{i=1}^s W_{i0} \bar{X}_i$  where  $\bar{X}_i$  is the mean of the observations in the  $i$ th stratum. We will derive the influence function of this estimator, and then show that it is, in fact, efficient. Let  $w_i(x) = 1_{A_i}(x)$  where  $A_1, \dots, A_s$  is the partition of  $\mathbf{X}$ , and let  $IP_n$  denote the empirical distribution of the  $X$ 's. Then we can write

$$\begin{aligned} \sum_{i=1}^s W_{i0} \bar{X}_i &= \sum_{i=1}^s W_{i0} \frac{\int x w_i(x) dIP_n(x)}{\int w_i(x) dIP_n(x)} \\ &= \sum_{i=1}^s W_{i0} E_P(X | I = i) \\ &\quad + \sum_{i=1}^s W_{i0} \left\{ \frac{\int x w_i(x) dIP_n(x)}{\int w_i(x) dIP_n(x)} - E_P(X | I = i) \right\} \\ &= E_Q X + \sum_{i=1}^s \frac{W_{i0}}{\lambda_i} \int (x - E_P(X | I = i)) w_i(x) d(IP_n - P)(x) \\ &\quad - \sum_{i=1}^s \frac{W_{i0}}{\lambda_i} \left\{ \frac{\int x w_i(x) dIP_n(x)}{\int w_i(x) dIP_n(x)} - E_P(X | I = i) \right\} \\ &\quad \cdot \left\{ \int w_i(x) dIP_n(x) - \lambda_i \right\}. \end{aligned}$$

It is easy to see that the last term on the right-hand side is  $O_p(n^{-1})$ . Hence the



influence function of  $\hat{\mu}_n$  is

$$\tilde{I}(x) = \sum_{i=1}^s \frac{W_{i0}}{\lambda_i} [x - E_P(X | I = i)] w_i(x).$$

Note that for this model  $I$  is a function of  $X$  since  $X \in A_i$  if and only if  $I = i$ . It is easy to check that  $\tilde{I} \perp [w_i] = [\bar{w}_i]$ , and hence that  $\tilde{I}$  is in  $\dot{P}_0$  (cf. (23)). By proposition 3.3.1,  $\tilde{I}$  is the efficient influence function.  $\square$

We conclude with a technical note.

**Theorem 3.**

A. Suppose that  $\sum_{i=1}^s \lambda_i w_i(\cdot)/W_i$  is bounded. Then

$$\dot{P} \supset \{a - E_P(a | I) : a \in \dot{Q}^c\}.$$

B. Suppose further that  $\sum_{i=1}^s \lambda_i w_i(\cdot)/W_i$  is bounded away from 0 and  $\int w_i(x) w_j(x) dP(x) > 0$ . Then

$$\dot{P} = \{a - E_P(a | I) : a \in \dot{Q}^c\}.$$

**Proof.** A. Since  $\dot{P} \subset L_2(P)$ , we may assume without loss of generality that  $\lambda \equiv \min\{\lambda_1, \dots, \lambda_s\} > 0$ . Suppose that  $Q \in \mathcal{Q}$  and that  $\{Q_\eta\}$  is a curve in  $\mathcal{Q}$  with tangent  $a \in L_2^0(Q)$ :

$$(a) \quad \left\| \frac{q_\eta^{1/2} - q^{1/2}}{\eta} - \frac{1}{2} a q^{1/2} \right\| \rightarrow 0 \quad \text{as } \eta \rightarrow 0.$$

Note that

$$\begin{aligned} (b) \quad & \left| \frac{p_\eta^{1/2} - p^{1/2}}{\eta} - \frac{1}{2} (a - E_P(a | I = i)) p^{1/2} \right| \\ & \leq \lambda_i^{1/2} \left| \eta^{-1} \left( \frac{w_i^{1/2} q_\eta^{1/2}}{\alpha_i^{1/2}} - \frac{w_i^{1/2} q^{1/2}}{\alpha_i^{1/2}} \right) - \frac{1}{2} a \frac{w_i^{1/2} q^{1/2}}{\alpha_i^{1/2}} \right| \\ & \quad + \lambda_i^{1/2} \left| \eta^{-1} \left( \frac{w_i^{1/2} q_\eta^{1/2}}{\alpha_{\eta i}^{1/2}} - \frac{w_i^{1/2} q_\eta^{1/2}}{\alpha_i^{1/2}} \right) + \frac{1}{2} E_P(a | I = i) \frac{w_i^{1/2} q_\eta^{1/2}}{\alpha_i^{1/2}} \right| \\ & \quad + \frac{1}{2} \lambda_i^{1/2} \left| E_P(a | I = i) \left( \frac{w_i^{1/2} q_\eta^{1/2}}{\alpha_i^{1/2}} - \frac{w_i^{1/2} q^{1/2}}{\alpha_i^{1/2}} \right) \right| \end{aligned}$$

$$\equiv A + B + C,$$

where  $\alpha_i = W_i(Q)$  and  $\alpha_{\eta i} = W_i(Q_\eta)$ . The integral with respect to counting measure  $\times \mu$  of the square of A converges to 0 by (a) and our boundedness assumption. Upon noting that

$$|E_P(a | I = i)| \leq \lambda^{-1} \sum_{i=1}^s \lambda_i \int \frac{w_i}{\alpha_i} |a| dQ$$

is finite for  $a \in L_2^0(Q)$ , the same is true for  $C$ .

To handle  $B$ , we need to show that

$$(c) \quad \frac{1}{\eta} \left( \frac{1}{\alpha_{\eta i}^{1/2}} - \frac{1}{\alpha_i^{1/2}} \right) + \frac{1}{2} \frac{E_P(a \mid I=i)}{\alpha^{1/2}} \rightarrow 0 \quad \text{as } \eta \rightarrow 0,$$

and this will follow from  $d\alpha_{\eta i}/d\eta|_{\eta=0} = \int a w_i dQ$ . But

$$\begin{aligned} & \left| \eta^{-1}(\alpha_{\eta i} - \alpha_i) - \int a w_i dQ \right| \\ &= \left| \int w_i \left( \frac{q_\eta - q}{\eta} - a q \right) d\mu \right| \\ &\leq \|w_i\|_\infty \left\{ \int \left| (q_\eta^{1/2} + q^{1/2}) \left( \frac{q_\eta^{1/2} - q^{1/2}}{\eta} - \frac{1}{2} a q^{1/2} \right) \right| d\mu \right. \\ &\quad \left. + \frac{1}{2} \int |a q^{1/2} (q_\eta^{1/2} - q^{1/2})| d\mu \right\} \\ &\rightarrow 0 \quad \text{as } \eta \rightarrow 0 \end{aligned}$$

by (a) and the Cauchy-Schwarz inequality. Note that  $\|w_i\|_\infty < \infty$  is implied by  $\lambda > 0$ .

B. Let  $\{P_\eta\}$  be a path in  $\mathbf{P}$  with a tangent at  $P = P_0$ . Let  $Q_\eta \longleftrightarrow P_\eta$ . For any measurable set  $A$ , let  $\bar{P}_\eta(A) \equiv \sum_{i=1}^s P_\eta(i, A)$ . Note that

$$(d) \quad Q_\eta(A) = \int_A \frac{1}{\sum_{i=1}^s (\lambda_{i\eta}/W_i(Q_\eta)) w_i(x)} d\bar{P}_\eta(x),$$

where  $\lambda_{i\eta} = P_\eta(I=i)$ ,  $i = 1, \dots, s$ . In particular,

$$(e) \quad W_i(Q_\eta) = \int \frac{w_i(x)}{\sum_{j=1}^s (\lambda_{j\eta}/W_j(Q_\eta)) w_j(x)} d\bar{P}_\eta(x).$$

Let  $V = \{v \in R^s : \sum v_i = 0\}$ . Define  $v(\eta) \in V$  by

$$v_i(\eta) = \log\left(\frac{\lambda_{i\eta}}{W_i(Q_\eta)}\right) - s^{-1} \sum_j \log\left(\frac{\lambda_{j\eta}}{W_j(Q_\eta)}\right). \quad \text{Then it follows from (e) that } f(\eta, v(\eta)) = 0, \text{ where}$$

$$(f) \quad f_i(\eta, v) = \lambda_{i\eta} - \frac{\int e^{v_i} w_i(x) d\bar{P}_\eta(x)}{\sum_j e^{v_j} w_j(x)}.$$

Now, for any  $u \in V$ ,  $|u| = 1$ ,

$$u^T \frac{\partial}{\partial \lambda} f(\eta, v + \lambda u) \Big|_{\lambda=0}$$

$$\begin{aligned}
 &= - \int \left\{ \frac{\sum_i e^{v_i} w_i(x) u_i^2}{\sum_i e^{v_i} w_i(x)} - \left( \frac{\sum_i e^{v_i} w_i(x) u_i}{\sum e^{v_i} w_i(x)} \right)^2 \right\} dP_\eta(x) \\
 &= - \int \sum_{i=1}^s e^{v_i} w_i(x) \left( u_i - \frac{\sum_j e^{v_j} w_j(x) u_j}{\sum_j e^{v_j} w_j(x)} \right)^2 \frac{dP_\eta(x)}{\sum_i e^{v_i} w_i(x)} \\
 &< 0
 \end{aligned}$$

by the assumptions of part B. Hence for any  $\eta$  the function  $-f(\eta, \cdot)$  is a gradient of a strictly convex function with invertible and continuous Hessian. An argument similar to that given in the proof of (c), shows that  $f(\cdot, v)$  has a continuous derivative. We conclude from these facts that  $v(\eta)$  has a continuous derivative at  $\eta = 0$ .

Write (d) as

$$dQ_\eta(x) = \frac{d\bar{P}_\eta(x)}{\gamma_\eta \sum e^{v_i(\eta)} w_i(x)},$$

where  $\gamma_\eta = \int (\sum e^{v_i(\eta)} w_i(x))^{-1} d\bar{P}_\eta(x)$ . So  $Q_\eta$  is a "biased" sample of  $\bar{P}_\eta$  (with one stratum and a bounded weight function that may depend on  $\eta$ ). An argument similar to the proof of part A shows that  $\{Q_\eta\}$  has a tangent  $a$ , say. But then it follows from part A that the tangent of  $\{P_\eta\}$  is  $a - E(a | I)$ . Part B follows.  $\square$

In section 6.4 we will prove that the conclusion of part B of the theorem is satisfied in the truncated regression model, although the assumption of B fails there.

### 4.5 MIXTURE MODELS

The semiparametric mixture models we shall discuss are described as follows. Begin with  $Q = \{Q_{(\theta, \eta)} : \theta \in \Theta, \eta \in H\}$ ,  $\Theta \subset R^r, H \subset R^q$  a regular parametric model dominated by  $\mu$  with densities  $f(\cdot, \theta, \eta)$ . Let  $G$  be a (big) subset of the set of all distributions on  $H$  and define

$$P = \{P_{(\theta, G)} : \theta \in \Theta, G \in G\},$$

where

$$(1) \quad P_{(\theta, G)} = \int Q_{(\theta, \eta)} dG(\eta).$$

Such models have arisen extensively in the econometric literature in the guise of regression models with errors in the independent variables or simultaneous equation models. Essentially, in these situations we start with a parametric regression model: Suppose that  $X' = (Z', Y)$ , where  $Z'$  has distribution  $G$  and given  $Z' = z', Y$  has density  $f(\cdot, r(z', \theta))$  with respect to the measure  $\lambda$ . But suppose that  $Z'$  can only be observed with error, so that we actually observe  $Z$  with the conditional density  $q(\cdot, z', \theta)$  of  $Z$  given  $Z' = z'$  (with respect to a

measure  $m$ ). Then the density of  $X = (Z, Y)$  with respect to  $m \times \lambda$  is

$$p(z, y, \theta, G) = \int f(y, r(z', \theta)) q(z, z', \theta) dG(z'),$$

a mixture model with  $z'$  playing the role of  $\eta$ . Mixture models are also closely related to models arising when we assume that we have independent nonidentically distributed observations governed by a "structural" parameter  $\theta$  and incidental parameters  $\eta_1, \dots, \eta_n$ . Specifically, suppose  $X_i$ ,  $1 \leq i \leq n$ , is a  $k$ -vector following a parametric model with density  $f(\cdot, \theta, \eta_i)$ . For instance, if  $X_i = (X_{i1}, \dots, X_{ik})^T$  then  $X_{ij}$ ,  $1 \leq j \leq k$ , could be observed from the  $i$ th of  $n$  strata of a population for which  $\theta$  is a parameter. If we suppose the  $\eta_i$  rather than being fixed unknown constants are themselves independent and identically distributed according to  $G$ , we obtain the mixture model (1). In this formulation, the model is an empirical Bayes model of the form studied by Robbins (1956), save that our interest is in estimating a common parameter  $\theta$  of interest rather than simultaneous estimation of the incidental parameters  $\eta_i$ . The connection between models with an increasing number of incidental parameters and mixture models is discussed further in Bickel and Klaassen (1986). See also Van der Vaart (1988a,b), Bhanja and Ghosh (1991), and Pfanzagl (1991).

Mixture models in the generality we consider were introduced by Kiefer and Wolfowitz (1956), following pioneering work by Neyman and Scott (1948). Kiefer and Wolfowitz considered the problem of nonparametric maximum likelihood estimation in this context. We shall touch on their work in chapter 7.

Identifiability of  $\theta$  and  $G$  in these models is usually not easy to settle. Some results in the Gaussian case and a good bibliography are in Bruni and Koch (1985).

### Examples

Lindsay, in an interesting series of papers (1980)–(1985), considers an important subclass of the mixture models based on exponential family models as defined in section 1. Lindsay's principal models are given by

$$(2) \quad f(x, \theta, \eta) = \exp\{\theta^T S(x) + \eta^T T(x) - b(\theta, \eta)\}$$

and

$$(3) \quad f(x, \theta, \eta) = \exp\{\eta^T T(x, \theta) - c(\theta, \eta)\},$$

where  $\theta$  and  $\eta$  are Euclidean. That is, in model (3) for  $\theta$  fixed,  $Q_2$  is the exponential family generated by  $\mu$  and the statistic  $T(X, \theta)$ . We study an extension which covers both cases, namely

$$(4) \quad f(x, \theta, \eta) = \exp\{\eta^T T(x, \theta) + S(x, \theta) - b(\theta, \eta)\}.$$

That is, we permit the measure generating  $Q_2$  to depend in an arbitrary way on  $\theta$ . Simple formulae for information bounds and influence functions are available essentially only for this class of mixture models, but see Robbins and Zhang (1989). These situations and enriched models in which we also have observations on the mixing distribution are our main concern in this section. Note that

in the mixture models resulting from (4),  $T(X, \theta)$  is a sufficient statistic for  $X$  with respect to  $G \in \mathcal{G}$ . For situations in which such a sufficient statistic exists, a finite sample inequality for the performance of estimators of  $\theta$  is given in Klaassen and Van Zwet (1985) and Klaassen, Van der Vaart, and Van Zwet (1988). Here is a classical example due to Neyman and Scott (1948).

**Example 1. The Neyman-Scott models.**

The basic model here is Gaussian location and scale for  $k$  independent, identically distributed observations. Write  $X = (X^{(1)}, \dots, X^{(k)})^T$ .

**Model L.** Let  $\theta = E_Q X^{(1)}$ ,  $\eta = -[2 \text{Var}_Q X^{(1)}]^{-1}$ . Then the conditional density of  $X$  given  $\eta$  at  $x = (x^{(1)}, \dots, x^{(k)})^T$  is

$$f(x, \theta, \eta) = \exp \left\{ \eta \sum_{j=1}^k (x^{(j)} - \theta)^2 - \frac{k}{2} \log \left( -\frac{\pi}{\eta} \right) \right\}.$$

This is a model of type (3) with

$$T(X, \theta) = \sum_{j=1}^k (X^{(j)} - \theta)^2.$$

**Model S.** Let  $\theta = -[2 \text{Var}_Q X^{(1)}]^{-1}$ ,  $\eta = E_Q X^{(1)} / \text{Var}_Q X^{(1)}$ . Then

$$(5) \quad f(x, \theta, \eta) = \exp \left\{ \eta \sum_{j=1}^k x^{(j)} + \theta \sum_{j=1}^k (x^{(j)})^2 - \frac{k}{2} \left( -\frac{\eta^2}{2\theta} + \log \left( -\frac{\pi}{\theta} \right) \right) \right\}.$$

Thus model S is of type (2). Model S can be thought of as a semiparametric ANOVA type II model. For instance, we may want to estimate the precision of an instrument by taking  $k$  measurements on each of a set of  $n$  samples of an unknown substance.

The parameter  $\theta$  is identifiable in model S for  $k \geq 2$  since  $(2n)^{-1} \sum_{i=1}^n (X_i^{(1)} - X_i^{(2)})^2$  is a consistent estimate of  $-(2\theta)^{-1}$ . For  $k = 1$ ,  $\theta$  is not identifiable. For instance, if  $Q_{(\theta, \eta)} = N(\eta\alpha, \alpha)$  and  $G = N(0, (1 - \alpha)/\alpha^2)$ , then  $P_{(\theta, G)} = N(0, 1)$  for all  $0 < \alpha \leq 1$ . Model L for  $k = 1$  can be thought of as a submodel of the symmetric location model since we are merely restricting the errors in that model to be scale mixtures of Gaussian distributions with mean 0. Since model L is a submodel of the symmetric location model for  $k = 1$ ,  $\theta$  is clearly identifiable. Also note that model L is a submodel of the elliptic model (example 4.2.3) for  $k \geq 1$ . □

**Example 2. Errors in variables.**

Errors in variables and simultaneous equation models have been the subjects of an enormous literature. A fairly recent reference with a good bibliography is Anderson (1984). In the simplest form of this model to fall under our framework, an observation  $X = (Z, Y)$ , where

$$\begin{aligned} Z &= Z' + \varepsilon_1, \\ Y &= \alpha + \beta Z' + \varepsilon_2. \end{aligned}$$

Here  $Z'$  plays the role of  $\eta$ . We will later identify  $Z'$  with  $U$  and write  $Z' \sim G$

or  $U \sim G$ . If we do not make any assumptions on the distributions of  $Z'$  and  $(\varepsilon_1, \varepsilon_2)$ , then neither  $\alpha$  nor  $\beta$  are identifiable. However,  $\beta$  is identifiable under various sets of assumptions. We consider the following two models. In both cases we make an assumption made throughout the literature that the errors are Gaussian, and, for simplicity, independent,  $\varepsilon_i \sim N(0, \sigma_i^2)$ ,  $i = 1, 2$ .

**Restricted model.** Here we require that  $\sigma_1^2, \sigma_2^2$  satisfy a restriction which makes  $\beta$  identifiable. A common one is

$$(6) \quad \sigma_1^2 = c_0 \sigma_2^2,$$

where  $c_0$  is known. Kendall and Stuart (1979) discuss situations in which (6) is reasonable. No restrictions are placed on the distribution  $G$  of  $Z'$ .

**Reiersøl model.** In 1950, Reiersøl showed that  $\beta$  is identifiable even if (6) does not hold provided that  $G$  does not contain distributions  $G$  with a Gaussian component; recall example 1.3.3.

Both of these models are of type (4). Thus if  $\theta = (\beta, \alpha, \sigma_1^2, \sigma_2^2)$ ,  $x = (z, y)$ ,

$$(7) \quad f(x, \theta, \eta) = \exp \left\{ \frac{-1}{2\sigma_1^2} (z - \eta)^2 - \frac{1}{2\sigma_2^2} (y - \alpha - \beta\eta)^2 - \log(2\pi\sigma_1\sigma_2) \right\},$$

leading to

$$(8) \quad T(x, \theta) = \frac{z}{\sigma_1^2} + \frac{\beta(y - \alpha)}{\sigma_2^2}$$

in the Reiersøl model and

$$(9) \quad T(x, \theta) = \frac{1}{\sigma_1^2} \{ z + c_0 \beta (y - \alpha) \}$$

in the restricted model. □

Here is a non-Gaussian example discussed in Lindsay (1980).

**Example 3. Bernoulli pairs with common odds ratio.**

Suppose  $X = (I_1, I_2)$  where  $I_1, I_2$  are independent indicator variables with  $P(I_j=1) = p_j = 1 - q_j$ ,  $j = 1, 2$ , and subject to  $\theta = \log(p_2 q_1 / (p_1 q_2))$ , the log odds ratio, being fixed  $(p_1, p_2)$  vary randomly from observation to observation. If we let  $\eta = \log(p_1/q_1)$ , then we can suppose  $\eta \sim G$  where  $G$  is any distribution on  $R$ . Here

$$f(i_1, i_2, \theta, \eta) = \exp \{ \theta i_2 + \eta (i_1 + i_2) - \log(1 + e^\eta) - \log(1 + e^{\theta + \eta}) \}$$

so that this is a model of type (2) with

$$T(X) = I_1 + I_2, \quad S(X) = I_2. \quad \square$$

*Formal Calculation of  $\hat{P}_i$ ,  $i = 1, 2$*

Suppose  $(U, X)$  is distributed so that  $U \sim G$  and given  $U = \eta$ ,  $X$  has density  $f(\cdot, \theta, \eta)$ . Then  $X \sim P_{(\theta, G)}$  and

$$(10) \quad l(x, \theta, G) = \log \int f(x, \theta, \eta) dG(\eta).$$

Formally, if  $P_0 = P_{(\theta, G)}$ ,

$$\begin{aligned}
 \dot{i}_1(X, \theta, G) &= \frac{\int \dot{f}(X, \theta, \eta) dG(\eta)}{\int f(X, \theta, \eta) dG(\eta)} \\
 (11) \qquad &= \frac{\int (\dot{f}/f)(X, \theta, \eta) f(X, \theta, \eta) dG(\eta)}{\int f(X, \theta, \eta) dG(\eta)} \\
 &= E(\dot{i}(X, \theta, U) | X),
 \end{aligned}$$

where

$$\dot{i}(X, \theta, \eta) \equiv \frac{\dot{f}(X, \theta, \eta)}{f(X, \theta, \eta)}.$$

Similarly, if  $\{P_{(\theta, G_\gamma)} : |\gamma| < 1\}$  is a parametric submodel of  $P_2$  with  $dG_\gamma/dG = g_\gamma, g_0 = 1$ , then

$$(12) \qquad \frac{\partial}{\partial \gamma} \ln(X, \theta, G_\gamma) \Big|_{\gamma=0} = \frac{\int f(X, \theta, \eta) a(\eta) dG(\eta)}{\int f(X, \theta, \eta) dG(\eta)} = E(a(U) | X),$$

where  $a(\eta) \equiv (\partial/\partial \gamma) \log g_\gamma |_{\gamma=0}$ . If  $G$  is rich enough, (12) suggests that

$$(13) \qquad \dot{P}_2 \supset \{E(a(U) | X) : a \in L_2(G), Ea(U) = 0\}.$$

Sufficient conditions for (13) are given in theorem 2 at the end of this section. While equality in (13) does not hold in general, under mild conditions on  $G$  it is often true that  $\dot{P}_2$  is the closure in  $L_2(P)$  of the right side. On the other hand, if  $G$  is discrete then it is easily seen that even the closure of the right side is strictly contained in  $\dot{P}_2$ ; see the discussion in section 6.5 after example 6.5.5, continued. Even if equality in (13) holds, there is no closed form for  $\Pi_0(\cdot | \dot{P}_2)$  in general since  $\dot{P}_2$  is then the closure of the range of the rather general linear operator  $\dot{i}_2 : L_2^0(G) \rightarrow L_2^0(P)$  given by

$$(14) \qquad \dot{i}_2 a = \int a(u) K(\cdot, u) dG(u),$$

where

$$(15) \qquad K(x, u) = \frac{f(x, \theta, u)}{\int f(x, \theta, u) dG(u)}.$$

However, Godambe (1976) (in a finite sample context) and Pfanzagl and Wefelmeyer (1982) showed that the following explicit characterization of  $\dot{P}_2$  holds for a class somewhat more general than that of the models (4).

**Definition 1.** Let  $\{Q_\eta : \eta \in H\}$  be a parametric model dominated by  $\mu$ . If  $X \sim Q_\eta$ , we say a statistic  $T(X)$  is  $G$ -strongly complete for a given probability distribution  $G$  if for any  $w$

$$(16) \quad G(\{\eta : E_{\eta} w(T(X)) = 0\}) = 1$$

implies  $Q_{\eta}(w(T(X)) = 0) = 1$  for all  $\eta \in H$ .

**Theorem 1.** Suppose that  $\mathbf{P}$  is the mixture model (1), that (13) holds, and  $T(X, \theta)$  is a statistic which satisfies:

- (i)  $T(X, \theta)$  is sufficient for  $\eta$  if  $X \sim Q_{(\theta, \eta)}$ ,  $\theta$  fixed,  $\eta \in H$ .
- (ii)  $T(X, \theta)$  is  $G$ -strongly complete for  $\{Q_{(\theta, \eta)} : \eta \in H\}$ .

Then

$$(17) \quad \dot{\mathbf{P}}_2 = \{w(X) \in L_2(P_0) : w(X) \text{ a function of } T(X, \theta), Ew(X) = 0\}$$

and

$$(18) \quad \Pi_0(h(X) | \dot{\mathbf{P}}_2) = E(h(X) | T(X, \theta)) - Eh(X).$$

**Proof.** By (i), the density of  $P \in \mathbf{P}$  may be written as (see (21) below)

$$(a) \quad p(x, \theta, G) = f(x, \theta | t) \int f_T(t, \theta, \eta) dG(\eta), \quad t = T(x, \theta).$$

Thus if  $p(x, \theta, G_{\gamma})$  is a regular parametric submodel of  $\mathbf{P}_2(\theta)$  with tangent  $a \in \dot{\mathbf{P}}_2$  then, with  $s_{\gamma} \equiv s(\theta, G_{\gamma})$  and  $s_0 \equiv s(\theta, G)$ ,

$$(b) \quad s_{\gamma}(x) - s_0(x) = \frac{1}{2} \gamma a(x) s_0(x) + o(\gamma) \quad \text{in } L_2(\mu).$$

Hence

$$(c) \quad 2 \frac{s_{\gamma}(x) - s_0(x)}{s_0(x)} = \gamma a(x) + o(\gamma) \quad \text{in } L_2(P_0),$$

where, by (a), the left-hand side is a function only of  $t$ . Thus  $a$  is, in fact, a function of  $t$  only, and hence  $\dot{\mathbf{P}}_2$  is contained in the right-hand side of (17).

To prove that the right-hand side of (17) is contained in  $\dot{\mathbf{P}}_2$ , let  $h(X) \in L_2(P_0)$  be such that:

$$(d) \quad h(X) \perp \dot{\mathbf{P}}_2,$$

$$(e) \quad h(X) = w(T(X, \theta)),$$

$$(f) \quad Eh(X) = 0.$$

By (13), (d), and (f) it follows that

$$(g) \quad E(h(X)E(v(U) | X)) = 0$$

for all  $v(U) \in L_2(P_0)$ . From (g)

$$(h) \quad Eh(X)v(U) = 0,$$

and hence,

$$E(h(X) | U) = 0 \quad \text{a.s.},$$

or, by (e),



$$G(\{\eta : E_{\eta} w(T(X, \theta)) = 0\}) = 1.$$

Therefore, by (ii)

$$h(X) = 0 \quad \text{a.s. } P_0,$$

and (17) follows. Claim (18) follows from (17) by (A.3.1).  $\square$

**Corollary 1.** Suppose that:

- (i) The mixture model is of type (4).
- (ii) The exponential family in (4) is of full rank.
- (iii) The maps  $\theta \rightarrow T(x, \theta)$ ,  $\theta \rightarrow S(x, \theta)$  are continuously differentiable with continuous derivatives  $\dot{T}(x, \theta)$ ,  $\dot{S}(x, \theta)$ .
- (iv)  $\int |\eta|^2 dG(\eta) < \infty$  and the maps

$$\theta \rightarrow \int |\dot{T}(x, \theta)|^2 f(x, \theta, G) d\mu(x) < \infty,$$

$$\theta \rightarrow \int |\dot{S}(x, \theta)|^2 f(x, \theta, G) d\mu(x) < \infty$$

are continuous.

- (v) The interior of the support of  $G$  is not empty.

Then (18) holds for this model, and

$$(19) \quad \Gamma^*(X, P_0 | \theta, P) = \{\dot{T}(X, \theta) - E(\dot{T}(X, \theta) | T)\} E(U | T) \\ + \dot{S}(X, \theta) - E(\dot{S}(X, \theta) | T),$$

where  $T = T(X, \theta)$ .

A more general result may be found in Van der Vaart (1988b).

**Proof.** It follows from proposition 2.1.1 and theorem 2 below that (i)–(iv) imply that (11) and (13) are valid. Moreover, (i)–(iii) and (v) imply that  $T(X, \theta)$  is sufficient for  $\eta$  in  $\{Q_{(\theta, \eta)} : \eta \in H\}$  and  $G$ -strongly complete; see Lehmann (1986, theorem 4.3.1, page 142). Now (11) yields, with  $\dot{b}$  the derivative of  $b$  with respect to  $\theta$ ,

$$(20) \quad \dot{I}_1(X, \theta, G) = \dot{T}(X, \theta)E(U | X) + \dot{S}(X, \theta) - E(\dot{b}(\theta, U) | X) \\ = \dot{T}(X, \theta)E(U | T) + \dot{S}(X, \theta) - E(\dot{b}(\theta, U) | T)$$

by the sufficiency of  $T(X, \theta)$ . Applying (18), we obtain (19).  $\square$

Suppose that the conditions of theorem 1 hold and  $T(X, \theta) = T(X)$  independent of  $\theta$ . Then, by the Neyman-Halmos-Savage factorization theorem we can write

$$(21) \quad f(X, \theta, \eta) = f(X, \theta | T(X))f_T(T(X), \theta, \eta) \quad \text{a.e. } \mu$$

where  $f_T(\cdot, \theta, \eta)$  is the marginal density of  $T(X)$  with respect to  $\mu^{T^{-1}}$  and  $f(\cdot, \theta | t)$  is the conditional density of  $X$  given  $T = t$  with respect to the conditional measure  $\mu$  given  $T = t$ . Call

$$(22) \quad l(X, \theta | T(X)) = \log f(X, \theta | T(X))$$

the conditional log-likelihood of  $X$  given  $T(X)$ . Then

$$(23) \quad \mathbf{l}(X, \theta, G) = \mathbf{l}(X, \theta | T(X)) + \log p_T(T(X), \theta, G),$$

where

$$p_T(t, \theta, G) = \int f_T(t, \theta, \eta) dG(\eta),$$

the marginal density of  $T(X)$  at  $P_{(\theta, G)}$ . Differentiating (23) formally with respect to  $\theta$  we obtain

$$(24) \quad \dot{\mathbf{l}}_1(X, \theta, G) = \dot{\mathbf{l}}(X, \theta | T(X)) + \frac{\dot{p}_T}{p_T}(T(X), \theta, G).$$

Moreover, by proposition A.3.4

$$\frac{\dot{p}_T}{p_T}(T(X), \theta, G) = E(\dot{\mathbf{l}}_1(X, \theta, G) | T(X)),$$

so that

$$(25) \quad E(\dot{\mathbf{l}}(X, \theta | T(X)) | T(X)) = 0.$$

Therefore, by theorem 1, (24), and (25),

$$(26) \quad \mathbf{l}^*(X, P_0 | \theta, P) = \dot{\mathbf{l}}_1(X, \theta, G) - \Pi_0(\dot{\mathbf{l}}_1(X, \theta, G) | P_2) \\ = \dot{\mathbf{l}}(X, \theta | T(X)).$$

Therefore we can expect that in such mixture models the way to arrive at efficient procedures is to condition on  $T(X)$ . The conditional distribution depends on  $\theta$  only and procedures efficient for these parametric conditional models will be efficient in the full semiparametric model. This idea has been exploited by Andersen (1973), Lindsay (1980), and others in the context of models of type (2). For such models, by lemma 2.7.8, page 58 of Lehmann (1986),

$$(27) \quad \dot{\mathbf{l}}(X, \theta | T(X)) = S(X) - E_\theta(S(X) | T(X))$$

in agreement with corollary 1 since  $S(X, \theta) = \theta^T S(X)$ . Since the conditional model corresponds to an exponential family in  $\theta$ , the conditional maximum likelihood estimate which solves

$$(28) \quad \sum_{i=1}^n S(X_i) = \sum_{i=1}^n E_\theta(S(X_i) | T(X_i))$$

is, in general, efficient.

Although the conditional likelihood approach fails if  $T$  depends on  $\theta$ , it is worth noting that the efficient score function (19) depends on  $G$  only via the marginal density  $p_T(t, \theta, G)$  of  $T$ . Note first that as a consequence of (21) (with  $T(X)$  replaced by  $T(X, \theta)$ ) the conditional distribution of  $X$  given  $T(X, \theta) = t$  has the density  $f(\cdot, \theta | t)$ , so that both  $E(\dot{T}(X, \theta) | T)$  and  $E(\dot{S}(X, \theta) | T)$  do not depend on  $G$ . Next, write

$$E(U | T = t) = \frac{\int \eta \exp\{\eta t - b(\theta, \eta)\} dG(\eta)}{\int \exp\{\eta t - b(\theta, \eta)\} dG(\eta)},$$

and note that the marginal density of  $T$  is of the form

$$p_T(t, \theta, G) = \int \exp\{\eta t - b(\theta, \eta)\} dG(\eta) v(t, \theta),$$

where  $v(t, \theta)$  is independent of  $G$ . Therefore

$$(29) \quad E(U | T = t) = \frac{\nabla_t p_T}{p_T}(t, \theta, G) - \frac{\nabla_t v}{v}(t, \theta),$$

where  $\nabla_t$  denotes the  $q \times 1$  gradient with respect to  $t$ .

**Example 1. The Neyman-Scott models, continued.**

Suppose that conditions (iv) and (v) of corollary 1 hold. Here (i)–(iii) hold automatically. We can therefore apply corollary 1 to calculate the efficient score function  $\Gamma^*$ .

**Model L.**

$$(30) \quad \begin{aligned} \Gamma^*(X, P_0 | \theta, P) &= E(-2U | \sum_{j=1}^k (X^{(j)} - \theta)^2) \sum_{j=1}^k (X^{(j)} - \theta) \\ &= \dot{\Gamma}_1(X, \theta, G) \end{aligned}$$

from (20), since

$$\begin{aligned} E\left(\sum_{j=1}^k (X^{(j)} - \theta) \mid \sum_{j=1}^k (X^{(j)} - \theta)^2\right) \\ = E\left(\sum_{j=1}^k -(X^{(j)} - \theta) \mid \sum_{j=1}^k (X^{(j)} - \theta)^2\right) = 0 \end{aligned}$$

by symmetry considerations.

Thus, fully adaptive estimation of  $\theta$  is possible in model L for any  $k$ . Such estimates can be constructed using the fact that model L is a submodel of the elliptic model of example 4.2.3.

**Model S.** Here  $T(X) = \sum_{j=1}^k X^{(j)}$  and  $S(X) = \sum_{j=1}^k [X^{(j)}]^2$ . Since

$$S(X) = \sum_{j=1}^k (X^{(j)} - k^{-1}T(X))^2 + k^{-1}T^2(X),$$

where  $-2\theta \sum_{j=1}^k (X^{(j)} - k^{-1}T(X))^2$  is independent of  $T(X)$  and has a  $\chi_{k-1}^2$  distribution,

$$(31) \quad E(S(X) | T(X)) = -\frac{k-1}{2\theta} + k^{-1}T^2(X).$$

Therefore, by (26), (27) and (31),

$$\Gamma^*(X, P_0 | \theta, P) = \sum_{j=1}^k (X^{(j)} - k^{-1}T(X))^2 + \frac{k-1}{2\theta}.$$

Equivalently, if  $\sigma^2 = -\frac{1}{2\theta}$ ,

$$(32) \quad I^*(X, P_0 | \sigma^2, \mathbf{P}) = \frac{\sigma^{-2}}{2} \left\{ \sum_{j=1}^k \left( \frac{X^{(j)} - k^{-1}T(X)}{\sigma} \right)^2 - (k-1) \right\},$$

so

$$(33) \quad I(P_0 | \sigma^2, \mathbf{P}) = \frac{k-1}{2\sigma^4}.$$

The conditional maximum likelihood estimate of  $\sigma^2$  which is just

$$\hat{\sigma}^2 = \frac{1}{n(k-1)} \sum_{i=1}^n \sum_{j=1}^k (X_i^{(j)} - k^{-1}T(X_i))^2$$

is efficient.

If  $G$  is  $N(0, \tau^2)$ , it is interesting to note that  $\{P_{(\theta, G_\gamma)}\}$  with  $G_\gamma = N(0, \tau^2 e^\gamma)$ , the normal ANOVA model II, is a least favorable submodel as was noted by Hammerstrom (1978).  $\square$

The following example has been studied by Lindsay (1980), (1982).

**Example 4. Paired exponential mixture model.**

Here,  $X = (X^{(1)}, X^{(2)})^T$  with density

$$f(x_1, x_2, v, \eta) = \exp\{-\eta x_1 + \eta v x_2 + \log(\eta^2 v)\}, \quad x_1, x_2 > 0,$$

so that given  $\eta$ ,  $X^{(1)}, X^{(2)}$  are independent exponential variables with the ratio of hazard rates,  $v$ , independent of  $\eta$ . Then, the model is of type (3),

$$T(X, v) = - (X^{(1)} + vX^{(2)}),$$

$$E\left(\frac{\partial T}{\partial v} \mid T\right) = - E(X^{(2)} \mid X^{(1)} + vX^{(2)}).$$

The marginal density of  $T$  is given by

$$p_T(t, G) = -t \int_0^\infty \eta^2 \exp(\eta t) dG(\eta) \quad \text{for } t < 0.$$

Applying (19) and (29) we find that

$$(34) \quad I^*(X, P \mid v, \mathbf{P}) = \{-X^{(2)} + E(X^{(2)} \mid T)\} E(U \mid T) \\ = -\{X^{(2)} - E(X^{(2)} \mid T)\} \left\{ \frac{p'_T}{p_T}(T) - \frac{1}{T} \right\}.$$

so that, since  $(v \mid T)^{-1} X^{(2)} \mid T \sim \text{Uniform}(0, 1)$ ,

$$I(P \mid v, \mathbf{P}) = \frac{1}{12v^2} E \left[ \left\{ 1 + T \frac{p'_T}{p_T}(T) \right\} - 2 \right]^2 \\ = \frac{1}{12v^2} \{4 + I_{\text{scale}}(p_T)\} \\ = \frac{1}{3v^2} + \frac{1}{12v^2} I_{\text{scale}}(p_T).$$

Here  $(3v^2)^{-1}$  is the information for  $v$  based on the maximal invariant  $Z \equiv X^{(1)}/X^{(2)}$ , which has density  $v(v+z)^{-2}$ ,  $z > 0$ .

Suppose  $G = G_0$  is fixed and  $G_\gamma \equiv G_0(\cdot e^{-\gamma})$ , the scale family generated by  $G_0$ . In the parametric submodel  $\{P_{(v,G_\gamma)}\}$

$$I(X, v, \gamma) = \log \left\{ \int \eta^2 v \exp\{-e^\gamma \eta (X^{(1)} + vX^{(2)})\} dG_0(\eta) \right\} + 2\gamma,$$

so that

$$\dot{i}_1 = v^{-1} (1 - vX^{(2)}E(U | T)), \quad \dot{i}_2 = 2 + TE(U | T).$$

Since  $E(X^{(1)} | T) = E(vX^{(2)} | T) = -\frac{1}{2}T$ , we may verify

$$\dot{i}_1 - \Pi_0(\dot{i}_1 | [\dot{i}_2]) = I^*(X, P | v, P).$$

That is, the scale family is a least favorable submodel. Nevertheless, as Lindsay (1980) points out, the efficient score function and corresponding information bound are not attained if we reduce by invariance and base inference on the marginal likelihood of  $X_i^{(1)} / X_i^{(2)}$ ,  $i = 1, \dots, n$ , since the resulting score function and information do not depend on  $G$ , while (34) patently does. We discuss estimation in this model further in example 7.8.4; see also example 7.5.1.  $\square$

Our treatment of this example has relied heavily on the exponential family structure of the conditional density of  $X$  given  $\eta$ . In fact, however, for this example, the exponential family and mixture structures are not essential. The key feature of the model is the exchangeability of the "core" distribution: with

$$f(x_1, x_2, 1) \equiv \int_0^\infty \eta^2 \exp\{-(\eta x_1 + \eta x_2)\} dG(\eta),$$

the density  $f(x_1, x_2, 1)$  is exchangeable:

$$f(x_1, x_2, 1) = f(x_2, x_1, 1).$$

Thus the result of proposition 4.2.1, with  $V \equiv \{\text{the scale group on } R^+\}$ , applies. It then becomes apparent that the present example 4 can be considerably extended as follows: instead of starting with the exponential density, suppose that  $f_0$  is a fixed known (Lebesgue) density on  $R^+$ , let

$$g(x_1, x_2; G) \equiv \int_0^\infty \eta^2 f_0(\eta x_1) f_0(\eta x_2) dG(\eta)$$

so that  $g$  is exchangeable, and let

$$f(x_1, x_2, v, G) \equiv v g(x_1, vx_2, G) \quad \text{for } x_1, x_2, v > 0.$$

This is a paired-sample model, of the type considered in proposition 4.2.1.

**Example 2. Errors in variables, continued: Restricted model.**

With  $T(X, \theta)$  as in (9), define

$$\begin{aligned} (35) \quad \hat{\eta}(X, \theta) &= \frac{T(X, \theta)}{1/\sigma_1^2 + \beta^2/\sigma_2^2} \\ &= T(X, \theta)\sigma_1^2(1 + c_0\beta^2)^{-1} = \frac{Z + c_0\beta(Y - \alpha)}{1 + c_0\beta^2} \end{aligned}$$

$$= Z' + \frac{\varepsilon_1 + c_0 \beta \varepsilon_2}{1 + c_0 \beta^2}.$$

Conditioning on  $\hat{\eta}$  and on  $T$  is the same. However,  $\hat{\eta}$  plays a distinguished role as the best unbiased estimator of  $\eta$  in  $\mathbf{Q}_2 = \{Q_{(\theta, \eta)} : \eta \in H\}$  where  $Q_{(\theta, \eta)}$  has density  $f$  given in (7). Let

$$R(X, \theta) = \frac{Y - \alpha - \beta Z}{\sigma(\theta)},$$

where

$$\sigma^2(\theta) = \text{Var}(Y - \alpha - \beta Z) = \sigma_1^2 \left( \frac{1}{c_0} + \beta^2 \right).$$

Note that, by Basu's theorem (e.g., Lehmann (1983, theorem 1.5.5, page 46)),  $R$  and  $\hat{\eta}$  are independent and  $R \sim N(0, 1)$ . For this model, by comparing (4) and (7) we find that

$$(36) \quad S(X, \theta) = -\frac{1}{2\sigma_1^2} (Z^2 + c_0(Y - \alpha)^2).$$

Then

$$(37) \quad \frac{\partial T}{\partial \beta} = \frac{c_0}{\sigma_1^2} (Y - \alpha), \quad \frac{\partial S}{\partial \beta} = 0,$$

$$\frac{\partial T}{\partial \alpha} = -\frac{\beta c_0}{\sigma_1^2}, \quad \frac{\partial S}{\partial \alpha} = \frac{c_0}{\sigma_1^2} (Y - \alpha),$$

and

$$\frac{\partial T}{\partial \sigma_1^2} = -\frac{T}{\sigma_1^2}, \quad \frac{\partial S}{\partial \sigma_1^2} = -\frac{S}{\sigma_1^2}.$$

Note that

$$(38) \quad Z = \hat{\eta} - \sigma^{-1}(\theta) \beta \sigma_1^2 R,$$

$$(39) \quad Y - \alpha = \beta \hat{\eta} + \sigma^{-1}(\theta) \sigma_1^2 c_0^{-1} R.$$

Further, the density of  $\hat{\eta}(X, \theta)$  is

$$w(t, \theta) = \tilde{\sigma}^{-1}(\theta) \int \phi\left(\frac{t - \eta}{\tilde{\sigma}(\theta)}\right) dG(\eta)$$

where

$$\tilde{\sigma}^2(\theta) = E(\hat{\eta}(X, \theta) - Z)^2 = \frac{\sigma_1^4}{c_0 \sigma^2(\theta)}.$$

Therefore, arguing as for (29),

$$(40) \quad E(U | T) = E(U | \hat{\eta}) = \hat{\eta} + \tilde{\sigma}^2(\theta) \frac{w'}{w}(\hat{\eta}, \theta)$$

since  $U = Z'$ .

Next, we write  $I^*$  given by (19) in this case, as  $(I_1^*, I_2^*, I_3^*)^T$ , corresponding to  $(\beta, \alpha, \sigma_1^2)^T$ . Substituting from (37)–(40) and using the independence of  $\hat{\eta}$  and  $R$  we obtain

$$(41) \quad \begin{aligned} I_1^* &= \sigma^{-1}(\theta)R(\hat{\eta} + \tilde{\sigma}^2(\theta)\frac{w'}{w}(\hat{\eta}, \theta)), \\ I_2^* &= \sigma^{-1}(\theta)R, \\ I_3^* &= \frac{1}{2}\sigma_1^{-2}(R^2 - 1). \end{aligned}$$

We can finally obtain  $I^*(\cdot, P_0 | \beta, \mathbf{P})$  from (3.4.22). Since

$$I_3^* \perp [I_1^*, I_2^*],$$

we obtain from (41)

$$(42) \quad \begin{aligned} I^*(\cdot, P_0 | \beta, \mathbf{P}) &= I_1^* - \Pi_0(I_1^* | [I_2^*]) \\ &= \sigma^{-1}(\theta)R(\hat{\eta} - EZ' + \tilde{\sigma}^2(\theta)\frac{w'}{w}(\hat{\eta}, \theta)), \end{aligned}$$

and, since  $\text{Var}(\hat{\eta}) = \text{Var}(Z') + \tilde{\sigma}^2(\theta)$ ,

$$(43) \quad \begin{aligned} I(P_0 | \beta, \mathbf{P}) &= \sigma^{-2}(\theta)\{\text{Var}(\hat{\eta}) + \tilde{\sigma}^4(\theta)I(w) - 2\tilde{\sigma}^2(\theta)\} \\ &= \sigma^{-2}(\theta)[\text{Var}(Z') + \tilde{\sigma}^2(\theta)(\tilde{\sigma}^2(\theta)I_w - 1)], \end{aligned}$$

where

$$I_w = \int \frac{[w']^2}{w}(t, \theta) dt,$$

the Fisher information for location of  $w$ . Formula (43) follows from  $\int t w'(t, \theta) dt = -1$  and  $\int w'(t, \theta) dt = 0$ .

We can draw some interesting qualitative consequences from (42) and (43).

(A1) If  $c_0 \rightarrow 0$  with  $\sigma_2^2$  fixed, then  $\tilde{\sigma}^2(\theta) \rightarrow 0$  and in the limit we are in the case where  $Z'$  is observed without error,  $\hat{\eta} = Z' = Z$ . In this case

$$I(P_0 | \beta, \mathbf{P}) = \frac{\text{Var}(Z)}{\sigma_2^2},$$

the reciprocal of the asymptotic variance of the least squares estimate divided by  $n$  as it should be.

(A2) The score function (42) can be compared to the score function when  $Z = Z'$  which follows from example 4.2.2, continued, with  $\psi(u) = u$ . The only difference is that in (42) the unobservable  $Z'$  is replaced by

$\hat{\eta} + \tilde{\sigma}^2(\theta)(w'/w)(\hat{\eta}, \theta)$ , the best predictor of  $Z'$  when  $(\theta, G)$  is known (cf. (40)).

(A3) If  $Z'$  is normal, then  $w$  is a normal density,  $\text{Var}(\hat{\eta}) = I^{-1}(w)$ , and (43) becomes, after some calculation using again  $\text{Var}(\hat{\eta}) = \text{Var}(Z') + \tilde{\sigma}^2(\theta)$ ,

$$\frac{\text{Var}^2(Z')}{\sigma^2(\theta) \text{Var}(\hat{\eta})}$$

If  $c_0 = 1$  this is just the reciprocal of the asymptotic variance of  $n^{1/2}(\hat{\beta}_p - \beta)$  where  $\hat{\beta}_p$  minimizes the sum of squares of perpendicular distances of  $(Z_i, Y_i)$  from  $y = \alpha + \beta z$ ,

$$\sum_{i=1}^n \frac{(Y_i - \alpha - \beta Z_i)^2}{1 + \beta^2}$$

Note that  $\hat{\beta}_p$  is the maximum likelihood estimate under the restricted model  $c_0 = 1$  and  $G$  Gaussian. The efficiency of  $\hat{\beta}_p$  at Gaussian  $G$  is to be expected since  $\hat{\beta}_p$  is locally regular in the full  $c_0 = 1$  restricted model; see, e.g., Gleser (1981). However  $\hat{\beta}_p$  is not efficient at any non-Gaussian  $G$ .

(A4) It is not hard to see, and shown in Bickel and Ritov (1987), that a least favorable parametric submodel at  $(\theta, G)$  is the location-scale submodel

$$\{P_{(\theta, G(m,s))} : G_{m,s}(z) = G\left(\frac{z-m}{s}\right), \quad s > 0, \quad m \text{ real}\}$$

Put another way, we can do as well not knowing  $G$  at all as knowing  $G$  up to a change of location and scale. Bickel and Ritov show how to construct efficient estimates of  $\beta$  in the model  $\mathbf{P}$ . The parameters  $\sigma_1^2$  and  $\alpha$  also admit efficient estimates.  $\square$

### Example 2. Errors in variables, continued: Reiersøl model.

If we permit  $\sigma_2^2$  to vary freely, the efficient score function for  $\theta$  now has four components,  $I_1^*, I_2^*, I_3^*, I_4^*$  corresponding to  $\beta, \alpha, \sigma_1^2, \sigma_2^2$ . The first two are as in (41). However, calculating as in the restricted model but with

$$S(X, \theta) = -\frac{1}{2} \left\{ \frac{Z^2}{\sigma_1^2} + \frac{(Y - \alpha)^2}{\sigma_2^2} \right\},$$

$$\frac{\partial T}{\partial \sigma_1^2} = -\frac{Z}{\sigma_1^4}, \quad \frac{\partial T}{\partial \sigma_2^2} = -\frac{\beta}{\sigma_2^4}(Y - \alpha),$$

we obtain

$$(44) \quad I_3^* = \sigma^{-2}(\theta) \frac{\beta^2}{2} (R^2 - 1) + \beta \tilde{\sigma}^2(\theta) \sigma^{-1}(\theta) \sigma_1^{-2} R \frac{w'}{w}(\hat{\eta}, \theta),$$

$$I_4^* = \frac{\sigma^{-2}(\theta)}{2} (R^2 - 1) - \beta \tilde{\sigma}^2(\theta) \sigma^{-1}(\theta) \sigma_2^{-2} R \frac{w'}{w}(\hat{\eta}, \theta),$$

so that, when  $\beta \neq 0$ ,

$$[I_2^*, I_3^*, I_4^*] = [R, (R^2 - 1), R \frac{w'}{w}(\hat{\eta}, \theta)].$$



Thus, when  $\beta \neq 0$  it follows from (41) and (44) by use of (3.4.22), after some calculation using again  $\int tw'(t, \theta) dt = -1$ , that

$$(45) \quad I^*(X, P_0 | \beta, \mathbf{P}) = \sigma^{-1}(\theta)R \left[ \hat{\eta} - E\hat{\eta} + I^{-1}(w) \frac{w'}{w}(\hat{\eta}, \theta) \right]$$

and

$$(46) \quad \begin{aligned} I(P_0 | \beta, \mathbf{P}) &= \sigma^{-2}(\theta)(Var(\hat{\eta}) - I^{-1}(w)) \\ &= \sigma^{-2}(\theta)(Var(Z') + \tilde{\sigma}^2(\theta) - I^{-1}(w)). \end{aligned}$$

We again draw some qualitative conclusions.

(B1) If  $G_m \rightarrow \delta_\mu$  weakly, then  $L(Z | P_{(\theta, G_m)}) \rightarrow L(\mu + \varepsilon_1) = N(\mu, \sigma_1^2)$  and, under regularity conditions,  $Var(\hat{\eta}) - I^{-1}(w) \rightarrow 0$ ; cf. (A3). Then by (46)

$$I(P_{(\theta, G_m)} | \beta, \mathbf{P}) \rightarrow 0.$$

It is clear the  $\beta$  cannot be identified if  $G$  is degenerate. A similar computation using (A3) and (46) yields  $I(P_{(\theta, G)} | \beta, \mathbf{P}) = 0$  for  $G$  Gaussian. This is appropriate since if  $G$  is Gaussian,  $\beta$  cannot be identified either.

(B2) If  $I_a$  is the information under the restricted model and  $I_b$  that under the Reiersøl model, then necessarily  $I_a \geq I_b$ . The inequality is always strict since, from (43) and (46),

$$\sigma^2(\theta)(I_a - I_b) = \frac{\{\tilde{\sigma}^2(\theta)I(w) - 1\}^2}{I(w)} > 0.$$

The last inequality follows since the Fisher information  $I(w)$  of  $w$ , which is the density of the convolution of  $G$  with  $N(0, \tilde{\sigma}^2(\theta))$ , is always smaller than the Fisher information of  $N(0, \tilde{\sigma}^2(\theta))$ , which is just  $\tilde{\sigma}^{-2}(\theta)$ , unless  $G$  is a point mass.

(B3) It again turns out that a least favorable model at  $(\theta, G)$  is the location-scale family generated by  $G$ .

(B4) The point  $\beta = 0$  requires separate discussion. When  $\beta = 0$  it follows immediately from (44) that the score function  $I_3^* = 0$ , which is reasonable since  $\sigma_1^2$  is, in this case, not identifiable. Furthermore,  $[I_3^*, I_4^*] = [R^2 - 1]$  and hence for the Reiersøl model  $\mathbf{P}_b$ ,  $I^*(X, P_0 | \beta, \mathbf{P}_b)$  becomes the score function (42) for the restricted model  $\mathbf{P}_a$ . Bickel and Ritov (1987) exhibit a locally regular estimate  $\hat{\beta}$  which has the corresponding influence function at  $\beta = 0$ . However, the point  $\beta = 0$  is *not* a regular point of the least favorable submodel discussed in (A4) and (B3), since the information matrix at  $\beta = 0$  is singular.

Further applications of these techniques to other errors in variables models may be found in Bickel and Ritov (1987). The more general situation in which  $\varepsilon_1, \varepsilon_2$  are independent but with unknown distributions has been treated by Spiegelman (1979). Efficient estimates are unknown for this model.  $\square$

### The Has'minskii-Ibragimov Model

In their 1983 paper, Has'minskii and Ibragimov consider an enriched version of the mixture model in which in addition to  $Y_i, i = 1, \dots, n$ , independent and

identically distributed according to  $P_{(\theta, G)}$  given by (1) we also observe independent  $U_1', \dots, U_n'$  distributed according to  $G$  and independent of the  $Y_i$ . In this model  $G$  is always consistently estimable and the existence of regular estimates of  $\theta$  depends essentially on the regularity of the parametric model  $P_1$ . This assertion is made more precise in corollary 2 below. If  $G \ll \lambda$ , the density of  $X = (U', Y)$  with respect to  $\lambda \times \mu$  is

$$(47) \quad p(x, \theta, G) = g(u') \int q(y, \theta, \eta) g(\eta) d\eta,$$

where  $x = (u', y)$ . Then

$$(48) \quad \dot{i}_1 = E \left( \frac{\dot{q}}{q}(Y, \theta, U) \mid Y \right),$$

where  $U \sim G$  (is unobservable), and given  $U, Y \sim Q_{(\theta, U)}$ . On the other hand, formally again,

$$(49) \quad \dot{P}_2^0 = \{a(U') + E[a(U) \mid Y] : a \in L_2(G), Ea(U') = 0\}.$$

Has'minskii and Ibragimov are led to construct their estimate in terms of the solution of a Fredholm equation. Alternatively, projections onto  $\dot{P}_2$  can be expressed in terms of conditional expectations and ACE as follows. We expect that this expression can be used for estimation as well.

Of course,  $(U, Y)$  are dependent with joint density  $q(y, \theta, u)g(u)$ . Define a new random variable  $U''$  so that  $L(U'' \mid Y) = L(U \mid Y)$  and  $U, U''$  are conditionally independent given  $Y$ .

**Proposition 1.** Suppose that  $\dot{P}_2$  is the closure of the set in (49). Then, for  $h = h(Y) \in L_2(P)$ ,

$$(50) \quad \Pi_0(h \mid \dot{P}_2) = (I^{U'} + E^Y)(I^U + E^U E^Y)^{-1} E^U h$$

$$(51) \quad = (I^{U'} + E^Y) ACE_1(h \mid U, U''),$$

where  $I^U a = a(U)$ ,  $I^{U'} a = a(U')$ ,  $E^Y a = E(a(U) \mid Y)$ ,  $E^U h = E(h(Y) \mid U)$ .

**Proof.** If  $h(Y) \in L_2(P)$ , then in order for

$$(a) \quad \Pi_0(h \mid \dot{P}_2) = a^*(U') + E^Y a^*(U) = (I^{U'} + E^Y) a^*$$

with  $E a^*(U') = 0$ , we must have, for every  $a \in L_2(G)$  with  $E a(U') = 0$ ,

$$\begin{aligned} 0 &= E[h(Y) - a^*(U') - E^Y a^*][a(U') + E^Y a] \\ &= E[h(Y)a(U)] - E[a^*(U')a(U')] \\ &\quad - E[E(a^*(U) \mid Y)E(a(U) \mid Y)] \end{aligned}$$

$$(b) \quad = E\{a(U)[E^U h - a^*(U) - E^U E^Y a^*]\}.$$

Hence  $a^*$  must satisfy

$$(c) \quad (I^U + E^U E^Y) a^* = E^U h.$$

Now  $E^U E^Y$ , and hence also  $I^U + E^U E^Y$ , are nonnegative definite, self-adjoint operators on  $L_2(G)$ , and for all  $a \in L_2(G)$

$$\begin{aligned} \|a\| \|(I^U + E^U E^Y)a\| &\geq \langle a, (I^U + E^U E^Y)a \rangle \\ &\geq \|a\|^2. \end{aligned}$$

Consequently, by corollary A.1.2,  $I^U + E^U E^Y$  has a bounded inverse, and (50) follows immediately from (c) and (a).

To see (51), first note that by the definition of ACE from section A.4, symmetry in  $U$ ,  $U''$ , and the conditional independence of  $U$ ,  $U''$ ,

$$ACE(h | U, U'') = a(U) + a(U''),$$

where  $a(U) \equiv ACE_1(h | U, U'')$  is uniquely defined by

$$E([h(Y) - a(U) - a(U'')]b(U)) = 0$$

for all  $b(U) \in L_2(P)$ . Equivalently

$$E(h(Y) - a(U) - a(U'') | U) = 0 \quad \text{a.s.}$$

or, since  $E(a(U'') | U) = E^U E(a(U'') | U, Y) = E^U E^Y a$ ,

$$(d) \quad E^U h - a(U) - E^U E^Y a = 0 \quad \text{a.s.}$$

Hence,  $(I^U + E^U E^Y)a = E^U h$  and, by (c),  $a = a^*$  and (51) follows.  $\square$

We conclude with:

**Corollary 2.** If  $I(P | \theta, P_1(G)) > 0$ , then  $I(P | \theta, P) > 0$ .

**Proof.** Suppose that  $I(P | \theta, P) = 0$ . Then, by (51) for  $a(U) = ACE_1(\hat{I}_1 | U, U'')$ ,  $E[\hat{I}_1(Y) - E^Y a - a(U')]^2 = 0$ , which, by independence of  $U'$  and  $Y$  implies  $E a^2(U') = 0$ . But then it follows that  $I(P | \theta, P_1(G)) = E \hat{I}_1^2(Y) = 0$ .  $\square$

### Other Information Bounds in Mixture Models

Lindsay (1982) and Kumon and Amari (1984) take a different approach to information in this context. Their point of departure is M-estimates  $\hat{\theta}_n$  solving

$$(52) \quad \sum_{i=1}^n \psi(X_i, \theta) = 0,$$

where  $\psi \neq 0$  is such that

$$(53) \quad \int \psi(x, \theta) dP_{(\theta, G)}(x) = 0 \quad \text{for all } \theta, G \in G.$$

For simplicity suppose  $k = 1$ . For general mixture models it is unclear whether there are any such  $\psi$  but for model (4) we can obtain such functions for instance by taking an arbitrary  $\psi_0(x, \theta)$  and defining

$$(54) \quad \psi(X, \theta) = \psi_0(X, \theta) - E_{\theta}(\psi_0(X, \theta) | T(X, \theta)),$$

where  $T(X, \theta)$  satisfies (i) and (ii) in theorem 1 for all  $G \in G$ .

As we will see in theorem 7.2.1, if (52) has a consistent solution then, under mild conditions,

$$L(\sqrt{n}(\hat{\theta}_n - \theta)) \rightarrow N(0, \frac{E\psi^2(X, \theta)}{[E\dot{\psi}(X, \theta)]^2}).$$

These authors now seek optimal estimates within the class of estimates defined by (52) or a subclass thereof. For instance, Lindsay concentrates on  $\psi$  such that, as for the usual score function,  $\int \psi^2(x, \theta) dP_{(\theta, G)}(x) = - \int \dot{\psi}(x, \theta) dP_{(\theta, G)}(x)$ . The information bounds obtained by these authors, essentially in the context of model (4) are typically larger than  $I^{-1}(P | \theta, P)$ . This is the case, for instance, in example 4 where Lindsay's approach leads to the information bound based on  $X_i^{(1)} / X_i^{(2)}, i = 1, \dots, n$ .

*Regularity Conditions*

Let  $P$  be a mixture model as in (1). Define

$$(55) \quad J(\theta, \eta) = \int \frac{|\dot{f}(x, \theta, \eta)|^2}{f(x, \theta, \eta)} d\mu(x)$$

where  $\dot{f}$  is the derivative of  $f$  with respect to  $\theta$ . Then,  $J$  is the trace of the Fisher information matrix for  $\theta$  with  $\eta$  fixed. Since  $Q$  is regular the map  $\theta \rightarrow J(\theta, \eta)$  is continuous.

**Theorem 2.** Suppose that:

- (i)  $Q$  is regular.
- (ii)  $\int J(\theta, \eta) dG(\eta) < \infty$  for all  $\theta$  and the map

$$(56) \quad \theta \rightarrow \int J(\theta, \eta) dG(\eta)$$

is continuous for all  $G \in G$ .

- (iii) The information matrix  $I(\theta) \equiv I(P_0 | \theta, P_1(G))$  for  $\theta$  in the model with  $G$  fixed,

$$(57) \quad I(\theta) = \int \frac{\dot{p}\dot{p}^T}{p}(x, \theta, G) d\mu(x),$$

is nonsingular.

Then  $P_1(G)$  is regular and (11) and (13) hold.

**Proof.** By (i) and (8.14.3), (8.7.7), and (7.6.2) of Dieudonné (1960)

$$(a) \quad \begin{aligned} &|\Delta|^{-1} \| f^{1/2}(\cdot, \theta + \Delta, \eta) - f^{1/2}(\cdot, \theta, \eta) \| \\ &= |\Delta|^{-1} \| \frac{1}{2} \int_0^1 \Delta^T \dot{f} f^{-1/2}(\cdot, \theta + \zeta\Delta, \eta) d\zeta \| \\ &\leq |\Delta|^{-1} \sup_{0 \leq \zeta \leq 1} \| \frac{1}{2} \Delta^T \dot{f} f^{-1/2}(\cdot, \theta + \zeta\Delta, \eta) \| \end{aligned}$$

$$\leq \frac{1}{2} \sup_{0 \leq \zeta \leq 1} J^{1/2}(\theta + \zeta \Delta, \eta).$$

Furthermore, as noted in (b) of the proof of proposition A.5.5, formula (2.1.2) of regularity is equivalent to

$$(b) \quad \|s(\theta + h) - s(\theta) - 2(s(\theta + h) + s(\theta))^{-1} \dot{s}^T(\theta) h s(\theta)\| = o(|h|).$$

Fix  $G \in \mathbf{G}$ . We will show (b) for  $\mathbf{P}_1(G)$ , i.e., with  $s(x, \theta)$  replaced by  $s(x, \theta, G) = \{\int f(x, \theta, \eta) dG(\eta)\}^{1/2}$ . Indeed, by lemma A.5.1 we have

$$(c) \quad \int \left\{ \lambda^{-1} (s(x, \theta + \lambda \Delta, G) - s(x, \theta, G)) - \frac{\Delta^T \int \dot{f}(x, \theta, \eta) dG(\eta)}{s(x, \theta + \lambda \Delta, G) + s(x, \theta, G)} \right\}^2 d\mu(x) \\ \leq \iint \left\{ \lambda^{-1} (f^{1/2}(x, \theta + \lambda \Delta, \eta) - f^{1/2}(x, \theta, \eta)) - \frac{\Delta^T \dot{f}(x, \theta, \eta)}{f^{1/2}(x, \theta + \lambda \Delta, \eta) + f^{1/2}(x, \theta, \eta)} \right\}^2 d\mu(x) dG(\eta).$$

The inner integral of the right-hand side of (c) converges to 0 as  $\lambda \rightarrow 0$  for every  $\eta$  by (i) and (b), and is, in view of (a) and (ii), bounded uniformly in  $\lambda \in [0, 1]$  by a  $G$ -integrable function of  $\eta$ . Consequently, both sides of (c) converge to 0 as  $\lambda \rightarrow 0$ , which yields (b) for  $\mathbf{P}_1(G)$ .

By the same argument as in (e) through (j) of the proof of proposition A.5.5, and by (iii) regularity of  $\mathbf{P}_1(G)$  is proved. Furthermore, (11) is valid.

To prove (13), let  $h$  be bounded and  $Eh(U) = 0$ , and, with  $g_\gamma \equiv dG_\gamma/dG_0$ ,  $g_0 = 1$ , for  $\gamma$  sufficiently small set

$$g_\gamma(\eta) = \left( \int e^{\gamma h} dG \right)^{-1} e^{\gamma h(\eta)}.$$

We leave it to the reader to show by arguing as above that  $P_{(\theta, G_\gamma)}$  is regular with tangent  $E(h(U) | X)$  at  $P_{(\theta, G_0)}$  so that (13) is proved.  $\square$

## 4.6 MISSING DATA MODELS

The mixture model discussed in the previous section can be thought of as being obtained by an information reducing state-space transformation as follows. Begin with a semiparametric model in which  $X^0 = (U, X)$  is distributed so that  $U \sim G$ , and given  $U = \eta$ ,  $X$  follows the parametric model  $\{Q_{(\theta, \eta)} : \theta \in \Theta\}$ . However, only  $X = T(X^0)$  is observed.

In the general missing data model we start out with some core or latent model  $\mathbf{Q}$  for  $X^0$  and a known map  $T: X^0 \rightarrow T$  and suppose that we observe only  $X = T(X^0)$ . That is

$$\mathbf{P} = \{QT^{-1} : Q \in \mathbf{Q}\}.$$

An important example of such a mapping is right or left *censoring*. That is,  $X^0 = (Z, Y, C)$  where  $Z$  is a vector of covariates,  $Y$  and  $C$  are real (usually positive), and we observe  $X = T(X^0) = (Z, 1_{[Y \leq C]}, Y \wedge C)$  (right censoring) or  $(Z, 1_{[Y \geq C]}, Y \vee C)$  (left censoring). More general schemes such as interval censoring can also be put in this framework.

As the mixture model results suggest, linking  $\dot{P}$  and  $\dot{Q}$  in such models is relatively easy, but computation of projections is, in general, difficult.

We obtain from Proposition A.5.5 that if  $b(X^0) \in \dot{Q}$  then  $E(b(X^0) | X) \in \dot{P}$ . Hence

$$(1) \quad \dot{P} \supset \{E(b(X^0) | X) : b \in \dot{Q}\}.$$

We may expect that in most cases of interest:

$$(2) \quad \dot{P} = \{E(b(X^0) | X) : b \in \dot{Q}\}.$$

**Example 1. Nonparametric core model.**

Suppose that the core model  $Q$  is nonparametric and consists of all  $Q$  dominated by some  $\sigma$ -finite measure  $\mu$ . Then  $\dot{Q} = \{b \in L_2(Q) : E b(X^0) = 0\}$ . In general,  $\dot{P}$  can be very small or very large depending on  $T$ . As a (trivial) example of the former, consider  $T(X^0) \equiv 0$ ; then  $\dot{P} = \{0\}$ . However,  $P = \{ \text{all } P \ll \mu T^{-1} \}$ , and it is clear that

$$(3) \quad \dot{P} = \{b(X) : E_P b(X) = 0\}.$$

Of course, (3) can hold even if  $P$  is considerably smaller; recall the Pfanzagl-Godambe theorem 4.5.1. We do not know a general theorem for deciding when (3) holds in terms of  $Q$  or  $\dot{Q}$  and  $T$ . However, whenever (3) holds, if we can construct a locally asymptotically linear, regular estimate of  $v(P)$  for the missing data model, the estimate is necessarily efficient by proposition 3.3.1. For example, if  $X = T(X^0) = |X^0|$  for  $X^0 \sim Q$  symmetric about 0 on  $R^k$ , then it is clearly true that  $P = \{ \text{all } P \}$  (by considering  $Q$  concentrating on one coordinate of  $X^0$ ), and hence (3) holds. Thus  $n^{-1} \sum_{i=1}^n 1_{\{|X_i^0| \leq r\}}$  is an efficient estimate of  $P(X \leq r)$ . We will return to this theme in section 6.6. □

**Example 2. Missing observations on a component of  $X^0$ .**

Let  $X^0 = (Y, Z, U)$  take values in  $R^m \times R^d \times \{0, 1\}$ , and let  $Q$  be all distributions of  $X^0$  such that if  $Q \in \mathcal{Q}$  then  $Q \ll \mu$  and  $U$  is independent of  $(Y, Z)$ .

Under the missing observation model we observe  $X = (Y, UZ, U)$ . That is, we assume that observations on the independent variable  $Z$  are missing ( $= 0$  with probability  $1 - p$ ) independently of what their values might have been. Generalizations of such models with both  $Y$  and  $Z$  missing jointly or separately and  $(Y, Z)$  Gaussian have been studied intensively; see Dempster, Laird, and Rubin (1977) and Little and Rubin (1987).

Here the tangent space for the model  $Q$  is

$$\dot{Q} = \{b_1(Y, Z) + b_2(U) : E_Q b_i = 0, b_i \in L_2(Q), i = 1, 2\},$$

and the tangent space  $\dot{P}$  satisfies (2). Let  $A : \dot{Q} \rightarrow \dot{P}$  be defined by

$$(4) \quad (Ab)(X) = E(b(X^0) | X) \\ = Ub_1(Y, Z) + (1 - U)E(b_1(Y, Z) | Y) + b_2(U)$$

for  $b = b_1 + b_2$  as above. We proceed to calculate the projection operator using theorem A.2.2, formula (A.2.13). Here, if  $h = h(Y, UZ, U)$  with  $Eh(Y, UZ, U) = 0$ , and  $p = P(U = 1)$ , then, by the independence of  $U$  and  $(Y, Z)$  and (A.1.9),

$$(5) \quad A^T h = ph(Y, Z, 1) + (1 - p)h(Y, 0, 0) \\ + UEh(Y, Z, 1) + (1 - U)Eh(Y, 0, 0).$$

Hence

$$(6) \quad A^T A(b) = pb_1(Y, Z) + (1 - p)E(b_1(Y, Z) | Y) + b_2(U),$$

and, by simple calculations,

$$(7) \quad (A^T A)^{-1}(b) = \frac{1}{p}b_1(Y, Z) - \frac{1-p}{p}E(b_1(Y, Z) | Y) + b_2(U).$$

Combining (4)–(7) we find for any function  $h$  as above,

$$\begin{aligned} \Pi_0(h | \dot{P}) &= A(A^T A)^{-1}A^T h \\ &= A(A^T A)^{-1}\{ph(Y, Z, 1) + (1 - p)h(Y, 0, 0) \\ &\quad + UEh(Y, Z, 1) + (1 - U)Eh(Y, 0, 0)\} \\ &= A\{h(Y, Z, 1) - (1 - p)E(h(Y, Z, 1) | Y) + (1 - p)h(Y, 0, 0) \\ &\quad + UEh(Y, Z, 1) + (1 - U)Eh(Y, 0, 0)\} \\ &= U\{h(Y, Z, 1) - (1 - p)E(h(Y, Z, 1) | Y) + (1 - p)h(Y, 0, 0)\} \\ &\quad + (1 - U)\{pE(h(Y, Z, 1) | Y) + (1 - p)h(Y, 0, 0)\} \\ &\quad + UEh(Y, Z, 1) + (1 - U)Eh(Y, 0, 0) \\ &= Uh(Y, Z, 1) - (U - p)E(h(Y, Z, 1) | Y) \\ &\quad + (1 - p)h(Y, 0, 0) - p^{-1}(U - p)Eh(Y, 0, 0), \end{aligned}$$

where the last equality follows from  $Eh(Y, UZ, U) = 0$ . In particular, for a function  $h$  such that  $h(Y, 0, 0) \equiv 0$ ,

$$(8) \quad \Pi_0(h | \dot{P}) = Uh(Y, Z, 1) - (U - p)E(h(Y, Z, 1) | Y).$$

For further treatment of projections and projection formulas, see sections 5.4 and 5.5.

We want to estimate functionals of the joint distribution of  $(Y, Z)$ , for instance, the regression coefficient

$$v = \Sigma_Z^{-1} E(Z - EZ)Y^T,$$

where

$$\Sigma_Z \equiv E(Z - EZ)(Z - EZ)^T.$$

A preliminary  $\sqrt{n}$ -consistent estimate is available by using the  $(Y_i, Z_i, 1)$  observations

$$(9) \quad \hat{v}_n = \left\{ \sum_{i=1}^n U_i (Z_i - \bar{Z})(Z_i - \bar{Z})^T \right\}^{-1} \left\{ \sum_{i=1}^n U_i (Z_i - \bar{Z}) Y_i^T \right\},$$

where  $\bar{Z} = (\sum_{i=1}^n U_i Z_i) / \sum_{i=1}^n U_i$ . The influence function of  $\hat{v}_n$  is easily seen to be

$$h_0(Y, UZ, U) = p^{-1} \Sigma_Z^{-1} U (Z - EZ) (Y^T - Z^T v),$$

and this is also a pathwise derivative for  $v$  on  $\mathbf{P}$ . By definition 3.3.2 and proposition 3.3.1, the efficient influence function for  $v$  is  $\Pi_0(h_0 | \mathbf{P})$  given by (8) which can be written

$$(10) \quad \Sigma_Z^{-1} \{ p^{-1} U (Z - EZ) (Y^T - Z^T v) \\ - (p^{-1} U - 1) [(E(Z | Y) - EZ) Y^T - E((Z - EZ) Z^T | Y) v] \}.$$

If  $m = d = 1$ , the variance of the second term in (10) is, we expect, the reduction in variance gained by using an efficient estimator. It is interesting to note that if  $Y$  is independent of  $Z$  and  $v = 0$  the gain is zero. That is, if the linear model  $Y = vZ + \varepsilon$  is assumed and we test  $H: v = 0$  then the naive test statistic is efficient. Asymptotically, using the observations with missing values is of no use. We will sketch the construction of efficient estimates for this example in example 7.8.2.  $\square$

**Example 3. Regression with missing observations on the covariates.**

Suppose again that  $X^0 = (Y, Z, U)$  takes values in  $R^m \times R^d \times \{0, 1\}$ ,  $U$  is independent of  $(Y, Z)$  and  $X = (Y, UZ, U)$ . This time we assume that  $(Y, Z)$  follows a regression model as in section 4.3. Then if  $\mu$  is a product measure on  $R^m \times R^d \times \{0, 1\}$  with the last component counting measure,  $Q \in \mathbf{Q}$  has density

$$(11) \quad \frac{dQ}{d\mu}(y, z, u) = h(z) f(y | \theta, z) \{ pu + (1-p)(1-u) \},$$

where  $f(y | \theta, z)$ ,  $\theta \in \Theta$ ,  $z \in R^d$ , is a fully parametric family. Suppose, for simplicity that  $p \equiv P(U = 1)$  is known. We compute the efficient score function for  $\theta$ . By (A.5.13) and (4) the score function for  $\theta$  when  $h$  is known is

$$(12) \quad \dot{\mathbf{i}}(X, P | \theta, \mathbf{P}_1) = E(\dot{\mathbf{i}} | X) = U \dot{\mathbf{i}} + (1-U) E(\dot{\mathbf{i}} | Y),$$

where

$$\dot{\mathbf{i}} = \frac{\dot{f}}{f}(Y | \theta, Z) = \dot{\mathbf{i}}(X, P | \theta, \mathbf{Q})$$

is the score function for  $\theta$  in the model  $\mathbf{Q}$  described by (11). Assume (2) holds. Then

$$\dot{\mathbf{Q}}_2 = \{ b(Z): b \in L_2^0(\mathbf{Q}) \}$$

and



$$(13) \quad \dot{P}_2 = \{E(b(Z) | X) : b \in \dot{Q}_2\}.$$

As in (4), define  $A : \dot{Q}_2 \rightarrow \dot{P}$  by

$$(14) \quad Ab = E(b(Z) | X) = Ub(Z) + (1-U)E_Q(b(Z) | Y).$$

Then  $A^T : \dot{P} \rightarrow \dot{Q}_2$  is given by

$$(15) \quad A^T w = pE_Q(w(Y, Z, 1) | Z) + (1-p)E_Q(w(Y, 0, 0) | Z).$$

We now introduce the operators  $E_Q^Y, E_Q^Z$  of conditional expectations given  $Y$  and  $Z$  respectively (in the model  $Q$ ). We obtain

$$(16) \quad \begin{aligned} A^T A b(Z) &= pb(Z) + (1-p)E_Q^Z E_Q^Y b(Z) \\ &= [pI + (1-p)E_Q^Z E_Q^Y] b(Z), \end{aligned}$$

where  $I$  is the identity. Now  $E_Q^Z E_Q^Y : \dot{Q}_2 \rightarrow \dot{Q}_2$  is a self-adjoint nonnegative definite operator. Hence, as in the proof of proposition 4.5.1 of the Has'minskii-Ibragimov model,  $A^T A$  has a bounded inverse,  $\|(A^T A)^{-1}\| \leq p^{-1}$ . We conclude that

$$(17) \quad \begin{aligned} \Pi_0(\dot{i}(\cdot, P | \theta, P_1) | \dot{P}_2) &= A(A^T A)^{-1} A^T \dot{i}(\cdot, P | \theta, P_1) \\ &= (1-p)A \{pI + (1-p)E_Q^Z E_Q^Y\}^{-1} E_Q^Z E_Q^Y \dot{i}, \end{aligned}$$

which yields  $I^*(\cdot, P | \theta, P)$  by (3.4.2). Note that the same operator enters here as in proposition 4.5.1 of the Has'minskii-Ibragimov model. Efficient estimation in this model is developed by Bajamonde (1991).

If the joint distribution of  $(Y, Z)$  is assumed to be normal, then the model  $Q$ , and hence also the model  $P$ , is simply a parametric family. The likelihood function for the parametric model  $P$  is not simple, but the maximum likelihood estimator can be found easily by the EM algorithm (Dempster, Laird, and Rubin (1977)). □

**Example 4. Linear regression with right censoring.**

This example deals with the estimation of the regression coefficients of a linear regression model when the observations on the dependent variable are censored. Let  $X^0 = (Y, Z, C)$  take values in  $R \times R^m \times R$  and  $\mu$  be a product of Lebesgue measure on the real line with a measure on  $R^{m+1}$ . The model  $Q$  is defined as the set of all distributions such that  $\varepsilon \equiv Y - v^T Z$  is independent of  $(Z, C)$ . Thus

$$Q = \{Q : \frac{dQ}{d\mu} = f(y - v^T z) h(z, c) \text{ for some } f, h, \text{ and } v \in R^m\}.$$

We observe  $X = (Z, Y \wedge C, 1_{[Y \leq C]}) \equiv (Z, V, \Delta)$ .

This model arises in survival analysis when  $Y$  is a log failure time (the "accelerated time model"). See, for example, Kalbfleisch and Prentice (1980) and Lawless (1982). Estimators of  $v$  suitable for different submodels have been suggested by Miller (1976), Buckley and James (1979), and Koul, Susarla, and Van Ryzin (1981).

We begin by calculating scores in the model  $Q$ . If  $f$  has finite Fisher information for location  $I_f = \int ((f')^2 / f) < \infty$ , then the score for  $v$  is

$$(18) \quad \dot{I}_1(y, z | v, \mathbf{Q}) = -z \frac{f'}{f}(y - v^T z) \equiv z \psi(\varepsilon)$$

with  $\psi \equiv -f'/f$ . Furthermore, the score (operator) for  $f$  with distribution function  $F$  (see sections 3.4, 5.4, and 5.5) is

$$(19) \quad \dot{I}_2 a(y, z | \mathbf{Q}) = a(y - v^T z) = a(\varepsilon) \quad \text{for } a \in L_2^0(F).$$

To compute scores in the induced model  $\mathbf{P} = \mathbf{Q}T^{-1}$ , we use proposition A.5.5 and compute them as

$$(20) \quad \begin{aligned} \dot{I}_1(X | v, \mathbf{P}) &= E(\dot{I}_1(X^0 | v, \mathbf{Q}) | X), \\ \dot{I}_2 a(X | \mathbf{P}) &= E(\dot{I}_2 a(X^0 | \mathbf{Q}) | X). \end{aligned}$$

These calculations are in turn done via the martingale facts developed in section A.3. The end result is stated most easily in terms of the two counting process martingales  $\mathbf{M}$  and  $\mathbf{M}_{uc}$  defined by

$$(21) \quad \mathbf{M}(t) = 1_{[\varepsilon \leq t]} - \int_{-\infty}^t 1_{[\varepsilon \geq s]} d\Lambda(s)$$

and

$$(22) \quad \mathbf{M}_{uc}(t) = 1_{[\varepsilon \wedge \delta \leq t, \Delta=1]} - \int_{-\infty}^t 1_{[\varepsilon \wedge \delta \geq s]} d\Lambda(s),$$

where  $\delta \equiv C - v^T Z$  is the natural variable which censors  $\varepsilon$ . Also recall the  $R$  and  $L$  operators introduced in section A.1 and A.3. Note that

$$(23) \quad \begin{aligned} R\psi(t) &= \psi(t) - E(\psi(\varepsilon) | \varepsilon > t) \\ &= \psi(t) - \lambda(t) \\ &= -\frac{f'}{f}(t) - \frac{f}{1-F}(t) = -\frac{\lambda'}{\lambda}(t), \end{aligned}$$

while

$$(24) \quad \begin{aligned} J(t) &= \int_{-\infty}^t (R\psi)^2 dF \\ &= \int_{-\infty}^t \left(\frac{f'}{f}\right)^2 dF + \frac{f^2}{1-F}(t) = \int_{-\infty}^t \left(\frac{\lambda'}{\lambda}\right)^2 dF \end{aligned}$$

is the Fisher information for location based on observation of  $\varepsilon \wedge t$ . For  $a \in L_2(F)$ , let

$$A(t) \equiv E(a(\varepsilon) | \varepsilon > t).$$

The following proposition, which gives the efficient influence function for  $v$  in the model  $\mathbf{P} = \mathbf{Q}T^{-1}$ , is due to Ritov (1984); see also Ritov and Wellner (1988).

**Proposition 1.** Suppose that  $|Z| \leq M$  a.s. and  $f$  has finite Fisher information for location. Then, with  $\psi = -f'/f$ ,

$$(25) \quad \dot{I}_1(X | v, \mathbf{P}) = Z \int R\psi d\mathbf{M}_{uc}$$

$$\begin{aligned}
&= Z \{ \Delta R \psi(V - v^T Z) + \lambda(V - v^T Z) \} \\
&= Z \{ \Delta \psi(V - v^T Z) + (1 - \Delta) \lambda(V - v^T Z) \}, \\
\dot{I}_2 a(X | P) &= \int Ra \, dM_{uc} \\
(26) \quad &= \Delta Ra(V - v^T Z) + A(V - v^T Z) \\
&= \Delta a(V - v^T Z) + (1 - \Delta)A(V - v^T Z),
\end{aligned}$$

and the efficient score function for estimation of  $v$  is

$$\begin{aligned}
(27) \quad I_1^*(X, P | v, P) &= \int (Z - E(Z | \delta \geq s)) R \psi(s) \, dM_{uc}(s) \\
&= \int (Z - E(Z | V - v^T Z \geq s)) R \psi(s) \, dM_{uc}(s).
\end{aligned}$$

Thus the information for estimation of  $v$  is

$$\begin{aligned}
(28) \quad I(P | v, P) &= \int E \{ 1_{[\delta \geq s]} (Z - E(Z | \delta \geq s)) \\
&\quad (Z - E(Z | \delta \geq s))^T (R \psi(s))^2 \, dF(s) \} \\
&= E Z Z^T J(C - v^T Z) - \int D(s) D^T(s) \bar{G}(s) \, dJ(s),
\end{aligned}$$

where  $D(s) \equiv E(Z | \delta \geq s)$  and  $\bar{G}(s) \equiv P(\delta \geq s)$ .

**Proof.** For  $t \in R$  let

$$(a) \quad \mathcal{F}_t \equiv \sigma\{Z, 1_{[e \leq s]} : s \leq t\} = \sigma\{Z, \varepsilon \wedge t, 1_{[e \leq t]}\}.$$

Then  $M$  defined in (21) is a counting process martingale with respect to  $\{\mathcal{F}_t\}$ . By proposition A.3.6,

$$(b) \quad E(\dot{I}_1(X^0 | v, Q) | \mathcal{F}_t) = E(Z \psi(\varepsilon) | \mathcal{F}_t) = Z \int_{-\infty}^t R \psi \, dM$$

and

$$(c) \quad E(\dot{I}_2 a(X^0 | Q) | \mathcal{F}_t) = E(a(\varepsilon) | \mathcal{F}_t) = \int_{-\infty}^t Ra \, dM.$$

Since  $\delta \equiv C - v^T Z$  is independent of  $\varepsilon$ , while  $V - v^T Z = Y \wedge C - v^T Z = \varepsilon \wedge \delta$  and  $\Delta = 1_{[Y \leq C]} = 1_{[e \leq \delta]}$ , it follows that

$$(d) \quad \dot{I}_1(X | v, P) = E(\dot{I}_1(X^0 | v, Q) | X) = Z \int_{-\infty}^{\delta} R \psi \, dM = Z \int R \psi \, dM_{uc}$$

and

$$(e) \quad \dot{I}_2 a(X | P) = E(\dot{I}_2 a(X^0 | Q) | X) = \int_{-\infty}^{\delta} Ra \, dM = \int Ra \, dM_{uc};$$

i.e., (25) and (26) hold. The second forms in both (25) and (26) follow from

$$(f) \quad \int_{-\infty}^t Ra \, d\Lambda = -A(t)$$

(which is, in turn, an immediate consequence of  $L \circ R = \text{identity}$ ) and (23).

To find the efficient score function  $I_1^*$  for  $v$ , we want to use the approach of section 3.4, method 2, to calculate  $\dot{I}_2 a^* = \Pi_0(\dot{I}_1 | \dot{P}_2)$  where

$\dot{P}_2 = \{[\dot{I}_2 a : a \in L_2^0(F)]\}$ , and thereby obtain  $I_1^* = \dot{I}_1 - \Pi_0(\dot{I}_1 | \dot{P}_2) = \dot{I}_1 - \dot{I}_2 a^*$ . Thus we want to find  $a^*$  with  $\int a^* dF = 0$  so that

$$\dot{I}_1 - \dot{I}_2 a^* \perp \dot{I}_2 a \quad \text{in } L_2(P)$$

for all  $a \in L_2(F)$  with  $\int a dF = 0$ ; i.e.

$$E(\dot{I}_1 - \dot{I}_2 a^*) \dot{I}_2 a = 0 \quad \text{for all } a \in L_2^0(F).$$

We do this by conditioning on  $Z$  and  $\delta$  and using the martingale calculus based on the martingale  $M$ . Recall that if  $Y_i \equiv \int f_i dM$ ,  $i = 1, 2$ , where  $M$  is as in (21), then

$$\begin{aligned} \text{(g)} \quad E Y_1 Y_2 &= E \langle Y_1, Y_2 \rangle = E \int f_1 f_2 d\langle M \rangle \\ &= E \int f_1(s) f_2(s) 1_{[\varepsilon \geq s]} d\Lambda(s). \end{aligned}$$

Note that

$$\text{(h)} \quad \dot{I}_1 - \dot{I}_2 a^* = \int 1_{[\delta \geq \cdot]} (ZR\psi - Ra^*) dM.$$

Hence, using (g),

$$\begin{aligned} &E(\dot{I}_1 - \dot{I}_2 a^*) \dot{I}_2 a \\ &= E E\{(\dot{I}_1 - \dot{I}_2 a^*) \dot{I}_2 a | Z, \delta\} \\ &= E E\left\{\int 1_{[\delta \geq s]} (ZR\psi(s) - Ra^*(s)) Ra(s) 1_{[\varepsilon \geq s]} d\Lambda(s) | Z, \delta\right\} \\ &= E\left(\int 1_{[\delta \geq s]} (ZR\psi(s) - Ra^*(s)) Ra(s) 1_{[\varepsilon \geq s]} d\Lambda(s)\right) \\ &= E \int 1_{[\delta \geq s]} (ZR\psi(s) - Ra^*(s)) Ra(s) dF(s) \\ &\quad \text{since } \varepsilon \text{ is independent of } Z \text{ and } \delta \\ \text{(i)} \quad &= \int \{E(Z 1_{[\delta \geq s]}) R\psi(s) - E 1_{[\delta \geq s]} Ra^*(s)\} Ra(s) dF(s). \end{aligned}$$

It is easily seen that if

$$a^* = L(E(Z | \delta \geq \cdot) R\psi),$$

then

$$\text{(j)} \quad Ra^*(s) = E(Z | \delta \geq s) R\psi(s) = \frac{E(Z 1_{[\delta \geq s]})}{E 1_{[\delta \geq s]}} R\psi(s),$$

and the right side of (i) equals zero for all  $a$ . Thus by substitution of (j) into (h), (27) holds:

$$\begin{aligned} \text{(k)} \quad I_1^*(X, P | \nu, P) &= \int 1_{[\delta \geq s]} (Z - E(Z | \delta \geq s)) R\psi(s) dM(s) \\ &= \int (Z - E(Z | V - \nu^T Z \geq s)) R\psi(s) dM_{uc}(s) \end{aligned}$$

since  $Z$  is independent of  $\varepsilon$  and  $[\varepsilon \geq s, \delta \geq s] = [V - \nu^T Z = \varepsilon \wedge \delta \geq s]$ .

Finally, to calculate  $I(P | \nu, P)$  we use the martingale  $M_{uc}$ :

$$\begin{aligned}
E(I_1^* I_1^{*T}) &= E\langle I_1^*, I_1^{*T} \rangle \\
&= E \int 1_{[\delta \geq s]} (Z - E(Z | \delta \geq s)) (Z - E(Z | \delta \geq s))^T \\
&\quad (R\Psi)^2 1_{[e \geq s]} d\Lambda(s) \\
&= E \int 1_{[\delta \geq s]} (Z - E(Z | \delta \geq s)) (Z - E(Z | \delta \geq s))^T (R\Psi)^2 dF(s)
\end{aligned}$$

by independence of  $\varepsilon$  and  $(Z, \delta)$ . □

The proposition implies that there is no loss of information due to not knowing the joint distribution of  $(Z, C)$ . (This follows because  $I_2 a^*(X | P)$  is the projection of  $I_1(X | v, P)$  on the subset of  $\dot{P}_2$  given by the closure in  $L_2(P)$  of  $\{E(g(\varepsilon) | X) : g \in \dot{Q}_2\}$ , i.e., the tangent space when  $h$  is known; cf. (19) and (20).) There is however loss of information when the distribution of  $\varepsilon$  is unknown and it is given by the second term in (28):

$$\int D(s) D^T(s) \bar{G}(s) dJ(s).$$

It is easy to show that if  $F$  is known up to a shift, then the efficient score function for  $v$  is given by

$$\left( Z - \frac{EZJ(\delta)}{EJ(\delta)} \right) \left[ 1_{[e \leq \delta]} \Psi(\varepsilon) + 1_{[e > \delta]} \frac{f}{F}(\delta) \right],$$

and accordingly the information for  $v$  is

$$(29) \quad EZZ^T J(\delta) - \frac{E(ZJ(\delta))E(Z^T J(\delta))}{EJ(\delta)}.$$

By Fubini

$$\begin{aligned}
E[ZJ(\delta)] &= E(E(Z | \delta)J(\delta)) \\
(30) \quad &= \int E(Z | \delta=v) J(v) dG(v) \\
&= \iint E(Z | \delta=v) 1_{[s \leq v]} dJ(s) dG(v) \\
&= \int D(s) \bar{G}(s) dJ(s).
\end{aligned}$$

In the same way

$$(31) \quad EJ(\delta) = \int J(s) dG(s) = \int \bar{G}(s) dJ(s).$$

Comparing (28) to (29)–(31), we see that the loss of information when  $f$  is unknown relative to the case when  $f$  is a member of a shift family is given by

$$\begin{aligned}
(32) \quad &\int D(s) D^T(s) \bar{G}(s) dJ(s) - \frac{\int D(s) \bar{G}(s) dJ(s) \int D^T(s) \bar{G}(s) dJ(s)}{\int \bar{G}(s) dJ(s)} \\
&= \int (D(s) - d)(D(s) - d)^T \bar{G}(s) dJ(s),
\end{aligned}$$

with

$$d = \frac{\int D(s) \bar{G}(s) dJ(s)}{\int \bar{G}(s) dJ(s)} = \frac{EZJ(\delta)}{EJ(\delta)}$$

Hence there is, in general, a loss of information. This loss is 0 if  $\delta = C - v^T Z$  is independent of  $Z$  (e.g.,  $v = 0$  and the censoring variable  $C$  is independent of the covariates) which is the assumption needed to ensure consistency of Miller's estimator; see, e.g., James and Smith (1984).

It is not too hard to modify either  $\hat{I}_1$  in (25) or  $I_1^*$  in (27) to obtain estimating equations which yield inefficient estimates of  $v$ . The first approach was taken by Buckley and James (1979), and Ritov (1984), (1990), while the second approach was followed by Tsiatis (1990) (by reasoning from the perspective of rank tests). We briefly describe these two approaches.

To get an estimating equation from  $\hat{I}_1$  in (25), fix  $\psi_0$  corresponding to some  $F_0$ ; Buckley and James (1979) chose  $\psi_0(x) = x$  corresponding to  $F_0$  standard normal. Then replace  $\psi$  by  $\psi_0$  in the first form of  $\hat{I}_1$  given in (25). This yields, by (f) of the proof of proposition 1,

$$\begin{aligned} W_0(v, F)(Z, V, \Delta) &\equiv Z \int R \psi_0 dM_{uc} \\ &= Z \left\{ \Delta \psi_0(V - v^T Z) + (1 - \Delta) \frac{\int_{V - v^T Z}^{\infty} \psi_0(u) dF(u)}{1 - F(V - v^T Z)} \right\}, \end{aligned}$$

which has expectation 0 under any  $P \in \mathbf{P}$ . Suppose  $(Z_i, V_i, \Delta_i)$ ,  $i = 1, \dots, n$ , are i.i.d.  $P \in \mathbf{P}$ . Since the efficient estimator of  $F$  if we knew  $v$  is the Kaplan-Meier estimator  $\hat{F}_n^v$  based on  $\{(V_i - v^T Z_i, \Delta_i) : i = 1, \dots, n\}$ , we are led naturally to solving the equation

$$(33) \quad \sum_{i=1}^n W_0(v, \hat{F}_n^v)(Z_i, V_i, \Delta_i) = 0$$

for  $v = \hat{v}_n$  as an estimator of  $v$ ; the resulting estimator is the Buckley-James (1979) estimator (when  $\psi_0(x) = x$ ).

The alternative approach of Tsiatis (1990) is essentially based on a score equation related to the efficient score  $I_1^*$  in (27). Fix  $F_0(x) = 1 - \exp(-e^x)$ , the double exponential or extreme value distribution, function of  $\varepsilon = \log(\tilde{\varepsilon})$  with  $\tilde{\varepsilon} \sim \text{exponential}$ . Then  $\psi_0(x) = e^x - 1$ ,  $\lambda_0(x) = e^x$ , and  $R_0 \psi_0(x) = \psi_0(x) - \lambda_0(x) = -1$ . If we replace  $R\psi$  in (27) by  $R_0 \psi_0$  and ignore the compensator of  $M_{uc}$ , we obtain

$$\begin{aligned} T_0(v, D)(Z, V, \Delta) &= \Delta(Z - D(V - v^T Z)) R_0 \psi_0(V - v^T Z) \\ &= -\Delta(Z - D(V - v^T Z)), \end{aligned}$$

which has expectation 0 under any  $P \in \mathbf{P}$ . Hence, as suggested by Tsiatis (1990), we are led to solving the (estimating) equations

$$(34) \quad \sum_{i=1}^n T_0(v, \hat{D})(Z_i, V_i, \Delta_i) = 0$$

for  $v = \hat{v}_n$  as yet another (family of) possible estimator(s) where  $\hat{D}(u) \equiv \sum_i Z_i 1_{[v_i - v^T Z_i \geq u]} / \sum_i 1_{[v_i - v^T Z_i \geq u]}$ . Ritov (1990) studies the relationships between the estimators resulting from (33) and (34). We return to these estimators in example 7.7.2. □

### 4.7 TRANSFORMATION MODELS

In this section we describe a large class of models, all of which involve one or more (monotone) transformations  $\tau: R \rightarrow R$  as parameters. Since this group  $T$  of transformations will arise repeatedly throughout the section, we define it carefully now.

**Definition 1.** Let  $T$  denote the group of all strictly increasing continuous functions from  $R$  to  $R$ . We call  $T$  the *transformation group*.

The simplest transformation model of the type we want to study here is a semiparametric generalization of the classical linear regression model different from the generalizations introduced in section 4.3. Now we allow a (monotone) transformation of the response instead of a (usually nonmonotone) transformation of one or more predictors as in section 4.3.

**Example 1. Linear regression-transformation model.**

Suppose that  $Z \sim H$  on  $R^k$  and  $\varepsilon \sim G_0$  are independent; we suppose (for now) that both  $H$  and  $G_0$  are known. (The assumption that  $G_0$  is known will be relaxed in example 6 below.)

Suppose that for some  $v \in R^k$

$$(1) \quad T = v^T Z + \varepsilon,$$

and we observe  $X = (Z, Y)$ , where  $Y = \tau^{-1}(T)$  for some  $\tau \in T$ . Equivalently,

$$(2) \quad \tau(Y) = v^T Z + \varepsilon,$$

with  $\tau \in T, v \in R^k$ .

If  $\varepsilon \sim N(0, 1)$ , this is the generalization of the Box-Cox (1964) model studied by Doksum (1987). If  $e^\varepsilon \sim \text{Pareto}(\eta), P(\varepsilon \geq t) = (1 + \eta e^t)^{-1/\eta}$  with  $\eta > 0$ , then (2) is the semiparametric Pareto regression-transformation model, which includes the Cox model as the limiting special case  $\eta = 0$ , introduced by Clayton and Cuzick (1985a), (1986). This model can be viewed as the Cox-model example 3.4.2 with an unobserved covariate, or *frailty*,  $W > 0$ , as follows. Suppose that the conditional (on  $Z$  and  $W$ ) hazard function is

$$\lambda(y|z, w) = w e^{v^T z} \lambda_0(y),$$

or

$$\Lambda(y|z, w) = w e^{v^T z} \Lambda_0(y).$$

Thus

$$\bar{F}(y|z, w) = \exp[-w e^{v^T z} \Lambda_0(y)]$$

and, if  $W \sim \text{Gamma}(1/\eta, 1/\eta)$  with density

$$f_W(w) = \frac{(1/\eta)^{1/\eta} w^{1/\eta-1} \exp(-w/\eta)}{\Gamma(1/\eta)}, \quad w > 0,$$

it follows that

$$\begin{aligned} \bar{F}(y | z) &= \int_0^\infty \exp[-w e^{v^T z} \Lambda_0(y)] \left(\frac{1}{\eta}\right)^{1/\eta} \frac{w^{1/\eta-1}}{\Gamma(1/\eta)} \exp\left(-\frac{w}{\eta}\right) dw \\ &= [1 + \eta e^{v^T z} \Lambda_0(y)]^{-1/\eta}. \end{aligned}$$

Hence,

$$(3) \quad \log \Lambda_0(Y) = -v^T Z + \varepsilon,$$

where  $e^\varepsilon \sim \text{Pareto}(\eta)$ . Note that, apart from the irrelevant minus sign, this is a particular case of (2) with  $\tau(y) = \log \Lambda_0(y) \uparrow$ . In the Cox model, (3) holds with  $e^\varepsilon \sim \text{exponential}(1)$ .  $\square$

Example 1 has many interesting variants, extensions, and relatives. A first natural generalization is to relax the linear regression structure required in (1).

### Example 2. Regression-transformation model.

Suppose that  $(T | Z = z) \sim P_\theta(\cdot | z)$  for some  $\theta \in \Theta$  where  $Z \sim H$  on  $R^k$  is known and, for each fixed  $z$ ,  $\{P_\theta(\cdot | z) : \theta \in \Theta\}$  is a regular parametric model on  $R$ . Then we observe  $X = (Z, Y) = (Z, \tau^{-1}(T))$  for some  $\tau \in \mathbf{T}$ . Of course, example 1 is the special case with  $P_\theta(T \in A | Z = z) = P_0(T - \theta^T z \in A | Z = z)$  with density  $g_0(t - \theta^T z)$  if  $G_0$  has density  $g_0$ . This model was studied by Bickel (1986) in the context of the two-sample problem with  $P(Z = 0) = 1 - P(Z = 1)$ .  $\square$

Note that the distribution  $H$  of  $Z$  in example 2 was assumed to be known, and therefore not a function of  $\theta$ . This is not the case in our second extension of example 1.

### Example 3. Joint distribution-transformation model.

Suppose that  $(Z, T) \sim P_\theta$  for some  $\theta \in \Theta \subset R^k$ , where  $P_0 \equiv \{P_\theta : \theta \in \Theta\}$  is a regular parametric model on  $R^{c+1}$ ; we call  $P_0$  the *core model*. Often  $\theta = (v, \eta)$  where  $v \in R^m$  is a vector of regression parameters and  $\eta \in R^{k-m}$  are additional (nuisance) parameters. We observe  $X = (Z, Y) \equiv (Z, \tau^{-1}(T)) \sim P_{\theta, \tau}$  for some  $\tau \in \mathbf{T}$ . Thus, if  $P_\theta$  has density  $p_0(z, t; \theta)$  with respect to  $\mu \times$  Lebesgue measure on  $R^{c+1}$ , where  $\mu$  is a  $\sigma$ -finite measure on  $R^c$ , then

$$\begin{aligned} P_{\theta, \tau}(Z \in A, Y \leq y) &= \int_A \int_{-\infty}^{\tau(y)} p_0(z, t; \theta) dt d\mu(z) \\ &= P_{\theta, t}(Z \in A, Y \leq \tau(y)) \\ &= P_\theta(Z \in A, T \leq \tau(y)) \end{aligned}$$

for Borel sets  $A$  in  $R^c$ ; or, if  $\tau$  is absolutely continuous with derivative  $\tau'$ , the density of  $P_{\theta, \tau}$  is



$$p(z, y; \theta, \tau) = p_0(z, \tau(y); \theta) \tau'(y).$$

An interesting special case of this example is when  $Z \in R$ ,  $\mu$  is Lebesgue measure, and  $(Z, T)$  has a distribution  $F_\theta(z, t)$  on  $R \times [0, 1]$  where the marginal distribution of  $T$  is *Uniform*(0, 1). Then the unknown marginal distribution  $G$  of  $Y \equiv \tau^{-1}(T)$  is just the transformation  $\tau$ .

This is a *copula model* in which the marginal distribution of  $Y$  is completely unknown, and the marginal distribution of  $Z$  as well as the dependence between  $Z$  and  $Y$  depends only on the finite-dimensional parameter  $\theta$ . In fact, example 1 with  $k = 1$  is an example of this.  $\square$

Our next example considers copula models more generally, and allows one or both marginal distributions to be completely unknown (or nonparametric), but with the dependence between the variables specified by a completely parametric family of distributions.

**Example 4. Copula models.**

Suppose that  $C_\theta$  is a distribution function on  $[0, 1]^2$  for  $\theta \in \Theta \subset R^k$ . Let  $I$  denote Lebesgue measure and the identity function  $I(u) = u$  on  $[0, 1]$ . We suppose that  $C_\theta$  has uniform marginals and density  $c_\theta$  with respect to Lebesgue measure  $I \times I$  on  $[0, 1]^2$ .

If  $(U, V) \sim C_\theta$  for some  $\theta \in \Theta$ , we observe  $X \equiv (S, T) \equiv (\sigma^{-1}(U), \tau^{-1}(V)) = (G^{-1}(U), H^{-1}(V))$ , where  $G$  and  $H$  are df's on  $R$  with densities  $g$  and  $h$  respectively. Thus, the joint distribution function and density function of  $X = (S, T)$  are given by

$$(4) \quad F_{S,T}(s, t) = C_\theta(G(s), H(t))$$

and

$$(5) \quad f_{S,T}(s, t) = f_{S,T}(s, t; \theta, G, H) = c_\theta(G(s), H(t)) g(s) h(t)$$

respectively.  $\square$

Because of its relative simplicity, we will study this example of a transformation model in some detail. To begin we will examine the special case of the model in which the marginal distribution  $H$  of  $T$  is known, and, without loss of generality, is *Uniform*(0, 1).

**Example 4.A. Copula model with one unknown marginal df.**

In this submodel, we suppose that the marginal df  $H$  of  $T$  is known; without loss we may assume that  $T \sim \text{Uniform}(0, 1)$ , or equivalently  $H(t) = t$ ,  $0 \leq t \leq 1$ . Thus the joint distribution and density functions are

$$(6) \quad F_{S,T}(s, t) = C_\theta(G(s), t)$$

and

$$(7) \quad f_{S,T}(s, t) = c_\theta(G(s), t) g(s)$$

respectively. The resulting model  $\mathbf{P}_A$  involves the finite-dimensional parameter  $\theta \in \Theta$  and the infinite-dimensional parameter  $\sigma = G$  or equivalently its density  $g$ .  $\square$

**Example 4.B. Copula model with two unknown marginals.**

In this version of example 4, both marginal distributions are completely unknown, and thus the joint distribution and density functions of  $X = (S, T)$  are given by (4) and (5). The resulting model  $\mathbf{P}_B$  involves the finite-dimensional parameter  $\theta \in \Theta$  and the two infinite-dimensional parameters  $\sigma = G$  and  $\tau = H$ , or equivalently the densities  $g, h$  of  $G, H$ .  $\square$

Example 4 contains many interesting special cases by considering different parametric families of df's  $C_\theta$  on  $[0, 1]^2$  with uniform marginals. We now give a few of these cases as illustrations without any attempt to be exhaustive. All of the following examples of families  $\{C_\theta\}$  involve a one-dimensional parameter  $\theta \in \Theta \subset R$ ; it would be interesting to consider more general parametric families  $C_\theta$  with  $\theta \in \Theta \subset R^k$ .

**Example 4.I. Archimedean copulas.**

Let  $\Psi$  denote the collection of functions  $\psi : [0, 1] \rightarrow [0, \infty]$  satisfying  $\psi(1) = 0, \psi(0) = \infty, \psi'(t) < 0, \psi''(t) > 0$ . For any  $\psi \in \Psi$ ,

$$(8) \quad C(u, v) \equiv \psi^{-1}(\psi(u) + \psi(v)), \quad (u, v) \in [0, 1]^2,$$

is a distribution function on  $[0, 1]^2$  with uniform marginals; see, e.g., Kimeldorf and Sampson (1975a,b), Schweizer and Sklar (1983), or Genest and MacKay (1986a,b). If  $\{\psi_\theta : \theta \in \Theta\}$  is a (smooth) parametric family of functions with  $\psi_\theta \in \Psi$  for each  $\theta \in \Theta$ , then  $\{C_\theta : \theta \in \Theta\}$  defined by (8), with  $\psi$  replaced by  $\psi_\theta$ , is a parametric family of copula functions. Here are some further special cases.

**I.1.**  $\psi_\theta(u) = (u^{-\theta} - 1)/\theta, 0 \leq u \leq 1, 0 < \theta$ . This choice of  $\psi_\theta$  yields the copula models studied by Clayton (1978), Oakes (1982). In parallel to our discussion of example 1, it can be derived via mixing of an independence model with a Gamma distribution of the common proportionality of hazards or "frailty." The resulting copula function is

$$(9) \quad C_\theta(u, v) = [u^{-\theta} + v^{-\theta} - 1]^{-1/\theta}.$$

If  $(U, V) \sim C_\theta$ , the corresponding distribution of  $(U^\#, V^\#) \equiv (1 - U, 1 - V)$  yields the copula function  $C_\theta^\#$  with

$$(10) \quad \bar{C}_\theta^\#(u, v) = [(1 - u)^{-\theta} + (1 - v)^{-\theta} - 1]^{-1/\theta},$$

since

$$\begin{aligned} \bar{C}_\theta^\#(u, v) &\equiv P_\theta(U^\# \geq u, V^\# \geq v) \\ &= P_\theta(U \leq 1 - u, V \leq 1 - v) \\ &= C_\theta(1 - u, 1 - v). \end{aligned}$$

This is actually the copula model studied by Clayton (1978), Oakes (1982), and Maguluri (1986), (1993); if  $\lambda^\#$  denotes a corresponding hazard rate, then this model has the property that

$$\lambda^\#(u | V^\# = v) = (1 + \theta)\lambda^\#(u | V^\# \geq v).$$

I.2.  $\psi_\theta(u) = (-\log u)^{1/\theta}$ ,  $0 \leq u \leq 1$ ,  $0 < \theta < 1$ . This family  $\psi_\theta$  corresponds to the strictly stable frailty distributions advocated by Hougaard (1986a), (1986b). The resulting family of copula functions is given by

$$(11) \quad C_\theta(u, v) = \exp\{-[(-\log u)^{1/\theta} + (-\log v)^{1/\theta}]^\theta\},$$

which is related to Gumbel's (1960) bivariate distribution (of type B) of extreme values.

I.3.  $\psi_\theta(u) = \log\left(\frac{1-\theta}{1-\theta^u}\right)$ ,  $0 \leq u \leq 1$ ,  $0 < \theta$ ,  $\theta \neq 1$ . The corresponding family of copula functions is

$$(12) \quad C_\theta(u, v) = \log\{1 - (1 - \theta^u)(1 - \theta^v)/(1 - \theta)\} / \log \theta, \quad (u, v) \in [0, 1]^2.$$

This is a family introduced by Frank (1979). It can be shown that this is the only Archimedean copula such that  $L_\theta(1-U, 1-V) = L_\theta(U, V)$ , i.e., such that  $1 - C_\theta(1-u, 1) - C_\theta(1, 1-v) + C_\theta(1-u, 1-v) = C_\theta(u, v)$ .  $\square$

**Example 4.II. Morgenstern distributions.**

There are many copula functions  $C_\theta$  that are *not* Archimedean—not of the form (8). The well-known Morgenstern (1956) distribution

$$(13) \quad C_\theta(u, v) = uv[1 + \theta(1-u)(1-v)], \quad (u, v) \in [0, 1]^2, \quad |\theta| < 1,$$

with density

$$c_\theta(u, v) = 1 + \theta(1-2u)(1-2v), \quad (u, v) \in [0, 1]^2,$$

is a particularly simple example.  $\square$

**Example 4.III. Bivariate normal copula function.**

Suppose that  $\Phi_\theta$  is the bivariate normal distribution function with zero means, unit variances, and correlation  $\theta \in (-1, 1)$ , let  $\phi_\theta$  be its density, and let  $\Phi$  denote the standard normal distribution function on  $R$  with density  $\phi = \Phi'$ . Thus  $\Phi_\theta(s, t)$  has marginals  $\Phi(s)$  and  $\Phi(t)$ , and

$$C_\theta(u, v) \equiv \Phi_\theta(\Phi^{-1}(u), \Phi^{-1}(v))$$

is a (non-Archimedean) copula function with density

$$c_\theta(u, v) = \phi_\theta(\Phi^{-1}(u), \Phi^{-1}(v)) \frac{1}{\phi(\Phi^{-1}(u))\phi(\Phi^{-1}(v))}. \quad \square$$

**Example 4.IV. Plackett's constant odds model.**

Suppose that  $\theta \in (0, \infty)$ , and  $(u, v) \in [0, 1]^2$ . Then we obtain a copula function  $C_\theta$  if we require that  $U$  and  $V$  have constant *odds ratio*  $\theta$ , i.e.,

$$\theta = \frac{C_\theta(u, v)(1-u-v+C_\theta(u, v))}{(u-C_\theta(u, v))(v-C_\theta(u, v))} \quad \text{for all } (u, v) \in [0, 1]^2.$$

Thus,  $C_\theta$  satisfies a quadratic equation. This is Plackett's (1965) copula function.  $\square$

For other examples of copula models, see Oakes (1986) and Marshall and Olkin (1988).

**Example 5. Regression-copula models.**

Suppose that  $(Z, T) \sim P_\theta$  for some  $\theta \in \Theta \subset R^k$  where  $\{P_\theta : \theta \in \Theta\}$  is a regular parametric model on  $R^{c+d}$ ; here  $Z \in R^c$  and  $T \in R^d$ . We observe  $X = (Z, Y)$  where  $T \equiv (\tau_1(Y_1), \dots, \tau_d(Y_d))$  for some  $(\tau_1, \dots, \tau_d) \in \mathbf{T}^d$ . This model contains all the preceding models as special cases; it also contains the multivariate regression models introduced by Clayton and Cuzick (1985a). Because of its complexity, it will not be analyzed in full detail here. Nonetheless, the analysis which we will give below of the simpler models 3 and 4 can be extended to this more general model.  $\square$

**Example 6. Linear regression-transformation model with unknown error distribution.**

Suppose, as in example 1, that

$$(14) \quad T = v^T Z + \varepsilon \quad \text{for some } v \in N \subset R^k,$$

where  $Z \sim H$  is known, but now suppose that  $\varepsilon$  has distribution function  $G$  with density  $g$  which is completely unknown. Again we observe  $X = (Z, Y)$  with  $\tau(Y) = T$  for some (unknown)  $\tau \in \mathbf{T}$ . Now the model has parameters  $v \in R^k$ ,  $\tau \in \mathbf{T}$ , and  $g \in \mathbf{G}$ , the collection of all densities on  $R$ . With some thought, it is clear that  $v$  is no longer identifiable without some restriction on the scale of the error distribution  $G$ . However, as implied by the results of Han (1987),  $v/|v|$  is identifiable in general.  $\square$

*Tangent Spaces, Efficient Score, and Information Calculations for Example 3*

In order to tie this example to example 4 we reparametrize so that the marginal distribution of  $Y$  under  $\tau = I$  is *Uniform*(0, 1) at  $\theta = \theta_0$ . We do this as follows: let  $I$  denote the identity function on  $R$  and identity element of the group  $\mathbf{T}$ , and let

$$F_\theta(t) \equiv P_\theta(T \leq t) = P_{(\theta, I)}(Y \leq t).$$

Define

$$\begin{aligned} \tilde{P}_\theta(Z \in A, T \leq u) &\equiv P_\theta(Z \in A, F_{\theta_0}(T) \leq u) \\ &= P_\theta(Z \in A, T \leq F_{\theta_0}^{-1}(u)) \end{aligned}$$

so that

$$\tilde{P}_\theta(T \leq u) = P_{\theta_0}(T \leq F_{\theta_0}^{-1}(u)) = u, \quad 0 \leq u \leq 1.$$

Then, with  $\tilde{\tau} \equiv F_{\theta_0}^{-1} \circ \tau$ ,

$$\begin{aligned} \tilde{P}_{(\theta, \tau)}(Z \in A, Y \leq y) &\equiv \tilde{P}_{(\theta, I)}(Z \in A, Y \leq \tau(y)) \\ &\equiv \tilde{P}_\theta(Z \in A, T \leq \tau(y)) \\ &= P_\theta(Z \in A, T \leq F_{\theta_0}^{-1}(\tau(y))), \\ &= P_{(\theta, \tilde{\tau})}(Z \in A, Y \leq y), \end{aligned}$$

and hence

$$\tilde{P}_{(\theta_0, \tau)}(Y \leq y) = P_{\theta_0}(T \leq F_{\theta_0}^{-1}(\tau(y))) = \tau(y).$$

Thus  $\tilde{P}_{(\theta, \tau)} = P_{(\theta, \tilde{\tau})}$  with  $\tilde{\tau} = F_{\theta_0}^{-1} \circ \tau$ , and  $\tau$  is the marginal df of  $Y$  under  $\tilde{P}_{(\theta_0, \tau)}$ .

**Convention:** We will henceforth assume that the reparametrization  $\tilde{\tau} = F_{\theta_0}^{-1} \circ \tau$  has been carried out, and we drop the tilde's. Thus, without loss of generality, the marginal distribution of  $T$  under  $P_{\theta_0}$  is *Uniform*(0, 1), the marginal distribution of  $Y$  under  $P_{(\theta_0, \tau)}$  is  $\tau \equiv G$ , and we identify  $\tau$  with the marginal distribution function  $G$  with density  $g$ . Furthermore, we denote  $T$  under  $P_{\theta_0}$  by  $U$  to stress its uniformity and to connect up with example 4. It is important to note that  $T$  need not be uniform under  $P_{\theta}$  for  $\theta \neq \theta_0$ .

Now we proceed with the calculation of the tangent spaces at  $\theta = \theta_0$ , suppressing the subscript 0. We will assume throughout that the following logarithmic derivatives of the density  $p_0(\cdot, \cdot; \theta)$  of the core model  $\mathbf{P}_0$  exist:

$$(15) \quad \dot{\mathbf{i}}_{\theta}(z, u; \theta) \equiv \dot{\mathbf{i}}_{\theta}(z, u) \equiv \nabla_{\theta} \log p_0(z, u; \theta),$$

$$(16) \quad \dot{\mathbf{i}}_u(z, u; \theta) \equiv \dot{\mathbf{i}}_u(z, u) \equiv \frac{\partial}{\partial u} \log p_0(z, u; \theta),$$

$$(17) \quad \ddot{\mathbf{i}}_{\theta u}(z, u; \theta) \equiv \ddot{\mathbf{i}}_{\theta u}(z, u) \equiv \frac{\partial}{\partial u} \dot{\mathbf{i}}_{\theta}(z, u).$$

In proving our results we need regularity conditions like

$$(18) \quad \int \frac{\partial}{\partial u} p_0(z, u; \theta) d\mu(z) = \frac{\partial}{\partial u} \int p_0(z, u; \theta) d\mu(z),$$

$$(19) \quad \int \nabla_{\theta} p_0(z, u; \theta) d\mu(z) = \nabla_{\theta} \int p_0(z, u; \theta) d\mu(z).$$

Note that these interchanges are valid trivially in the important special case that the support of  $\mu$  can be chosen to be finite. Furthermore, we define the functions  $\alpha$ ,  $\gamma$ , and  $\delta$ , when they exist, as follows:

$$(20) \quad \alpha(u) \equiv E_{\theta}[\dot{\mathbf{i}}_u^2(Z, U; \theta) \mid U = u] = \int \dot{\mathbf{i}}_u^2(z, u; \theta) p_0(z, u; \theta) d\mu(z),$$

$$(21) \quad \gamma(u) \equiv E_{\theta}[\dot{\mathbf{i}}_{\theta}(Z, U; \theta) \dot{\mathbf{i}}_u(Z, U; \theta) \mid U = u] \\ = \int \dot{\mathbf{i}}_{\theta}(z, u; \theta) \dot{\mathbf{i}}_u(z, u; \theta) p_0(z, u; \theta) d\mu(z),$$

$$(22) \quad \delta(u) \equiv E_{\theta}[\ddot{\mathbf{i}}_{\theta u}(Z, U; \theta) \mid U = u] = \int \ddot{\mathbf{i}}_{\theta u}(z, u; \theta) p_0(z, u; \theta) d\mu(z).$$

Recall that  $I$  denotes the identity function  $I(u) = u$ , the identity element of the transformation group  $\mathbf{T}$ , and Lebesgue measure on  $[0, 1]$ . In view of invariance of model 3 under monotone transformations of the  $Y$  axis, it is clear that the tangent spaces at different points  $P_{\theta, \tau} = P_{\theta, G}$  of the model  $\mathbf{P}$  with  $\theta$  the same are related by simple monotone transformations. The following proposition formalizes this and relates the tangent space at a general point  $P_{\theta, G}$  to the tangent space at  $P_{\theta, I}$ .

**Proposition 1.** (Tangent spaces) Suppose that  $\log p_0(z, u; \theta)$  is continuously differentiable with respect to  $u$  and  $\theta$  with  $\dot{\mathbf{i}}_u, \dot{\mathbf{i}}_{\theta} \in L_2(P_{\theta})$ , that  $E_{\theta} \dot{\mathbf{i}}_{\theta} \dot{\mathbf{i}}_{\theta}^T$  is

nonsingular, that (18) holds and that  $\alpha(u) \leq M u^{-2}$  (or  $\alpha(u) \leq M(1-u)^{-2}$ ),  $0 < u < 1$  for some  $M < \infty$ . Then the tangent space of the model  $\mathbf{P}$  is given by

$$(23) \quad \dot{\mathbf{P}}(P_{(\theta, G)}) = \{h(z, G(y)) : h \in \dot{\mathbf{P}}(P_{(\theta, I)})\},$$

where

$$(24) \quad \dot{\mathbf{P}}(P_{(\theta, I)}) \supset \dot{\mathbf{P}}_{\theta} + \dot{\mathbf{P}}_g,$$

$$(25) \quad \dot{\mathbf{P}}_{\theta} = [\dot{\mathbf{i}}_{\theta}],$$

$$(26) \quad \dot{\mathbf{P}}_g \supset [\dot{\mathbf{i}}_g a : a \in L_2^0(I)],$$

and

$$(27) \quad \dot{\mathbf{i}}_g a(z, u) \equiv a(u) + \dot{\mathbf{i}}_u(z, u; \theta) \int_0^u a \, dI.$$

**Proof.** Let  $\chi : R \rightarrow R$  be a bounded function with a continuous bounded derivative  $\chi'$  such that  $\chi \geq \frac{1}{2}$ ,  $\chi' \geq 0$ ,  $\chi(0) = \chi'(0) = 1$ . For example, take  $\chi(x) = \frac{1}{2} + (1 + e^{-4x})^{-1}$ . Let  $a \in L_2^0(I)$  be extended to  $[0, \infty)$  and equal 0 on  $(1, \infty)$ . Define

$$(a) \quad G_{\eta}(u) = \min \left\{ 1, \int_0^u \chi(\eta a(v)) \, dv \right\},$$

$0 \leq u, \eta \in R$ , with derivative  $g_{\eta}(u)$ , and consider the parametric model with densities

$$(b) \quad p(z, u; \theta, \eta) = p_0(z, G_{\eta}(u); \theta) g_{\eta}(u), \quad 0 \leq u.$$

Using proposition 2.1.1 and part A of the next proposition and its proof, one may verify that this parametric model is regular in  $\theta$  and that the score function with respect to  $\eta$  at  $\eta = 0$  is well defined and equals  $\dot{\mathbf{i}}_g a(z, u)$ . The rest of the proof is standard.  $\square$

To compute efficient score functions, efficient influence functions and information for  $\theta$  in the model 3, in view of theorem 3.4.1 we need to compute the projection of  $\dot{\mathbf{i}}_{\theta}$  onto the orthocomplement of  $[\dot{\mathbf{i}}_g a]$ . Note that we have not shown equality in (24) or (26), and that hence projection on  $[\dot{\mathbf{i}}_g a]$  might yield an inefficient score function; recall the discussion in section 3.4. Typically, however, we will get the efficient one. We first show, under the simple boundedness hypothesis of proposition 1, that  $\dot{\mathbf{i}}_g$  is bounded and  $\mathbf{R}(\dot{\mathbf{i}}_g)$  is closed. This implies that the projection on  $[\dot{\mathbf{i}}_g a] = \mathbf{R}(\dot{\mathbf{i}}_g)$  exists and can be computed in terms of the projection operator

$$\Pi_0(\cdot | [\dot{\mathbf{i}}_g]) = \dot{\mathbf{i}}_g (\dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_g)^{-1} \dot{\mathbf{i}}_g^T.$$

This gives a qualitative "solution" to the projection problem in the model of example 3. To compute projections in example 3 we will use classical Sturm-Liouville theory as presented in e.g., Tricomi (1957) or Hille (1969).

**Proposition 2.** Suppose that the assumptions of proposition 1 are satisfied. Then:

- A.  $\dot{\mathbf{i}}_g$  is bounded;  $\|\dot{\mathbf{i}}_g\| \leq (1 + 4M)^{1/2}$ .
- B.  $\mathbf{R}(\dot{\mathbf{i}}_g)$  is closed and  $\mathbf{N}(\dot{\mathbf{i}}_g) = \{0\}$ .
- C.  $(\dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_g)^{-1}$  exists and is bounded.
- D. The projection operator  $\Pi_0(\cdot \mid [\dot{\mathbf{i}}_g])$  equals  $\dot{\mathbf{i}}_g(\dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_g)^{-1} \dot{\mathbf{i}}_g^T$ .

**Proof.** To prove A, we first note that since  $T = U \sim \text{Uniform}(0, 1)$  marginally,

$$\int p_0(z, u; \theta) d\mu(z) = 1,$$

and hence, by (18),

$$E[\dot{\mathbf{i}}_u(Z, U) \mid U = u] = 0 \quad \text{for almost all } u \in [0, 1].$$

Thus

$$(a) \quad E[\dot{\mathbf{i}}_g a(Z, U) \mid U] = a(U) \quad \text{a.s.}$$

and

$$\begin{aligned} E[(\dot{\mathbf{i}}_g a)^2] &= \text{Var}[\dot{\mathbf{i}}_g a] \\ &= \text{Var}[E(\dot{\mathbf{i}}_g a \mid U)] + E[\text{Var}(\dot{\mathbf{i}}_g a \mid U)] \\ &= \text{Var}[a(U)] + E\{E[(\dot{\mathbf{i}}_u(Z, U) \int_0^U a dI)^2 \mid U]\} \quad \text{by (a)} \end{aligned}$$

$$(b) \quad = \int_0^1 a^2 dI + E\{(\frac{1}{U} \int_0^U a dI)^2 U^2 \alpha(U)\}$$

$$(c) \quad \leq (1 + 4M) \int_0^1 a^2 dI$$

by the hypothesized bound on  $\alpha$  and by Hardy's inequality (A.1.23). Assertion A follows from (c).

To prove B, note that (a) implies

$$(d) \quad \dot{\mathbf{i}}_g^{-1} h(u) = E[h(Z, U) \mid U = u]$$

for  $h \in \mathbf{R}(\dot{\mathbf{i}}_g)$ . Since  $\dot{\mathbf{i}}_g^{-1}$  is a bounded operator, corollary A.1.3 applies, and we conclude that  $\mathbf{R}(\dot{\mathbf{i}}_g)$  is closed. That  $\mathbf{N}(\dot{\mathbf{i}}_g) = \{0\}$  follows immediately from (a).

Statement C follows from B (see corollary 5.4.2), and D is a consequence of C. □

**Proposition 3.** Suppose that the hypotheses of proposition 1 hold and that  $\dot{\mathbf{i}}_\theta(z, u)$  is continuously differentiable with respect to  $u$  with derivative  $\dot{\mathbf{i}}_{\theta u} \in L_2^k(P_\theta)$ . Then

$$(28) \quad \dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_g a(u) = a(u) + \int_0^1 (1_{[u \leq s]} - s) \left\{ \int_0^s a dI \right\} \alpha(s) ds$$

and

$$(29) \quad \mathbf{i}_g^T \mathbf{i}_\theta(u) = - \int_0^1 (1_{[u \leq s]} - s) \delta(s) ds,$$

where  $\alpha$  and  $\delta$  are given by (20) and (22).

**Proof.** Applying proposition A.1.5 on (27) and noting that the range of  $\mathbf{i}_g^T$  is in  $L_2^0(I)$ , we arrive at (28) and (29).  $\square$

To actually compute the projection of  $\mathbf{i}_\theta$  onto  $[\mathbf{i}_g a]$ , we need to solve

$$(30) \quad \mathbf{i}_g^T \mathbf{i}_g a_* = \mathbf{i}_g^T \mathbf{i}_\theta$$

for  $a_* \in L_2^0(I)$  (cf. theorems A.2.1 and A.2.2). In view of proposition 3, (30) becomes

$$(31) \quad a_*(u) + \int_0^1 (1_{[u \leq s]} - s) \left( \int_0^s a_* dl \right) \alpha(s) ds = - \int_0^1 (1_{[u \leq s]} - s) \delta(s) ds,$$

where  $\alpha$  and  $\delta$  are the fixed functions given in (20) and (22). Integrating with respect to  $u$  from 0 to  $t$  across (31), and letting

$$A_*(t) = \int_0^t a_* dl$$

yields

$$(32) \quad A_*(t) + \int_0^1 (s \wedge t - st) A_*(s) \alpha(s) ds = - \int_0^1 (s \wedge t - st) \delta(s) ds.$$

We let

$$K(s, t) \equiv s \wedge t - st$$

denote the Brownian bridge covariance kernel appearing in (32). If  $a'_* = A''_*$  exists, note that  $A_*$  satisfies the differential equation (differentiate (32) twice or (31) once)

$$(33) \quad A''_*(t) - \alpha(t) A_*(t) = \delta(t)$$

subject to the boundary conditions  $A_*(0) = A_*(1) = 0$  (since  $a_* \in L_2^0(I)$ ). This is a Sturm-Liouville differential equation; to solve it, we will draw upon classical Sturm-Liouville theory as given in Tricomi (1957, pages 127–136), or Hille (1969, chapter 8).

To solve (32) for  $A_*$ , we define an operator  $T$  by

$$TA(t) \equiv A(t) + \int_0^1 K(s, t) A(s) \alpha(s) ds \equiv (I + R)A(t);$$

this is the operator on the left side in (32). Multiplying across (32) by  $\sqrt{\alpha(t)}$  yields

$$\begin{aligned} \sqrt{\alpha(t)} A_*(t) + \int_0^1 K(s, t) \sqrt{\alpha(s) \alpha(t)} \sqrt{\alpha(s)} A_*(s) ds \\ = - \sqrt{\alpha(t)} \int_0^1 K(s, t) \delta(s) ds \end{aligned}$$

or

$$\tilde{A}_*(t) + \int_0^1 \tilde{K}(s, t) \tilde{A}_*(s) ds = - \sqrt{\alpha(t)} \int_0^1 K(s, t) \delta(s) ds$$



with

$$\tilde{A}_* \equiv \sqrt{\alpha} A_*, \quad \tilde{K}(s, t) \equiv K(s, t) \sqrt{\alpha(s)\alpha(t)}.$$

This is just as in Tricomi (1957, page 3). Therefore we define

$$(34) \quad \begin{aligned} \tilde{T}\tilde{A}(t) &\equiv \tilde{A}(t) + \int_0^1 \tilde{K}(s, t)\tilde{A}(s) ds \\ &= (I + \tilde{R})\tilde{A}(t), \end{aligned}$$

where  $\tilde{R}$  is self-adjoint and nonnegative definite since the kernel  $K$  is. Consequently

$$\begin{aligned} \|\tilde{A}\| \|\tilde{T}\tilde{A}\| &\geq \langle \tilde{A}, (I + \tilde{R})\tilde{A} \rangle = \|\tilde{A}\|^2 + \langle \tilde{A}, \tilde{R}\tilde{A} \rangle \\ &\geq \|\tilde{A}\|^2, \end{aligned}$$

and by corollary A.1.2  $\tilde{T}$  has a bounded inverse  $\tilde{T}^{-1}$  with  $\|\tilde{T}^{-1}\| \leq 1$ .

Hence,

$$A_*(t) = (\alpha(t))^{-1/2} \tilde{T}^{-1}(-\sqrt{\alpha(t)} \int_0^1 K(s, t) \delta(s) ds)$$

exists if

$$\sqrt{\alpha(\cdot)} \int_0^1 K(s, \cdot) \delta(s) ds \in L_2(I).$$

Note that, for  $\alpha$  as in proposition 1, this holds if for some  $M < \infty, \varepsilon > 0$ ,

$$(35) \quad |\delta(s)| \leq M[s(1-s)]^{\varepsilon-3/2}.$$

As in Tricomi (1957) and Hille (1969), we introduce a Green's function  $\Delta(u, v)$  solving

$$(36) \quad \Delta(u, v) + \int_0^1 K(u, s)\Delta(s, v)\alpha(s) ds = K(u, v);$$

or, in terms of  $\tilde{\Delta}(u, v) \equiv \sqrt{\alpha(u)} \Delta(u, v)$ ,

$$\tilde{\Delta}(u, v) + \int_0^1 \tilde{K}(u, s)\tilde{\Delta}(s, v) ds = \sqrt{\alpha(u)}K(u, v).$$

Since  $\sqrt{\alpha(\cdot)}K(\cdot, v) \in L_2(I)$  for each  $v \in [0, 1]$ ,  $\tilde{\Delta}(\cdot, v)$  exists and  $\tilde{\Delta}(\cdot, v) \in L_2(I)$  for each  $v \in [0, 1]$ . Thus

$$\Delta(u, v) \equiv \frac{1}{\sqrt{\alpha(u)}} \tilde{\Delta}(u, v)$$

exists.

By classical Sturm-Liouville theory, it then follows, e.g., under (35), that

$$(37) \quad A_*(t) = - \int_0^1 \Delta(t, v) \delta(v) dv$$

and

$$(38) \quad a_*(t) = - \int_0^1 \frac{\partial}{\partial t} \Delta(t, v) \delta(v) dv = - \int_0^1 \Delta'(t, v) \delta(v) dv,$$

where  $\Delta(t, v)$  satisfies (36), and hence

$$(39) \quad \frac{\partial}{\partial t} \Delta(t, v) + \int_0^1 (1_{[t \leq s]} - s) \Delta(s, v) \alpha(s) ds, = 1_{[t \leq v]} - v.$$

This formula will have consequences for estimation of  $G$  which we exploit in chapter 6. Now we can calculate the information for  $\theta$  in terms of  $\Delta$  in the case of the model  $\mathbf{P}$ . The efficient score function for  $\theta$  is

$$(40) \quad \begin{aligned} \mathbf{I}_\theta^*(z, u) &= \dot{\mathbf{I}}_\theta(z, u) - \dot{\mathbf{I}}_g a_*(z, u) \\ &= \dot{\mathbf{I}}_\theta(z, u) - a_*(u) - \dot{\mathbf{I}}_u(z, u; \theta) A_*(u). \end{aligned}$$

Therefore, by (21),

$$(41) \quad \begin{aligned} I_\theta^* &\equiv I(P_{(\theta, G)} | \theta, \mathbf{P}) \\ &= E[\mathbf{I}_\theta^* (\mathbf{I}_\theta^*)^T] = E[\mathbf{I}_\theta^* \dot{\mathbf{I}}_\theta^T] \\ &= E[\dot{\mathbf{I}}_\theta \dot{\mathbf{I}}_\theta^T(Z, U)] - E[a_*(U) \dot{\mathbf{I}}_\theta^T(Z, U)] - E[A_*(U) \dot{\mathbf{I}}_\theta^T(Z, U) \dot{\mathbf{I}}_u(Z, U)] \\ &\equiv I_\theta + \int \int_0^1 A_*(u) \frac{\partial}{\partial u} \nabla_\theta p_0(z, u; \theta) du d\mu(z) \\ &\quad - E\{A_*(U) E[\dot{\mathbf{I}}_\theta^T(Z, U) \dot{\mathbf{I}}_u(Z, U) | U]\} \\ &= I_\theta + E_\theta[A_*(U) \delta^T(U)] \\ &= I_\theta - \int_0^1 \int_0^1 \Delta(u, v) \delta(v) \delta^T(u) du dv. \end{aligned}$$

These results are essentially as in Bickel (1986), but here we have relaxed Bickel's hypothesis that  $\alpha$  be bounded to just  $\alpha(u) \leq Mu^{-2}$ .

We now give an alternative view of equation (33) that will be very useful in the construction of efficient estimators of  $\theta$ . Suppose that  $\log p_0(z, u; \theta)$  is twice differentiable with respect to both  $\theta$  and  $u$ . Then, since  $p_0$  has uniform second marginal (at  $\theta = \theta_0$ )

$$(42) \quad \int_0^1 p_0(z, u; \theta) d\mu(z) = 1 \quad \text{for } 0 \leq u \leq 1 \text{ and } \theta = \theta_0.$$

Differentiation of (42) twice with respect to  $u$  yields, under regularity conditions and as in the proof of proposition 2, both

$$(43) \quad E_\theta[\dot{\mathbf{I}}_u(Z, U) | U] = 0 \quad \text{a.s.}$$

and

$$(44) \quad \alpha(u) = E_\theta[\dot{\mathbf{I}}_u^2(Z, U) | U = u] = -E_\theta[\ddot{\mathbf{I}}_{uu}(Z, U) | U = u],$$

where

$$\ddot{\mathbf{I}}_{uu}(z, u) \equiv \frac{\partial}{\partial u} \dot{\mathbf{I}}_u(z, u; \theta) = \frac{\partial^2}{\partial u^2} \log p_0(z, u; \theta).$$

Differentiation of (42) under (19) with respect to  $\theta$  does not yield similar results

since for  $\theta \neq \theta_0$ ,  $T$  need not be uniformly distributed. By using (44) to reexpress  $\alpha$ , (33) can be rewritten as

$$\begin{aligned}
 0 &= \delta(u) - A''_*(u) + \alpha(u)A_*(u) \\
 &= E_\theta[\ddot{I}_{\theta u}(Z,U) | U = u] - A''_*(u) - E_\theta\left[\frac{\partial}{\partial u} \dot{I}_u(Z,U) | U = u\right]A_*(u) \\
 &= E_\theta\left\{\frac{\partial}{\partial u} \dot{I}_\theta(Z,U) - a'_*(U) \right. \\
 &\quad \left. - \left[\frac{\partial}{\partial u} \dot{I}_u(Z,U)\right]A_*(U) - \dot{I}_u(Z,U)a_*(U) | U = u\right\} \text{ by (43)} \\
 (45) \quad &= E_\theta\left\{\frac{\partial}{\partial u} I_\theta^*(Z,U) | U = u\right\}
 \end{aligned}$$

where  $I_\theta^*$  is given by (40).

Finally, here is the key lemma from Sturm-Liouville differential equation theory that shows how to calculate the important Green's function  $\Delta$ .

**Lemma 1.**

- A. Under the conditions of proposition 2,  $\Delta(\cdot, v)$  is continuously differentiable on  $(0, v)$  and  $(v, 1)$ ;  $(\partial/\partial u)\Delta(u, v)$  has a jump discontinuity of  $-1$  at  $u = v$ . Moreover,  $\Delta(\cdot, v)$  is the unique solution of the homogeneous Sturm-Liouville equation

$$(46) \quad y''(u) - \alpha(u)y(u) = 0$$

everywhere except at  $u = v$  which satisfies the boundary conditions  $y(0) = 0 = y(1)$ .

- B. Suppose that  $y_1, y_2$  are fundamental solutions of (46) satisfying  $y_1(0) = 0, y'_1(0) \neq 0, y_2(1) = 0, y'_2(1) \neq 0$ . Then  $y_1, y_2$  are linearly independent and

$$(47) \quad \Delta(u, v) = \frac{1}{D}y_1(u \wedge v)y_2(u \vee v),$$

where the constant  $D \equiv y'_1y_2 - y_1y'_2$  is the Wronskian.

**Proof.** See Tricomi (1957, pages 131–133), or Hille (1969, theorem 3.2.1 and lemma 8.5.1). □

*Tangent Spaces, Efficient Score, and Information Calculation for Example 4*

Now we calculate the tangent spaces for examples 4.A and 4.B. Much of this parallels the calculations for example 3. To do this we will assume throughout that the following logarithmic derivatives of the density  $c_\theta$  of the copula function  $C_\theta$  exist:

$$\begin{aligned}
 \dot{I}_\theta(u, v; \theta) &\equiv \dot{I}_\theta(u, v) \equiv \nabla_\theta \log c_\theta(u, v), \\
 \dot{I}_u(u, v; \theta) &\equiv \dot{I}_u(u, v) \equiv \frac{\partial}{\partial u} \log c_\theta(u, v),
 \end{aligned}$$

(48)

$$\dot{\mathbf{i}}_v(u, v; \theta) \equiv \dot{\mathbf{i}}_v(u, v) \equiv \frac{\partial}{\partial v} \log c_\theta(u, v),$$

and the following mixed derivatives

$$(49) \quad \begin{aligned} \ddot{\mathbf{i}}_{\theta u}(u, v; \theta) &\equiv \ddot{\mathbf{i}}_{\theta u}(u, v) \equiv \frac{\partial}{\partial u} \dot{\mathbf{i}}_\theta(u, v), \\ \ddot{\mathbf{i}}_{\theta v}(u, v; \theta) &\equiv \ddot{\mathbf{i}}_{\theta v}(u, v) \equiv \frac{\partial}{\partial v} \dot{\mathbf{i}}_\theta(u, v). \end{aligned}$$

Again we will need regularity conditions like

$$(50) \quad \begin{aligned} \int \frac{\partial}{\partial u} c_\theta(u, v) dv &= \frac{\partial}{\partial u} \int c_\theta(u, v) dv = 0, \\ \int \frac{\partial}{\partial v} c_\theta(u, v) du &= \frac{\partial}{\partial v} \int c_\theta(u, v) du = 0, \end{aligned}$$

and

$$(51) \quad \begin{aligned} \int \frac{\partial}{\partial u} \nabla_\theta c_\theta(u, v) dv &= \frac{\partial}{\partial u} \nabla_\theta \int c_\theta(u, v) dv = 0, \\ \int \frac{\partial}{\partial v} \nabla_\theta c_\theta(u, v) du &= \frac{\partial}{\partial v} \nabla_\theta \int c_\theta(u, v) du = 0. \end{aligned}$$

We define the functions  $\alpha$ ,  $\beta$ , and  $\gamma$ , when they exist, as follows:

$$(52) \quad \alpha(u) \equiv E_\theta[\dot{\mathbf{i}}_u^2(U, V; \theta) \mid U = u] = \int_0^1 \dot{\mathbf{i}}_u^2(u, v; \theta) c_\theta(u, v) dv,$$

$$(53) \quad \beta(v) \equiv E_\theta[\dot{\mathbf{i}}_v^2(U, V; \theta) \mid V = v] = \int_0^1 \dot{\mathbf{i}}_v^2(u, v; \theta) c_\theta(u, v) du,$$

$$(54) \quad \begin{aligned} \gamma(u) &\equiv E_\theta[\dot{\mathbf{i}}_\theta(U, V; \theta) \dot{\mathbf{i}}_u(U, V; \theta) \mid U = u] \\ &= \int_0^1 \dot{\mathbf{i}}_\theta(u, v; \theta) \dot{\mathbf{i}}_u(u, v; \theta) c_\theta(u, v) dv. \end{aligned}$$

In view of invariance of the models 4.A and 4.B under monotone transformations of the  $S$  axis, and  $S$  and  $T$  axes, respectively, it is clear that the tangent spaces with respect to  $G$  and  $H$  at different points  $P_{(\theta, G, H)}$  but with the same  $\theta$  of the model  $\mathbf{P}_B$  of example 4.B are also related by simple monotone transformations of the axes. The following proposition formalizes this and relates the tangent space at a general point  $P_{(\theta, G, H)}$  to the tangent space at  $P_{(\theta, I, I)}$ .

**Proposition 4.** (Tangent spaces) Suppose that the continuous derivatives  $\dot{\mathbf{i}}_\theta$ ,  $\dot{\mathbf{i}}_u$ , and  $\dot{\mathbf{i}}_v$  in (48) exist, that  $\|\dot{\mathbf{i}}_\theta\|$ ,  $\dot{\mathbf{i}}_u$ ,  $\dot{\mathbf{i}}_v \in L_2(C_\theta)$ , that  $E_\theta \dot{\mathbf{i}}_\theta \dot{\mathbf{i}}_\theta^T$  is nonsingular, that (49) and (50) hold, and that  $\alpha(u) \leq Mu^{-2}$  (or  $\alpha(u) \leq M[u(1-u)]^{-2}$ ), for  $0 < u < 1$ , and  $\beta(v) \leq Mv^{-2}$  (or  $\beta(v) \leq M[v(1-v)]^{-2}$ ), for  $0 < v < 1$ , for some  $M < \infty$ . Then the tangent space of the model  $\mathbf{P}_B$ , example 4.B, is

$$(55) \quad \dot{\mathbf{P}}_B(P_{(\theta, G, H)}) = \{h(G(s), H(t)) : h \in \dot{\mathbf{P}}(P_{(\theta, I, I)})\},$$

where

$$(56) \quad \dot{\mathbf{P}}(P_{(\theta, I, I)}) \supset [\dot{\mathbf{P}}_\theta + \dot{\mathbf{P}}_g + \dot{\mathbf{P}}_h],$$

$$(57) \quad \dot{P}_\theta = [\dot{I}_\theta],$$

$$(58) \quad \dot{P}_g \supset [\dot{I}_g a : a \in L_2^0(I)], \quad \dot{P}_h \supset [\dot{I}_h b : b \in L_2^0(I)],$$

$$(59) \quad \dot{I}_g a(u, v) \equiv a(u) + \dot{I}_u(u, v; \theta) \int_0^u a \, dI,$$

$$(60) \quad \dot{I}_h b(u, v) \equiv b(v) + \dot{I}_v(u, v; \theta) \int_0^v b \, dI.$$

**Proof.** We choose  $\chi$  and  $G_\eta$  as in the proof of proposition 2 and we define densities  $c_{\theta, \eta}$  on  $[0, \infty) \times [0, 1]$  by

$$(a) \quad c_{(\theta, \eta)}(u, v) = c_\theta(G_\eta(u), v) g_\eta(u), \quad 0 \leq u, \quad 0 \leq v \leq 1.$$

Using proposition 2.1.1 and part A (and its proof) of proposition 5 below, we can verify that the model is regular in  $\theta$  with  $\dot{I}_\theta \in \dot{P}^k(P_{(\theta, I, I)})$  and that for any  $a \in L_2(I)$ ,  $\dot{I}_g a$  belongs to  $\dot{P}(P_{(\theta, I, I)})$ . The same technique may be applied to  $\dot{I}_h b$ . □

To compute efficient score functions,, efficient influence functions and information for  $\theta$  in the models 4.A and 4.B, in view of theorem 3.4.1, we need to compute the projection of  $\dot{I}_\theta$  onto the orthocomplement of  $[\dot{I}_g a]$  or  $[\dot{I}_g a] + [\dot{I}_h b]$  respectively. Again, note that we have not shown equality in (56). We first examine the model 4.A with just one unknown marginal distribution and show, under a simple boundedness hypothesis, that  $\dot{I}_g$  is bounded and  $R(\dot{I}_g)$  is closed. This implies that the projection on  $[\dot{I}_g a] = R(\dot{I}_g)$  exists and can be computed in terms of the projection operator

$$\Pi_0(\cdot \mid [\dot{I}_g]) = \dot{I}_g(\dot{I}_g^T \dot{I}_g)^{-1} \dot{I}_g^T.$$

This gives a qualitative “solution” to the projection problem in model 4.A. We then show that the sum space  $[\dot{I}_g] + [\dot{I}_h]$  is closed under the same basic hypotheses, and hence we can compute the projection of  $\dot{I}_\theta$  onto  $[\dot{I}_g] + [\dot{I}_h]$  by means of the alternating projection methods of section A.4—with a geometric convergence rate guaranteed by theorem A.4.2. This “solves” the projection problem for model 4.B.

Of course, we also need to actually calculate the projections in models 4.A and 4.B. If this can be done in model 4.A, then, indeed, we can use the iterative methods of section A.4 to find projections in the model 4.B, in principle. As in example 3, to compute projections in the copula model 4.A with one unknown marginal df we will use classical Sturm-Liouville theory.

**Proposition 5.** Suppose that the conditions of proposition 4 regarding  $\dot{I}_\theta$ ,  $\dot{I}_u$ , and  $\alpha$  hold. Then:

- A.  $\dot{I}_g$  is bounded;  $\|\dot{I}_g\| \leq (1 + 4M)^{1/2}$ .
- B.  $R(\dot{I}_g)$  is closed and  $N(\dot{I}_g) = \{0\}$ .
- C.  $(\dot{I}_g^T \dot{I}_g)^{-1}$  exists and is bounded.
- D. The projection operator  $\Pi_0(\cdot \mid [\dot{I}_g]) = \dot{I}_g(\dot{I}_g^T \dot{I}_g)^{-1} \dot{I}_g^T$ .

**Proof.** Exactly the same as the proof of proposition 2. □

Now consider the model 4.B when both marginals are unknown. Under the conditions of proposition 4, just as in proposition 5, the score operator  $\dot{i}_h$  given in (60) is bounded and has closed range; note that  $[\dot{i}_g a]$  and  $[\dot{i}_h b]$  denote the (closed) ranges of  $\dot{i}_g$  and  $\dot{i}_h$  respectively. To find the efficient score for  $\theta$  in the model 4.B, assuming that the inclusions in proposition 4 are in fact equalities, we now need to find the projection of  $\dot{i}_\theta$  onto  $[\dot{P}_g + \dot{P}_h] = [\dot{i}_g a] + [\dot{i}_h b]$ . In fact, the sum space  $[\dot{i}_g] + [\dot{i}_h]$  is closed; we show this by bounding the maximal cosine of angles between  $[\dot{i}_g]$  and  $[\dot{i}_h]$  away from 1, then by theorem A.4.2.B,  $[\dot{i}_g] + [\dot{i}_h]$  is closed.

**Proposition 6.** Suppose that  $\alpha(u) \leq Mu^{-2}$  and  $\beta(v) \leq Mv^{-2}$  for some  $0 < M < \infty$ . Then

$$(61) \quad \rho([\dot{i}_g], [\dot{i}_h]) \equiv \sup\{ \langle \dot{i}_g a, \dot{i}_h b \rangle_0 : \|\dot{i}_g a\|_0 = \|\dot{i}_h b\|_0 = 1 \} \\ \leq \left( \frac{4M}{1+4M} \right)^{1/2} < 1,$$

and hence  $[\dot{i}_g] + [\dot{i}_h]$  is closed.

**Proof.** Note that, by definition of  $\dot{i}_g a$ ,

$$\begin{aligned} \langle \dot{i}_g a, \dot{i}_h b \rangle_0 &= E(\dot{i}_g a \dot{i}_h b) \\ &= E[a(U)\dot{i}_h b(U, V)] + E[\dot{i}_u(U, V) \int_0^U a \, dI \dot{i}_h b(U, V)] \\ (a) \quad &= E[\dot{i}_u(U, V) \int_0^U a \, dI \dot{i}_h b(U, V)], \end{aligned}$$

since, by Fubini for  $b \in L_2^0(I)$ ,

$$\begin{aligned} E[a(U)\dot{i}_h b(U, V)] &= E\{a(U)E[\dot{i}_h b(U, V) | U]\} \\ &= E[a(U)0] = 0. \end{aligned}$$

From (a) it follows that, for  $a, b$  with  $\|\dot{i}_g a\|_0 = 1 = \|\dot{i}_h b\|_0$ ,

$$(b) \quad |\langle \dot{i}_g a, \dot{i}_h b \rangle_0| \leq \{E[\dot{i}_u^2(U, V) (\int_0^U a \, dI)^2]\}^{1/2} \equiv \sqrt{C^2},$$

where, by (b) and (c) of the proof of proposition 2,

$$(c) \quad 1 = \|\dot{i}_g a\|_0^2 = E[(\dot{i}_g a)^2] = \int_0^1 a^2 \, dI + C^2 \\ \leq (1 + 4M) \int_0^1 a^2 \, dI,$$

or

$$(d) \quad \int_0^1 a^2 \, dI \geq \frac{1}{1+4M}.$$

Thus, from (c),

$$(e) \quad C^2 = 1 - \int_0^1 a^2 \, dI \leq 1 - \frac{1}{1+4M} = \frac{4M}{1+4M},$$

and the claim follows from (b) and (e).

Now we return to the model  $P_A$  of example 4.A with just one unknown marginal df, and study the projection operator

$$\Pi_0(\cdot | [\dot{\mathbf{i}}_g]) = \dot{\mathbf{i}}_g (\dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_g)^{-1} \dot{\mathbf{i}}_g^T .$$

First we calculate  $\dot{\mathbf{i}}_g^T$  and  $\dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_g$ .

**Proposition 7.** Suppose that the hypotheses of proposition 4 hold, that  $\dot{\mathbf{i}}_\theta(u, v)$  is continuously differentiable with respect to  $u$  with derivative  $\ddot{\mathbf{i}}_{\theta u} \in L_2^k(P_\theta)$ , and that (51) holds. Then

$$(62) \quad \dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_g a(u) = a(u) + \int_0^1 (1_{[u \leq s]} - s) \left\{ \int_0^s a \, dl \right\} \alpha(s) \, ds$$

and

$$(63) \quad \dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_\theta(u) = \int_0^1 (1_{[u \leq s]} - s) \gamma(s) \, ds,$$

where  $\alpha$  and  $\gamma$  are given by (52) and (54).

**Proof.** By straightforward calculation from the form (59) for  $\dot{\mathbf{i}}_g$  using proposition A.1.5 and remembering that the range of  $\dot{\mathbf{i}}_g^T$  is in  $L_2^0(I)$ , we arrive at (62) and

$$\dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_\theta(u) = - \int_0^1 (1_{[u \leq s]} - s) \tilde{\delta}(s) \, ds,$$

where

$$\begin{aligned} \tilde{\delta}(u) &\equiv \int \ddot{\mathbf{i}}_{\theta u}(u, v; \theta) c_\theta(u, v) \, dv \\ (a) \quad &= \int \frac{\partial}{\partial u} \nabla_\theta c_\theta(u, v) \, dv - \gamma(u) \\ &= -\gamma(u) \end{aligned}$$

by (51). □

As in (36), we introduce a Green's function  $\Delta(u, v)$  solving

$$(64) \quad \Delta(u, v) + \int_0^1 K(u, s) \Delta(s, v) \alpha(s) \, ds = K(u, v);$$

or, in terms of  $\tilde{\Delta}(u, v) \equiv \sqrt{\alpha(u)} \Delta(u, v)$ ,

$$\tilde{\Delta}(u, v) + \int_0^1 \tilde{K}(u, s) \tilde{\Delta}(s, v) \, ds = \sqrt{\alpha(u)} K(u, v).$$

Since  $\sqrt{\alpha(\cdot)} K(\cdot, v) \in L_2(I)$  for each  $v \in [0, 1]$ ,  $\tilde{\Delta}(\cdot, v)$  exists and  $\tilde{\Delta}(\cdot, v) \in L_2(I)$  for each  $v \in [0, 1]$ . Thus

$$\Delta(u, v) \equiv \frac{1}{\sqrt{\alpha(u)}} \tilde{\Delta}(u, v)$$

exists.

By classical Sturm-Liouville theory, it then follows, for example under the condition

$$(65) \quad |\gamma(u)| \leq M[u(1-u)]^{\varepsilon-3/2}, \quad 0 < u < 1,$$

for some  $M < \infty$ ,  $\varepsilon > 0$  that

$$(66) \quad A_*(t) = \int_0^1 \Delta(t, v) \gamma(v) dv$$

and

$$(67) \quad a_*(t) = \int_0^1 \frac{\partial}{\partial t} \Delta(t, v) \gamma(v) dv = \int_0^1 \Delta'(t, v) \gamma(v) dv,$$

where  $\Delta(t, v)$  satisfies (64), and hence

$$(68) \quad \frac{\partial}{\partial t} \Delta(t, v) + \int_0^1 (1_{[t \leq s]} - s) \Delta(s, v) \alpha(s) ds = 1_{[t \leq v]} - v.$$

Equivalently, as in example 3,  $A_*$  satisfies the differential equation (33) with  $\delta$  replaced by  $-\gamma(t)$ . The formula (68) will have consequences for estimation of  $G$  which we exploit in chapter 6. Now we can calculate the information for  $\theta$  in terms of  $\Delta$  in the case of the model  $P_A$ . The efficient score function for  $\theta$  is

$$(69) \quad \begin{aligned} I_\theta^*(u, v) &= \dot{I}_\theta(u, v) - \dot{I}_g a_*(u, v) \\ &= \dot{I}_\theta(u, v) - a_*(u) - \dot{I}_u(u, v; \theta) A_*(u). \end{aligned}$$

Therefore, as for example 3, but using now (a) of the proof of proposition 7, we obtain

$$(70) \quad \begin{aligned} I_\theta^* &\equiv I(P_{(\theta, G)} | \theta, P_A) \\ &= E [I_\theta^*(I_\theta^*)^T] = E [I_\theta^* i_\theta^T] \\ &= I_\theta - \int_0^1 \int_0^1 \Delta(u, v) \gamma(v) \gamma^T(u) du dv. \end{aligned}$$

Just as for example 3, we now give an alternative view of the equation (33) that will be very useful in the construction of efficient estimators of  $\theta$ . Suppose now that  $\log c_\theta(u, v)$  is twice differentiable with respect to both  $\theta$  and  $u$ . Then, since  $c_\theta$  has uniform marginals (for all  $\theta \in \Theta$ ),

$$(71) \quad \int_0^1 c_\theta(u, v) dv = 1 \quad \text{for } 0 \leq u \leq 1 \text{ and } \theta \in \Theta.$$

Differentiation of (71) twice with respect to  $u$  as in example 3 yields

$$(72) \quad 0 = E_\theta \left\{ \frac{\partial}{\partial u} I_\theta^*(U, V) | U = u \right\},$$

where  $I_\theta^*$  is given by (69).

### More Calculations for Special Cases of Example 3

We now calculate more explicitly for three of the special cases of example 3, which are in fact instances of example 1 too.

#### Example 3, continuation 1.

We consider the Clayton-Cuzick special case of example 1, with core model

$$p_0(z, t; \theta) = e^{\theta z} (1 + \eta e^{\theta z} t)^{-1-1/\eta} 1_{(0, \infty)}(t)$$

with  $\eta > 0$ . We regard  $\eta$  as known.



If  $|Z| \leq C$  a.s., then the functions  $\alpha$  and  $\delta$  of propositions 2 and 3 (see (20) and (22)) satisfy

$$\alpha(u) \leq M(1 - u)^{2\eta-2}$$

and

$$|\delta(u)| \leq M(1 - u)^{\eta-1}$$

for some  $M > 0$ . Thus the hypothesis of proposition 1 and (35) hold (even for  $\eta = 0$ , the limiting case of the Cox model). Note that in this sense the Cox model is the "most singular" model in this family since its  $\alpha$  grow most rapidly as  $u \uparrow 1$ . For the proofs of these claims, see Klaassen (1993).  $\square$

**Example 3, continuation 2.**

We study the Cox model with core density

$$(73) \quad p(z, t; \theta) = e^{\theta z} \exp(-e^{\theta z} t) h(z)$$

with respect to  $\mu \times$  Lebesgue measure on  $R \times [0, \infty)$ ,  $\theta \in R$ . Fix  $\theta_0$  and transform to the unit interval by the transformation  $F_{\theta_0}$  yielding

$$(74) \quad p_0(z, u; \theta) = \frac{p(z, F_{\theta_0}^{-1}(u); \theta)}{f_{\theta_0}(F_{\theta_0}^{-1}(u))}$$

Straightforward but tedious computations lead to (33) with (suppressing the subscript 0)

$$\alpha(u) = \frac{\text{Var}(e^{\theta Z} | U = u)}{f_{\theta}^2(F_{\theta}^{-1}(u))},$$

$$\delta(u) = - \frac{E(Ze^{\theta Z} | U = u)}{f_{\theta}(F_{\theta}^{-1}(u))},$$

for which

$$(75) \quad A_*(u) = f_{\theta}(F_{\theta}^{-1}(u)) \int_0^u \frac{E(Z | U = s)}{f_{\theta}(F_{\theta}^{-1}(s))} ds$$

is the solution. In terms of the original core model this leads to the efficient score function

$$(76) \quad I_{\theta}^*(Z, T) = Z - E(Z | T) - e^{\theta Z} \int_0^T [Z - E(Z | T = s)] ds$$

and the information

$$(77) \quad I(P_{\theta} | \theta, P) = E[(I_{\theta}^*)^2] = E \text{Var}(Z | T),$$

as can be seen by partial integration. Note that in the parametric model with the transformation  $\tau = G$  known,

$$I(P_{\theta} | \theta, P_1(G)) = E(Z(1 - e^{\theta Z T}))^2 = E Z^2.$$

Deriving (33) within the original core model of (73) leads to an equivalent first order differential equation; see Klaassen (1989). There it is also shown that

the Cox partial likelihood estimator is efficient in that it is asymptotically linear in the efficient influence function; cf. Greenwood and Wefelmeyer (1991).  $\square$

**Example 3, continuation 3.**

This is the generalized Box-Cox special case of examples 1 (and 3). Thus

$$(78) \quad p_0(z, t; \theta) = \phi(t - \theta z)h(z).$$

It follows that if  $|Z| \leq C$  a.s., then the functions  $\alpha$  and  $\delta$  of propositions 2 and 3 satisfy

$$\alpha(u) \leq \frac{M(u(1-u))^{-2}}{-\log[u(1-u)]}, \quad 0 < u < 1,$$

and

$$|\delta(u)| \leq \frac{M(u(1-u))^{-1}}{\{-\log[u(1-u)]\}^{1/2}}, \quad 0 < u < 1.$$

In the special case of  $h$  normal, explicit computation is possible as suggested to us by Jack Cuzick. Choose  $Z \sim N(0, \sigma^2)$  and note that the marginal distribution of  $T = \theta Z + \varepsilon$  is  $N(0, 1 + \theta^2\sigma^2)$ . Transforming  $T$  to the unit interval by

$$F_\theta(t) = \Phi((1 + \theta^2\sigma^2)^{-1/2}t)$$

leads to the Sturm-Liouville problem

$$(79) \quad \begin{aligned} A''_*(u) - \frac{\theta^2\sigma^2}{\phi^2(\Phi^{-1}(u))}A_*(u) - \frac{\theta\sigma^2}{\phi(\Phi^{-1}(u))}\Phi^{-1}(u) &= 0, \\ A_*(0) = A_*(1) &= 0, \end{aligned}$$

for which

$$(80) \quad A_*(u) = -\frac{\theta\sigma^2}{2 + \theta^2\sigma^2}\Phi^{-1}(u)\phi(\Phi^{-1}(u))$$

is the solution. In terms of the original core model (78) (with  $h(z) = \phi(z/\sigma)/\sigma$ ) this yields the efficient score function

$$(81) \quad I_\theta^*(Z, T) = \frac{1}{2 + \theta^2\sigma^2}\{2(T - \theta Z)Z + \theta\sigma^2[1 - (T - \theta Z)^2]\}$$

and the information bound

$$(82) \quad I^{-1}(P_\theta | \theta, P) = \frac{1}{\sigma^2} + \frac{1}{2}\theta^2.$$

If the transformation  $\tau = G$  is known, the information bound, of course equals

$$I^{-1}(P_\theta | \theta, P_1(G)) = \{E[Z(T - \theta Z)]^2\}^{-1} = \frac{1}{\sigma^2}. \quad \square$$

*More Calculations for Special Cases of Example 4*

We now calculate more explicitly for several of the special cases of example 4.

**Example 4.I.1, continued.**

We call this special case the Clayton-Oakes model. Tedious but straightforward calculation yields

$$(83) \quad \alpha(u) = c(\theta) u^{-2} \quad \text{with} \quad c(\theta) \equiv \frac{\theta^2 (\theta + 1)}{3\theta + 1}.$$

The homogeneous equation (46) can be solved explicitly here and two solutions are

$$\begin{aligned} y_1(u) &= u^{1/2+\delta}, \\ y_2(u) &= u^{1/2-\delta} - u^{1/2+\delta}, \end{aligned}$$

with

$$\delta \equiv \left(\frac{1}{4} + c(\theta)\right)^{1/2} > \frac{1}{2}.$$

Note that these functions  $y_1$  and  $y_2$  do not satisfy the conditions of lemma 1.B, since  $y'_1(0) = 0$ . Nevertheless  $\Delta(u, v)$  defined by (47) with  $D = y'_1 y_2 - y_1 y'_2 = 2\delta$  and hence

$$(84) \quad \Delta(u, v) = \begin{cases} u^{1/2+\delta} [v^{1/2-\delta} - v^{1/2+\delta}] / 2\delta & 0 \leq u \leq v, \\ [u^{1/2-\delta} - u^{1/2+\delta}] v^{1/2+\delta} / 2\delta & v \leq u \leq 1, \end{cases}$$

satisfies the conditions of lemma 1.A. In principle, the solution of (33) with  $\delta = -\gamma$  can be obtained now by using (66) and (84) after computing

$$(85) \quad \gamma(u) = \frac{u^{-1}}{1+3\theta} \left\{ (1+\theta)(1+2\theta) \sum_{k=0}^{\infty} \frac{(1-u^\theta)^k}{1+(2+k)\theta} + \theta(1+\theta) \log u + \frac{2\theta^2}{1+2\theta} \right\}.$$

□

**Example 4.I.3, continued.**

For the Frank family (12), tedious but straightforward calculation yields

$$(86) \quad \alpha(u) = \frac{1}{3}(\log \theta)^2 \equiv d^2(\theta) \equiv d^2, \quad 0 \leq u \leq 1.$$

In this case too the homogeneous Sturm-Liouville equations of lemma 1 can be solved for  $y_1$  and  $y_2$  to obtain the Green's function  $\Delta(u, v)$  explicitly. We find that

$$(87) \quad y_1(u) = \sinh(du), \quad y_2(u) = \frac{\cosh(du)}{\cosh(d)} - \frac{\sinh(du)}{\sinh(d)},$$

are solutions satisfying lemma 1.B, so that the Wronskian is

$$D \equiv y'_1 y_2 - y_1 y'_2 = \frac{d}{\cosh(d)},$$

and hence

$$(88) \quad \Delta(u, v) = \frac{\sinh(d(u \wedge v)) \sinh(d(1 - u \vee v))}{2d \sinh(d)}.$$

However, computation of  $\gamma$ , and hence  $I_\theta^*$ , is awkward. □

**Example 4.II, continued.**

For the Morgenstern distributions (13), it is easy to calculate

$$\begin{aligned} \alpha(u) &= 4\theta^2 \int_0^1 \frac{(1-2v)^2}{1+\theta(1-2u)(1-2v)} dv \\ (89) \quad &= \frac{2}{(1-2u)^2} \left\{ \frac{1}{\theta(1-2u)} \log\left(\frac{1+\theta(1-2u)}{1-\theta(1-2u)}\right) - 2 \right\}, \end{aligned}$$

which is, in contrast to  $\alpha$  for the Clayton-Oakes model given in (83), a bounded function of  $u$ . This can be seen most easily by bounding the integral expression of (89) for  $\alpha$  by

$$\frac{4\theta^2}{1-|\theta|} \int_0^1 (1-2v)^2 dv = \frac{4\theta^2}{2(1-|\theta|)}.$$

Anyway, the hypotheses of proposition 4 hold. However, we do not know an explicit solution of the homogeneous Sturm-Liouville equation (46) for this  $\alpha$ , nor for the inhomogeneous equation (33) with  $\delta = -\gamma$ .

We can calculate

$$\begin{aligned} I_\theta &= E_\theta[I_\theta^2] \\ (90) \quad &= \int_0^1 \int_0^1 \frac{(1-2u)^2(1-2v)^2}{1+\theta(1-2u)(1-2v)} du dv \\ &= \frac{1}{\theta^3} \int_0^\theta \left\{ \frac{1}{2x} \log\left(\frac{1+x}{1-x}\right) - 1 \right\} dx \end{aligned}$$

and

$$\begin{aligned} (91) \quad \gamma(u) &= \frac{2}{\theta(1-2u)} - \frac{1}{\theta^2(1-2u)^2} \log\left(\frac{1+\theta(1-2u)}{1-\theta(1-2u)}\right) \\ &= -\frac{1-2u}{2\theta} \alpha(u). \end{aligned}$$

□

**Example 4.III, continued.**

Let  $(U, V)$  on  $[0, 1]^2$  have the distribution of the bivariate normal copula model with parameter  $\theta$  and define  $(Y, Z)$  by  $Y = \Phi^{-1}(U)$ ,  $Z = \Phi^{-1}(V)$ . Then  $(Y, Z)$  has a bivariate normal distribution with means 0, variances 1, and correlation  $\theta$ , the conditional distribution of  $Z$  given  $Y$  is  $N(\theta Y, 1 - \theta^2)$ , and

$$(92) \quad \alpha(u) = E\left\{ \left[ \frac{\theta(Z - \theta Y)}{(1 - \theta^2)\phi(Y)} \right]^2 \mid Y = \Phi^{-1}(u) \right\} = \frac{\theta^2}{1 - \theta^2} \frac{1}{\phi^2(\Phi^{-1}(u))},$$

$$\begin{aligned} (93) \quad \gamma(u) &= \text{Cov}\left( \frac{\theta Z}{(1 - \theta^2)\phi(Y)}, \frac{(1 + \theta^2)YZ - \theta Z^2}{(1 - \theta^2)^2} \mid Y = \Phi^{-1}(u) \right) \\ &= \frac{\theta}{1 - \theta^2} \frac{\Phi^{-1}(u)}{\phi(\Phi^{-1}(u))}. \end{aligned}$$

As in (79) and (80),  $A_*(u) = c\Phi^{-1}(u)\phi(\Phi^{-1}(u))$  is a solution for (33) with  $\delta = -\gamma$  for  $c = \theta/(2 - \theta^2)$ . In terms of  $(Y, Z)$  this yields the efficient score function in the case of one unknown marginal distribution

$$(94) \quad I_{\theta}^*(Y, Z) = \frac{1}{(1 - \theta^2)^2(2 - \theta^2)} \{-\theta Y^2 + 2YZ - \theta(2 - \theta^2)Z^2 + \theta(1 - \theta^2)\}.$$

Another straightforward but tedious computation yields

$$(95) \quad I_*(\theta) = \frac{2}{(1 - \theta^2)^2(2 - \theta^2)},$$

whereas

$$(96) \quad I(\theta) = \frac{1 + \theta^2}{(1 - \theta^2)^2},$$

so that

$$\frac{I_*(\theta)}{I(\theta)} = \frac{2}{(1 + \theta^2)(2 - \theta^2)} \geq \frac{8}{9}$$

for all  $\theta \in [-1, 1]$ . □

# 5 | Information Bounds for Infinite-Dimensional Parameters

## 5.1 INTRODUCTION

Our goal in this chapter is to extend the theory and methods developed in chapters 2 and 3 to the case of “general” infinite-dimensional parameters in non-parametric and semiparametric models. We now consider a general  $\mathbf{P} \subset \mathbf{M}$ , a “parameter”  $v(P)$  such as a distribution function, or cumulative hazard function, or perhaps  $P$  itself, which is to be estimated, and a fixed  $P_0 \in \mathbf{P}$ . Thus  $v: \mathbf{P} \rightarrow \mathbf{B}$  where  $\mathbf{B}$  is a set of functions which, in our examples, will invariably be a (subset of a) Banach space. The theorems we present here are simplified versions of the more elaborate results of Hájek (1970), (1972), LeCam (1972), (1979), and Millar (1979), (1983), (1985).

We first extend the notions of (local) regularity, and pathwise differentiability of chapters 2 and 3. To avoid measurability hypotheses, we define regularity in terms of the theory of weak convergence developed by Hoffmann-Jørgensen (1984), Dudley (1985), and Van der Vaart and Wellner (1990), and presented here in appendix A.8. We then give two versions of the Hájek convolution theorem for  $\mathbf{B}$ -valued parameters. These results improve somewhat on those of Millar (1985) in that  $\mathbf{B}$  is not required to be separable.

A number of simple examples are given in section 3.

The geometry of scores developed in section 2.4 is further developed in section 5.4. In this section we introduce score operators and develop a view of information and inverse information in terms of efficient score and influence operators. We also present a recent result due to Van der Vaart (1991) giving a necessary and sufficient condition for pathwise differentiability in a differentiable model.

In fact we will not handle completely arbitrary general parameters, but only those which *can* be estimated at a  $\sqrt{n}$ -rate, primarily df's and probability measures indexed by relatively small classes of sets. While we make some effort to distinguish cases in which the  $\sqrt{n}$ -rates do and do not hold, our treatment is not complete. Recall example 3.1.1 which showed  $I = 0$  for estimation of a density at a fixed point. For further results about parameters which typically *cannot* be estimated at a  $\sqrt{n}$ -rate, e.g., density functions and nonparametric regression functions, the reader will have to look elsewhere: for example Ibragimov and

Has'minskii (1981), Devroye and Györfi (1985), Stone (1980), or Donoho and Liu (1987), (1991a,b).

Further, more involved examples, applying this theory to the infinite-dimensional components of the models in chapter 4, are given in chapter 6.

## 5.2 CONVOLUTION THEOREMS FOR REGULAR ESTIMATES OF INFINITE-DIMENSIONAL PARAMETERS

Suppose that  $\mathbf{P}$  is a model (a collection of probability measures  $P$  on some measurable space  $(\mathbf{X}, \mathcal{B})$ ), and that  $v : \mathbf{P} \rightarrow \mathbf{B}$  where  $\mathbf{B}$  is a Banach space with norm  $\|\cdot\|_{\mathbf{B}}$ . We call  $v(P)$  a *general parameter* or  *$\mathbf{B}$ -valued parameter*. For the extensions of the convolution theorem to estimation of  $\mathbf{B}$ -valued parameters which we develop in this section, we first need to give suitable definitions of pathwise derivatives and regular estimators.

The reader should think of this section as an extension of the results for  $\mathbf{B} = R^m$  developed in section 3.3. Indeed, theorem 2 below uses (a slight extension of) theorem 3.3.2 as a basic building block. The approach of section 3.4, using scores and score operators, is developed further in sections 5.4 and 5.5.

### *Differentiable Functions $v$ and Influence Operators*

Let  $\mathbf{B}^*$  denote the dual space of  $\mathbf{B}$  (the space of all real-valued bounded linear functionals defined on  $\mathbf{B}$ ), and let  $\langle b, b^* \rangle_{\mathbf{B}} \equiv b^*(b)$  denote the value of  $b^*$  at  $b \in \mathbf{B}$ ; see appendix A.1 for basic facts and further references. We begin with an extension of definition 3.3.1 from  $\mathbf{B} = R$  to a general Banach space  $\mathbf{B}$ .

**Definition 1.** The parameter  $v$  with values in  $\mathbf{B}$  is *pathwise differentiable* at  $P_0 \in \mathbf{P}$  if there exists a bounded linear function (operator)  $\dot{v}(P_0) \equiv \dot{v} : \dot{\mathbf{P}} \rightarrow \mathbf{B}$  such that for any curve  $\{P_\eta\} \subset \mathbf{P}$  with tangent  $h \in \dot{\mathbf{P}}^0$ , (i.e.,  $s_\eta = s_0 + \eta h s_0/2 + o(\eta)$ ) we have

$$(1) \quad \left\| \frac{v(P_\eta) - v(P_0)}{\eta} - \dot{v}(h) \right\|_{\mathbf{B}} = o(1).$$

It follows immediately that pathwise differentiability of  $v$  with derivative  $\dot{v}$  implies

$$(2) \quad \left\langle \frac{v(P_\eta) - v(P_0)}{\eta} - \dot{v}(h), b^* \right\rangle_{\mathbf{B}} = o(1)$$

for all  $b^* \in \mathbf{B}^*$ . Thus pathwise differentiability implies differentiability in the weak topology of  $\mathbf{B}$ , and we say that  $v$  is *pathwise weak-differentiable* if (2) holds. By linearity of  $b^*$ , (2) is equivalent to

$$(3) \quad \langle v(P_\eta), b^* \rangle_{\mathbf{B}} = \langle v(P_0), b^* \rangle_{\mathbf{B}} + \eta b^* \dot{v}(h) + o(\eta)$$

for all  $b^* \in \mathbf{B}^*$ , where, by definition of the adjoint  $\dot{v}^T : \mathbf{B}^* \rightarrow \dot{\mathbf{P}}$ ,

$$(4) \quad b^* \dot{v}(h) = \langle \dot{v}(h), b^* \rangle_{\mathbf{B}} = \langle h, \dot{v}^T b^* \rangle_0.$$

Note that for a fixed  $b^* \in \mathbf{B}^*$ ,  $\langle v(P), b^* \rangle_{\mathbf{B}} = b^* v(P)$  is a real-valued parameter defined on  $\mathbf{P}$ . In view of (3), pathwise weak-differentiability of  $v$  implies that  $b^* v(P)$  is pathwise differentiable at  $P_0$  in the sense of definition 3.3.1 with derivative  $b^* \dot{v} : \dot{\mathbf{P}} \rightarrow R$ . Since  $b^* \dot{v}$  is a bounded linear functional on the Hilbert space  $\dot{\mathbf{P}}$ , it can be represented as an inner product (by the Riesz representation theorem given in example A.1.8)

$$(5) \quad b^* \dot{v}(h) = \langle \dot{v}_{b^*}, h \rangle_0$$

for a unique function  $\dot{v}_{b^*} \in \dot{\mathbf{P}}$ . This is exactly as in (3.3.5). Comparison of (4) and (5) shows that, in fact,

$$(6) \quad \dot{v}_{b^*} = \dot{v}^T b^*.$$

The function  $\dot{v}_{b^*}$  is sometimes called the *canonical gradient* of  $v$  in the direction  $b^*$  at  $P_0$ .

Just as in section 3.3, the derivative map  $\dot{v}$  is uniquely defined on  $\dot{\mathbf{P}}^0$  and hence on  $\dot{\mathbf{P}}$ . And, again as in section 3.3, the parameter  $v$  is often described most naturally as the restriction to  $\mathbf{P}$  of a parameter  $v_e$  defined on a larger model  $\mathbf{M}_0 \supset \mathbf{P}$ . Suppose  $v_e$  is pathwise differentiable on  $\mathbf{M}_0$  with derivative  $\dot{v}_e : \dot{\mathbf{M}}_0 \rightarrow \mathbf{B}$ . Then necessarily

$$\dot{v}_e|_{\dot{\mathbf{P}}} = \dot{v},$$

and hence

$$(7) \quad \dot{v}^T b^* = \Pi_0(\dot{v}_e^T b^* | \dot{\mathbf{P}}).$$

Since  $\dot{v}^T b^*$  plays the role of  $\dot{v}$  on the right side of (3.3.5), in view of definition 3.3.2 we define the efficient influence operator for estimation of  $v$  as follows.

**Definition 2.** If  $v$  is pathwise (weak-) differentiable with derivative  $\dot{v}$ , then the *efficient influence operator*  $\tilde{\mathbf{I}}(P_0 | v, \mathbf{P}) \equiv \tilde{\mathbf{I}}_v : \mathbf{B}^* \rightarrow \dot{\mathbf{P}}$  for estimation of  $v$  in  $\mathbf{P}$  is defined by

$$(8) \quad \tilde{\mathbf{I}}_v(b^*) \equiv \dot{v}^T b^* = \dot{v}_{b^*}.$$

Thus the range of  $\tilde{\mathbf{I}}_v$  is just the collection of influence functions for estimating the real parameters  $b^* v$ . Note that for  $h \in \dot{\mathbf{P}}$  and any extension  $v_e$  of  $v$

$$(9) \quad \langle \dot{v}_e(h), b^* \rangle_{\mathbf{B}} = \langle h, \dot{v}_e^T b^* \rangle_0 = \langle h, \tilde{\mathbf{I}}_v(b^*) \rangle_0$$

since  $\dot{v}_e^T b^* - \Pi_0(\dot{v}_e^T b^* | \dot{\mathbf{P}}) \perp h \in \dot{\mathbf{P}}$ .

In analogy with the definition of information bound in (2.3.1) and of information in definition 3.3.2, we define the *inverse information covariance functional* for  $v$ ,  $I^{-1}(P_0 | v, \mathbf{P}) \equiv I_v^{-1} : \mathbf{B}^* \times \mathbf{B}^* \rightarrow R$ , by

$$(10) \quad I_v^{-1}(b_1^*, b_2^*) \equiv E_0[\tilde{\mathbf{I}}_v(b_1^*) \tilde{\mathbf{I}}_v(b_2^*)].$$



That is,  $I_v^{-1}$  gives the covariance structure of the stochastic process  $b^* \rightarrow \tilde{I}_v(b^*)$ . Further,

$$I_v^{-1}(b_1^*, b_2^*) = \langle \tilde{I}_v(b_1^*), \tilde{I}_v(b_2^*) \rangle_0 = \langle b_1^*, \tilde{I}_v^T \tilde{I}_v(b_2^*) \rangle_{B^*},$$

where  $\tilde{I}_v^T : \dot{P} \rightarrow B^{**}$  so that  $\tilde{I}_v^T \tilde{I}_v : B^* \rightarrow B^{**}$ . Motivated by the Euclidean case in which, according to (2.3.1) and (2.3.2),  $E\tilde{I}\tilde{I}^T$  equals the information bound  $I^{-1}$ , we call  $\tilde{I}_v^T \tilde{I}_v : B^* \rightarrow B^{**}$  the *information bound operator* for  $v$  in the model  $P$ .

If  $B = R^m$  (or more generally if  $B$  is Hilbert),  $B, B^*$ , and  $B^{**}$  can all be identified. Then the information bound operator is just the matrix  $I^{-1}(P_0 | v, P)$  of (3.3.24) viewed as an operator from  $R^m$  to  $R^m$  in the usual way. The inverse information covariance functional is the same matrix viewed as a bilinear functional,  $(x, y) \rightarrow x^T I^{-1} y$  for  $x, y \in R^m$ .

**Definition 3.** If there exists a map  $\tilde{I} : X \rightarrow B$  such that for all  $b^* \in B^*$

$$(11) \quad b^* \tilde{I} = \dot{v}_{b^*} = \tilde{I}_v(b^*) \quad P \text{ - a.s. ,}$$

then, as in definition 3.3.2, we call  $\tilde{I}$  the *efficient influence function*.

If  $v$  is pathwise (weak-) differentiable, such an influence function does not always exist, but if  $B = R^m$  with the Euclidean norm, it clearly does. Furthermore, if  $\dot{P}$  is finite-dimensional with  $\{h_1, \dots, h_k\}$  as an orthonormal basis, then  $\tilde{I} = \sum_{j=1}^k \dot{v}(h_j) h_j$  will do, since, for all  $x \in X, \tilde{I}(x) \in B$  and

$$(12) \quad b^* \tilde{I}(x) = \sum_{j=1}^k b^* \dot{v}(h_j) h_j(x) = \sum_{j=1}^k \langle \dot{v}_{b^*}, h_j \rangle_0 h_j(x) = \dot{v}_{b^*}(x)$$

and (11) is satisfied. Usually, however, we will encounter models with  $\dot{P}$  infinite-dimensional. Then we typically can establish existence of  $\tilde{I}$  directly in connection with regularity and asymptotic linearity, to which we now turn.

### Regular Estimators of a B-Valued Parameter $v$

Now let  $X_1, \dots, X_n, \dots$  be i.i.d. in  $(X, \mathcal{B})$  with distribution  $P \in \mathcal{P}$ . An "estimator" of  $v(P)$  is a map  $T_n = t_n(X_1, \dots, X_n)$  from  $X^n$  to  $B$ ; here we use quotes around "estimator" because no measurability assumptions will be imposed on  $t_n$ . This approach, which is motivated by recent developments in empirical process theory, avoids measurability restrictions and allows "estimators" which may not be measurable in the Borel  $\sigma$ -field of  $B$ . This difficulty already arises in the theory of the empirical distribution function of real-valued random variables treated as an element of  $D[0,1]$  or  $D[-\infty, \infty]$  with the supremum metric (see, e.g., Billingsley (1968, section 18, pages 150-153)), and becomes much more severe when dealing with general empirical processes indexed by classes of sets or functions. Thus our definition of (locally) regular "estimators" will be formulated in terms of the Hoffmann-Jørgensen and Dudley theory of weak convergence as described in Hoffmann-Jørgensen (1984), Dudley (1985), and appendix 8. A sketch of another approach to measurability issues

may be found in appendix 8 as well. We will suppress the quotes around “estimator.”

**Definition 4.** The estimator sequence  $\{T_n\}$  of  $v(P)$  is said to be *regular* at  $P_0 \in \mathbf{P}$  if there exists a tight Borel measurable random element  $\mathbb{Z}$  in  $\mathbf{B}$  such that for every curve  $\{P_{\eta}\} \subset \mathbf{P}$  through  $P_0$  and every sequence  $\{\eta_n\}$  with  $\eta_n = O(n^{-1/2})$ ,

$$(13) \quad \sqrt{n}(T_n - v(P_{\eta_n})) \Rightarrow \mathbb{Z} \quad \text{as } n \rightarrow \infty$$

under  $P_n \equiv P_{\eta_n}$ .

Call  $\mathbb{Z}$  *separable* if there exists a separable Borel set  $C \subset \mathbf{B}$  with  $P(\mathbb{Z} \in C) = 1$ . With the possible exception of set-theoretic pathological cases, it is no loss of generality to assume that a random element in the Borel  $\sigma$ -field is separable. For a discussion, see, e.g., Dudley (1985, pages 148–149).

Several definitions of asymptotic linearity are possible.

**Definition 5.** A linear operator  $\Psi_v(\cdot, P_0) \equiv \Psi_v : \mathbf{B}^* \rightarrow L_2^0(P_0)$  is called an *influence operator* at  $P_0$ . The estimator sequence  $\{T_n\}$  is said to be *weakly asymptotically linear* at  $P_0$  with influence operator  $\Psi_v$  if for every  $b^* \in \mathbf{B}^*$

$$(14) \quad \sqrt{n}b^*(T_n - v(P_0)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \Psi_v(b^*)(X_i) \Rightarrow 0$$

as  $n \rightarrow \infty$ , under  $P_0$ . Similarly, the estimator sequence  $\{T_n\}$  is called *asymptotically linear* at  $P_0$  with influence function  $\psi(\cdot, P_0) \equiv \psi : \mathbf{X} \rightarrow \mathbf{B}$  if

$$(15) \quad \sqrt{n}(T_n - v(P_0)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi(X_i) \Rightarrow 0$$

as  $n \rightarrow \infty$ , under  $P_0$ , and it is called *weakly asymptotically linear* at  $P_0$  with influence function  $\psi$  if for all  $b^* \in \mathbf{B}^*$

$$(16) \quad \sqrt{n}b^*(T_n - v(P_0)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n b^*\psi(X_i) \Rightarrow 0$$

as  $n \rightarrow \infty$ , under  $P_0$ . We define the uniform and local versions of these linearity concepts as in definitions 2.2.6 and 2.2.7. Weaker versions of asymptotic linearity, involving a proper subset  $\mathbf{B}_0^*$  of  $\mathbf{B}^*$ , are given in definition 7.6.1.

Now we can formulate a convolution theorem for estimators with values in  $\mathbf{B}$ . The following theorem is adapted from Van der Vaart and Wellner (1990). Their proof proceeds by representing the entire limit process  $\mathbb{Z}$  as a sum of two independent processes for finite-dimensional submodels of the full model, and then finding the “optimal representation” as a limit within this collection of representations (or decompositions). Our present proof differs from theirs, however, building up from the finite-dimensional theorem 2 below.

**Theorem 1.** (Convolution theorem for regular estimators) Suppose that:

- (i)  $v$  is pathwise weak-differentiable at  $P_0 \in \mathbf{P}$  with derivative  $\dot{v}$ .

(ii)  $\{T_n\}$  is regular with limit  $\mathbb{Z}$ .

(iii)  $\dot{P}^0$  is linear; thus  $\dot{P} = \overline{P}^0$ .

Then there exist tight Borel measurable random elements  $\mathbb{Z}_0$  and  $\Delta_0$  in  $\mathbf{B}$  such that:

A.  $L(\mathbb{Z}) = L(\mathbb{Z}_0 + \Delta_0)$ .

B.  $\mathbb{Z}_0$  and  $\Delta_0$  are independent.

C.  $\mathbb{Z}_0$  is mean 0 Gaussian with

$$(17) \quad \text{Cov}(b_1^* \mathbb{Z}_0, b_2^* \mathbb{Z}_0) = I_v^{-1}(b_1^*, b_2^*),$$

where  $I_v^{-1}$  is the inverse information covariance functional for  $v$ .

D. Suppose that the efficient influence function  $\tilde{I} : \mathbf{X} \rightarrow \mathbf{B}$  exists and satisfies

$$(18) \quad \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i) \Rightarrow \mathbb{Z}_0.$$

Then

$$(19) \quad \begin{pmatrix} \sqrt{n}(T_n - v(P_0)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i) \\ \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}(X_i) \end{pmatrix} \Rightarrow \begin{pmatrix} \Delta_0 \\ \mathbb{Z}_0 \end{pmatrix}$$

with  $\Delta_0$  and  $\mathbb{Z}_0$  independent.

E. If (18) holds, then  $\Delta_0$  is degenerate at  $0 \in \mathbf{B}$  if and only if  $\{T_n\}$  is asymptotically linear at  $P_0$  with influence function  $\tilde{I}$ .

Note that (19) is of the same type as the convolution results (3.3.26) and (2.3.3). It implies conclusions A through C, but only under the assumption of existence of  $\tilde{I}$  and the asymptotic normality (18). The proof of theorem 1, as well as the proofs of the other theorems in this section, is given at the end of the section. Our approach here is to first find optimal representations of all the finite-dimensional distributions of the process  $\mathbb{Z}$ ; this is carried out in theorem 2 below. We then combine these (via tightness) to obtain a representation (or decomposition) of the entire process. Theorem 2 itself involves a weaker notion of regularity as follows:

**Definition 6.** The estimator sequence  $\{T_n\}$  of  $v(P)$  is said to be *weakly regular* at  $P_0 \in \mathbf{P}$  if there exists a process  $\{b^* \mathbb{Z} : b^* \in \mathbf{B}^*\}$  on  $(R^{\mathbf{B}}, \mathcal{B}^{\mathbf{B}})$  such that for every curve  $\{P_\eta\} \subset \mathbf{P}$  through  $P_0$ , every sequence  $\{\eta_n\}$  with  $\eta_n = O(n^{-1/2})$ , and every  $b^* \in \mathbf{B}^*$

$$(20) \quad \sqrt{n}(b^* T_n - b^* v(P_n)) \Rightarrow b^* \mathbb{Z} \quad \text{as } n \rightarrow \infty$$

under  $P_n \equiv P_{\eta_n}$ .

With this definition we can formulate a second convolution theorem for estimates with values in  $\mathbf{B}$ . The following theorem separates the issues of "finite-dimensional" convergence and tightness. The proof, which is adapted from Van

der Vaart (1988a) and is similar in spirit to those of Millar (1985), proceeds by decomposing the finite dimensional laws first.

**Theorem 2.** (Convolution theorem for weakly-regular estimators)

Suppose that:

- (i)  $v$  is pathwise weak-differentiable at  $P_0 \in \mathbf{P}$ .
- (ii)  $\{T_n\}$  is weakly regular with limit  $\mathbf{Z}$ .
- (iii)  $\dot{\mathbf{P}}^0$  is linear; thus  $\dot{\mathbf{P}} = \overline{\dot{\mathbf{P}}^0}$ .

Then there exist processes  $\{b^* \mathbf{Z}_0 : b^* \in \mathbf{B}^*\}$  and  $\{b^* \Delta_0 : b^* \in \mathbf{B}^*\}$  on  $(R^{\mathbf{B}^*}, \mathcal{B}^{\mathbf{B}^*})$  such that:

- A.  $L(b^* \mathbf{Z}) = L(b^* \mathbf{Z}_0 + b^* \Delta_0)$  for all  $b^* \in \mathbf{B}^*$ .
- B.  $\mathbf{Z}_0$  and  $\Delta_0$  are independent.
- C.  $\mathbf{Z}_0$  is mean 0 Gaussian satisfying (17).
- D. For every  $h \in \dot{\mathbf{P}}$  and  $b^* \in \mathbf{B}^*$

$$(21) \quad \left( \begin{array}{c} \sqrt{n} b^* (T_n - v(P_0)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{l}_v(b^*)(X_i) \\ \frac{1}{\sqrt{n}} \sum_{i=1}^n h(X_i) \end{array} \right) \Rightarrow \begin{pmatrix} b^* \Delta_0 \\ S_h \end{pmatrix}$$

as  $n \rightarrow \infty$  with  $\Delta_0$  and  $S_h$  independent.

- E.  $\Delta_0$  is degenerate at  $0 \in R^{\mathbf{B}^*}$  if and only if for every  $b^* \in \mathbf{B}^*$ ,  $\{b^* T_n\}$  is asymptotically linear with influence function  $\tilde{l}_v(b^*)$  in the sense of weak convergence ( $\Rightarrow$ ).

The convolution theorems 1 and 2 for (weakly) regular estimators make it possible to define (weakly) efficient estimators. Just as in the case of  $\mathbf{B} = R^k$  in chapters 2 and 3, (weakly) regular estimators  $\{T_n\}$  which have limit  $\mathbf{Z} = \mathbf{Z}_0$  in (13) (or (20)) so that  $\Delta_0 = 0$  in  $\mathbf{B}$ , have an optimality property.

**Definition 7.** If  $\{T_n\}$  is a regular estimator of  $v(P)$  satisfying (13) with  $\mathbf{Z} = \mathbf{Z}_0$ , we say that  $\{T_n\}$  is *efficient at  $P_0$* . If  $\{T_n\}$  is an estimator sequence of  $v(P)$  satisfying (20) with  $\mathbf{Z} = \mathbf{Z}_0$ , we say that  $\{T_n\}$  is *weakly efficient at  $P_0$* . If  $\{T_n\}$  is (weakly) efficient at all regular  $P_0$ , then we call  $\{T_n\}$  (weakly) *efficient*.

We define a class of loss functions  $l : \mathbf{B} \rightarrow R^+$ ; call  $l$  *subconvex* if

$$l(0) = 0 \leq l(b) \text{ for every } b \in \mathbf{B},$$

$$l(b) = l(-b),$$

$$\{b : l(b) \leq c\} \text{ is convex and closed for every } c \in R.$$

If  $w : [0, \infty) \rightarrow [0, \infty)$  is continuous (or even just lower semicontinuous) and nondecreasing with  $w(0) = 0$ , then  $l(b) = w(\|b\|)$  is subconvex. Thus both  $l(b) = 1_{\{|b| > t\}}$ ,  $t \in R^+$  and  $l(b) = \|b\|^r$  with  $r > 0$  are subconvex.

**Proposition 1.** (Asymptotic optimality theorem) Suppose that the hypotheses of theorem 1 hold and  $l : \mathbf{B} \rightarrow R^+$  is subconvex. Then

$$\liminf_{n \rightarrow \infty} E_* I(\sqrt{n}(T_n - v(P))) \geq El(\mathbb{Z}) \geq El(\mathbb{Z}_0).$$

For asymptotic minimax theorems—which apply to all estimators, not just regular estimators—see Millar (1983) or Van der Vaart and Wellner (1990). As a consequence of our convolution theorems, we obtain the following generalization of proposition 3.3.1.

**Corollary 1.**

- A. Suppose that the conditions of theorem 1 hold, and that  $\{T_n\}$  is asymptotically linear at  $P_0$  with influence function  $\psi$ . Then  $\{T_n\}$  is efficient at  $P_0$  if and only if  $b^* \psi \in \dot{P}$  for all  $b^* \in \mathbf{B}^*$  and then the efficient influence function  $\tilde{l}$  exists and equals  $\psi$ .
- B. Suppose that the conditions of theorem 2 hold, and that  $\{T_n\}$  is weakly asymptotically linear at  $P_0$  with influence function  $\psi$ . Then  $\{T_n\}$  is weakly efficient at  $P_0$  if and only if  $b^* \psi \in \dot{P}$  for all  $b^* \in \mathbf{B}^*$  and then the efficient influence function  $\tilde{l}$  exists and equals  $\psi$ .
- C. Suppose that the conditions of theorem 2 hold, and that  $\{T_n\}$  is weakly asymptotically linear at  $P_0$  with influence operator  $\psi_v$ . Then  $\{T_n\}$  is weakly efficient at  $P_0$  if and only if  $\psi_v(b^*) \in \dot{P}$  for all  $b^* \in \mathbf{B}^*$  and then  $\psi_v = \tilde{l}_v$ .

Now we give a theorem which shows connections between the hypotheses in the convolution theorem 1. Roughly speaking, in the presence of joint convergence of the estimator and the scores, regularity of the estimator implies pathwise differentiability of the function  $v$ , and vice versa. Again, our treatment owes much to Van der Vaart (1988a) and (1988c).

**Theorem 3.** (Equivalence of regularity and differentiability) Suppose that (iii) of theorem 1 holds (i.e.  $\dot{P}^0$  is linear), and that for every  $h \in \dot{P}^0$

$$(22) \quad \left( \begin{array}{c} \sqrt{n}(T_n - v(P_0)) \\ \frac{1}{\sqrt{n}} \sum_{i=1}^n h(X_i) \end{array} \right) \Rightarrow \begin{pmatrix} \mathbb{Z} \\ S \end{pmatrix} \quad \text{as } n \rightarrow \infty \text{ under } P_0.$$

- A. If  $\{T_n\}$  is regular with limit  $\mathbb{Z}$ , then  $v$  is pathwise differentiable at  $P_0$  with derivative  $\dot{v}$ .
- B. If  $v$  is pathwise differentiable at  $P_0$  with derivative  $\dot{v}$  and if  $(\mathbb{Z}, S) = (\bar{\mathbb{Z}} + \bar{\Delta}, S)$  with  $\bar{\Delta}$  and  $(\bar{\mathbb{Z}}, S)$  independent and  $(\bar{\mathbb{Z}}, S)$  Gaussian with  $E(S\bar{\mathbb{Z}}) = \dot{v}(h)$ , then  $\{T_n\}$  is regular with limit  $\mathbb{Z}$ .

A somewhat simpler version of theorem 3 would replace the hypothesis (22) in theorem 3 by asymptotic linearity of  $\{T_n\}$  as in proposition 3.3.1. Note that we do not have strict equivalence in the formulation of this theorem; in part B we need additional structure on  $(\mathbb{Z}, S)$ . However, in view of the convolution theorem, it is a reasonable assumption.

*Some Special Cases*

In most of the examples in the sequel,  $\mathbf{B}$  is a Banach space of real-valued functions defined on a separable metric space  $(\mathbf{T}, d)$ . Usually

$$\mathbf{B} = l^\infty(\mathbf{T}) \equiv \{b : \mathbf{T} \rightarrow R : \sup_{t \in \mathbf{T}} |b(t)| < \infty\}$$

with the supremum norm  $\|\cdot\|_{\mathbf{B}} = \|\cdot\|_\infty$ .

If the sample space is a Euclidean space  $R^k$ , a frequent choice of  $(\mathbf{T}, d)$  is  $(R^k, |\cdot|)$  where  $|\cdot|$  is the usual Euclidean metric. Another choice for  $(\mathbf{T}, d)$  is  $(R^k, \tilde{d})$  where  $\tilde{d}^2(s, t) = P(X \in (-\infty, s] \Delta (-\infty, t])$  and  $\Delta$  denotes symmetric difference; this metric  $\tilde{d}$  is exactly the  $L_2(P)$  distance between functions in the collection  $\mathbf{F} = \{1_{(-\infty, t]} : t \in R^k\}$ . Still more generally, a frequent choice is  $(\mathbf{T}, d) = (\mathbf{F}, \tilde{d})$ , where  $\mathbf{F} \subset L_2(P)$  and  $\tilde{d}$  is the usual  $L_2(P)$ -metric  $\tilde{d}(f, g) \equiv \|f - g\| \equiv e_P(f, g)$ . For more on this case, see example 5.3.8.

If  $\mathbf{T}$  is uncountable,  $\mathbf{B} = l^\infty(\mathbf{T})$  is an inseparable Banach space, but the limit processes  $\mathbf{Z}$  of regular estimators are (without loss) concentrated on a complete separable subspace  $\mathbf{B}_0$ , such as the space  $UC(\mathbf{T})$  of uniformly continuous (with respect to  $d$ ) real-valued functions on  $\mathbf{T}$ . This is just analogous to (but more general than) the situation in Billingsley (1968):  $D[0, 1]$  viewed as a subset of  $l^\infty([0, 1])$  with the supremum metric is inseparable, but the limit processes concentrate on complete separable subspaces such as  $C[0, 1] = UC[0, 1]$ .

For  $t \in \mathbf{T}$ , let  $\pi_t$  denote the projection or evaluation map from  $\mathbf{B} = l^\infty(\mathbf{T})$  to  $R$ ,  $\pi_t(b) = b(t)$  for  $b \in \mathbf{B}$ . Since  $\mathbf{B} = l^\infty(\mathbf{T})$  is endowed with the supremum norm,  $\|\cdot\|_\infty \equiv \|\cdot\|_{\mathbf{T}}$ ,  $\pi_t \in (l^\infty(\mathbf{T}))^* = \mathbf{B}^*$  for each  $t \in \mathbf{T}$ . Moreover, the tight limit process  $\mathbf{Z}_0$  is (without loss) concentrated on a separable subspace  $\mathbf{B}_0$  of  $\mathbf{B}$ , and hence its law is determined by its finite-dimensional distributions. Thus the Gaussian distribution  $L(\mathbf{Z}_0)$  on  $(\mathbf{B}_0, \mathcal{B}_0)$  is determined by its inverse information covariance function  $I_V^{-1}(P_0 | \nu, \mathbf{P}) \equiv I_V^{-1} : \mathbf{T} \times \mathbf{T} \rightarrow R$  defined by

$$(23) \quad I_V^{-1}(s, t) \equiv E_0(\tilde{I}_s \tilde{I}_t) \equiv E_0(\tilde{I}_V(\pi_s) \tilde{I}_V(\pi_t)) = I_V^{-1}(\pi_s, \pi_t),$$

where  $I_V^{-1}$  on the right side is the inverse information covariance functional defined in (10).

Also note that for  $b^* = \pi_t$ , (5) becomes

$$(24) \quad \dot{v}(h)(t) = \pi_t \dot{v}(h) = \langle \dot{v}_{\pi_t}, h \rangle_0$$

for some function  $\dot{v}_{\pi_t} \in \dot{\mathbf{P}} \subset L_2^0(P_0)$ . By (6) we have

$$(25) \quad \dot{v}^T \pi_t = \dot{v}_{\pi_t} \equiv \dot{v}_t \quad \text{for } t \in \mathbf{T},$$

and hence we can rewrite (24) as

$$(26) \quad \dot{v}(h)(t) = \int \dot{v}_t h dP_0$$

for  $h \in \dot{\mathbf{P}}$  and  $t \in \mathbf{T}$ . Of course,  $\dot{v}_t$  is just the pathwise derivative of the real-valued parameter  $\nu_t \equiv \nu(P)(t) \equiv \pi_t \nu(P)$  in the sense of definition 3.3.1.

As we will see, the functions  $\{\dot{v}_t : t \in \mathbf{T}\}$  appear naturally in all our examples, and together they form the efficient influence function.

At this point the reader may wish to proceed to the examples in section 5.3 before returning to the proofs.

*Proofs*

We first prove theorem 2; our proof of theorem 1 will build thereon.

**Proof of theorem 2.** Let  $b^* \in \mathbf{B}^*$ . Then (ii) implies that

$$(a) \quad b^* \mathbf{Z}_n \equiv b^* \sqrt{n}(T_n - v(P_n)) \Rightarrow b^* \mathbf{Z}.$$

Consequently, by corollary A.8.1, for  $\pi \equiv (b_1^*, \dots, b_q^*) \in (\mathbf{B}^*)^q$

$$(b) \quad (b_1^* \mathbf{Z}_n, \dots, b_q^* \mathbf{Z}_n) \Rightarrow (b_1^* \mathbf{Z}, \dots, b_q^* \mathbf{Z}) \equiv \mathbf{Z}_\pi.$$

By the same argument as in the proof of theorem 3.3.2, which is based on theorem 2.3.1, but with lemmas A.8.4 and A.8.5 and corollary A.8.1 replacing the corresponding standard weak convergence results from A.7 used in the proofs of these theorems,

$$(c) \quad \mathbf{L}(\mathbf{Z}_\pi) = \mathbf{L}(Z_{0\pi} + \Delta_\pi) = Q_\pi^\# * R_\pi^\#,$$

where

$$Q_\pi^\# \equiv \mathbf{L}(Z_{0\pi}) = N_q(0, (I_v^{-1}(b_i^*, b_j^*))),$$

$$(d) \quad R_\pi^\# \equiv \mathbf{L}(\Delta_\pi).$$

In fact, this is just a generalization of theorem 3.3.2 to allow for the possibility that the  $b_i^* T_n$  may not be measurable. Furthermore, (c) can be generalized (for  $q = 1$ ) to

$$\left( \begin{array}{c} \sqrt{n} b^*(T_n - v(P_0)) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{I}_v(b^*)(X_i) \\ \frac{1}{\sqrt{n}} \sum_{i=1}^n h(X_i) \end{array} \right) \Rightarrow \left( \begin{array}{c} \Delta_\pi \\ \mathbf{W} \end{array} \right)$$

as  $n \rightarrow \infty$ , with  $\Delta_\pi$  and  $\mathbf{W}$  independent, and  $h \in \dot{\mathbf{P}}$ .

The system of finite-dimensional distributions  $\{Q_\pi^\# : \pi \in (\mathbf{B}^*)^q, q = 1, 2, \dots\}$  is consistent, and hence yields a cylinder probability measure  $Q^\#$  on the field

$$\mathcal{B}_*^0 \equiv \bigcup_{\pi} \mathcal{B}_*(\pi)$$

where  $\mathcal{B}_*(\pi)$  is the smallest  $\sigma$ -field making  $\pi \equiv (b_1^*, \dots, b_q^*)$  measurable. From this and consistency of

$$\{\mathbf{L}(b_1^* \mathbf{Z}, \dots, b_q^* \mathbf{Z}) : b_i^* \in \mathbf{B}^*, q = 1, 2, \dots\},$$

consistency of

$$\{L(\Delta_\pi) : \pi \in (\mathbf{B}^*)^q, q = 1, 2, \dots\}$$

follows. Thus, by the Kolmogorov existence theorem,  $L(\mathbf{Z}) = L(\mathbf{Z}_0 + \Delta_0)$  where  $\mathbf{Z}_0$  and  $\Delta_0$  are independent in the sense of cylinder probability measures on the field  $\mathcal{B}_*^0$ , and A through E hold.  $\square$

**Proof of theorem 1.** We will use the notation developed in the proof of theorem 2. Denote the cylinder measures corresponding to  $\mathbf{Z}_0$  and  $\Delta_0$  by  $Q^\#$  and  $R^\#$  respectively. Since  $\{T_n\}$  is regular,  $\mathbf{Z}$  is tight and Borel measurable. Hence, for every  $\varepsilon > 0$  there is a compact set  $K \in \mathcal{B}_{Borel}$  such that  $P(\mathbf{Z} \in K) > 1 - \varepsilon$ . Then, for any  $\pi = (b_1^*, \dots, b_q^*) \in (\mathbf{B}^*)^q$ ,

$$(a) \quad P(\pi\mathbf{Z} \in \pi K) \geq P(\mathbf{Z} \in K) > 1 - \varepsilon.$$

But then by (c) of the proof of theorem 2,

$$1 - \varepsilon < P(\pi\mathbf{Z} \in \pi K) = \int Q_\pi^\#(\pi K - w) dR_\pi^\#(w).$$

Hence there is a  $w_0 \in R^q$  such that

$$Q_\pi^\#(\pi K - w_0) > 1 - \varepsilon.$$

By symmetry of  $Q_\pi^\#$

$$Q_\pi^\#(w_0 - \pi K) > 1 - \varepsilon,$$

and hence

$$Q_\pi^\#(\frac{1}{2}(\pi K - \pi K)) > 1 - 2\varepsilon$$

by noting that  $A \cap B \subset \frac{1}{2}(A + B)$  and taking  $A = \pi K - w_0, B = -A$ . Note that  $\tilde{K} \equiv (K - K)/2$  is compact since it is the image of  $K \times K$  under the continuous map  $(x, y) \rightarrow (x - y)/2$ . It therefore follows from theorem 1.4.23 of Araujo and Giné (1980) that the cylinder measure  $Q^\#$  can be extended to a tight Borel measure  $Q$  on  $(\mathbf{B}, \mathcal{B}_{Borel})$ ; denote the corresponding random element by  $\mathbf{Z}_0$ .

Now for  $R^\#$  we have

$$\begin{aligned} 1 - \varepsilon < P(\pi\mathbf{Z} \in \pi K) &= \int R_\pi^\#(\pi K - z) dQ_\pi^\#(z) \\ &< \int_{\tilde{\pi K}} R_\pi^\#(\pi K - z) dQ_\pi^\#(z) + 2\varepsilon. \end{aligned}$$

Hence there is a  $z_0 \in \tilde{\pi K}$  such that

$$R_\pi^\#(\pi K - z_0) > 1 - 3\varepsilon.$$

Then, again by theorem 1.4.23 of Araujo and Giné (1980),  $R^\#$  extends to a tight Borel measure  $R$  on  $(\mathbf{B}, \mathcal{B}_{Borel})$ , and A-C hold.

It remains only to prove D and E. Since  $\sqrt{n}(T_n - v(P_0)) \Rightarrow \mathbf{Z}$  and  $n^{-1/2} \sum_{i=1}^n \tilde{I}(X_i) \Rightarrow \mathbf{Z}_0$  by hypothesis, it follows from corollary A.8.1 that for every subsequence  $\{n'\}$  there is a further subsequence  $\{n''\}$  such that

$$(b) \quad (\sqrt{n''}(T_{n''} - v(P_0)), \frac{1}{\sqrt{n''}} \sum_{i=1}^{n''} \tilde{I}(X_i)) \Rightarrow (\mathbf{Z}, \mathbf{Z}_0)$$



for some joint law of  $\mathbf{Z}$  and  $\mathbf{Z}_0$ . Hence by the continuous mapping theorem (proposition A.8.1)

$$(c) \quad \left( \begin{array}{c} \sqrt{n''}(T_{n''} - v(P_0)) - \frac{1}{\sqrt{n''}} \sum_{i=1}^{n''} \tilde{I}(X_i) \\ \frac{1}{\sqrt{n''}} \sum_{i=1}^{n''} \tilde{I}(X_i) \end{array} \right) \Rightarrow \left( \begin{array}{c} \mathbf{Z} - \mathbf{Z}_0 \\ \mathbf{Z}_0 \end{array} \right).$$

By D of theorem 2, it follows from the convergence in (c) that the finite-dimensional distributions of  $\mathbf{Z} - \mathbf{Z}_0$  and  $\mathbf{Z}_0$  are independent. Hence  $\mathbf{Z} - \mathbf{Z}_0$  and  $\mathbf{Z}_0$  are independent. Again by D of theorem 2, the limit is the same for every subsequence, and hence convergence holds for the full sequence. Thus D holds. E is an immediate consequence of D.  $\square$

**Proof of proposition 1.** Since  $l$  is lower semicontinuous, the first inequality follows from regularity of  $\{T_n\}$  with limit  $\mathbf{Z}$  and the portmanteau theorem lemma A.8.3(vi). To prove the second inequality, first note that for any subconvex function  $l$  there exist functions  $l_k$  of the form

$$l_k(b) = \underline{l}_k(b_1^*(b), \dots, b_k^*(b)) \leq l(b)$$

where  $b_i^* \in \mathbf{B}^*$ ,  $l_k \uparrow l$  a.s.  $\mathbf{L}(\mathbf{Z}_0)$ , and the  $\underline{l}_k$ 's are bowl-shaped (i.e., subconvex functions on  $R^k$ ); see, e.g., Van der Vaart (1988a, lemma 3.10, page 61). This step uses separability of  $\mathbf{Z}_0$ . Thus

$$\begin{aligned} (a) \quad E l(\mathbf{Z}) &\geq E l_k(\mathbf{Z}) \\ &= E \underline{l}_k(b_1^* \mathbf{Z}, \dots, b_k^* \mathbf{Z}) \\ &\geq E \underline{l}_k(b_1^* \mathbf{Z}_0, \dots, b_k^* \mathbf{Z}_0) \end{aligned}$$

by corollary A.6.2 of Anderson's inequality

$$= E l_k(\mathbf{Z}_0).$$

Letting  $k \rightarrow \infty$  in (a) yields the second inequality by monotone convergence.  $\square$

**Proof of corollary 1.** First we note that proposition 3.3.1 is valid when the weak convergence ( $\rightarrow$ ) is replaced by  $\Rightarrow$ -convergence, since in its proof lemma A.9.3 and theorem 3.3.2.B can be replaced by lemma A.8.6 and theorem 1.E (with  $\mathbf{B} = R^m$ ) respectively. Now we prove A. If  $T_n$  is efficient at  $P_0$ , then  $\Delta_0$  is degenerate at 0 in theorem 1. Consequently  $b^* T_n$  is an efficient linear estimator of  $b^* v(P)$  with influence function  $b^* \psi$  and proposition 3.3.1 implies  $b^* \psi \in \dot{\mathbf{P}}$  and  $b^* \underline{\psi} = \tilde{I}_v(b^*)$ . Since this holds for all  $b^* \in \mathbf{B}^*$ , the efficient influence function  $\mathbf{I}$  exists and equals  $\psi$ . On the other hand, if  $b^* \psi \in \dot{\mathbf{P}}$ , then proposition 3.3.1 and theorem 1 show that  $b^* \Delta_0$  is degenerate at 0. Consequently  $\Delta_0$  is degenerate at 0 by proposition 1.4.15 of Araujo and Giné (1980), and  $\{T_n\}$  is efficient at  $P_0$ .

B holds by the proof of A with theorem 1 replaced by theorem 2. The proof of C is similar to the one for B.  $\square$

To prove theorem 3 we need the following lemma:

**Lemma 1.** (Van der Vaart) Suppose that  $\mathbf{B} = R$  and that  $\{T_n\}$  is a regular estimator of  $v(P) \in R$  (in the sense of definition 4) with limit  $\mathbb{Z}$ . Let  $\{P_\eta\} \subset \mathbf{P}$  be a curve in  $\mathbf{P}$  with tangent  $h \in \dot{\mathbf{P}}^0$ .

A.

$$(27) \quad \eta^{-1}(v(P_\eta) - v(P)) = O(1) \quad \text{as } \eta \rightarrow 0.$$

B. Let  $S_{h,n} \equiv n^{-1/2} \sum_{j=1}^n h(X_j)$ . Then

$$(28) \quad \lim_{\eta \rightarrow 0} \frac{v(P_\eta) - v(P)}{\eta}$$

exists if and only if

$$(29) \quad \left( \sqrt{n}(T_n - v(P)), S_{h,n} \right) \Rightarrow (\mathbb{Z}, S) \quad \text{under } P = P_0,$$

and then the limit is

$$(30) \quad \lim_{\eta \rightarrow 0} \frac{v(P_\eta) - v(P)}{\eta} = \begin{cases} \frac{ES e^{it\mathbb{Z}}}{it E e^{it\mathbb{Z}}} & \text{for every } t \neq 0 \text{ with } E e^{it\mathbb{Z}} \neq 0, \\ ES\mathbb{Z} & \text{if } E|S\mathbb{Z}| < \infty. \end{cases}$$

**Proof.** We first prove A. Fix  $c > 0$  and a curve  $\{P_\eta\}$  in  $\mathbf{P}$  with tangent  $h$ . We want to show that

$$\eta^{-1}(v(P_{c\eta}) - v(P)) = O(1).$$

Let  $\eta_m \rightarrow 0$  be arbitrary. Let  $n_m$  be the integer part of  $\eta_m^{-2}$  and set  $c_m \equiv c \eta_m \sqrt{n_m}$ ,  $P_{n_m} \equiv P_{c_m/\sqrt{n_m}}$ . Note that  $c_m \rightarrow c$  as  $m \rightarrow \infty$ , and

$$(a) \quad \eta_m^{-1}(v(P_{c\eta_m}) - v(P)) = \sqrt{n_m}(v(P_{c_m/\sqrt{n_m}}) - v(P))(1 + o(1)).$$

By corollary A.8.1 there is a further subsequence  $\{n'\} \subset \{n_m\}$  such that

$$(b) \quad \left( \sqrt{n'}(T_{n'} - v(P)), S_{h,n'} \right) \Rightarrow (\mathbb{Z}, S), \quad \text{under } P,$$

where  $L(S) = N(0, \|h\|_0^2) \equiv N(0, \sigma^2)$ . Thus if  $\Lambda_{n'} \equiv \Lambda_{n'}(P_{n'}, P)$  denotes the log-likelihood ratio of the product of  $n'$  copies of  $P_{n'}$  and  $P$ ,

$$(c) \quad (\sqrt{n'}(T_{n'} - v(P)), \Lambda_{n'}) \Rightarrow (\mathbb{Z}, cS - \frac{1}{2}c^2\sigma^2) \quad \text{under } P.$$

By contiguity, (c) implies that

$$(d) \quad \left\{ \sqrt{n'}(T_{n'} - v(P)) \right\} \text{ is asymptotically tight under } P_{n'}.$$

But by regularity

$$(e) \quad \left\{ \sqrt{n'}(T_{n'} - v(P_{n'})) \right\} \text{ is also asymptotically tight under } P_{n'}.$$

Combining (d) and (e) shows that

$$(f) \quad \left\{ \sqrt{n'} (v(P_{n'}) - v(P)) \right\} \text{ is tight.}$$

Let  $\{n''\} \subset \{n'\}$  be a subsequence such that  $\sqrt{n''} (v(P_{n''}) - v(P))$  converges to a limit  $d(c)$ . Since this holds for any initial sequence  $\{\eta_m\}$ , (27) follows.

Note that (c), (f), and Le Cam's third lemma A.8.6 imply that under  $P_{n''}$

$$\sqrt{n''} (T_{n''} - v(P_{n''})) \Rightarrow \mathbb{Z}_c,$$

where

$$(g) \quad P(\mathbb{Z}_c \in A) = E 1_A(\mathbb{Z} - d(c)) \exp(cS - \frac{1}{2}c^2\sigma^2).$$

By regularity of  $\{T_n\}$ ,  $L(\mathbb{Z}_c) = L(\mathbb{Z})$  for all  $c$ , and hence (g) yields

$$(h) \quad e^{itd(c)} E \exp\{it\mathbb{Z}\} = E \exp\{it\mathbb{Z} + cS - \frac{1}{2}c^2\sigma^2\}.$$

Together with (29) this shows that  $d(c)$  does not depend on the particular sequence  $\{n''\}$ , and consequently that (28) holds.

By (21) of theorem 2.D with  $\mathbf{P} = \{P_\eta\}$ ,  $b^* = 1$ ,  $\tilde{l}_v(1) = d(1)\sigma^{-2}h$ , we see that (28) implies (29). Finally, note that

$$(i) \quad \begin{aligned} d(c) &\equiv \lim_{\eta \rightarrow 0} \eta^{-1} (v(P_{c\eta}) - v(P)) = c \lim_{\eta \rightarrow 0} (c\eta)^{-1} (v(P_{c\eta}) - v(P)) \\ &= c d(1). \end{aligned}$$

Substitution of (i) into (h) and differentiation with respect to  $c$  at  $c = 0$  yield the first equality of (30). The second one is immediate by subtracting  $ES$  from the numerator of the first expression in the right-hand side of (30), taking the limit as  $t \rightarrow 0$ , and applying the dominated convergence theorem.  $\square$

**Proof of theorem 3.** First consider the case  $\mathbf{B} = R$ . Given  $h \in \dot{\mathbf{P}}^0$ , set

$$(a) \quad \dot{v}(h) \equiv \lim_{\eta \downarrow 0} \frac{v(P_\eta) - v(P_0)}{\eta},$$

where  $\{P_\eta\}$  has tangent  $h$ . This is well defined, by lemma 1, and if  $L(\mathbb{Z}, S_h)$  is the limiting law of  $(\sqrt{n}(T_n - v(P)), S_{h,n})$ , then

$$(b) \quad \dot{v}(h) = \frac{E S_h e^{itZ}}{it E e^{itZ}} \quad \text{for all } t \neq 0 \text{ with } E e^{itZ} \neq 0.$$

For the proof of A with  $\mathbf{B} = R$ , it remains only to show that  $\dot{v}$  is linear and continuous. Let  $\{h_0, h_1, \dots\}$  be a sequence in  $\dot{\mathbf{P}}^0$ , and let  $L(\mathbb{Z}, S_0, S_1, \dots)$  be a limit distribution in  $R \times R^\infty$  of

$$(c) \quad (\sqrt{n}(T_n - v(P_0)), S_{h_0,n}, S_{h_1,n}, \dots).$$

Then, for  $a_1, a_2 \in R$ , by applying the representation in (b) to both  $h \equiv a_1 h_1 + a_2 h_2$  and  $h_i$ ,  $i = 1, 2$ , it follows that  $\dot{v}$  is linear:

$$(d) \quad \dot{v}(a_1 h_1 + a_2 h_2) = a_1 \dot{v}(h_1) + a_2 \dot{v}(h_2).$$

If  $h_i \rightarrow h_0$  in  $L_2(P_0)$ , then since this implies that  $E(S_i - S_0)^2 = \|h_i - h_0\|_0^2 \rightarrow 0$ , by using (b) again it follows that

$$(e) \quad |\dot{v}(h_i) - \dot{v}(h_0)|^2 \rightarrow 0,$$

so  $\dot{v}$  is continuous.

Now consider the case of a general  $\mathbf{B}$ . By regularity of  $\{T_n\}$  and lemma A.8.4,  $\{\mathbf{Z}_n \equiv \sqrt{n}(T_n - v(P_0))\}$  converges and is asymptotically tight under  $P_0$ . Then by (22) and the generalization of Le Cam's third lemma, lemma A.8.6,  $\{\mathbf{Z}_n\}$  converges and hence, by lemma A.8.4, is asymptotically tight under  $P_n \equiv P_{c_n/\sqrt{n}}$  with  $c_n \rightarrow c$ . Combining this with regularity of  $\{T_n\}$  shows that  $\{\sqrt{n}(v(P_n) - v(P_0))\}$  is relatively compact.

Now let  $b^* \in \mathbf{B}^*$ . Then  $\{b^* T_n\}$  satisfies the conditions of the theorem as an estimator of the real-valued parameter  $b^* v(P)$ . Thus by the first part of the proof,  $b^* v(P)$  is pathwise differentiable at  $P_0$ , and

$$(f) \quad \lim_{\eta \downarrow 0} b^* \frac{v(P_\eta) - v(P_0)}{\eta}$$

exists. But since  $\mathbf{B}^*$  separates points of  $\mathbf{B}$ , (f) identifies the limit points of  $\sqrt{n}(v(P_{c_n/\sqrt{n}}) - v(P_0))$  and we have

$$\sqrt{n}(v(P_{c_n/\sqrt{n}}) - v(P_0)) \rightarrow c \dot{v}(h).$$

For A, it remains only to show that this defines a continuous linear map from  $\dot{\mathbf{P}}^0$  to  $\mathbf{B}$ .

By the first part of the proof,  $b^* \dot{v}$  is a continuous linear map for every  $b^* \in \mathbf{B}^*$  because it is the derivative of  $b^* v$ . But then continuity and linearity of  $\dot{v}$  follow from lemma A.1.1 applied to  $\mathbf{X} = \dot{\mathbf{P}}^0$ ,  $\mathbf{Y} = \mathbf{B}$ , and  $\mathbf{Y}^* \equiv \mathbf{B}^*$ .

For the proof of B, let  $P_n$  be as above. By (22) and corollary A.9.2 of Le Cam's second lemma, it follows that

$$(g) \quad \left( \sqrt{n}(T_n - v(P_0)), \Lambda_n \right) \Rightarrow \left( \mathbf{Z}, cS - \frac{1}{2}c^2 \|h\|_0^2 \right) \equiv (\mathbf{Z}, \Lambda)$$

under  $P_0$ . By Le Cam's third lemma A.8.6,

$$(h) \quad \sqrt{n}(T_n - v(P_0)) \Rightarrow \tilde{\mathbf{Z}} \quad \text{under } P_n$$

where

$$P(\tilde{\mathbf{Z}} \in B) = E1_B(\mathbf{Z})e^\Lambda.$$

Since  $(\tilde{\mathbf{Z}}, S)$  is Gaussian with  $ES\tilde{\mathbf{Z}} = \dot{v}(h)$ , this implies that

$$(i) \quad L(\tilde{\mathbf{Z}}) = L(\mathbf{Z} + c\dot{v}(h)).$$

Here is a calculation:

$$\begin{aligned}
 \text{(j)} \quad E \exp[itb^* \tilde{Z}] &= E \exp[itb^* Z + cS - c^2 \|h\|_0^2 / 2] \\
 &= E \exp[itb^* \bar{\Delta}] E \exp[itb^* \bar{Z} + cS - c^2 \|h\|_0^2 / 2] \\
 &= E \exp[itb^* \bar{\Delta}] E \exp[itb^* [\bar{Z} + cES\bar{Z}]] \\
 &= E \exp[itb^* [Z + c\dot{v}(h)]],
 \end{aligned}$$

where the third equality is an instance of (c) of the proof of lemma A.9.3. On the other hand, pathwise differentiability implies

$$\text{(k)} \quad \sqrt{n}(v(P_n) - v(P_0)) \rightarrow c \dot{v}(h).$$

Hence, by (h), (i), and (k), under  $\{P_n\}$ ,

$$\begin{aligned}
 \text{(l)} \quad \sqrt{n}(T_n - v(P_n)) &= \sqrt{n}(T_n - v(P_0)) - \sqrt{n}(v(P_n) - v(P_0)) \\
 &\Rightarrow Z + c\dot{v}(h) - c\dot{v}(h) \\
 &= Z;
 \end{aligned}$$

that is,  $\{T_n\}$  is regular. □

### 5.3 EXAMPLES

Now we give several simple examples of the theory in section 2; for more examples, see sections 4, 5, and chapter 6.

In examples 1–4 we consider estimation of a df  $F$  on  $R$  in various models  $\mathbf{P}$  ranging from completely nonparametric (example 1) to parametric (example 4).

**Example 1. Estimation of a df  $F$  on  $R$ .**

Suppose that  $(\mathbf{X}, \mathcal{B}) = (R, \mathcal{B})$  and  $\mathbf{P} = \{P \text{ on } R, P \ll \mu \equiv \text{Lebesgue measure}\}$ . For any probability measure  $P$  on  $R$ , let  $v(P)$  be the corresponding distribution function

$$v(P)(t) \equiv P((-\infty, t]) \equiv F(t) \quad \text{for } t \in \mathbf{T} = R.$$

Thus  $v: \mathbf{P} \rightarrow \mathbf{B}$  with  $\mathbf{B} \equiv D[-\infty, \infty]$ , the space of right-continuous functions with left limits with the supremum norm  $\|\cdot\|_{\mathbf{B}} = \|\cdot\|_{\infty}$ . Then  $v$  is pathwise differentiable at every  $P_0 \in \mathbf{P}$  with derivative  $\dot{v}(P_0) \equiv \dot{v}: \dot{\mathbf{P}} \rightarrow \mathbf{B}$  given by

$$\text{(1)} \quad \dot{v}(h)(t) = \int [1_{(-\infty, t]}(x) - F_0(t)] h(x) dP_0(x) \quad \text{for } t \in R.$$

Indeed, if  $\{P_{\eta}\} \subset \mathbf{P}$  and  $s(\eta) = s(0) + \eta hs(0)/2 + o(\eta)$ , then

$$\begin{aligned}
 v(P_{\eta})(t) - v(P_0)(t) - \eta \dot{v}(h)(t) &= \int_{(-\infty, t]} s^2(\eta) d\mu - \int_{(-\infty, t]} s^2(0) d\mu - \eta \int_{(-\infty, t]} h s^2(0) d\mu \\
 &= \int_{(-\infty, t]} \left\{ (s(\eta) - s(0))(s(\eta) + s(0)) - 2\eta \left( \frac{1}{2} h s(0) \right) s(0) \right\} d\mu
 \end{aligned}$$

$$\begin{aligned}
 &= \int_{(-\infty, t]} \left\{ (s(\eta) - s(0)) - \eta \left( \frac{1}{2} h s(0) \right) \right\} 2s(0) d\mu \\
 &\quad + \int_{(-\infty, t]} (s(\eta) - s(0))^2 d\mu \\
 &\equiv A(t) + B(t)
 \end{aligned}$$

where

$$\|A\|_\infty \leq 2 \|s(\eta) - s(0) - \eta \left( \frac{1}{2} h s(0) \right)\| = o(\eta)$$

and

$$\|B\|_\infty = \|s(\eta) - s(0)\|^2 = o(\eta).$$

Thus (5.2.1) holds, and the derivative  $\dot{v}$  is given by (1).

From (1) and (5.2.24) through (5.2.26), it follows immediately that for  $b^* = \pi_t$ , the projection or evaluation maps, gradients  $\dot{v}_t = \dot{v}^T \pi_t$  are given by

$$(2) \quad \dot{v}_t(x) = 1_{(-\infty, t]}(x) - F_0(t), \quad x \in R,$$

for any  $t \in R$ , and the efficient influence function is  $\tilde{l}(x)(\cdot) = 1_{[x, \infty)}(\cdot) - F_0(\cdot)$ , provided  $\dot{v}_t \in \dot{P}$ . Indeed, by example 3.2.1,  $\dot{P} = L_2^0(P_0)$ , and consequently  $\dot{v}_t \in \dot{P}$  and  $\tilde{l}_t = \tilde{l}_v(\pi_t) = \dot{v}_t$ . Hence the inverse information covariance function of (5.2.23) is, for  $s, t \in R$ , just

$$(3) \quad I_v^{-1}(s, t) = F_0(s) \wedge F_0(t) - F_0(s)F_0(t),$$

and the Gaussian process  $Z_0$  of theorems 5.2.1 and 5.2.2 is the Brownian bridge process  $B_0$  on  $[0, 1]$  composed with  $v(P_0) \equiv F_0$ :

$$(4) \quad Z_0(t) = B_0(F_0(t)), \quad t \in R.$$

Hence theorems 5.2.1 and 5.2.2 assert that the estimation error

$$Z_n \equiv \sqrt{n}(T_n - v) = \sqrt{n}(T_n - F)$$

of any regular estimator  $T = \{T_n\}$  of  $v = F$  must have a limit law which is at least as dispersed as  $L(Z_0)$ .

If  $T_n = F_n = n^{-1} \sum_{i=1}^n 1_{[X_i, \cdot]}$  denotes the usual empirical df, so that  $Z_n = \sqrt{n}(F_n - F)$ , it is well known that  $F_n$  is an asymptotically linear (in fact, exactly linear), regular estimator of  $v = F$  and  $Z_n \Rightarrow Z_0$  (see, e.g., Shorack and Wellner (1986, chapter 3) or Pollard (1984, theorem V.2.11, page 97)). Therefore the empirical df  $F_n$  is efficient. The influence function  $\psi$  of  $F_n$  is (trivially) equal to the efficient influence function  $\tilde{l}$ . Theorem 5.2.1 as applied to (a continuous modification of) this example was proved by Beran (1977a). Of course, the original proof of optimality of  $F_n$  in a very strong global asymptotic minimax sense is due to Dvoretzky, Kiefer, and Wolfowitz (1956).  $\square$

**Example 2. Estimation of a df  $F$  with  $\gamma = \gamma(F)$  known.**

Now suppose that  $(\mathbf{X}, \mathcal{B}) = (R, \mathcal{B})$  as in example 1, but

$$\mathbf{P} = \{P \text{ on } R; P \ll \mu, \gamma(P) = 0\}$$

where  $\gamma: \mathbf{P} \rightarrow R^k$  is a pathwise differentiable function with derivative  $\dot{\gamma} \equiv \dot{\gamma}(P_0)$  satisfying  $\langle \dot{\gamma}, \dot{\gamma}^T \rangle_0$  nonsingular and  $P \rightarrow \dot{\gamma}(P)$  continuous at  $P_0$ . Then, as in example 3.2.3,

$$\dot{\mathbf{P}} = \{h \in L_2^0(P_0): \langle h, \dot{\gamma} \rangle_0 = 0\}.$$

Consider estimation of  $v(P) = F$  as in example 1. Let  $\mathbf{M}_0$  be the larger model as in example 1 *without* the constraint  $\gamma(P) = 0$ , and let  $v_e(P) = F$  be the  $v$  of example 1. Then,  $\dot{v}_e$  is as given in example 1, and, by (5.2.25), (5.2.7), and (A.2.16)

$$\begin{aligned} \dot{v}_t &= \Pi_0(\dot{v}_e^T \pi_t | \dot{\mathbf{P}}) \\ &= \dot{v}_e^T \pi_t - \langle \dot{v}_e^T \pi_t, \dot{\gamma}^T \rangle_0 \langle \dot{\gamma}, \dot{\gamma}^T \rangle_0^{-1} \dot{\gamma} \end{aligned}$$

and (cf. (5.2.23), (2) and (3))

$$\begin{aligned} I_v^{-1}(s, t) &= F_0(s) \wedge F_0(t) - F_0(s)F_0(t) \\ &\quad - \langle \dot{v}_e^T \pi_s, \dot{\gamma}^T \rangle_0 \langle \dot{\gamma}, \dot{\gamma}^T \rangle_0^{-1} \langle \dot{\gamma}, \dot{v}_e^T \pi_t \rangle_0. \end{aligned}$$

(Note that  $T$  is used two ways here:  $\dot{v}_e^T$  denotes the adjoint of the operator  $\dot{v}_e$ , while  $\dot{\gamma}^T$  is the row vector (of functions) corresponding to the column vector  $\dot{\gamma}$ .)

For example, if  $\gamma(P) = \int x dP(x) - \mu_0$  where  $\mu_0$  is fixed, then we know the mean of  $P$ ,  $\dot{\gamma}(x) = x - \mu_0$ , and, assuming  $\text{Var}_0(X) > 0$ , from (2)

$$\begin{aligned} I_v^{-1}(s, t) &= F_0(s) \wedge F_0(t) - F_0(s)F_0(t) \\ &\quad - \frac{(\int_{-\infty}^s (x - \mu_0) dF_0(x))(\int_{-\infty}^t (x - \mu_0) dF_0(x))}{\text{Var}_0(X)}. \end{aligned}$$

□

**Example 3. Estimation of a symmetric df.**

Suppose that  $\mathbf{X}, \mathcal{B}, \mu$ , and  $v = F$  are as in example 1, but now suppose that

$$\mathbf{P} = \left\{ P \ll \mu : f = \frac{dP}{d\mu} \text{ is symmetric about } 0 \right\}.$$

Hence  $F(t) = 1 - F(-t)$  for all  $t \in R$ . Let  $\mathbf{M}_0$  denote the model  $\mathbf{P}$  of example 1, let  $v_e$  be  $v = F$  as in example 1. Then  $\dot{v}_e$  is exactly as in example 1, but the tangent space of the model is

$$\dot{\mathbf{P}} = \{h \in L_2^0(P_0): h(x) = h(-x) \text{ a.s. } P_0\}$$

by example 3.2.2. Therefore  $\tilde{I}_t = \Pi_0(\dot{v}_e^T \pi_t | \dot{\mathbf{P}})$  is given by

$$\tilde{I}_t(x) = \frac{1}{2} \{1_{(-\infty, t]}(x) + 1_{(-\infty, t]}(-x)\} - F(t)$$

$$\begin{aligned}
 &= \frac{1}{2} \{1_{(-\infty, t]}(x) + 1_{[-t, \infty)}(x)\} - F(t) \\
 (5) \quad &= \frac{1}{2} \operatorname{sign}(t) \{1_{[0, |t|]}(|x|) - (2F(|t|) - 1)\}
 \end{aligned}$$

and the efficient influence function in estimating a symmetric df is  $\tilde{I}(x)(\cdot) = \frac{1}{2} \operatorname{sign}(\cdot) \{1_{[|x|, \infty)}(|\cdot|) - (2F(|\cdot|) - 1)\}$ . Here and henceforth we suppress the subscript 0. Thus the inverse information covariance function is easily computed to be

$$(6) \quad I_v^{-1}(s, t) = \frac{1}{4} \operatorname{sign}(st) \{G(s) \wedge G(t) - G(s)G(t)\},$$

where

$$G(t) \equiv |2F(t) - 1| = 2F(|t|) - 1 = P(|X| \leq |t|).$$

Hence the Gaussian process  $\mathbb{Z}_0$  of theorem 5.2.1 is given by

$$(7) \quad \mathbb{Z}_0(t) = \frac{1}{2} B_0(2F(t) - 1)$$

where  $B_0$  is a Brownian bridge process on  $[0, 1]$  and  $B_0(-t) \equiv -B_0(t)$  for  $0 \leq t \leq 1$ . Note from (5) that

$$n^{-1/2} \sum_{i=1}^n \tilde{I}_t(X_i) = \frac{1}{2} \{\sqrt{n}(F_n(t) - F(t)) - \sqrt{n}(F_n(-t-) - F(-t-))\},$$

where  $F_n$  is the empirical df and  $F_n(t-)$  denotes the left-limit of  $F_n$  at  $t$ , and hence from example 1,

$$(8) \quad L(\mathbb{Z}_0) = L\left(\frac{1}{2}(B_0(F) - B_0(\bar{F}))\right),$$

with  $\bar{F} = 1 - F$ .

It follows from (7) that

$$(9) \quad L(\|\mathbb{Z}_0\|_\infty) = L\left(\frac{1}{2}\|B_0\|_\infty\right)$$

where  $\|B_0\|_\infty$  has the well-known Kolmogorov-Smirnov distribution. Hence by proposition 5.2.1 under the assumption that  $f$  is symmetric about 0, efficient estimators yield confidence bands which for a continuous df  $F$  are exactly half as wide as the classical bands for  $F$  without any assumption.

If  $\bar{F}_n(t) \equiv 1 - F_n(-t-)$  and  $T_n \equiv (F_n + \bar{F}_n)/2$ , then  $T_n$  is an asymptotically linear (in fact exactly linear), regular estimator of  $v = F$ : with

$$\mathbb{Z}_n \equiv \sqrt{n}(T_n - v) = \sqrt{n}(T_n - F) = n^{-1/2} \sum_{i=1}^n \tilde{I}(X_i),$$

$$(10) \quad \mathbb{Z}_n \Rightarrow \mathbb{Z}_0.$$

Hence the symmetrized empirical df is efficient. Of course the empirical distribution function  $F_n$  is a regular but inefficient estimator of  $F$  in this model  $\mathcal{P}$  with limit process  $\mathbb{Z} = B_0(F)$  as in example 1; the decomposition of  $\mathbb{Z}$  given by theorem 5.2.1 and (8) is



$$\begin{aligned} \mathbf{Z} = B_0(F) &= \frac{B_0(F) - B_0(\bar{F})}{2} + \frac{B_0(F) + B_0(\bar{F})}{2} \\ &\equiv \mathbf{Z}_0 + \Delta_0. \end{aligned}$$

Indeed,  $\mathbf{Z}_0$  and  $\Delta_0$  are independent in view of

$$\begin{aligned} \text{Cov}(B_0(F(s)), B_0(F(t)) + B_0(\bar{F}(t))) \\ = \text{Cov}(B_0(\bar{F}(s)), B_0(F(t)) + B_0(\bar{F}(t))). \end{aligned}$$

Thus,

$$L(\Delta_0) = L\left(\frac{1}{2}(B_0(F) + B_0(\bar{F}))\right)$$

for this particular inefficient estimator of  $v = F$ .

It is also instructive to compute, for  $t \geq 0$ ,

$$\text{Var}(\mathbf{Z}_0(t)) = I_v^{-1}(t, t) = (F(t) - \frac{1}{2})(1 - F(t)),$$

whereas the variance of the limit process  $\mathbf{Z} = B_0(F)$  of the empirical df is

$$\text{Var}(\mathbf{Z}(t)) = F(t)(1 - F(t)),$$

and hence

$$(11) \quad \frac{\text{Var}(\mathbf{Z}_0(t))}{\text{Var}(\mathbf{Z}(t))} = \frac{F(|t|) - \frac{1}{2}}{F(|t|)} = 1 - \frac{1}{2F(|t|)}, \quad t \in R,$$

which equals 0 for  $t = 0$  and approaches 1/2 as  $|t| \rightarrow \infty$ .

Theorem 5.2.1 as applied to this example was given by Millar (1979) in the form of an asymptotic minimax theorem. It has been part of the "folklore" of nonparametric estimation for many years; see, e.g., Schuster (1973), (1975). This example is a special case of a nonparametric family  $\mathbf{P}$  defined by invariance with respect to a group of transformations; a more general treatment of this type of family  $\mathbf{P}$  will be given in section 6.3.  $\square$

**Example 4. Estimation of a df  $F$  in a regular parametric model.**

Now suppose that  $(\mathbf{X}, \mathcal{B}) = (R, \mathcal{B})$  as in example 1, but

$$\mathbf{P} = \{P_\theta : \theta \in \Theta \subset R^k\}$$

is a regular parametric model as in chapter 2. Then

$$\dot{\mathbf{P}} = [\dot{l}_i : i = 1, \dots, k] = [\dot{l}]$$

is a finite-dimensional subspace of  $L_2^0(P_0)$ . Consider estimation of  $v(P) = F$  as in examples 1-3. As in examples 2 and 3, let  $\mathbf{M}_0$  denote the model of example 1, and let  $v_e(P) = F$  on the larger model. Then

$$\begin{aligned}\dot{v}_t &= \Pi_0(\dot{v}_e^T \pi_t \mid \dot{\mathbf{P}}) = \langle \dot{v}_e^T \pi_t, \dot{\mathbf{i}} \rangle_0^T I^{-1}(\theta) \dot{\mathbf{i}} \\ &= \int_{-\infty}^t \dot{\mathbf{i}}^T(x) dP_0(x) \tilde{\mathbf{I}} \\ &= \sum_{i=1}^k \left( \int_{-\infty}^t \dot{\mathbf{i}}_i(x) dP_0(x) \right) \tilde{\mathbf{I}}_i\end{aligned}$$

and

$$\begin{aligned}I_V^{-1}(s, t) &= \langle \dot{v}_e^T \pi_s, \dot{\mathbf{i}}^T \rangle_0 I^{-1}(\theta) \langle \dot{\mathbf{i}}, \dot{v}_e^T \pi_t \rangle_0 \\ &= \int_{-\infty}^s \int_{-\infty}^t \dot{\mathbf{i}}^T(x) I^{-1}(\theta) \dot{\mathbf{i}}(y) dP_0(y) dP_0(x).\end{aligned}$$

□

Now we consider functions other than the df  $F$  corresponding to  $P$ : the cumulative hazard function  $\Lambda$  is considered in example 5, while the mean residual life function  $e$  is considered in examples 6 and 7.

**Example 5. Estimation of a cumulative hazard function.**

Suppose that  $\mathbf{X} = R^+$  and  $\mathbf{P}$  the class of distributions on  $\mathbf{X}$  dominated by Lebesgue measure as in example 1, but now suppose that  $v(P)$  is the cumulative hazard function  $\Lambda$  corresponding to  $P$  with df  $F$ ; more precisely

$$(12) \quad v(P)(t) = \Lambda(t), \quad t \in [0, \tau] \equiv \mathbf{T}$$

where

$$\Lambda(t) \equiv \int_{[0, t]} \frac{1}{1 - F(s-)} dF(s)$$

and  $F(\tau) < 1$ ; here we have written  $F(s-)$  in the denominator to indicate clearly the correct definition of  $\Lambda$  when  $F$  is not continuous, even though  $F(s-) = F(s)$  for  $P \in \mathbf{P}$ . This helps in defining natural estimators correctly both here and in examples 6.4.1 and 6.6.1. Since  $\dot{\mathbf{P}} = L_2^0(P)$  and, by example 1,

$$\begin{aligned}v(P_\eta)(t) - v(P)(t) &= -\log(1 - F(t) - \eta \int (1_{[0, t]} - F(t)) h dF + o(\eta)) \\ &\quad + \log(1 - F(t)) \\ &= \frac{\eta}{1 - F(t)} \int (1_{[0, t]} - F(t)) h dF + o(\eta)\end{aligned}$$

in  $\mathbf{B} = l^\infty(\mathbf{T})$ , it is straightforward to verify that  $v$  is pathwise differentiable with  $\dot{v}_t$  of (5.2.25) given by

$$(13) \quad \dot{v}_t(x) = \frac{1_{[0, t]}(x) - F(t)}{1 - F(t)}, \quad t \in [0, \tau].$$

Thus the efficient influence function  $\tilde{\mathbf{I}}$  for  $v(P)$  is given by

$$\tilde{\mathbf{I}}(x) = \frac{1_{[x, \infty)}(\cdot) - F(\cdot)}{1 - F(\cdot)},$$

and the inverse information covariance function (5.2.23) is simply

$$(14) \quad I_V^{-1}(s, t) = \frac{F(s) \wedge F(t) - F(s)F(t)}{(1 - F(s))(1 - F(t))} = C(s) \wedge C(t),$$

with  $C(t) \equiv F(t)/(1 - F(t))$ . Hence the Gaussian process  $Z_0$  of theorem 5.2.1 is

$$(15) \quad Z_0(t) = B(C(t)) \quad \text{for } t \in [0, \tau]$$

where  $B$  is standard Brownian motion. The natural cumulative hazard function estimator

$$\Delta_n(t) \equiv \int_{[0, t]} \frac{1}{1 - F_n(s-)} dF_n(s)$$

is a regular estimator with limit process  $Z_0$ , and hence is efficient. See, e.g., Shorack and Wellner (1986, chapter 6). □

**Example 6. Estimation of mean residual life.**

Suppose that  $X = R^+$ ,  $B = I^\infty(T)$ , and  $T = [0, \tau]$  as in the last example, but now suppose that  $P = \{P : P \ll \text{Lebesgue measure}, E_P X^2 \leq M\}$  for some  $M < \infty$  and that  $v(P)$  is the mean residual life function

$$(16) \quad v(P)(t) = E(X - t | X > t) \\ = \frac{\int_t^\infty (x - t) dF(x)}{\bar{F}(t)} \equiv e(t), \quad t \in T,$$

where  $\bar{F}(t) \equiv 1 - F(t) = P(X > t)$  and  $\bar{F}(\tau) > 0$ . By (28) of example 8 below with  $F = \{(\cdot - t) \vee 0 : t \geq 0\}$ , we see that for  $s(\eta) = s(0)$

$$+ \eta h s(0)/2 + o(\eta), F_\eta(x) = \int_0^x s^2(\eta) d\mu,$$

$$\int_t^\infty (x - t) dF_\eta(x) = \int_t^\infty (x - t) dF(x) \\ + \eta \int_t^\infty (x - t) h(x) dF(x) + o(\eta)$$

uniformly in  $t$ , and hence, by example 1,

$$v(P_\eta)(t) - v(P)(t) \\ = \int_t^\infty (x - t) dF(x) \left\{ \frac{1}{\bar{F}_\eta(t)} - \frac{1}{\bar{F}(t)} \right\} \\ + \frac{\eta}{\bar{F}_\eta(t)} \int_t^\infty (x - t) h(x) dF(x) + o(\eta) \\ = -\eta \frac{e(t)}{\bar{F}(t)} \int_t^\infty h(x) dF(x) + \frac{\eta}{\bar{F}(t)} \int_t^\infty (x - t) h(x) dF(x) + o(\eta)$$

uniformly in  $t \leq \tau$ . Consequently  $v$  is pathwise differentiable and, since  $\dot{P} = L_2^0(P)$ ,  $\dot{v}_t$  of (5.2.25) is given by

$$(17) \quad \dot{v}_t(x) = \{(x - t) - e(t)\} \frac{1_{[x > t]}}{\bar{F}(t)}, \quad t \in [0, \tau].$$

Hence the efficient influence function  $\tilde{I}$  from (5.2.11) exists and is given by  $\pi_t \tilde{I}(x) = \dot{v}_t(x)$ ,  $t \in \mathbf{T}$ . The inverse information covariance function (5.2.23) is

$$(18) \quad I_V^{-1}(s, t) = \frac{E\{[(X - t) - e(t)]^2 1_{[X > t]}\}}{\bar{F}(s)\bar{F}(t)}$$

$$= \frac{\text{Var}(X | X > t)}{\bar{F}(s)}$$

$$\equiv \frac{\sigma^2(t)\bar{F}(t)}{\bar{F}(s)\bar{F}(t)}, \quad \text{for } 0 \leq s \leq t \leq \tau,$$

and hence the Gaussian limit process  $\mathbf{Z}_0$  of theorem 5.2.1 can be represented as

$$(19) \quad \mathbf{Z}_0(t) = \frac{\sigma}{\bar{F}(t)} B(\bar{G}(t))$$

where  $B$  is standard Brownian motion on  $[0, 1]$ ,  $\sigma^2 \equiv \text{Var}(X)$ , and

$$\bar{G}(t) \equiv \frac{\sigma^2(t)\bar{F}(t)}{\sigma^2}$$

is a survival function on  $R^+$  (so that  $G \equiv 1 - \bar{G}$  is a df); see, e.g., Csörgő, Csörgő, and Horváth (1986). If  $\hat{e}_n \equiv v(P_n)$  denotes the empirical mean residual life function, then

$$(20) \quad \mathbf{Z}_n \equiv \sqrt{n}(\hat{e}_n - e) \Rightarrow \mathbf{Z}_0 \quad \text{in } D[0, \tau]$$

was proved by Yang (1978), Hall and Wellner (1980a), and Csörgő, Csörgő, and Horváth (1986).  $\square$

**Example 7. Estimation of mean residual life in a family  $\mathbf{P}$  with known mean.**

Now suppose that  $\mathbf{X} = R^+$  and  $v(P)(t) = e(t)$  as in example 6, but that the model  $\mathbf{P}$  is restricted to

$$(21) \quad \mathbf{P} = \{P : P \ll \text{Lebesgue measure}, E_P X^2 \leq M, E_P X = \mu_0\}$$

where  $M < \infty$  and  $\mu_0$  is fixed and known as in example 3.2.3. Then

$$(22) \quad \dot{\mathbf{P}} = \{h \in L_2^0(P_0) : \langle h, X - \mu_0 \rangle_0 = 0\}.$$

Let  $\mathbf{M}_0$  be the larger model as in example 6, and let  $v_e$  denote  $v$ , mean residual life, as defined on the larger model  $\mathbf{M}_0$ . Then  $\dot{v}_e$  is exactly as calculated in example 6, and

$$(23) \quad \dot{v}_t(x) \equiv \tilde{I}_t(x) = \Pi_0(\dot{v}_e^T \pi_t | \dot{\mathbf{P}})(x)$$

$$= \frac{c(x, t) - \bar{G}(t)c(x, 0)}{\bar{F}(t)},$$

where  $c(x, t) \equiv \{x - t - e(t)\} 1_{[x > t]}$  and  $\bar{G}$  as in example 6. By straightforward calculation, the inverse information covariance function is

$$(24) \quad \Gamma_v^{-1}(s, t) = \frac{\sigma^2}{\bar{F}(s)\bar{F}(t)} \{ \bar{G}(s) \wedge \bar{G}(t) - \bar{G}(s)\bar{G}(t) \},$$

and hence

$$(25) \quad \mathbf{Z}_0(t) = \frac{\sigma}{\bar{F}(t)} B_0(\bar{G}(t)),$$

where  $B_0$  is a Brownian bridge process.

For this submodel of example 6, the empirical mean residual life function  $\hat{e}_n = v_e(\hat{P}_n)$  is inefficient: it is a regular estimator with limit process  $\mathbf{Z} = \sigma B(\bar{G})/\bar{F}$  where  $B$  is Brownian motion. Hence the decomposition of the limit process  $\mathbf{Z}$  of  $\hat{e}_n$  is just

$$(26) \quad \begin{aligned} \mathbf{Z} &= \frac{\sigma}{\bar{F}} (B(\bar{G}) - \bar{G}B(1)) + \frac{\sigma}{\bar{F}} \bar{G}B(1) \\ &\equiv \mathbf{Z}_0 + \Delta_0, \end{aligned}$$

where  $\mathbf{L}(\Delta_0) = \mathbf{L}(\frac{\sigma \bar{G}Z}{\bar{F}})$  and  $\mathbf{L}(Z) = \mathbf{L}(B(1)) = N(0, 1)$ . Note that the independence of  $\mathbf{Z}_0$  and  $\Delta_0$  follows from

$$\text{Cov}(B(\bar{G}(s)) - \bar{G}(s)B(1), B(1)) = 0.$$

Of course this neat and clean decomposition is a result of the simple relationship between the constraint  $E_P X = \mu_0$  and the function  $v = e$  being estimated; other constraints will lead to somewhat more complicated results.  $\square$

To conclude this section we generalize example 1 considerably; in the following example we consider estimation of a measure  $P$  on an arbitrary sample space  $\mathbf{X}$ .

**Example 8. Estimation of a probability measure  $P$ .**

For a general sample space  $(\mathbf{X}, \mathcal{B})$ , let  $\mathbf{P} = \mathbf{M}$ , the collection of all distributions  $P$  on  $(\mathbf{X}, \mathcal{B})$ , and fix  $P \in \mathbf{P}$ . Suppose that  $\mathbf{F} \subset L_2(P)$ , and let  $v(P)$  be the function from  $\mathbf{F}$  to  $R$  defined by

$$(27) \quad v(P)(f) = \int f dP \equiv P(f) \quad \text{for } f \in \mathbf{F}.$$

Thus we are taking  $\mathbf{T} = \mathbf{F}$  and  $\mathbf{B} = l^\infty(\mathbf{F})$  with the supremum norm  $\|b\|_{\mathbf{B}} \equiv \sup_{f \in \mathbf{F}} |b(f)| \equiv \|b\|_{\mathbf{F}}$  for  $b \in \mathbf{B} = l^\infty(\mathbf{F})$ . Furthermore we use the  $L_2(P)$  pseudometric  $d \equiv e_P$  on  $\mathbf{F}$  given by  $d(f, g) = \|f - g\|_0 = e_P(f, g)$  for  $f, g \in \mathbf{F}$ .

Now suppose there is an envelope function  $f_e \in L_2(P)$  for  $\mathbf{F}$ ; i.e.,  $|f| \leq f_e$  pointwise for all  $f \in \mathbf{F}$ , and that  $f_e$  has bounded second moment over  $\mathbf{P}$ ,

$$\sup_{P \in \mathbf{P}} \int f_e^2 dP < \infty.$$

As we will show below,  $v$  is pathwise differentiable at any  $P \in \mathbf{P}$  with derivative  $\dot{v}(P) \equiv \dot{v} : \dot{\mathbf{P}} \rightarrow l^\infty(\mathbf{F})$  given by

$$(28) \quad \begin{aligned} \dot{v}(h)(f) &= \int (f(x) - \int f dP) h(x) dP(x) \\ &= \int \dot{v}_f(x) h(x) dP(x). \end{aligned}$$

Consequently, the efficient influence function  $\tilde{I}$  is given by

$$(29) \quad \tilde{I}(x)(f) = \dot{v}_f(x) = f(x) - \int f dP.$$

Then the inverse information covariance function (5.2.23) becomes

$$(30) \quad \begin{aligned} I_v^{-1}(f, g) &= E(f - Ef)(g - Eg) = P(fg) - P(f)P(g) \\ &= \text{Cov}(\mathbf{Z}_0(f), \mathbf{Z}_0(g)), \end{aligned}$$

where  $\mathbf{Z}_0$  is a  $P$ -Brownian bridge process.

For several particular classes of functions  $\mathbf{F}$  the empirical measure  $v(P_n)$  (i.e.  $v(P_n)(f) = P_n(f) = n^{-1} \sum_{i=1}^n f(X_i)$  for  $f \in \mathbf{F}$ ) is a regular estimator which attains this bound; note that  $P_n - P - n^{-1} \sum_{i=1}^n \tilde{I}(X_i) = 0$  and see, e.g., Pollard (1982), Giné and Zinn (1984), Ossiander (1987), and Dudley (1978), (1984), (1987). Wellner (1992) shows that if  $\mathbf{F}$  is a Donsker class for a fixed  $P$  (i.e., the CLT holds:  $\mathbf{Z}_0 \in UC(\mathbf{F}, e_P)$  a.s. and  $\mathbf{Z}_n \Rightarrow \mathbf{Z}_0$ ), then boundedness of the second moment of  $f_e$  over  $\mathbf{P}$  suffices for (local) regularity of  $v(P_n)$ .

Note that example 1 is the special case with  $\mathbf{X} = R$  and  $\mathbf{F} = \{1_{(-\infty, t]} : t \in R\}$ . Note, however that in the present setup we are using the pseudometric  $d$  given by  $d^2(s, t) = F(s \vee t) - F(s \wedge t) \equiv \int \{1_{(-\infty, s]} - 1_{(-\infty, t]}\}^2 dF$  for  $s, t \in R$ ; recall the discussion of  $\mathbf{B} = l^\infty(R^k)$  in section 5.2. This makes the limit process (uniformly) continuous even if  $F$  is discontinuous, whereas in example 1 the metric on  $\mathbf{T} = R$  is just the usual Euclidean metric, and the limit process is discontinuous if  $F$  is discontinuous.

More generally, when  $\mathbf{X} = R^d$ , we will frequently consider  $\mathbf{F} = \{1_C : C \in \mathbf{C}\}$  for some class of sets  $\mathbf{C}$ . Of particular interest are Vapnik-Chervonenkis classes  $\mathbf{C}$ ; see, e.g., Pollard (1984), Dudley (1984), or Shorack and Wellner (1986, chapter 26).

To prove (28), suppose  $s(\eta) = s(0) + \frac{1}{2}\eta hs(0) + o(\eta)$  and write

$$\begin{aligned} v(P_\eta)(f) - v(P_0)(f) - \eta \dot{v}(h)(f) &= \int f s^2(\eta) d\mu - \int f s^2(0) d\mu - \eta \int f h s^2(0) d\mu \\ &= \int f [(s(\eta) - s(0)) - \eta(\frac{1}{2}hs(0))](s(\eta) + s(0)) d\mu \\ &\quad + \frac{\eta}{2} \int f h (s(\eta) - s(0)) s(0) d\mu \\ &\equiv A_\eta(f) + B_\eta(f) \end{aligned}$$

where

$$|A_\eta(f)| \leq \{2 \int f_e^2 dP_\eta + 2 \int f_e^2 dP_0\}^{1/2} \|(s(\eta) - s(0)) - \frac{1}{2} \eta h s(0)\|$$

$$= o(\eta) \quad \text{as } \eta \rightarrow 0$$

uniformly in  $f \in \mathbf{F}$ , and

$$\begin{aligned} \left| \frac{2}{\eta} B_\eta(f) \right| &\leq \left| \int_{|f| \leq \eta^{-1/2}} f h(s(\eta) - s(0)) s(0) d\mu \right| \\ &\quad + \left| \int_{|f| > \eta^{-1/2}} f h(s(\eta) - s(0)) s(0) d\mu \right| \\ &\leq \eta^{-1/2} \|h\|_0 \|s(\eta) - s(0)\| \\ &\quad + \left\{ \int h^2 1_{|f| > \eta^{-1/2}} dP(P_\eta(f_e^2) + P(f_e^2)) \right\}^{1/2} \\ &\rightarrow 0 \quad \text{as } \eta \rightarrow 0, \end{aligned}$$

and hence  $v(P)$  is differentiable with derivative  $\dot{v}$  given by (28). Note that this is basically Van der Vaart's (1988a, lemma 5.21, page 167), but using the envelope function  $f_e$  to bound functions  $f \in \mathbf{F}$ .  $\square$

### 5.4 DIFFERENTIABILITY OF FUNCTIONS

A key hypothesis in both sections 3.3 and 5.2 was pathwise differentiability of  $v: \mathbf{P} \rightarrow \mathbf{R}$  or  $v: \mathbf{P} \rightarrow \mathbf{B}$ ; since the former is a special case of the latter, we consider the general case here. Given the derivative  $\dot{v}: \dot{\mathbf{P}} \rightarrow \mathbf{B}$ , we established convolution theorems for regular estimators of  $v$ . The optimal limit process  $\mathbf{Z}_0$  on  $\mathbf{B}$  is uniquely determined (if it exists) by its finite-dimensional distributions which are given in terms of the efficient influence operator  $\tilde{\mathbf{I}}_v: \mathbf{B}^* \rightarrow \dot{\mathbf{P}}$  by

$$(1) \quad b^* \mathbf{Z}_0 \sim N(0, E[\tilde{\mathbf{I}}_v(b^*)]^2).$$

On the other hand, in many models parametrized by a subset  $\mathbf{G}$  of a Hilbert space  $(\mathbf{H}, \langle \cdot, \cdot \rangle)$ , the tangent space  $\dot{\mathbf{P}}$  can be described in terms of the closure of the range of a bounded linear operator  $\dot{\mathbf{I}}, \dot{\mathbf{P}} = \overline{\mathbf{R}(\dot{\mathbf{I}})}$ . Many of the examples in chapter 4 had this structure. We now make this more precise, and give conditions for differentiability of certain parameters in terms of the operator  $\dot{\mathbf{I}}$ .

Suppose that  $\mathbf{P} = \{P_g: g \in \mathbf{G}\}$  where  $\mathbf{G}$  is a subset of  $\mathbf{H}$ . For a fixed  $g \in \mathbf{G}$ , suppose that  $\mathbf{G}_g$  is a collection of curves or paths  $\{g_\eta\}$  through  $g$  such that

$$(2) \quad g_\eta = g + \eta h + o(\eta) \quad \text{in } \mathbf{H}$$

for some  $h \in \mathbf{H}$ . Let  $\dot{\mathbf{G}}^0$  be the collection of  $h$ 's corresponding to  $\mathbf{G}_g$ . We assume that:

(A1)  $\dot{\mathbf{G}}^0$  is a closed and linear subspace of  $\mathbf{H}$ ; i.e.,  $\dot{\mathbf{G}} = \dot{\mathbf{G}}^0$ .

(A2) There is a bounded linear operator  $\dot{\mathbf{I}}_g \equiv \dot{\mathbf{I}} : \dot{\mathbf{G}}^0 \rightarrow L_2^0(P_g)$  such that, for every  $\{g_\eta\}$  in  $\mathbf{G}_g$  satisfying (2),

$$(3) \quad s(g_\eta) - s(g) - \eta \left(\frac{1}{2} \dot{\mathbf{I}} h\right) s(g) = o(\eta) \quad \text{in } L_2(\mu).$$

When (3) holds we say that  $\mathbf{P}$  is *Hellinger-differentiable at  $P_g$  along the curve  $\{P_{g_\eta}\}$* , and we call  $\dot{\mathbf{I}}$  the *score operator for  $g$* . Note that  $\mathbf{R}(\dot{\mathbf{I}}) \subset \dot{\mathbf{P}}^0$ .

(A3) The parameter  $v : \mathbf{P} \rightarrow \mathbf{B}$  is of the form

$$(4) \quad v(P_g) = \psi(g),$$

where  $\psi : \mathbf{G} \rightarrow \mathbf{B}$  is pathwise differentiable at  $g$  in the following sense. (A1) holds and there is a bounded linear operator  $\dot{\psi}_g \equiv \dot{\psi} : \dot{\mathbf{G}}^0 \rightarrow \mathbf{B}$  such that, for any family  $\{g_\eta\}$  satisfying (2),

$$(5) \quad \|\eta^{-1}(\psi(g_\eta) - \psi(g)) - \dot{\psi}(h)\|_{\mathbf{B}} = o(1).$$

The following theorem is due to Van der Vaart (1991).

**Theorem 1.** Suppose that (A1)–(A3) hold.

A. If  $v : \mathbf{P} \rightarrow \mathbf{B}$  given by (4) is pathwise differentiable at  $P_g \in \mathbf{P}$ , then

$$(6) \quad \mathbf{R}(\dot{\psi}^T) \subset \mathbf{R}(\dot{\mathbf{I}}^T),$$

where  $\dot{\psi}^T : \mathbf{B}^* \rightarrow \dot{\mathbf{G}}$  and  $\dot{\mathbf{I}}^T : L_2^0(P_g) \rightarrow \dot{\mathbf{G}}$  are the adjoints of  $\dot{\psi}$  and  $\dot{\mathbf{I}}$  respectively, and  $\mathbf{R}$  denotes the range of an operator (see (A.1.12)).

B. If (6) holds and  $\dot{\mathbf{P}} = \overline{\mathbf{R}(\dot{\mathbf{I}})}$ , then  $v$  is pathwise differentiable and the efficient influence operator  $\dot{\mathbf{I}}_v$  for estimation of  $v$  is the unique operator from  $\mathbf{B}^*$  to  $\dot{\mathbf{P}}$  satisfying

$$(7) \quad \dot{\psi}^T = \dot{\mathbf{I}}^T \dot{\mathbf{I}}_v.$$

**Remark 1.** Note that  $\mathbf{R}(\dot{\psi}^T)^\perp = \mathbf{N}(\dot{\psi})$  and  $\mathbf{R}(\dot{\mathbf{I}}^T)^\perp = \mathbf{N}(\dot{\mathbf{I}})$  by (A.1.18), so (6) implies

$$(8) \quad \mathbf{N}(\dot{\mathbf{I}}) \subset \mathbf{N}(\dot{\psi}).$$

Necessity of (8) is clear by reasoning as follows. If  $h \in \mathbf{N}(\dot{\mathbf{I}})$ , then by (3),  $\dot{\mathbf{I}}h = 0$ , and (5.2.1),

$$\psi(g_\eta) - \psi(g) = \eta \dot{\psi}(h) + o(\eta) = o(\eta),$$

and hence  $\dot{\psi}(h) = 0$ , i.e.  $h \in \mathbf{N}(\dot{\psi})$ .

**Remark 2.** The operator  $\dot{\mathbf{I}}^T \dot{\mathbf{I}}$  is a positive operator, and hence has a square root  $(\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{1/2}$ . Moreover,  $\mathbf{R}(\dot{\mathbf{I}}^T) = \mathbf{R}((\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{1/2})$  (see proposition A.1.6), so that (6) can be restated as

$$(9) \quad \mathbf{R}(\dot{\psi}^T) \subset \mathbf{R}((\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{1/2}),$$



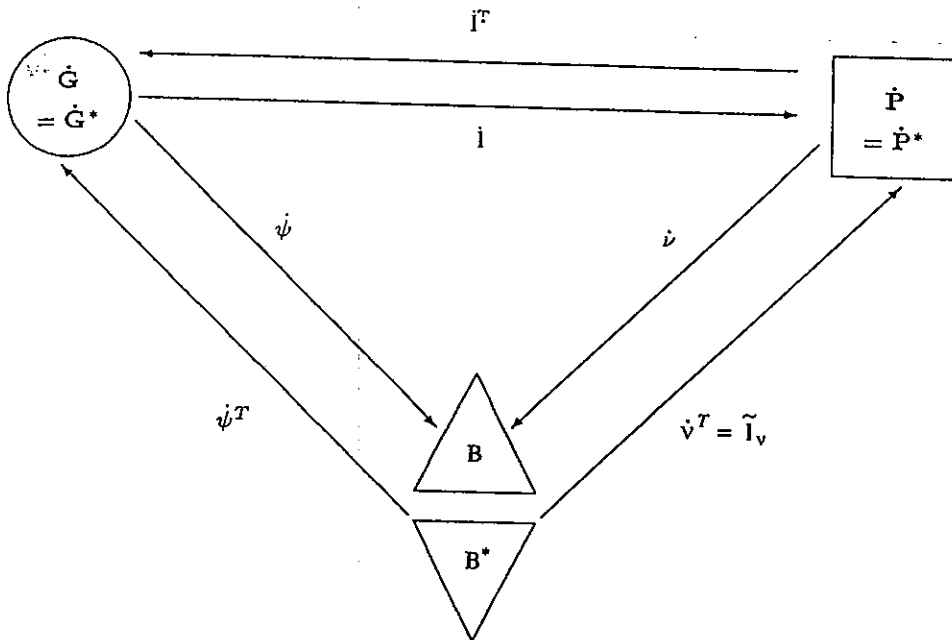


FIGURE 1. The function diagram for theorem 1.

and this is sometimes easier to check, especially if the spectral decomposition of  $\dot{\mathbf{i}}^T \dot{\mathbf{i}}$  is available.

**Remark 3.** There may be an extension  $v_e$  of  $v$  to  $P_e \supset P$ . Then  $\tilde{I}_v$  solves (7) (with the convention  $\dot{\mathbf{i}}^T h = 0$  for  $h \in \dot{P}^\perp$ ) and corresponds to inefficient influence functions.

**Proof of theorem 1.** Suppose that  $v$  is pathwise differentiable. Then by (5.2.1) and (3) through (5)

$$(a) \quad \dot{v}(\dot{h}) = \lim_{\eta \rightarrow 0} \frac{v(P_{g_\eta}) - v(P_g)}{\eta} = \dot{\psi}(h)$$

for every  $h \in \dot{G}$ . Therefore we have continuous linear maps  $\dot{\psi} : \dot{G} \rightarrow B$ ,  $\dot{\mathbf{i}} : \dot{G} \rightarrow \dot{P}$ ,  $\dot{v} : \dot{P} \rightarrow B$ , such that

$$(b) \quad \dot{\psi} = \dot{v} \dot{\mathbf{i}}.$$

Hence

$$(c) \quad \dot{\psi}^T = \dot{\mathbf{i}}^T \dot{v}^T,$$

and (6) follows.

Conversely, if (6) holds, define  $\dot{v} : R(\dot{\mathbf{i}}) \rightarrow B$  by

$$(d) \quad \dot{v}(\dot{h}) = \dot{\psi}(h) \quad \text{for } h \in \dot{G}.$$

This is well defined, since, if  $\dot{h}_1 = \dot{h}_2$ , then  $h_1 - h_2 \in N(\dot{\mathbf{i}}) \subset N(\dot{\psi})$  by (8), and hence  $\dot{\psi}(h_1) = \dot{\psi}(h_2)$ . Furthermore  $\dot{v}$  is a linear map on  $R(\dot{\mathbf{i}})$ .

It remains only to show that  $\dot{v}$  is continuous (i.e., bounded) and extendable to  $\dot{P}$ . For any  $b^* \in B^*$  it follows from (d) that

$$\begin{aligned} b^* \dot{v}(\dot{h}) &= \langle \dot{\psi}h, b^* \rangle_B = \langle h, \dot{\psi}^T b^* \rangle \\ &= \langle h, \dot{I}^T a \rangle \quad \text{by (6) for some } a = a(b^*) \in L_2^0(P_g) \\ &= \langle \dot{h}, a \rangle_0. \end{aligned}$$

Thus  $b^* \dot{v}$  is a continuous linear (real-valued) functional on  $R(\dot{I})$  for every  $b^* \in B^*$ . Boundedness of  $\dot{v}: R(\dot{I}) \rightarrow B$  follows from lemma A.1.1 (with  $X \equiv R(\dot{I})$  and  $Y \equiv B$ ), and by continuity  $\dot{v}$  can be extended to  $\dot{P} = \overline{R(\dot{I})}$ .

Finally, (7) is implied by (c) and the definition of  $\tilde{I}_v$ ;  $\tilde{I}_v b^* = \dot{v}^T b^*$ .

It remains only to prove uniqueness. Suppose that  $\tilde{I}_v^{(1)}, \tilde{I}_v^{(2)}$  are both solutions with  $d \equiv \tilde{I}_v^{(1)}(b^*) - \tilde{I}_v^{(2)}(b^*) \neq 0$  for some  $b^* \in B^*$ . Note that  $d \in \dot{P}$ . Then  $\dot{I}^T d = 0$ , and  $\langle \dot{I}^T d, h \rangle = 0 = \langle \dot{h}, d \rangle_0$  for all  $h \in \dot{G}$ . Thus  $d \perp \dot{P}$ , a contradiction. We conclude that  $\tilde{I}_v$  solving (7) is unique.  $\square$

### Corollaries and Consequences

The hypothesis (6) of theorem 1 can be strengthened in a variety of ways in order to imply pathwise differentiability of  $v$ . The following corollary was obtained by Van der Vaart (1991).

**Corollary 1.** Suppose that (A1)–(A3) hold and  $\overline{R(\dot{I})} = \dot{P}$ . Then  $v: P \rightarrow B$  is pathwise differentiable at  $P_g \in P$  if

$$(10) \quad R(\dot{I}) \text{ is closed}$$

and (8) holds, i.e.

$$N(\dot{I}) \subset N(\dot{\psi}).$$

**Proof.** By theorem 1.B it suffices to show that (8) and (10) imply (6). Since  $h \in N(\dot{\psi})$  implies  $\langle h, \dot{\psi}^T b^* \rangle = \langle \dot{\psi}h, b^* \rangle_B = 0$ , we have  $\dot{\psi}^T b^* \perp N(\dot{\psi})$ . Thus by (8) and (A.1.18),

$$\dot{\psi}^T b^* \in N(\dot{\psi})^\perp \subset N(\dot{I})^\perp = \overline{R(\dot{I}^T)} = R(\dot{I}^T),$$

where the last equality follows since  $R(\dot{I})$  is closed if and only if  $R(\dot{I}^T)$  is closed; see proposition A.1.7.D. Hence (6) holds.  $\square$

For a simple example in which (10) fails, see the indicator censoring example at the end of this section. Another class of examples for which (10) typically fails are the mixture models of section 4.5. Then, as in (4.5.13),

$$\dot{h}(x) = E(h(U) | X = x),$$

for  $h \in L_2^0(G)$ , and

$$\dot{I}^T a(u) = E(a(X) | U = u) - E a(X),$$

for  $a \in L_2(P)$ . We will return to this example in section 6.5. For examples in which (10) holds, recall the transformation models of section 4.7. (10) was proved for the joint distribution-transformation model in proposition 4.7.2 and for the copula model in proposition 4.7.5. Unfortunately, (10) fails quite frequently as we will see in the examples in chapter 6.

Corollary 1 replaces (6) by its implication (8) and closedness of the range of  $\dot{\mathbf{i}}$ ; this condition (10) is clearly not necessary since it does not hold in e.g. example 1 below, although there exist pathwise differentiable parameters there. Note that (10) does not involve the particular function  $\psi$ . The following corollary goes further, and replaces both (8) and (10) by conditions independent of the particular function  $\psi$ . It is essentially in the spirit of the hypotheses of section 4 of Begun, Hall, Huang, and Wellner (1983).

**Corollary 2.** Suppose that (A1)–(A3) hold. Then:

A.  $v : \mathbf{P} \rightarrow \mathbf{B}$  is pathwise differentiable at  $P_g \in \mathbf{P}$  if  $\overline{\mathbf{R}(\dot{\mathbf{i}})} = \dot{\mathbf{P}}$ ,

$$(11) \quad \mathbf{N}(\dot{\mathbf{i}}) = \{0\}$$

and (10) holds; i.e.,

$$\mathbf{R}(\dot{\mathbf{i}}) \text{ is closed.}$$

B.  $\dot{\mathbf{i}}^T \dot{\mathbf{i}}$  is one-to-one and onto if and only if (10) and (11) hold.

C. The inclusion

$$(12) \quad \mathbf{R}(\dot{\mathbf{i}}^T \dot{\mathbf{i}}) \subset \mathbf{R}(\dot{\mathbf{i}}^T)$$

holds with equality if and only if (10) holds.

When (10) fails, we can work instead with the condition (6) and “solve” (7): i.e. given  $b^* \in \mathbf{B}^*$ , find a solution  $\tilde{\mathbf{i}}_v b^* \in \dot{\mathbf{P}} \subset L_2^0(P_g)$  of

$$(13) \quad \dot{\psi}^T b^* = \dot{\mathbf{i}}^T (\tilde{\mathbf{i}}_v b^*).$$

Alternatively, in view of corollary 2.C, we can replace (6) by the sufficient condition

$$(14) \quad \mathbf{R}(\dot{\psi}^T) \subset \mathbf{R}(\dot{\mathbf{i}}^T \dot{\mathbf{i}})$$

and find  $h^*$  to solve, for fixed  $b^* \in \mathbf{B}^*$ ,

$$(15) \quad \dot{\psi}^T b^* = \dot{\mathbf{i}}^T \dot{\mathbf{i}} h^*$$

for  $h^* \in \dot{\mathbf{G}}$ . Note that (15) is similar to the usual “normal equations” of linear regression theory,  $X^T X \beta = X^T Y$  with  $\dot{\mathbf{i}}$  playing the role of  $X$ ,  $h^*$  replacing  $\beta$ , and  $\dot{\psi}^T b^*$  replacing  $X^T Y$ .

Write  $h^*$  satisfying (15) as

$$h^* = (\dot{\mathbf{i}}^T \dot{\mathbf{i}})^- \dot{\psi}^T b^*.$$

Then

$$(16) \quad \tilde{\mathbf{i}}_v b^* = \dot{\mathbf{i}} (\dot{\mathbf{i}}^T \dot{\mathbf{i}})^- \dot{\psi}^T b^* \in \mathbf{R}(\dot{\mathbf{i}}) \subset \dot{\mathbf{P}}.$$

**Proof of corollary 2.** A follows immediately from corollary 1 since (11) implies (8) trivially.

To prove B, note that if  $\dot{I}^T \dot{I}$  is one-to-one, then so is  $\dot{I}$ . If  $\dot{I}^T \dot{I}$  is onto, then  $\dot{I}^T$  is onto and  $\mathbf{R}(\dot{I}^T) = \dot{G}^0$  is closed by (A1). Thus  $\mathbf{R}(\dot{I})$  is closed too, as in the proof of corollary 1. Conversely, if  $\dot{I}^T \dot{I}h = 0$ , then  $\dot{I}h \in \mathbf{N}(\dot{I}^T) = \mathbf{R}(\dot{I})^\perp$  by (A.1.17). Since  $\dot{I}h \in \mathbf{R}(\dot{I})$ , this implies  $\dot{I}h = 0$ , and hence  $h = 0$  if  $\mathbf{N}(\dot{I}) = \{0\}$ . That  $\dot{I}^T \dot{I}$  is onto is a consequence of C since (10) then implies  $\mathbf{R}(\dot{I}^T \dot{I}) = \mathbf{R}(\dot{I}^T)$ , where by (A.1.18)

$$\mathbf{R}(\dot{I}^T) = \overline{\mathbf{R}(\dot{I}^T)} = \mathbf{N}(\dot{I})^\perp = \{0\}^\perp.$$

The inclusion in C is obvious. If  $\mathbf{R}(\dot{I})$  is closed, suppose  $h = \dot{I}^T \alpha \in \mathbf{R}(\dot{I}^T)$ . Then

$$\begin{aligned} h &= \dot{I}^T \alpha = \dot{I}^T (\Pi(\alpha | \mathbf{R}(\dot{I})) + \Pi(\alpha | \mathbf{R}(\dot{I})^\perp)) \\ &= \dot{I}^T (\dot{I}h_1 + \Pi(\alpha | \mathbf{N}(\dot{I}^T))) \quad \text{by (A.1.17)} \\ &= \dot{I}^T \dot{I}h_1 \quad \text{for some } h_1 \in \dot{G}, \end{aligned}$$

so  $h \in \mathbf{R}(\dot{I}^T \dot{I})$ .

Conversely, let  $\alpha \in \overline{\mathbf{R}(\dot{I})}$ . Consider  $\dot{I}^T \alpha$ ; by  $\mathbf{R}(\dot{I}^T) = \mathbf{R}(\dot{I}^T \dot{I})$  there exists  $h_0$  such that  $\dot{I}^T \alpha = \dot{I}^T \dot{I}h_0$ . Hence  $\alpha - \dot{I}h_0 \in \mathbf{N}(\dot{I}^T) = \mathbf{R}(\dot{I})^\perp$ . On the other hand,  $\alpha - \dot{I}h_0 \in \overline{\mathbf{R}(\dot{I})}$ . Thus,  $\alpha - \dot{I}h_0 = 0$ , or  $\alpha \in \mathbf{R}(\dot{I})$ ; hence  $\mathbf{R}(\dot{I})$  is closed.  $\square$

The following proposition gives relationships between various ranges and closures of ranges of the score and information operators; it will be used in connection with mixture models in section 6.5. We give it here because its proof uses the same type of arguments as the proofs of the preceding corollaries. Again corollary A.1.1 is the key tool.

- Proposition 1.** Suppose that  $\mathbf{P}$  is any model with score operator  $\dot{I}$ . Then:
- $\mathbf{R}(\dot{I})$  is dense in  $L_2^0(\mathbf{P})$  if and only if  $\dot{I}^T$  is one-to-one (i.e.,  $\mathbf{N}(\dot{I}^T) = \{0\}$ ).
  - If  $\dot{I}^T$  is one-to-one, then  $\overline{\mathbf{R}(\dot{I})} = L_2^0(\mathbf{P}) = \dot{M}$ .
  - If  $\dot{I}$  is one-to-one, then  $\mathbf{R}(\dot{I}^T \dot{I}) = L_2^0(\mathbf{G})$  if and only if  $(\dot{I}^T \dot{I})^{-1}$  is bounded.

**Proof.** A follows from corollary A.1.1:  $\mathbf{N}(\dot{I}^T) = \mathbf{R}(\dot{I})^\perp$ ; or see, e.g., Rudin (1973, corollary 4.12.(b), page 94). Part B is an immediate consequence of A.

To prove part C, first note that  $\mathbf{N}(\dot{I}^T \dot{I}) = \mathbf{N}(\dot{I})$ :  $\mathbf{N}(\dot{I}) \subset \mathbf{N}(\dot{I}^T \dot{I})$  is obvious; if  $h \in \mathbf{N}(\dot{I}^T \dot{I})$ , then

$$0 = \langle \dot{I}^T \dot{I}h, h \rangle_G = \langle \dot{I}h, \dot{I}h \rangle_0 = \|\dot{I}h\|_0^2,$$

so that  $\dot{I}h = 0$  a.e.  $P_0$ , and hence  $h \in \mathbf{N}(\dot{I})$ . Since  $\dot{I}$  is one-to-one by hypothesis,  $\dot{I}^T \dot{I}$  is also one-to-one. Now by corollary A.1.1 again, the self-adjoint operator  $\dot{I}^T \dot{I}$  has range which is dense in  $L_2^0(\mathbf{G})$ . By Banach's theorem, proposition A.1.7.B,  $(\dot{I}^T \dot{I})^{-1}$  is bounded if and only if the range of  $\dot{I}^T \dot{I}$  is equal to  $L_2^0(\mathbf{G})$ .  $\square$

### An Example

To illustrate the application of theorem 1, we give just one example here; more examples are given in chapter 6. The point of this example is to show concretely that when  $\mathbf{R}(\dot{\mathbf{I}})$  is not closed (as is frequently true), then some functions of  $G$  may be pathwise differentiable as a function of the distribution  $P$  of the observations (and then theorems 3.3.2, 5.2.1, and 5.2.2 provide information bounds for estimation of such functions), while other functions fail to be pathwise differentiable as a function of the distribution  $P$  of the observations (so that no information bound is available from the theory we present here, and rates slower than  $\sqrt{n}$  are likely). Theorem 1 says that if  $\dot{\mathbf{P}} = \overline{\mathbf{R}(\dot{\mathbf{I}})}$ , then (6) is *exactly* the condition which separates these two distinct situations.

The following example is a well-known problem in bioassay; see Ayer, Brunk, Ewing, Reid, and Silverman (1955) and the references therein. It has also been studied in considerable detail by Groeneboom (1988), (1991), who shows that the "right rate" for estimation of the df  $G$  is  $n^{1/3}$ ; also see Groeneboom and Wellner (1992).

#### Example 1. Indicator censoring model.

Suppose  $Y \sim G$  and  $Z \sim H$  on  $R$  are independent. We suppose that  $G$  and  $H$  are absolutely continuous with respect to Lebesgue measure on  $R$  with densities  $g$  and  $h$  respectively; without loss of generality we suppose that  $G$  and  $H$  are concentrated on  $[0, 1]$ . We observe  $X \equiv (Z, 1_{\{Y \leq Z\}}) \equiv (Z, \Delta)$ . Thus the density of  $X$  is

$$(17) \quad p(z, \Delta) = G(z)^\Delta \overline{G}(z)^{1-\Delta} h(z) \quad \text{for } \Delta \in \{0, 1\}, z \in [0, 1],$$

where  $\overline{G} \equiv 1 - G$ . In view of proposition A.5.5, this model is Hellinger-differentiable with respect to  $g^{1/2}$  and  $h^{1/2}$  in the sense of (3) with score operators  $\dot{\mathbf{I}}_1$  and  $\dot{\mathbf{I}}_2$  given by

$$(18) \quad \begin{aligned} \dot{\mathbf{I}}_1 a(z, \Delta) &= E(a(Y) \mid X = (z, \Delta)) \\ &= \Delta \frac{\int_0^z a dG}{G(z)} + (1 - \Delta) \frac{\int_z^1 a dG}{\overline{G}(z)}, \\ \dot{\mathbf{I}}_2 b(z, \Delta) &= b(z) \end{aligned}$$

for  $a \in \dot{\mathbf{G}} = L_2^0(G)$  and  $b \in \dot{\mathbf{H}} = L_2^0(H)$ . Note that (h) of the proof of proposition A.1.8 yields the boundedness of  $\dot{\mathbf{I}}_1$ . Thus,  $\dot{\mathbf{I}}: \dot{\mathbf{G}} \times \dot{\mathbf{H}} \rightarrow \dot{\mathbf{P}}$  defined via  $\dot{\mathbf{I}}(a, b) = \dot{\mathbf{I}}_1 a + \dot{\mathbf{I}}_2 b$  is bounded. By straightforward calculation

$$(19) \quad \begin{aligned} \dot{\mathbf{I}}_1^T \alpha(\dot{y}) &= E(\alpha(Z, \Delta) \mid Y = y) \\ &= \int_y^1 \alpha(z, 1) dH(z) + \int_0^y \alpha(z, 0) dH(z) \\ &= \int_0^1 \alpha(z, 1) dH(z) + \int_0^y [\alpha(z, 0) - \alpha(z, 1)] dH(z) \end{aligned}$$

since  $\alpha \in L_2^0(P)$ , and furthermore

$$(20) \quad \begin{aligned} \dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1 a(y) &= \int_0^1 K(y,t) a(t) dG(t), \\ \dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_2 b &= 0, \end{aligned}$$

where

$$(21) \quad K(y,t) = \int_{t \vee y}^1 \frac{1}{G(z)} dH(z) + \int_0^{t \wedge y} \frac{1}{\overline{G}(z)} dH(z) - 1.$$

Then  $\int K(y,t) dG(y) = 0$ , and straightforward but tedious computation shows that

$$\begin{aligned} &\int_0^1 \int_0^1 K^2(s,t) dG(s) dG(t) \\ &= 2 \int_0^1 \int_0^v \frac{G(u)(1-G(v))}{G(v)(1-G(u))} dH(u) dH(v) \leq 1. \end{aligned}$$

Thus  $\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1$  is a Hilbert-Schmidt operator. In particular it is a compact operator, and it follows that  $(\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1)^{-1}$  does not exist as a bounded operator; see, e.g., Rudin (1973, theorem 4.18(e), page 98). By corollary 2.B this shows that  $\mathbf{R}(\dot{\mathbf{I}}_1)$  is not closed since  $\mathbf{N}(\dot{\mathbf{I}}_1) = \{0\}$ . Because, by (20),  $\mathbf{R}(\dot{\mathbf{I}}_1) \perp \mathbf{R}(\dot{\mathbf{I}}_2)$ , this implies that  $\mathbf{R}(\dot{\mathbf{I}})$  is not closed and hence that corollary 2.A does not apply.

We therefore consider estimation of the specific function  $G$ , the distribution function corresponding to  $g$ :

$$(22) \quad v(P_{(g,h)}) = \psi(g) = \int_0^1 g dI \equiv G,$$

where  $I$  is Lebesgue measure on  $[0,1]$ . Then  $\dot{\psi}: \dot{\mathbf{G}} \times \dot{\mathbf{H}} \rightarrow \mathbf{B} \equiv C[0,1]$  is given by

$$(23) \quad \dot{\psi}(a,b)(t) = \int_0^1 [1_{[0,t]}(y) - G(t)] a(y) dG(y)$$

for  $a \in \dot{\mathbf{G}}_T = L_2^0(G)$ ,  $b \in \dot{\mathbf{H}} = L_2^0(H)$ ; recall example 5.3.1. Hence for  $b^* = \pi_t$ ,  $\dot{\psi}^T b^* \in \dot{\mathbf{G}} \times \dot{\mathbf{H}}$  is given by

$$\dot{\psi}^T \pi_t(y,z) = (1_{[0,t]}(y) - G(t), 0)$$

a discontinuous function of  $y$ . On the other hand, all the functions in  $\mathbf{R}(\dot{\mathbf{I}}_1) \subset \dot{\mathbf{G}}$  are continuous (since  $H$  is continuous) as can be easily seen from (19). Therefore  $\dot{\psi}^T \pi_t \notin \mathbf{R}(\dot{\mathbf{I}}^T)$ , that is, (6) fails to hold for this  $\psi$ , and hence, by theorem 1,  $v$  defined by (22) is *not* pathwise differentiable at any  $P \in \mathbf{P}$ ; see also corollary 5.5.1.

What functions  $v$  are differentiable in this model? It is possible to give a characterization for linear functionals  $\int f(y) dG(y)$  where  $f$  is a fixed given function. Differentiability essentially corresponds to the existence of  $\beta(Z, \Delta) \in L_2(P)$  such that

$$f(Y) = E(\beta(Z, \Delta) | Y) \quad \text{a.s.}$$

Consider such a functional

$$(24) \quad v(P_g) \equiv \psi(g) = \int_0^1 f(y)g(y) dy = E\beta(Z, \Delta).$$

Then  $\dot{\psi}_1 : L_2^0(G) \rightarrow R$  is given by

$$(25) \quad \dot{\psi}_1 a = \int_0^1 f(y)a(y) dG(y),$$

so  $\dot{\psi}_1^T : R \rightarrow L_2^0(G)$  becomes

$$(\dot{\psi}_1^T 1)(y) = f(y) - Ef(Y),$$

and it is easily seen that  $\dot{\psi}_1^T 1$  is in the range of  $\dot{\mathbf{i}}_1^T$  since

$$\begin{aligned} f(y) - Ef(Y) &= E(\beta(Z, \Delta) | Y=y) - E\beta(Z, \Delta) \\ &= \dot{\mathbf{i}}_1^T(\beta - E\beta)(y). \end{aligned}$$

Consequently, by theorem 1.B,  $v$  is pathwise differentiable. In view of

$$E\left(\frac{1-\Delta}{h(Z)} | Y=y\right) = y,$$

the mean functional with  $f$  the identity is pathwise differentiable if  $(1-\Delta)/h(Z) \in L_2(P)$ , i.e., if

$$(26) \quad \int_0^1 \frac{\bar{G}(z)}{h(z)} dz < \infty.$$

Hence if  $h$  is known and (26) holds,

$$\hat{v}_n \equiv \frac{1}{n} \sum_{i=1}^n \frac{1-\Delta_i}{h(Z_i)}$$

is an asymptotically linear estimator of  $v(P)$ , the mean of  $G$ , with influence function

$$\psi(z, \Delta) = \frac{1-\Delta}{h(z)} - v(P).$$

Furthermore, by the last expression of (19),

$$\begin{aligned} (27) \quad \mathbf{N}(\dot{\mathbf{i}}_1^T) &= \{\alpha(z, \Delta) \in L_2^0(P) : \alpha(z, 1) = \alpha(z, 0) \text{ a.e. } H, \\ &\quad \int \alpha(z, 1) dH(z) = 0\} \\ &= L_2^0(H). \end{aligned}$$

Hence by (A.1.17)

$$(28) \quad \overline{\mathbf{R}(\dot{\mathbf{i}}_1^T)} = \mathbf{N}(\dot{\mathbf{i}}_1^T)^\perp = \{\alpha \in L_2^0(P) : E(\alpha(Z, \Delta) | Z) = 0 \text{ a.s.}\},$$

and to calculate the efficient influence function  $\tilde{\mathbf{I}}$  we can simply project  $\psi$  above onto  $\overline{\mathbf{R}(\dot{\mathbf{i}}_1^T)} = \mathbf{N}(\dot{\mathbf{i}}_1^T)^\perp$  given in (28). A simple calculation yields

$$(29) \quad \tilde{\mathbf{I}}(z, \Delta) = \psi(z, \Delta) - E(\psi(Z, \Delta) | Z=z) = \frac{G(z) - \Delta}{h(z)},$$

which satisfies  $\dot{\mathbf{I}}_1^T \tilde{\mathbf{I}}(y) = (\dot{\mathbf{I}}_1^T \tilde{\mathbf{I}}_v 1)(y) = y - EY = (\dot{\psi}_1^T 1)(y)$ , i.e., (7) under  $h$  known; see also (13). Thus under (26) the mean functional  $v$  of (24) with  $f$  the identity is differentiable, and

$$(30) \quad E(\tilde{\mathbf{I}}^2) = \int_0^1 \frac{G(z)(1-G(z))}{h(z)} dz.$$

This information bound is achieved by the mean of the Ayer, et al. (1955) estimator of  $G$ ; see Groeneboom (1988), (1991), and Groeneboom and Wellner (1992).  $\square$

## 5.5 THE "CALCULUS" OF EFFICIENT SCORE AND INFLUENCE OPERATORS

With the methods and techniques of sections 5.2 and 5.4 now in hand, we can systematically develop a "calculus" of efficient score and influence operators. By keeping the geometric perspective of section 2.4, we extend the formulas developed there for finite-dimensional parameter spaces to the infinite-dimensional parameter spaces which occur in nonparametric and semiparametric models. The basic results of Van der Vaart in section 4 will be used repeatedly. Some of the formulas obtained extend and elaborate on those of Begun, et al. (1983). All the proofs are given together at the end of the section.

### *Score, Influence, and Information Operators for Parametrizations via a Hilbert Space*

As in section 4, suppose that  $\mathbf{P} = \{P_g : g \in \mathbf{G}\}$  where  $\mathbf{G}$  is a subset of a Hilbert space  $(\mathbf{H}, \langle \cdot, \cdot \rangle)$ . We begin with some terminology and definitions. If the model  $\mathbf{P}$  is Hellinger-differentiable in the sense of (5.4.3), we called  $\dot{\mathbf{I}} : \dot{\mathbf{G}} \rightarrow L_2^0(P_g)$  the *score operator for  $g$* . Thus  $\dot{\mathbf{I}}^T : L_2^0(P_g) \rightarrow \dot{\mathbf{G}}$ , and  $\dot{\mathbf{I}}^T \dot{\mathbf{I}} : \dot{\mathbf{G}} \rightarrow \dot{\mathbf{G}}$ . In analogy with the parametric case which we explore in more detail below),  $\dot{\mathbf{I}}^T \dot{\mathbf{I}}$  is called the *information operator for  $g$* .

When  $\mathbf{N}(\dot{\mathbf{I}}) = \{0\}$  and  $\mathbf{R}(\dot{\mathbf{I}})$  is closed, so that  $(\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{-1}$  exists by corollary 5.4.2.B, we can express the efficient influence operator  $\tilde{\mathbf{I}}_v : \mathbf{B}^* \rightarrow \dot{\mathbf{P}}$  from (5.2.8) for any function  $v(P_g) = \psi(g)$  with  $\psi : \mathbf{G} \rightarrow \mathbf{B}$  differentiable, as (cf. (5.4.16))

$$(1) \quad \tilde{\mathbf{I}}_v b^* = \dot{\mathbf{I}}(\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{-1} \dot{\psi}^T b^* \equiv \tilde{\mathbf{I}}_g \dot{\psi}^T b^*, \quad \text{for } b^* \in \mathbf{B}^*,$$

where

$$(2) \quad \tilde{\mathbf{I}}_g \equiv \dot{\mathbf{I}}(\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{-1} : \dot{\mathbf{G}} \rightarrow \dot{\mathbf{P}}$$

is called the *efficient influence operator for  $g$* . Thus we can view the calculation of the efficient influence operator  $\tilde{\mathbf{I}}_v$  for  $v$  as proceeding in two steps: (1) calculation of the efficient influence operator  $\tilde{\mathbf{I}}_g$  for  $g$ ; and (2) calculation of  $\dot{\psi}$  and  $\dot{\psi}^T$  for the specific differentiable function  $\psi$ .

Finally, since the information bound operator for  $v$  is  $\tilde{\mathbf{I}}_v^T \tilde{\mathbf{I}}_v : \mathbf{B}^* \rightarrow \mathbf{B}^{**}$  with  $\tilde{\mathbf{I}}_v$  given by (1), we have



$$(3) \quad \tilde{I}_v^T \tilde{I}_v = \dot{\psi} \tilde{I}_g^T \tilde{I}_g \dot{\psi}^T,$$

where we call

$$(4) \quad \tilde{I}_g^T \tilde{I}_g = (\dot{I}^T \dot{I})^{-1} : \dot{G} \rightarrow \dot{G}$$

the *inverse information operator* for  $g$ . As we will see more explicitly below, formula (1) is a generalization of the parametric model formula (2.3.2), while (3) similarly generalizes (2.3.1).

To fix ideas, we apply theorem 5.4.1 and the preceding notation to some cases of interest beginning with:

*Parametric Models*

Here  $G \subset R^k$ . To make the correspondence with the notation of chapter 2, we let  $G \equiv \Theta$  and  $g \equiv \theta$ . We suppose that  $P$  is a regular parametric model,

$$P = \{P_\theta : \theta \in \Theta\}.$$

Fix  $\theta_0 \in \Theta \subset R^k$ . Then  $\dot{G}^0 = \dot{G} = R^k$  and  $\dot{P} = [\dot{i}_1, \dots, \dot{i}_k]$  where  $\dot{i}_1, \dots, \dot{i}_k$  are the derivatives of  $I$  with respect to the components of  $\theta$  as in section 2.1.

**Notational correspondence:** We now *change* notation from chapters 2 to 4, and let  $\underline{i} \equiv (\dot{i}_1, \dots, \dot{i}_k)$ , a *row vector* (recall that in chapters 2–4  $\dot{I}$  and  $\tilde{I}$  were column vectors of  $L_2^0(P_0)$  functions). With this minor change, the general formulas in section 5.4 and (1)–(4) above specialize nicely in this case to yield formulas which are familiar from chapters 2 and 3.

**Reminder:** In the following, we use the superscript  $T$  to denote *both* the transpose of a vector or vector of functions in  $R^k$  or  $R^m$  and the adjoint of an operator. These two distinct uses should be kept clearly in mind in reading the rest of this example.

With these conventions,  $\dot{I} : \dot{G} \rightarrow \dot{P} \subset L_2^0(P_0)$  defined in (5.4.3) is given by

$$(5) \quad \dot{I}(t) = \underline{i} t = \sum_{i=1}^k t_i \dot{i}_i \quad \text{for } t \in \dot{G} = R^k,$$

$\dot{I}^T : L_2^0(P_0) \rightarrow \dot{G} = R^k$  is given by

$$(6) \quad \dot{I}^T h = \langle \underline{i}^T, h \rangle_0 = E_0(\underline{i}^T h) \quad \text{for } h \in L_2^0(P_0),$$

and  $\dot{I}^T \dot{I} : \dot{G} \rightarrow \dot{G}$  is

$$(7) \quad \dot{I}^T \dot{I}(t) = \langle \underline{i}^T, \underline{i} \rangle_0 t = I(\theta_0) t.$$

Also note that  $\tilde{I}_g : \dot{G} \rightarrow \dot{P}$  as defined in (2) is given by

$$(8) \quad \begin{aligned} \tilde{I}_g(t) &= \dot{I}(\dot{I}^T \dot{I})^{-1}(t) = \underline{i}_1 \langle \underline{i}^T, \underline{i} \rangle_0^{-1} t \\ &= \underline{i} I^{-1}(\theta_0) t = \tilde{I} t \end{aligned}$$

for  $t \in \dot{\mathbf{G}} = R^k$  where  $\tilde{\mathbf{I}} \equiv \dot{\mathbf{I}}^{-1}(\theta_0)$  is the efficient influence function for  $\theta$  from (2.3.2), but now viewed as a row vector of  $L_2^0(P_0)$  functions. Thus (8) is just the efficient influence function for estimation of  $t^T v(P_\theta) = t^T \theta$ . More generally, consider estimation of  $v(P_\theta) = \psi(g) = q(\theta)$  with  $q : \Theta \rightarrow R^m$  differentiable as in section 2.3. Then  $\dot{\psi} : \dot{\mathbf{G}} = R^k \rightarrow R^m$  is given by  $\dot{\psi} b = \dot{q}(\theta_0) b$  and  $\dot{\psi}^T : R^m \rightarrow R^k$  by  $\dot{\psi}^T a = \dot{q}^T(\theta_0) a$ . Consequently, the efficient influence operator (1) becomes

$$\tilde{\mathbf{I}}_v a = \tilde{\mathbf{I}}_g \dot{\psi}^T a = \dot{\mathbf{I}}^{-1}(\theta_0) \dot{q}^T(\theta_0) a$$

and  $\dot{\mathbf{I}}^{-1}(\theta_0) \dot{q}^T(\theta_0)$  is the efficient influence function for  $v$  from (2.3.2), but written as a row vector now. Similarly, the information bound operator (3) from  $R^m$  to  $R^m$  becomes

$$\tilde{\mathbf{I}}_v^T \tilde{\mathbf{I}}_v = \dot{\psi} \tilde{\mathbf{I}}_g^T \tilde{\mathbf{I}}_g \dot{\psi}^T = \dot{q}(\theta_0) \dot{\mathbf{I}}^{-1}(\theta_0) \dot{q}^T(\theta_0),$$

the information bound for  $v$  from (2.3.1). The inverse information covariance function from (5.2.23) for estimation of  $v(P_\theta) = \theta$  becomes

$$(9) \quad \dot{I}^{-1}(s, t) = \langle s^T, (\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{-1} t \rangle = s^T \dot{I}^{-1}(\theta_0) t$$

for  $(s, t) \in \dot{\mathbf{G}} \times \dot{\mathbf{G}}$ , which is just the covariance between efficient estimates of  $s^T \theta$  and  $t^T \theta$ .

We now interpret these formulas from the point of view of theorem 5.4.1 and its corollaries. The range of  $\dot{\mathbf{I}}$  given by (5) is finite-dimensional, and hence closed. Thus corollary 5.4.2.C applies, and  $\mathbf{R}(\dot{\mathbf{I}}^T) = \mathbf{R}(\dot{\mathbf{I}}^T \dot{\mathbf{I}})$ . From this it follows that

$$(10) \quad \mathbf{R}(\dot{\mathbf{I}}^T) = \mathbf{R}(\dot{\mathbf{I}}^T \dot{\mathbf{I}}) = I(\theta_0) R^k \subset R^k$$

with equality if  $I(\theta_0)$  is nonsingular. Thus, since in a regular parametric model  $I(\theta_0)$  is nonsingular, every  $v(P_\theta) = \psi(g) = q(\theta)$  as above is pathwise differentiable as it should be. Note that if the map  $\theta \rightarrow I(\theta)$  is smooth, but  $I(\theta_0)$  is singular, then by theorem 5.4.1,  $v(P_\theta)$  is pathwise differentiable at  $\theta_0$  and hence still potentially regularly estimable if and only if

$$\dot{q}^T(\theta_0) R^m = \mathbf{R}(\dot{\psi}^T) \subset \mathbf{R}(\dot{\mathbf{I}}^T) = \mathbf{R}(\dot{\mathbf{I}}^T \dot{\mathbf{I}}) = I(\theta_0) R^k,$$

i.e., if and only if the rows of the matrix  $\dot{q}(\theta_0)$  are in the column space of  $I(\theta_0)$ .

### Models with Composite $g$

In the rest of this section we consider models with  $\mathbf{G} = \mathbf{G}_1 \times \mathbf{G}_2$  with  $\mathbf{G}_i \subset \mathbf{H}_i$  Hilbert spaces with inner products  $\langle \cdot, \cdot \rangle_i$ ,  $i = 1, 2$ . These are the generalization to abstract parameter spaces of the parametric models we considered in section 2.4.

We suppose that

$$\mathbf{P} = \{P_g : g = (g_1, g_2), g_i \in \mathbf{G}_i \subset \mathbf{H}_i, i = 1, 2\},$$

where both  $\mathbf{G}_1$  and  $\mathbf{G}_2$  may be infinite-dimensional. We assume that there are bounded linear operators  $\dot{\mathbf{I}}_1, \dot{\mathbf{I}}_2$  with  $\dot{\mathbf{I}}_i : \dot{\mathbf{G}}_i \rightarrow L_2^0(P_0)$  such that

$$(11) \quad s(g_1(\eta), g_2(\eta)) - s(g_1, g_2) - \eta \frac{1}{2} (\dot{\mathbf{i}}_1 h_1 + \dot{\mathbf{i}}_2 h_2) s(g) = o(\eta)$$

in  $L_2(\mu)$  whenever  $\eta \rightarrow 0$  and

$$g_i(\eta) = g_i + \eta h_i + o(\eta) \quad \text{in } \mathbf{H}_i, \quad i = 1, 2.$$

An example of such a model with both  $\mathbf{G}_1$  and  $\mathbf{G}_2$  infinite-dimensional is the copula model of example 4.7.4 with  $\theta$  known and  $g$  and  $h$  unknown; as shown in proposition 4.7.6,  $\dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2$  is closed in this example under mild hypotheses.

Now we are in the framework of section 5.4 with  $\dot{\mathbf{G}} = \dot{\mathbf{G}}_1 \times \dot{\mathbf{G}}_2$  and

$$(12) \quad \dot{\mathbf{i}}(h_1, h_2) = \dot{\mathbf{i}}_1 h_1 + \dot{\mathbf{i}}_2 h_2.$$

Since

$$\begin{aligned} \langle \dot{\mathbf{i}}(h_1, h_2), \alpha \rangle_0 &= \langle \dot{\mathbf{i}}_1 h_1, \alpha \rangle_0 + \langle \dot{\mathbf{i}}_2 h_2, \alpha \rangle_0 \\ &= \langle h_1, \dot{\mathbf{i}}_1^T \alpha \rangle_1 + \langle h_2, \dot{\mathbf{i}}_2^T \alpha \rangle_2, \end{aligned}$$

for  $\alpha \in L_2^0(P_0)$ , it follows that

$$(13) \quad \dot{\mathbf{i}}^T \alpha = (\dot{\mathbf{i}}_1^T \alpha, \dot{\mathbf{i}}_2^T \alpha)^T \in \dot{\mathbf{G}}_1 \times \dot{\mathbf{G}}_2 \quad \text{for } \alpha \in L_2^0(P).$$

In analogy with (3.4.2) we define, for  $i = 1, 2$ , and  $j \neq i$ ,

$$(14) \quad \mathbf{I}_i^* \equiv (I - P_j) \dot{\mathbf{i}}_i: \dot{\mathbf{G}}_i \rightarrow \dot{\mathbf{P}},$$

where  $P_j \equiv \Pi_0(\cdot | \dot{\mathbf{P}}_j)$  and  $I$  is the identity.

**Corollary 1.** Let (A1) and (A2) of section 5.4 hold, and  $\overline{\mathbf{R}(\dot{\mathbf{i}})} = \dot{\mathbf{P}}$ ,  $\mathbf{R}(\dot{\mathbf{i}}_2) = \dot{\mathbf{P}}_2$ . Suppose that  $v(P_g) = \chi(g_1)$  where  $\chi: \mathbf{G}_1 \rightarrow \mathbf{B}$  is pathwise differentiable such that (A3) of section 5.4 holds.

A. If

$$(15) \quad \mathbf{R}(\dot{\chi}^T) \subset \mathbf{R}(\mathbf{I}_1^{*T}),$$

then  $v(P_g)$  is pathwise differentiable. In this case, the efficient influence operator  $\tilde{\mathbf{I}}_v$  for  $v$  is uniquely determined by the system of equations

$$(16) \quad \begin{aligned} \dot{\chi}^T b^* &= \mathbf{I}_1^{*T} \tilde{\mathbf{I}}_v b^* \\ 0 &= \dot{\mathbf{i}}_2^{T \sim} \tilde{\mathbf{I}}_v b^* \end{aligned} \quad \text{for } b^* \in \mathbf{B}^*.$$

B. If

$$(17) \quad \mathbf{R}(\dot{\chi}^T) \subset \mathbf{R}(\dot{\mathbf{i}}_1^T)$$

and

$$(18) \quad \mathbf{R}(\mathbf{I}_2^{*T}) = \mathbf{R}(\dot{\mathbf{i}}_2^T),$$

then (15) and the conclusions of A hold.

**Remark 1.** Note that  $\mathbf{R}(\mathbf{I}_1^{*T}) \subset \mathbf{R}(\dot{\mathbf{i}}_1^T)$  since  $\mathbf{I}_1^{*T} = \dot{\mathbf{i}}_1^T (I - P_2)$ , and similarly with 1 and 2 interchanged.

**Remark 2.** If both  $N(I_2^*) = \{0\}$  and  $R(I_2^*)$  is closed, then  $R(I_2^{*T}) = R(\dot{i}_2^T)$  is closed. To see this, note that

$$\{0\} = N(I_2^*) \supset N(\dot{i}_2)$$

since  $\dot{i}_2^* \equiv (I - P_1)\dot{i}_2$ . Since  $R(I_2^{*T})$  is closed if and only if  $R(\dot{i}_2^*)$  is closed (see proposition A.1.7.D), we have

$$R(I_2^{*T}) = \overline{R(\dot{i}_2^*)} = N(I_2^*)^\perp = \{0\}^\perp = \dot{G}_2$$

and hence, by remark 1,

$$\dot{G}_2 = R(I_2^{*T}) \subset R(\dot{i}_2^*) \subset \overline{R(\dot{i}_2^*)} \subset \dot{G}_2,$$

so  $R(I_2^{*T}) = R(\dot{i}_2^*) = \dot{G}_2$ .

Now we specialize these formulas by assuming that  $G_1 \equiv \Theta \subset R^k$ . This is the situation for the semiparametric models we have primarily considered. Suppose that

$$P = \{P_g : g = (g_1, g_2), g_i \in G_i \subset H_i, i = 1, 2\},$$

where  $G_1 \equiv \Theta \subset R^k$  and  $G_2 \subset H_2$  is infinite-dimensional. Assume that (11) holds with  $\dot{i}_1 h_1 = \underline{\dot{i}}_1 h_1$  where  $\underline{\dot{i}}_1 \in (L_2^0(P_g))^k$  and  $h_1 \in \overline{R^k}$  as in (5). In this situation we often have (and henceforth assume)  $R(\dot{i}) = \dot{P}$  and hence  $\dot{P} = \dot{P}_1 + \dot{P}_2 = [\underline{\dot{i}}_1] + \overline{R(\dot{i}_2)}$ ; note that we have used the latter condition in section 3.4. This fits into the framework of section 5.4 with  $\dot{G} = \dot{G}_1 \times \dot{G}_2 = R^k \times \dot{G}_2$  and

$$(19) \quad \dot{i}(h_1, h_2) = \underline{\dot{i}}_1 h_1 + \dot{i}_2 h_2.$$

Since by (6) and the definition of inner product in  $H_1 = R^k$ ,

$$\begin{aligned} \langle \dot{i}(h_1, h_2), \alpha \rangle_0 &= \langle \underline{\dot{i}}_1 h_1, \alpha \rangle_0 + \langle \dot{i}_2 h_2, \alpha \rangle_0 \\ &= h_1^T \langle \underline{\dot{i}}_1, \alpha \rangle_0 + \langle h_2, \dot{i}_2^T \alpha \rangle_2, \end{aligned}$$

for  $\alpha \in L_2^0(P_g)$ , it follows that

$$(20) \quad \dot{i}^T \alpha = (\langle \underline{\dot{i}}_1, \alpha \rangle_0, \dot{i}_2^T \alpha)^T \in \dot{G}_1 \times \dot{G}_2 = R^k \times H_2.$$

In this setting, as in (3.4.2) and (3.4.8) and in analogy with (14), it is natural to define the efficient score function  $\underline{\dot{i}}_1^*$  for estimation of  $g_1 = \theta$  by

$$(21) \quad \underline{\dot{i}}_1^* = \underline{\dot{i}}_1 - \Pi_0(\underline{\dot{i}}_1 | \dot{P}_2);$$

then the (efficient) information matrix from (3.4.11) for estimation of  $g_1 = \theta$  is the  $k \times k$  matrix

$$(22) \quad I_*(\theta) = E_0(\underline{\dot{i}}_1^{*T} \underline{\dot{i}}_1^*).$$

**Corollary 2.** Assume that the conditions of corollary 1 are satisfied and that the efficient information matrix  $I_*(\theta)$  defined in (22) is nonsingular.

A. Suppose that  $v(P_g) \equiv g_1 = \theta \in R^k \equiv B$ . Then  $v(P_g)$  is pathwise differentiable at  $g = (g_1, g_2)$  and  $\tilde{I}_v : R^k \rightarrow \dot{P}$  is the unique operator satisfying

$$(23) \quad \tilde{I}_v t = \underline{I}_1^* I_*^{-1}(\theta) t \quad \text{for } t \in \dot{G}_1^* = \dot{G}_1 = R^k.$$

B. Suppose  $v(P_g) = \chi(g_2)$  where  $\chi : G_2 \rightarrow B$  is pathwise differentiable with derivative  $\dot{\chi} : G_2 \rightarrow B$  such that (A3) of section 5.4 holds. Then  $v(P_g)$  is pathwise differentiable at  $g$  if and only if

$$(24) \quad R(\dot{\chi}^T) \subset R(\dot{I}_2^T).$$

When (24) holds, the efficient influence operator  $\tilde{I}_v : B^* \rightarrow \dot{P}$  for estimation of  $v$  is uniquely determined by the system of equations

$$(25) \quad \begin{aligned} 0 &= \langle \dot{I}_1^T, \tilde{I}_v b^* \rangle_0, \\ \dot{\chi}^T b^* &= \dot{I}_2^T \tilde{I}_v b^*. \end{aligned}$$

Note that the condition  $\overline{R(\dot{I}_1)} = \dot{P}$  of corollaries 1 and 2, which, in view of

$$R(\dot{I}_1) = R(\dot{I}_1) + R(\dot{I}_2) \subset \dot{P}_1 + \dot{P}_2 \subset \dot{P},$$

implies  $\dot{P} = \dot{P}_1 + \dot{P}_2$ , is not automatic, but in the situation of corollary 2,  $\overline{R(\dot{I}_1)} = \dot{P}_1$  is. Even the weaker condition  $\overline{R(\dot{I}_2)} = \dot{P}_2$  is not always known to hold, for instance for the Cox model and other transformation models of section 4.7.

Further interpretation of the corollary may be helpful: Part A says that  $\theta$  is potentially regularly estimable if  $I_*$  is nonsingular. We have noted in chapter 4 examples where this condition is violated and indeed  $\theta$  is unidentifiable even though  $R(\dot{I}_1^T) = R^k$  so that  $\theta$  is pathwise differentiable if  $g_2$  is assumed known. For instance, in the case control example 4.4.1 the parameters  $\mu_i$  are unidentifiable if the  $W_i$  are not known but became regularly estimable if they are known. Note that since A of corollary 1 does not hold in these cases, we must have  $R(\dot{I}_1^{*T}) \neq R(\dot{I}_2^T)$ . That is,  $P_2 \dot{I}_1 \neq \{0\}$  and there are also functions of  $g_2$  which are not pathwise differentiable if  $\theta$  is unknown (but still may be pathwise differentiable if  $\theta$  is known). To see this interplay, consider the general location model,  $X = \theta + \varepsilon$  where  $\varepsilon \sim G$  with  $G(|x| \leq M) = 1$  but  $G$  otherwise arbitrary. Then  $\theta$  is estimable if  $G$  is known but not if  $G$  is unknown, while  $E_G \varepsilon$  is estimable if  $\theta$  is known, but not otherwise; see example 1 below.

If we strengthen the hypothesis (24) in part B of corollary 2, an analogue of formula (5.4.16) holds.

**Corollary 3.** Let the assumptions of corollary 2.B be satisfied, and suppose that

$$(26) \quad R(\dot{\chi}^T) \subset R(\dot{I}_2^T \dot{I}_2).$$

Then, writing  $h_0^* = (\dot{I}_2^T \dot{I}_2)^{-1} \dot{\chi}^T b^*$  as a solution of  $\dot{\chi}^T b^* = \dot{I}_2^T \dot{I}_2 h_0^*$ , a solution of (25) is

$$(27) \quad \tilde{\mathbf{i}}_v b^* = \dot{\mathbf{i}}_2 (\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1} \dot{\chi}^T b^* - \dot{\mathbf{I}}_1^* \Gamma^{-1}(\theta) \langle \dot{\mathbf{I}}_1^T, \dot{\mathbf{i}}_2 (\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1} \dot{\chi}^T b^* \rangle_0.$$

If  $\mathbf{R}(\dot{\mathbf{i}}_2)$  is closed and  $\mathbf{N}(\dot{\mathbf{i}}_2) = \{0\}$  so that  $\mathbf{R}(\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2) = \mathbf{R}(\dot{\mathbf{i}}_2^T)$  and  $(\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1}$  exists, then  $(\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1}$  in (27) can be replaced by  $(\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1}$ , and the resulting formula was given by Begun et al. (1983). Since  $\dot{\mathbf{I}}_1^* \perp \mathbf{R}(\dot{\mathbf{i}}_2)$  in  $L_2^0(P_g)$ , the resulting inverse information covariance functional for  $v$  from (5.2.10) is calculated using adjoints and is given by

$$(28) \quad \Gamma_v^{-1}(b_1^*, b_2^*) = \langle \dot{\chi}^T b_1^*, (\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1} \dot{\chi}^T b_2^* \rangle_2 + \langle \dot{\chi}^T b_1^*, h_2^* \rangle_2 \Gamma^{-1}(\theta) \langle \dot{\chi}^T b_2^*, h_2^* \rangle_2,$$

where, by definition,

$$(29) \quad h_2^* \equiv (\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1} \dot{\mathbf{i}}_2^T \dot{\mathbf{I}}_1^*$$

is the unique  $k$ -column vector of elements of  $\dot{\mathbf{G}}_2$  such that  $\dot{\mathbf{I}}_1^* \equiv \dot{\mathbf{I}}_1 - \dot{\mathbf{i}}_2 h_2^{*T} \perp \dot{\mathbf{i}}_2 h_2$  for all  $h_2 \in \dot{\mathbf{G}}_2$  (componentwise). Note that  $(\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1}$  in (28) is the inverse information operator for  $g_2$  if  $g_1 \equiv \theta$  is known, and the second term is the increase in variance due to not knowing the parametric part  $g_1 = \theta$  of the model. Of course, (27) is essentially an extension of (2.4.12) (with 1 and 2 interchanged) to the semiparametric model setting.

Here are two examples to illustrate the formulas obtained in corollaries 2 and 3.

**Example 1. Estimation of a distribution function up to its mean.**

Suppose that the model is

$$\mathbf{P} = \{P \text{ on } (R, \mathcal{B}) : p(x; \theta, g) = g(x - \theta), \theta \in R, \\ g \text{ with mean } 0, I(G) < \infty\},$$

where the  $g$ 's are densities with respect to Lebesgue measure, and consider estimation of the distribution function  $G$  corresponding to  $g \equiv g_2$ ,  $v(P_{\theta, g}) \equiv \chi(g) \equiv \int_{-\infty}^{\cdot} g \equiv G$ . In this situation, it follows from example 3.2.3

with  $\gamma(G) \equiv \int x dG(x)$  that we may take

$$\dot{\mathbf{G}}_2 = \{h \in L_2^0(G) : h \perp \dot{\gamma}\},$$

where  $\dot{\gamma}(x) = x - \gamma(G)$ . Furthermore, (11) holds with

$$\dot{\mathbf{i}}_1 \equiv -\frac{g'}{g}(\cdot - \theta) \equiv \psi(\cdot - \theta)$$

and

$$\dot{\mathbf{i}}_2 h_2 = h_2(\cdot - \theta) \quad \text{for } h_2 \in \dot{\mathbf{G}}_2.$$

Then

$$\dot{\mathbf{i}}_2^T b = b(\cdot + \theta) - \frac{\langle b(\cdot + \theta), \dot{\gamma} \rangle_2}{\|\dot{\gamma}\|_2^2} \dot{\gamma}$$

$$= b(\cdot + \theta) - \frac{\text{Cov}_G(b(X + \theta), X)}{\text{Var}_G(X)} \dot{\gamma}$$

and

$$\dot{I}_2^T \dot{I}_2 = \text{identity on } \dot{G}_2.$$

We therefore compute from (29)

$$h_2^* = (\dot{I}_2^T \dot{I}_2)^{-1} \dot{I}_2^T \dot{I}_1 = \psi - \frac{\langle \psi, \dot{\gamma} \rangle_2}{\text{Var}_G(X)} \dot{\gamma} = \psi - \frac{1}{\text{Var}_G(X)} \dot{\gamma}$$

since  $\langle \psi, \dot{\gamma} \rangle_2 = - \int x g'(x) dx = 1$ ,

$$\dot{I}_1^* = \dot{I}_1 - \dot{I}_2 h_2^* = - \frac{1}{\text{Var}_G(X)} \dot{\gamma},$$

and from (22)

$$I_*^{-1}(\theta) = \text{Var}_G(X).$$

Since

$$\dot{\chi}^T \pi_t = 1_{(-\infty, t]} - G(t) - \frac{\langle 1_{(-\infty, t]}, \dot{\gamma} \rangle_2}{\|\dot{\gamma}\|_2^2} \dot{\gamma}$$

by example 5.3.1 and (5.2.25) and (5.2.7), it follows that the inverse information covariance function for estimation of  $v(P_{\theta, g}) = \chi(g) = G$  is, from (28), for  $s, t \in R$ ,

$$\begin{aligned} I_V^{-1}(s, t) &= \langle \dot{\chi}^T \pi_s, \dot{\chi}^T \pi_t \rangle_2 + \langle \dot{\chi}^T \pi_s, h_2^* \rangle_2 \text{Var}_G(X) \langle \dot{\chi}^T \pi_t, h_2^* \rangle_2 \\ (30) \quad &= G(s) \wedge G(t) - G(s)G(t) + g(s)g(t) \text{Var}_G(X) \\ &\quad + g(s) \text{Cov}_G(1_{(-\infty, t]}(X), X) + g(t) \text{Cov}_G(1_{(-\infty, s]}(X), X) \end{aligned}$$

since  $\int_{-\infty}^t \psi dG = -g(t)$ . This fits with our intuition, since the natural estimator of  $\bar{G}$  in this example is simply the empirical df  $F_n$  shifted by the sample mean  $\bar{X}$ :

$$\hat{G}_n \equiv F_n(\cdot + \bar{X}) = n^{-1} \sum_{i=1}^n 1_{[X_i - \bar{X} \leq \cdot]}.$$

This estimator has influence function

$$1_{[x - \theta \leq t]} - G(t) + g(t)(x - \theta),$$

and hence achieves the bound provided by (30). □

**Example 2. Estimation of the distribution function  $G$  corresponding to the baseline cumulative hazard function  $\Lambda$  in the Cox model.**

Now consider the Cox proportional hazards model of examples 3.4.2, 4.7.1, and 6.7.1.A. In example 6.7.1.A we will show that the inverse information

covariance function for estimation of  $\kappa(P_{\theta, g}) \equiv G$  when  $\theta$  is known is given by

$$I_{\kappa}^{-1}(s, t; \mathbf{P}_2) \equiv \bar{G}(s)\bar{G}(t) \int_0^{s \wedge t} \frac{1}{S_0} d\Lambda.$$

This is just the first term in (28). When  $\theta \equiv \mathbf{v}$  is *also* unknown, we need to include the second term in (28). From the  $a^*$  calculated in 3.4.2, we have  $h_2^* = a^* = L(S_1/S_0)$ . Since  $R \circ L = \text{identity}$  by proposition A.1.8.B,

$$\begin{aligned} \frac{S_1}{S_0}(t) &= R \circ L\left(\frac{S_1}{S_0}\right)(t) = Rh_2^*(t) = h_2^*(t) - \frac{\int_t^{\infty} h_2^* dG}{\bar{G}(t)} \\ &= \frac{S_1}{S_0}(t) - \int_0^t \frac{S_1}{S_0} d\Lambda - \frac{\int_t^{\infty} h_2^* dG}{\bar{G}(t)}, \end{aligned}$$

so that  $\int_0^{\infty} h_2^* dG = -\bar{G}(t) \int_0^t (S_1/S_0) d\Lambda$ , and hence, using  $\int_0^{\infty} h_2^* dG = 0$ ,  $\dot{\chi}^T \pi_t = \dot{1}_{(0, t]} - G(t)$ ,

$$\langle \dot{\chi}^T \pi_t, h_2^* \rangle_2 = \int_0^t h_2^* dG = \bar{G}(t) \int_0^t \frac{S_1}{S_0} d\Lambda.$$

Thus from (28), the inverse information covariance function for estimation of  $G$  in the full model is

$$I_{\kappa}^{-1}(s, t; \mathbf{P}) = \bar{G}(s)\bar{G}(t) \left\{ \int_0^{s \wedge t} \frac{1}{S_0} d\Lambda + I_*^{-1} \int_0^s \frac{S_1}{S_0} d\Lambda \int_0^t \frac{S_1}{S_0} d\Lambda \right\},$$

where  $I_*$  is as given in (3.4.55). □

### Proofs

**Proof of corollary 1.** A. Here  $\psi(g_1, g_2) = \chi(g_1)$ , so  $\dot{\psi}(h_1, h_2) = \dot{\chi}(h_1)$  and  $\dot{\psi}^T b^* = (\dot{\chi}^T b^*, 0) \in \dot{\mathbf{G}}_1 \times \dot{\mathbf{G}}_2$ . Hence

$$(a) \quad \mathbf{R}(\dot{\psi}^T) = \mathbf{R}(\dot{\chi}^T) \times \{0\}.$$

On the other hand,  $\dot{\mathbf{P}}_2^{\perp} = \mathbf{R}(\dot{\mathbf{i}}_2)^{\perp} = \mathbf{N}(\dot{\mathbf{i}}_2^T)$  by  $\overline{\mathbf{R}(\dot{\mathbf{i}}_2)} = \dot{\mathbf{P}}_2$  and (A.1.17), so

$$\begin{aligned} (b) \quad \mathbf{R}(\dot{\mathbf{i}}^T) &\supset \{(\dot{\mathbf{i}}_1^T \alpha, 0) : \alpha \in \dot{\mathbf{P}}_2^{\perp}\} \\ &= \{(\dot{\mathbf{i}}_1^T (I - P_2) \alpha, 0) : \alpha \in \dot{\mathbf{P}}_2^{\perp}\} \\ &= \mathbf{R}(\dot{\mathbf{i}}_1^{*T}) \times \{0\} \end{aligned}$$

since  $\dot{\mathbf{i}}_1^{*T} = \dot{\mathbf{i}}_1^T (I - P_2)$  by (14). But (a), (b), and (15) imply that  $\mathbf{R}(\dot{\psi}^T) \subset \mathbf{R}(\dot{\mathbf{i}}^T)$ , and therefore pathwise differentiability follows from theorem 5.4.1.B.

To see that  $\tilde{\mathbf{I}}_{\mathbf{v}}$  satisfies (16), note that by (13) and  $\dot{\psi}^T b^* = (\dot{\chi}^T b^*, 0)$ , equation (5.4.7) is equivalent to

$$(c) \quad \dot{\chi}^T b^* = \dot{\mathbf{i}}_1^T \tilde{\mathbf{I}}_{\mathbf{v}} b^*$$



and

$$(d) \quad 0 = \dot{\mathbf{i}}_2^T \tilde{\mathbf{l}}_v b^* .$$

Thus  $\tilde{\mathbf{l}}_v b^* \in N(\dot{\mathbf{i}}_2^T) = R(\dot{\mathbf{i}}_2^T)^\perp = \dot{\mathbf{P}}_2^\perp$  by (d), and hence  $(I - P_2)\tilde{\mathbf{l}}_v b^* = \tilde{\mathbf{l}}_v b^*$ . Substituting this into (c) yields (16) by definition of  $\dot{\mathbf{l}}_1^* = (I - P_2)\dot{\mathbf{l}}_1$ .

B. First, we will show

$$(e) \quad R(\dot{\mathbf{l}}^T) = R(\dot{\mathbf{l}}_1^T) \times R(\dot{\mathbf{i}}_2^T) .$$

Since the left-hand side is trivially contained in the right-hand side, it suffices to prove that the right-hand side is contained in the left. Let  $(a_1, a_2) \in R(\dot{\mathbf{l}}_1^T) \times R(\dot{\mathbf{i}}_2^T)$ . Then there exist  $\alpha_1, \alpha_2 \in L_2^0(P_g)$  such that  $a_i = \dot{\mathbf{l}}_i^T \alpha_i$ ,  $i = 1, 2$ . Note that by (18) and  $N(\dot{\mathbf{l}}_2^{*T}) \supset \dot{\mathbf{P}}_1$ , there exists a  $\beta_1 \in \dot{\mathbf{P}}_1^\perp$  satisfying

$$\dot{\mathbf{l}}_2^{*T} \beta_1 = a_2 - \dot{\mathbf{i}}_2^T \alpha_1 = \dot{\mathbf{i}}_2^T (\alpha_2 - \alpha_1) \in R(\dot{\mathbf{i}}_2^T) = R(\dot{\mathbf{l}}_2^{*T}) .$$

Since

$$\dot{\mathbf{l}}_2^{*T} \beta_1 = \dot{\mathbf{i}}_2^T (I - P_1) \beta_1 = \dot{\mathbf{i}}_2^T \beta_1 ,$$

we have both

$$a_2 = \dot{\mathbf{i}}_2^T (\alpha_1 + \beta_1)$$

and

$$a_1 = \dot{\mathbf{l}}_1^T (\alpha_1 + \beta_1)$$

using  $\beta_1 \in \dot{\mathbf{P}}_1^\perp \subset R(\dot{\mathbf{l}}_1^T)^\perp = N(\dot{\mathbf{l}}_1^T)$ . Hence  $(a_1, a_2) \in R(\dot{\mathbf{l}}^T)$ , and (e) holds.

Since  $R(\dot{\psi}^T) = R(\dot{\chi}^T) \times \{0\}$  by (a), the differentiability part of the claim follows from (17), (e), and theorem 5.4.1.B. Then equation (5.4.7) again yields (c) and (d), and hence, by the same argument as in the proof of A, (16) holds. But (16) implies (15).  $\square$

**Proof of corollary 2.** A. In this case  $\psi(g_1, g_2) = g_1 \equiv \theta$ , so  $\dot{\psi}(h_1, h_2) = h_1$  and  $\dot{\psi}^T(h_1^*) = (h_1^*, 0) \in \dot{\mathbf{G}}_1 \times \dot{\mathbf{G}}_2$ . Thus  $R(\dot{\psi}^T) = R^k \times \{0\}$ . Hence by (b) of the proof of corollary 1, by (14), (5), (6), (21), and (22), and by the non-singularity of  $I_*(\theta)$

$$\begin{aligned} (a) \quad R(\dot{\mathbf{l}}^T) &\supset R(\dot{\mathbf{l}}_1^{*T}) \times \{0\} \\ &= \{\dot{\mathbf{l}}_1^T (I - P_2) \alpha : \alpha \in \dot{\mathbf{P}}\} \times \{0\} \\ &\supset \{\dot{\mathbf{l}}_1^T (I - P_2) (I - P_2) \dot{\mathbf{l}}_1 t : t \in R^k\} \times \{0\} \\ &= \{I_*(\theta) t : t \in R^k\} \times \{0\} \\ &= R^k \times \{0\} = R(\dot{\psi}^T) . \end{aligned}$$

Hence  $\psi(P_g) = g_1 \equiv \theta$  is differentiable by theorem 5.4.1.B. It is easy to verify that  $\tilde{\mathbf{l}}_v$  of (23) satisfies (16) and maps  $R^k$  into  $\dot{\mathbf{P}}$ .

B. Now  $\psi(g_1, g_2) = \chi(g_2)$ , so  $\dot{\psi}(h_1, h_2) = \dot{\chi}(h_2)$  and

$$(b) \quad \dot{\psi}^T b^* = (0, \dot{\chi}^T b^*) \in \dot{G}_1 \times \dot{G}_2.$$

Thus

$$(c) \quad \mathbf{R}(\dot{\psi}^T) = \{0\} \times \mathbf{R}(\dot{\chi}^T).$$

In (a) we have shown that  $\mathbf{R}(I_1^{*T}) \supset R^k \supset \mathbf{R}(\dot{i}_1^T)$ . Consequently,  $\mathbf{R}(I_1^{*T}) = \mathbf{R}(\dot{i}_1^T)$ , an analogue of (18), holds and (e) of the proof of corollary 1 is valid:

$$(d) \quad \mathbf{R}(\dot{i}^T) = R^k \times \mathbf{R}(\dot{i}_2^T).$$

The first conclusion follows from (c), (d), and an application of theorem 5.4.1: thus  $v(P_g) = \chi(g_2)$  is pathwise differentiable if and only if (24) holds. The two formulas (25) follow immediately from (5.4.7), (b), and (20).  $\square$

**Proof of corollary 3.** It suffices to show that  $\tilde{I}_v$  given in (27) satisfies (25). But this follows easily using  $\langle \dot{i}_1^T, \underline{1}_1^* \rangle_0 = I_*(\theta)$  and  $\underline{1}_1^* \in (\dot{P}_2^\perp)^k = (\mathbf{R}(\dot{i}_2^\perp))^k = (\mathbf{N}(\dot{i}_2^T))^k$ .  $\square$

# 6 | Infinite-Dimensional Parameters: Further Examples

## 6.1 INTRODUCTION

The information bounds for infinite-dimensional parameters developed in chapter 5 were illustrated there with calculations for relatively simple models. In this chapter our goal is to find information bounds for the infinite-dimensional components of some of the models investigated in chapter 4, or at least explain why such bounds are not feasible in cases for which such bounds cannot exist.

Thus the situation in this chapter is the reverse of that in chapter 4. In the examples treated in chapter 4 the parametric part of the model was the parameter of interest and the infinite-dimensional parameter played the role of a nuisance parameter. In this chapter interest is focused on the infinite-dimensional parameters (or some pathwise differentiable functional thereof) of the same models—with the finite-dimensional parameters (if any) playing the role of nuisance parameters. When finite-dimensional nuisance parameters are present, the key decompositions presented in (5.5.27) and (5.5.28)—extending the basic decompositions given in (2.4.12) and (2.4.13) for the parametric case—are used repeatedly.

We begin in section 2 with models defined by constraints—finite or infinite-dimensional. Group models are treated in section 3. We derive information bounds for biased sampling models—including truncation and truncated regression models—in section 4. Mixture models present a special challenge for estimation of the infinite-dimensional (mixing distribution) component. Usually this is possible at  $\sqrt{n}$ -rates only for very particular special functionals of the mixing distribution. An explanation for this phenomena is given in section 5.

Missing data and censoring models are treated in section 6, with special attention to random censorship and generalizations thereof. Finally, information bounds are calculated for the infinite-dimensional (transformation) component of transformation models in section 7.

## 6.2 CONSTRAINED FAMILIES

Many interesting families  $\mathbf{P}$  are defined in terms of a system of constraints. The simple case of finitely many constraints was already introduced, and the

tangent spaces calculated, in section 3.2. Other very natural families are defined by infinite-dimensional systems of constraints. Examples of both types of such families will be given in this section.

### *Finitely Many Constraints*

Suppose that

$$(1) \quad \mathbf{P} \equiv \{P \ll \mu : \gamma_i(P) = 0, i = 1, \dots, r\},$$

where  $\gamma_i$  is pathwise differentiable with derivative  $\dot{\gamma}_i$  as a map from  $\mathbf{P}$  to  $R$  as in definition 3.3.1. Frequently the constraint functions  $\gamma_i(P)$  are defined by

$$(2) \quad \gamma_i(P) = E_P(T_i) - t_i = \int T_i dP - t_i,$$

where  $T_i$  is a measurable function of  $X$  and  $t_i$  is a fixed number. Then  $\gamma_i(P) = 0$  just means that the moments  $E(T_i)$  are assumed to be known and equal to  $t_i$ ,  $i = 1, \dots, r$ , in the constrained model (1).

If the functions  $P \rightarrow \dot{\gamma}_i(P)$  are continuous at  $P_0$ , then, as shown in example 3.2.3, the tangent space  $\dot{\mathbf{P}}$  of the model  $\mathbf{P}$  in (1) is

$$(3) \quad \dot{\mathbf{P}} = \{h \in L_2(P_0) : E_0(h) = 0, E_0(h \dot{\gamma}_i) = 0, i = 1, \dots, r\},$$

and the projection operator  $\Pi_0(\cdot | \dot{\mathbf{P}})$  is given, for  $h \in L_2(P_0)$ , by

$$(4) \quad \Pi_0(h | \dot{\mathbf{P}}) = h - \int h dP_0 - Cov_0(h, \dot{\gamma}^T) \Sigma_{\dot{\gamma}}^{-1} \dot{\gamma},$$

where we have assumed that  $\dot{\gamma} = (\dot{\gamma}_1, \dots, \dot{\gamma}_r)^T$  and

$$\Sigma_{\dot{\gamma}} \equiv Cov_0(\dot{\gamma}, \dot{\gamma}^T)$$

is nonsingular. The projection formula (4) makes the calculations required in section 5.2 straightforward, and, in fact, a calculation of exactly this kind was already carried out in example 5.3.2. Here is a simple example of such a calculation based on a general differentiable function  $\gamma$  with derivative  $\dot{\gamma}$ , and example 5.3.8.

#### **Example 1. Estimation of a constrained distribution $P$ .**

For a general sample space  $\mathbf{X}$ , let  $v(P)(f) \equiv \int f dP \equiv P(f)$  for  $f \in \mathbf{F}$  and  $P \in \mathbf{P}$  just as in example 5.3.8, but where now

$$(5) \quad \mathbf{P} \equiv \{P \ll \mu : \gamma(P) = 0\}$$

and  $\gamma : \mathbf{P} \rightarrow R^r$  is pathwise differentiable with derivative  $\dot{\gamma}$ . Just as before,  $v$  is pathwise differentiable in the unconstrained model with  $\dot{v}(h)$  given in (5.3.28), and by (5.2.8), (5.2.7), and (5.2.25) it suffices to calculate

$$\bar{\Gamma}_f \equiv \bar{\Gamma}(\pi_f) = \Pi_0(\dot{v}_f | \dot{\mathbf{P}})$$

where  $\dot{v}_f$  is given by (5.3.29). Thus by (4) the efficient influence function  $\tilde{\Gamma}$  from (5.2.11) is given by

$$(6) \quad \pi_f \tilde{\Gamma}(x) = \bar{\Gamma}_f(x) = f(x) - \int f dP_0 - Cov_0(f(X), \dot{\gamma}^T(X)) \Sigma_{\dot{\gamma}}^{-1} \dot{\gamma}(x).$$

It follows that the inverse information covariance function  $I_V^{-1} : F \times F \rightarrow R$  of (5.2.23) is

$$(7) \quad I_V^{-1}(f, g) = Cov_0(f(X), g(X)) - Cov_0(f(X), \dot{\gamma}^T(X)) \Sigma_V^{-1} Cov_0(\dot{\gamma}(X), g(X)),$$

and the limit process  $Z_0$  of theorem 5.2.1 is

$$(8) \quad Z_0(f) = \tilde{Z}_0(\tilde{I}_f) \quad \text{for } f \in F,$$

where  $\tilde{Z}_0$  is the  $P$ -Brownian bridge process  $Z_0$  of example 5.3.8 and  $\tilde{I}_f$  is given by (6).

In particular, if  $X = R^r$  and  $\gamma(P) = E_P X$ , the mean, then  $\dot{\gamma}(x) = x - E_P X$  and

$$I_V^{-1}(f, g) = Cov_0(f(X), g(X)) - Cov_0(f(X), X^T) \Sigma_X^{-1} Cov_0(X, g(X)).$$

The general version of this example with  $r$  (linear of the form (2)) constraints, together with the construction of estimates which achieve the bounds are contained in Sheehy (1987), (1988). Related calculations have been given by Haberman (1984), Levit (1975), and Koshevnik and Levit (1976).  $\square$

### Infinitely Many Constraints

Many interesting models can be viewed as constrained models with an infinite family of constraint functionals  $\{\gamma_u : u \in U\}$  rather than the finite constraint family  $\{\gamma_i : i = 1, \dots, r\}$  in (1). Here are several interesting examples:

**Example 2. Estimation of a bivariate distribution  $P$  with one known marginal.**

Suppose that  $X = R^2$  and write  $X = (U, V)$  for an observation from  $P \in \mathbf{P}$ ,

$$\mathbf{P} \equiv \{P \ll \mu : \gamma_u(P) = 0 \text{ for all } u \in R\},$$

where  $\mu$  is Lebesgue measure on  $R^2$  and

$$(9) \quad \gamma_u(P) \equiv E_P 1_{(-\infty, u]}(U) - F_0(u) = P(U \leq u) - F_0(u)$$

and  $F_0$  is a fixed, known df on  $R$ . Thus the family  $\mathbf{P}$  consists of all  $P$  on  $R^2$  with known marginal df  $F_0$  for  $U$ . As in example 3.2.3 it is plausible that

$$(10) \quad \dot{\mathbf{P}} = \{h \in L_2(P_0) : E_0(h) = 0, E_0(h \dot{\gamma}_u) = 0 \text{ for all } u \in R\}$$

with  $\dot{\gamma}_u \equiv 1_{(-\infty, u] \times R} - F_0(u)$ . That  $\dot{\mathbf{P}}$  is contained in the right side is easy, but further argument is required to prove the reverse inclusion. Now the space spanned by  $\{\dot{\gamma}_u : u \in R\}$  in  $L_2^0(P_0)$  is the subspace

$$L_2^{0U}(P_0) \equiv \{h \in L_2^0(P_0) : h \text{ is a function of } U \text{ only}\}.$$

Hence by proposition A.3.1

$$(11) \quad \Pi_0(h | \dot{\mathbf{P}}) = h - E(h | U) \quad \text{for } h \in L_2(P_0).$$

Now consider estimation of the bivariate df of  $(U, V)$ ,

$$(12) \quad v(P)(s, t) \equiv P(U \leq s, V \leq t) \equiv H(s, t), \quad (s, t) \in R^2.$$

It follows from example 5.3.8 that  $v$  is pathwise differentiable in the unconstrained model (consider the class of functions  $F \equiv \{1_{(-\infty, s] \times (-\infty, t]} : (s, t) \in R^2\}$ ) with

$$\dot{v}_{(s,t)}(u, v) = 1_{(-\infty, s] \times (-\infty, t]}(u, v) - H(s, t).$$

Hence by (5.2.8), (5.2.7), and (5.2.25) it suffices to calculate

$$\tilde{I}_{(s,t)} = \tilde{I}_v(\pi_{(s,t)}) = \Pi_0(\dot{v}_{(s,t)} | \dot{P}).$$

But by (11) the efficient influence function  $\tilde{I}$  from (5.2.11) is given by

$$(13) \quad \begin{aligned} \pi_{(s,t)} \tilde{I}(u, v) &= \tilde{I}_{(s,t)}(u, v) \\ &= 1_{(-\infty, s] \times (-\infty, t]}(u, v) - P(V \leq t | U = u) 1_{(-\infty, s]}(u), \end{aligned}$$

and the covariance function  $I_v^{-1}$  of (5.2.23) is, with  $x = (u, v)$ ,  $y = (s, t)$ ,

$$(14) \quad \begin{aligned} I_v^{-1}(x, y) &= E_0 \left( \{1_{(-\infty, u] \times (-\infty, v]}(U, V) - P(V \leq v | U) 1_{(-\infty, u]}(U)\} \right. \\ &\quad \cdot \left. \{1_{(-\infty, s] \times (-\infty, t]}(U, V) - P(V \leq t | U) 1_{(-\infty, s]}(U)\} \right) \\ &= H(u \wedge s, v \wedge t) \\ &\quad - E_0 \left( P(V \leq v | U) P(V \leq t | U) 1_{[U \leq u \wedge s]} \right). \end{aligned}$$

With no knowledge of the marginal df  $H(u, \infty)$  of  $U$ , it follows from (5.3.30) of example 5.3.8 that the limit process  $Z_0$  of an efficient estimator has covariance function

$$(15) \quad H(u \wedge s, v \wedge t) - H(u, v)H(s, t),$$

so that the difference in variance at a point  $x = (u, v)$  due to knowledge of the marginal distribution of  $U$  is

$$\begin{aligned} E_0 \{ P(V \leq v | U)^2 1_{[U \leq u]} \} - H^2(u, v), \\ = E_0 \{ E_0(1_{[V \leq v]} | U)^2 1_{[U \leq u]} \} - H^2(u, v), \end{aligned}$$

which takes values in  $[0, \frac{1}{4}]$  since

$$(16) \quad E \{ P^2(V \leq v | U) 1_{[U \leq u]} \} \leq E \{ P(V \leq v | U) 1_{[U \leq u]} \} \leq H(u, v)$$

and

$$(17) \quad E \{ (P(V \leq v | U) 1_{[U \leq u]})^2 \} \geq H^2(u, v).$$

Equality holds in (16) for all  $u$  and  $v$  if and only if the conditional distribution of  $V$  given  $U$  is degenerate, i.e., if and only if  $V$  is a function of  $U$ . Equality holds in (17) if and only if

$$(18) \quad P(V \leq v | U)1_{[U \leq u]} = H(u, v) \quad \text{a.s.,}$$

which in turn holds for all  $u$  and  $v$  if and only if  $U$  is a.s. constant. Thus if  $U$  is not degenerate, the reduction in variance achievable by taking advantage of knowledge of  $H(u, \infty) = F_0(u)$  is strictly positive. In view of (16) and (17), the ratio of (14) to (15) at a point  $(s, t) = (u, v)$ , the asymptotic relative efficiency, can range between 0 and 1, and the ratio may be zero if  $V$  is a function of  $U$ .

Note that if  $U$  and  $V$  are independent, then the difference between (15) and (14) at  $x = y = (u, v)$  is

$$\begin{aligned} P(U \leq u)P^2(V \leq v) - P^2(U \leq u)P^2(V \leq v) \\ = F_0(u)(1 - F_0(u))P^2(V \leq v) \\ = 0 \end{aligned}$$

only if  $U$  is degenerate; cf. (18).

More generally, for  $P$  indexed by a class of functions  $F$  with  $v(P)(f) \equiv \int f dP \equiv P(f)$  for  $f \in F$ ,

$$\dot{v}_f(u, v) = f(u, v) - \int f dP_0,$$

in the unconstrained model, and it follows that

$$\tilde{I}_f(u, v) = \dot{v}_f(u, v) - E_0(\dot{v}_f | U = u) = f(u, v) - E_0(f | U = u)$$

so that

$$(19) \quad I_v^{-1}(f, g) = E_0(f(X)g(X)) - E_0(E_0(f | U)E_0(g | U)).$$

It follows that the gain in variance at a "point"  $f$  is

$$Var_0(f(X)) - E_0 Var_0(f(X) | U) = Var_0(E_0(f(X) | U)),$$

which may vary between 0 and  $Var_0(f(X))$ .

One natural estimator for this model will be studied in example 7.5.8. □

**Example 3. Estimation of a bivariate distribution  $P$  with two known marginals.**

Now suppose that  $X = R^2$ , and consider estimation of

$$(20) \quad v(P)(s, t) \equiv P(U \leq s, V \leq t) = H(s, t), \quad (s, t) \in R^2,$$

as in example 2, but now  $P \in \mathbf{P}$  with

$$(21) \quad \mathbf{P} \equiv \left\{ P \ll \mu : \begin{array}{l} \gamma_u(P) = 0 \text{ for all } u \in R \\ \delta_v(P) = 0 \text{ for all } v \in R \end{array} \right\},$$

where  $\gamma_u$  is given by (9) and

$$(22) \quad \delta_v(P) \equiv E_P 1_{(-\infty, v]}(V) - G_0(v) = P(V \leq v) - G_0(v),$$

and  $G_0$  is a fixed, known df on  $R$ . Thus the family  $\mathbf{P}$  consists of all  $P$  on  $R^2$  with known marginal distributions  $F_0$  and  $G_0$  for  $U$  and  $V$  respectively. By

considerations similar to those of example 2, it is reasonable to suppose that the tangent space is

$$(23) \quad \begin{aligned} \dot{P} &= \{ h \in L_2^0(P_0) : \langle h, \dot{\gamma}_s \rangle = 0 \\ &= \langle h, \dot{\delta}_t \rangle \text{ for all } s, t \in R \}, \end{aligned}$$

where

$$\dot{\gamma}_s(u, v) = 1_{(-\infty, s] \times R}(u, v) - F_0(s),$$

and

$$\dot{\delta}_t(u, v) \equiv 1_{R \times (-\infty, t]}(u, v) - G_0(t).$$

Again, containment of  $\dot{P}$  in the right side is easily proved, but the reverse inclusion requires further argument. Bickel, Ritov, and Wellner (1991) prove the reverse inclusion under the further hypothesis of a strictly positive lower bound for the density of  $P$  with respect to the product measure  $F_0 \times G_0$ . We assume that equality holds in (23) and note that the functions  $\{\dot{\gamma}_s : s \in R\}$  and  $\{\dot{\delta}_t : t \in R\}$  span the subspaces

$$L_U \equiv \{ h \in L_2^0(P_0) : h \text{ is a function of } U \text{ only} \},$$

and

$$L_V \equiv \{ h \in L_2^0(P_0) : h \text{ is a function of } V \text{ only} \}.$$

Hence

$$(24) \quad \Pi_0(h | \dot{P}) = h - \Pi_0(h | L_U + L_V)$$

where  $L_U$  and  $L_V$  are not, in general, orthogonal subspaces of  $L_2^0(P_0)$ . This projection can be computed by the ACE algorithm of section A.4. If

$$\Pi_0(h | L_U + L_V) \equiv a^*(U) + b^*(V) \equiv ACE_U(h) + ACE_V(h),$$

it follows from propositions A.4.1 and A.3.1 that

$$a^*(U) = a_0(U) - E_0(b^*(V) | U),$$

$$b^*(V) = b_0(V) - E_0(a^*(U) | V),$$

where  $a_0(U) \equiv E_0(h | U)$  and  $b_0(V) \equiv E_0(h | V)$ . It follows that the efficient influence function  $\tilde{l}$  for estimation of  $v$  is given by

$$(25) \quad \pi_{(s,t)} \tilde{l} = \tilde{l}_{(s,t)} = 1_{(-\infty, s] \times (-\infty, t]} - \Pi_0(1_{(-\infty, s] \times (-\infty, t]} | L_U + L_V).$$

Alternatively, for estimation of  $P$  indexed by a class of functions  $F$ ,  $v(P)(f)$  for  $f \in F$ ,

$$(26) \quad \begin{aligned} \pi_f \tilde{l}(u, v) &= \tilde{l}_f(u, v) \\ &= f(u, v) - E_0 f - \Pi_0(f | L_U + L_V)(u, v) \\ &= f(u, v) - E_0 f - ACE_U(f)(u) - ACE_V(f)(v). \end{aligned}$$



Thus the inverse information covariance function  $\Gamma_V^{-1}$  of (5.2.23) is (see also proposition A.2.2.B)

$$\begin{aligned}\Gamma_V^{-1}(f, g) &= E_0 \{ (f - E_0 f - ACE_U(f) - ACE_V(f)) \\ &\quad \cdot (g - E_0 g - ACE_U(g) - ACE_V(g)) \} \\ &= Cov_0(f, g) - E_0 \{ (ACE_U(f) + ACE_V(f))(ACE_U(g) + ACE_V(g)) \}\end{aligned}$$

and

$$\begin{aligned}\Gamma_V^{-1}(f, f) &= E_0 (f - E_0 f - \Pi_0(f | L_U + L_V))^2 \\ &= Var_0(f) - E_0 (ACE_U(f) + ACE_V(f))^2.\end{aligned}$$

While it is not possible to give a more explicit expression for the covariance function  $\Gamma_V^{-1}$  in general, it is clear that the process  $\mathbf{Z}_0$  of theorem 5.2.1 is

$$\mathbf{Z}_0(f) \equiv \tilde{\mathbf{Z}}_0(\tilde{\mathbf{I}}_f), \quad f \in \mathbf{F},$$

where  $\tilde{\mathbf{Z}}_0$  is the  $P$ -Brownian bridge process  $\mathbf{Z}_0$  of example 5.3.8 and  $\tilde{\mathbf{I}}_f$  is given by (26).

This example is the continuous version of the problem of estimating cell probabilities in a two-way contingency table with known marginal distributions, a model which is often fitted by use of the "iterative proportional fitting" algorithm introduced by Deming and Stephan (1940). As is now well known, this algorithm converges to the minimum Kullback divergence estimator in this model; see, e.g., Haberman (1979), (1984), Ireland and Kullback (1968), Mosteller (1968), Causey (1972), and Gokhale and Kullback (1978). Bickel, Ritov, and Wellner (1991) studied the modified minimum chi-square estimates of Deming and Stephan (1940) with the number of cells increasing to infinity with sample size  $n$ , and showed that the related natural estimate of  $v(P)(f)$  for a fixed  $f$  is efficient.  $\square$

A number of generalizations of examples 2 and 3 are of interest and importance. One of them would be to allow one or both marginal distributions to be parametric families, such as normal or exponential, rather than completely known as in examples 2 and 3. It would also be interesting to explore the continuous analogues of some of the discrete models of Goodman (1985).

**Example 4. Estimation of a distribution with independence.**

Suppose that  $\mathbf{X} = \mathbf{X}_1 \times \dots \times \mathbf{X}_d$ , write  $X = (X_1, \dots, X_d)$  and suppose that

$$(27) \quad \mathbf{P} = \{ P : \gamma_u(P) = 0 \text{ for all } u \in \mathbf{U} \},$$

where

$$\gamma_u(P) \equiv E_P \left( \prod_{i=1}^d u_i(X_i) \right) - \prod_{i=1}^d E_P \left( u_i(X_i) \right)$$

and

$$\mathbf{U} \equiv \{ u : u(x) = u_1(x_1) \dots u_d(x_d), \quad u_i \in \mathbf{C}_i \}$$

where the families of functions  $C_j$  are sufficiently rich; e.g.,  $C_i$  containing all functions  $1_C$  for  $C$  in a  $\pi$ -system  $\tilde{C}_i$  (cf. Billingsley (1986, page 36)). It is clear that  $\mathbf{P}$  in (27) is just the family of independent distributions on the product space  $\mathbf{X}$ . Thus we may also write

$$(28) \quad \{\mathbf{P} = P : P = P_1 \times \dots \times P_d, P_i \text{ a measure on } \mathbf{X}_i, i = 1, \dots, d\}.$$

By reasoning as in example 3.2.1 it follows that

$$(29) \quad \dot{\mathbf{P}} = \left\{ h \in L_2 : h(x) = h_1(x_1) + \dots + h_d(x_d), \int h_i dP_i = 0, i = 1, \dots, d \right\},$$

and hence by (A.2.9) and proposition A.3.2

$$(30) \quad \Pi_0(h | \dot{\mathbf{P}}) = \sum_{i=1}^d [E_0(h | X_i = x_i) - E_0(h)].$$

Thus, for estimation of  $v(P)$  with  $v(P)(f) \equiv \int f dP$  for  $f \in \mathbf{F}$ ,  $\dot{v}_f = f - \int f dP$  by (5.3.29) of example 5.3.8 in the unconstrained model  $\mathbf{M}$ , and hence the efficient influence function  $\tilde{I}$  is given by

$$(31) \quad \begin{aligned} \pi_f \tilde{I}(x) &= \tilde{I}_f(x) \\ &= \Pi_0(\dot{v}_f | \dot{\mathbf{P}}) \\ &= \sum_{i=1}^d [E_0(f | X_i = x_i) - E_0(f)] \\ &= \sum_{i=1}^d [(f_i(x_i) - E_0(f_i)) \prod_{j \neq i} E_0(f_j)] \end{aligned}$$

if

$$(32) \quad f(x) = f_1(x_1) \dots f_d(x_d).$$

It follows that the inverse information covariance function  $I_V^{-1}$  of (5.2.23) becomes for  $f, g$  of the form (32),

$$I_V^{-1}(f, g | \mathbf{P}) = \sum_{i=1}^d \left( \prod_{j \neq i} E_0(f_j) E_0(g_j) \right) Cov_0(f_i, g_i).$$

Note that for  $f, g$  of the product form given in (32) the covariance function  $I_V^{-1}$  of the completely nonparametric empirical measure estimator given in example 5.3.8 becomes

$$\begin{aligned} I_V^{-1}(f, g | \mathbf{M}) &= Cov_0(f, g) \\ &= \prod_{i=1}^d E_0(f_i g_i) - \prod_{i=1}^d E_0(f_i) E_0(g_i) \\ &= \sum_{i=1}^d \left\{ \prod_{j=1}^{i-1} E_0(f_j g_j) \prod_{j=i+1}^d E_0(f_j) E_0(g_j) \right\} Cov_0(f_i, g_i), \end{aligned}$$

and hence the decrease in variance of efficient estimates at the point  $f = g = \prod_1^d f_i$  is

$$\begin{aligned} I_V^{-1}(f, f | \mathbf{M}) - I_V^{-1}(f, f | \mathbf{P}) &= \sum_{i=1}^d \text{Var}_0(f_i) \left\{ \prod_{j=1}^{i-1} \frac{E_0(f_j^2)}{(E_0 f_j)^2} - 1 \right\} \prod_{j \neq i} (E_0 f_j)^2. \end{aligned}$$

Note that examples 3 and 4 (in the case  $d = 2$ ) are, in a sense, dual to each other: In example 3 the marginal distributions are known while the joint or dependence structure is unknown, whereas in example 4 the situation is just the reverse. Also note that the copula models of section 4.7 are in a sense dual to the distributions with parametric marginals and unknown joint structure mentioned after example 4. □

### 6.3 GROUP MODELS

Our goal in this section is to provide bounds for estimation of the non-parametric components of the group models considered in section 4.2. We first treat the relatively simple special case in which the family  $\mathbf{P}$  is defined directly in terms of a group of transformations  $\mathbf{T}$  on the sample space  $\mathbf{X}$ , and we call these models *nonparametric group models*. We then later add the complicating factor of finite-dimensional nuisance parameters, usually in the form of a finite-dimensionally parametrized group as in section 4.2, and we call these models *semiparametric group models*.

#### Nonparametric Group Models

Let  $\mathbf{T}$  be a group of transformations on  $\mathbf{X}$ , and let

$$(1) \quad \mathbf{P} = \{ P \ll \mu : PT^{-1} = P \text{ for all } T \in \mathbf{T} \}.$$

Here are several examples.

**Example 1. Spherically symmetric distributions in  $R^d$ .**

Suppose that  $\mathbf{X} = R^d$  with  $d \geq 1$  and  $\mathbf{T} = \{ T : T(x) = Ox, O \text{ an orthogonal matrix} \}$ . Then  $\mathbf{P}$  in (1) is the family of all distributions which are invariant under orthogonal transformations, or equivalently, the family of all spherically symmetric distributions. Note that example 5.3.3 is the special case of this example with  $d = 1$ . □

**Example 2. Coordinatewise symmetric distributions.**

Suppose that  $\mathbf{X} = R^d$  and

$$\mathbf{T} = \{ T : T(x) = (s_1 x_1, \dots, s_d x_d) \text{ where } s_i = \pm 1, i = 1, \dots, d \}.$$

Then the model  $\mathbf{P}$  in (1) is the family of distributions which are symmetric in each coordinate. □

**Example 3. Exchangeable distributions.**

Let  $\Pi$  denote the collection of all permutations of the integers  $\{ 1, \dots, d \}$ , and for  $\pi \in \Pi$  and  $x \in R^d$  let  $\pi x \equiv (x_{\pi(1)}, \dots, x_{\pi(d)})$ . Suppose that

$X = R^d$  and  $T = \{T: T(x) = \pi x, \pi \in \Pi\}$ . Then  $P$  in (1) is the family of exchangeable or permutation symmetric distributions.  $\square$

**Example 4. Cyclically symmetric distributions.**

If  $X = R^d$  and  $T = \{T: T(x) = \pi x, \pi \in C\}$  where  $C \subset \Pi$  is the cyclic-subgroup of the permutation group  $\Pi$ , then  $P$  in (1) is the family of cyclically symmetric distributions.  $\square$

**Example 5. Dihedrally symmetric distributions.**

Suppose that  $X = R^d$  and

$$T = \{T: T(x) = s\pi x, s = \pm 1, \pi \in C\}$$

where  $C \subset \Pi$  is the cyclic subgroup of the permutation group  $\Pi$  as in example 4. (For the definition of the dihedral group, see, e.g., Herstein (1964, problem 2.6.17, page 46). This particular representation of the dihedral group yields one type of "dihedrally symmetric distributions." Another type of "dihedral symmetry" is obtained in terms of the matrices  $C$  and  $D$  defined by

$$C = \begin{pmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & & \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & 1 \\ 1 & 0 & 0 & \dots & 0 \end{pmatrix}, \quad \text{and} \quad D = \begin{pmatrix} 0 & 0 & \dots & 1 \\ 0 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ \dots & 1 & \dots & \dots \\ 1 & 0 & \dots & 0 \end{pmatrix}.$$

Note that  $C^d = D^2$  is the identity matrix and that  $C = \{C, C^2, \dots, C^d\}$ . Then another dihedral group is given by

$$T = \{T: T(x) = \pi x, \pi \in \{C, C^2, \dots, C^d, DC, \dots, DC^d\}\},$$

and this in turn leads to another type of "dihedrally symmetric distributions."  $\square$

**Example 6. Rotationally symmetric distributions on the sphere.**

Suppose that  $X = S^{d-1}$ , the unit sphere in  $R^d$ , and  $T = \{T: T(x) = Ox, O \text{ rotation about } [a]\}$ ,  $a \in R^d$ , so  $[a]$  is a fixed axis. Then  $P$  in (1) consists of all the distributions on  $S^{d-1}$  for which the conditional distributions given the angle between  $X$  and  $[a]$  are uniform; see, e.g., Watson (1983).  $\square$

From the results of appendix A.3 it follows that the tangent space  $\dot{P}$  of the model  $P$  in (1) is

$$(2) \quad \dot{P} = \{h \in L_2^0(P_0): h \text{ is almost invariant under } T\},$$

and hence, by proposition A.3.2, that

$$(3) \quad \Pi_0(h | \dot{P}) = E_0(h | \mathcal{B}_T) - E_0 h \quad \text{for } h \in L_2(P_0),$$

where  $\mathcal{B}_T$  is the invariant  $\sigma$ -field given in (A.3.5)

$$\mathcal{B}_T = \{A \in \mathcal{B}: T(A) = A \text{ for all } T \in T\}.$$

Of course, (3) is exactly the projection formula needed to carry out the calculations of section 5.3. Equivalently, if  $Y$  is a maximal invariant with respect to the

group  $T$  (so that  $Y$  is invariant,  $Y(Tx) = Y(x)$  for all  $T \in T$ , and  $Y(x_1) = Y(x_2)$  implies that  $x_2 = T(x_1)$  for some  $T \in T$ ), then  $\mathcal{B}_T = \sigma(Y) \equiv$  the  $\sigma$ -field generated by  $Y$ , and (3) becomes

$$(4) \quad \Pi_0(h | \dot{P}) = E_0(h | Y) - E_0 h.$$

When the group  $T$  is finite or compact, as is the case in all of the examples above, the conditional expectations in (3) or (4) can be calculated via proposition A.3.3:

$$(5) \quad E_0(h | \mathcal{B}_T)(x) = \int_T h(Tx) dm(T) \quad \text{a.s.}$$

where  $m$  is right-invariant Haar probability measure on  $T$ .

Now consider estimation of  $\nu(P)$  with

$$\nu(P)(f) \equiv \int f dP \equiv P(f), \quad f \in F,$$

as in example 5.3.8. From (2), (3) and (5.3.28)–(5.3.29), it follows that the efficient influence function  $\tilde{I}$  for estimation of  $\nu$  in the model (1) is given by

$$(6) \quad \pi_f \tilde{I}(x) = \tilde{I}_f(x) \equiv \tilde{I}(\pi_f)(x) = \Pi_0(\dot{\nu}_f | \dot{P}) = E_0(f | \mathcal{B}_T) - E_0 f.$$

Hence the covariance function  $I_\nu^{-1}$  of (5.2.23) is given by

$$(7) \quad I_\nu^{-1}(f_1, f_2) = \text{Cov}_0(E_0(f_1 | \mathcal{B}_T), E_0(f_2 | \mathcal{B}_T)),$$

and the decrease in variance at a point  $f \in F$  due to knowledge of invariance ( $P \in \mathcal{P}$  given by (1)) is

$$(8) \quad \text{Var}_0(f) - \text{Var}_0(E_0(f | \mathcal{B}_T)) = E_0 \text{Var}_0(f | \mathcal{B}_T).$$

Note that

$$(9) \quad 0 \leq E_0 \text{Var}_0(f | \mathcal{B}_T) \leq \text{Var}_0 f.$$

No decrease in variance occurs when equality holds in the first inequality of (9), and this is true if and only if

$$(10) \quad \text{Var}_0(f | \mathcal{B}_T) = 0 \quad \text{a.s.},$$

or, in other words, if and only if  $f$  is  $\mathcal{B}_T$ -measurable. On the other hand, the maximum decrease in variance occurs when equality holds in the second inequality of (9), and this is true if and only if (cf. (8))  $E_0(f | \mathcal{B}_T)$  is degenerate, i.e.,

$$(11) \quad E_0(f | \mathcal{B}_T) = E_0 f \quad \text{a.s.}$$

From (9) it follows that the ratio of (7) to (5.3.30) at a point  $f = f_1 = f_2 = g \in F$ , the asymptotic relative efficiency for estimation of  $\nu(P)(f) = \int f dP$ , ranges from 0 to 1, and is 0 if equality holds in the first inequality of (9), 1 if equality holds in the second inequality of (9).

Now we return to examples 1–5 and carry out the above calculations explicitly in these cases.

**Example 1. Spherically symmetric distributions in  $R^d$ , continued.**

Note that there is a one-to-one correspondence between  $X$  and  $(U, V)$  where  $V \equiv |X|$  is a maximal invariant with respect to  $\mathbf{T}$  and  $U \equiv X/|X|$  is distributed uniformly on the unit sphere  $S^{d-1}$  in  $R^d$  and is independent of  $V$ . The only unknown in this model is the one-dimensional distribution of the random variable  $V = |X|$ . If

$$(12) \quad \frac{dP}{d\mu}(x) = \tilde{g}(|x|^2),$$

where  $\tilde{g}: R^+ \rightarrow R^+$  and  $\mu$  is Lebesgue measure on  $R^d$ , then  $V$  has density  $c(d)v^{d-1}\tilde{g}(v^2)$ ,  $v \geq 0$ , where  $c(d) = 2\pi^{d/2}/\Gamma(d/2)$  is the area of the unit sphere  $S^{d-1}$  in  $R^d$ ,  $d \geq 2$ ,  $c(1) = 2$  (the "area" of the unit sphere in  $R$  with respect to counting measure). Thus

$$E_0(f | \mathcal{B}_T) = E_0(f | V),$$

and for functions  $f$  of the form

$$(13) \quad f(x) = g\left(\frac{x}{|x|}\right)h(|x|) = g(u)h(v)$$

with  $g \in L_2(U)$ ,  $h \in L_2(V)$ , it follows that

$$E_0(f | \mathcal{B}_T) = h(V)Eg(U)$$

and

$$\text{Var}_0(f | \mathcal{B}_T) = h^2(V)\text{Var} g(U),$$

where  $U \sim \text{Uniform}(S^{d-1})$ . Thus for functions  $f_i$  of the form (13), (7) becomes

$$(14) \quad I_V^{-1}(f_1, f_2) = \{Eg_1(U)Eg_2(U)\} \text{Cov}_0(h_1(V), h_2(V)),$$

and the difference (8) is

$$E_0 \text{Var}_0(f | V) = \text{Var} g(U)E_0h^2(V).$$

For example, if  $f(x) = g(u)h(v) = 1_C(u)1_{[0,r]}(v)$ , where  $C$  is a spherical cap on  $S^{d-1}$  of normalized area  $A \equiv (\text{area of } C)/(\text{area of } S^{d-1})$ , then the decrease in variance for estimation of

$$\begin{aligned} v(P) &= \int f dP \\ &= P(X \in \text{truncated cone determined by } r \text{ and } C) \\ &= AP(|X| \leq r), \end{aligned}$$

is

$$\text{Var} g(U)E_0h^2(V) = A(1-A)P_0(|X| \leq r) = (1-A)v(P),$$

which is bounded by  $\frac{1}{4}$  with the maximum attained when  $A = \frac{1}{2}$  and  $p \equiv P_0(V \leq r) = 1$ ; or the asymptotic relative efficiency is

$$\begin{aligned} \frac{Var_0(E_0(f|V))}{Var_0 f} &= \frac{A^2 p(1-p)}{Ap(1-Ap)} \\ &= \frac{A(1-p)}{1-Ap} \\ &\geq 0 \quad \text{with } = \text{ if } A = 0 \text{ or } p = 1, \\ &\leq 1 \quad \text{with } = \text{ if } A = 1. \end{aligned}$$

In this example the group  $T$  is big enough to essentially reduce a  $d$ -dimensional model (all  $P$  on  $R^d \ll \mu$ ) to a one-dimensional model: all distributions of  $V$  on  $R^+$ . □

**Example 2. Coordinatewise symmetric distributions, continued.**

In this example the maximal invariant is  $Y = (|X_1|, \dots, |X_d|)$  and  $X$  is equivalent to  $(S, Y)$  where  $S \equiv (S_1, \dots, S_d)$ , the vector of signs  $S_i \equiv X_i/|X_i|$ , is independent of  $Y$ . Therefore the conditional expectation in (3) or (4) becomes

$$(15) \quad E_0(f|Y) = \frac{1}{2^d} \sum_{s_i = \pm 1} f(s_1|X_1|, \dots, s_d|X_d|) \quad \text{a.s.}$$

Hence for any function  $f$  of the form

$$(16) \quad f(x) = g(s)h(y)$$

with  $h \in L_2(Y)$ , it follows that

$$(17) \quad \Gamma_v^{-1}(f_1, f_2) = \{Eg_1(S)Eg_2(S)\} Cov_0(h_1(Y), h_2(Y)),$$

and the difference (8) is

$$(18) \quad E_0(Var_0(f|Y)) = Var g(S)E_0h^2(Y),$$

where, again, the first two factors in (17) and the first factor in the right-hand side of (18) do not depend on  $P \in \mathbf{P}$ . □

**Example 3. Exchangeable distributions, continued.**

In this example the maximal invariant is  $Y \equiv (X_{(1)}, \dots, X_{(d)})$  where  $X_{(1)} \leq \dots \leq X_{(d)}$  are the order statistics, and  $X$  is equivalent to  $(R, Y)$  where  $R \equiv (R_1, \dots, R_d)$ , the vector of ranks  $R_j \equiv \# \text{ of } X_i\text{'s } \leq X_j$ , is independent of  $Y$ . It is well-known that

$$(19) \quad E_0(f|Y) = \frac{1}{d!} \sum_{\pi \in \Pi} f(X_{\pi(1)}, \dots, X_{\pi(d)});$$

see, e.g., Hájek and Šidák (1967, problem II.1.7, page 79); of course, this also follows easily from (5). In particular, for any function  $f$  of the form

$$(20) \quad f(x) = g(r)h(y)$$

with  $h \in L_2(Y)$ , it follows that

$$E_0(f|Y) = h(Y)Eg(R),$$

where  $R$  is distributed uniformly over the  $d!$  permutations of  $\{1, \dots, d\}$ . Thus for functions  $f_i$  of the form (20), (7) becomes

$$(21) \quad I_v^{-1}(f_1, f_2) = \{Eg_1(R)Eg_2(R)\} Cov_0(h_1(Y), h_2(Y)),$$

and the difference (8) is

$$(22) \quad E_0 Var_0(f | Y) = Var g(R)E_0 h^2(Y)$$

where the first two factors in (21) and the first factor in the right-hand side of (22) do not depend on  $P \in \mathbf{P}$ .  $\square$

**Example 4. Cyclically symmetric distributions, continued.**

In this example a maximal invariant  $Y = (Y_1, \dots, Y_d)$  can be chosen to be  $Y = C^h X$  with the random integer  $h \in \{0, 1, \dots, d-1\}$  such that  $Y_1 = X_{(1)}$  and  $C$  as in example 5. Then the conditional expectation in (4) becomes

$$(23) \quad E_0(f | Y)(x) = \frac{1}{d} \sum_{\pi \in C} f(\pi x). \quad \square$$

**Example 5. Dihedrally symmetric distributions, continued.**

For the first transformation group of this example the conditional expectation in (3) becomes

$$(24) \quad E_0(f | \mathcal{B}_T)(x) = \frac{1}{2d} \sum_{s=\pm 1, \pi \in C} f(s\pi x). \quad \square$$

Examples 4 and 5 seem to merit further exploration. Many other symmetry groups  $T$  are possible and potentially interesting; see, e.g., Serre (1977, chapter 5).

### Semiparametric Group Models

We now consider the situation discussed in section 4.2, but derive bounds for estimation of the "shape parameter"  $G$ . Suppose  $\mathbf{X} = R^d$ ,  $\mu$  is Lebesgue measure, and we are given

$$(25) \quad \mathbf{P} = \{Ga_\theta^{-1} : G \in G, \theta \in \Theta\},$$

where  $\Theta$  open  $\subset R^k$ ,

$$(26) \quad a_\theta(x) = a(x, \theta) = \Delta(\theta) + S(\theta)x,$$

and

$$\mathbf{A} \equiv \{a(\cdot, \theta) : \theta \in \Theta\}$$

is a Lie subgroup of the affine group as in assumptions (iii), (iv) of section 4.2. That is,  $\theta \rightarrow (\Delta(\theta), S(\theta))$  is continuously differentiable. Now suppose  $T$  is another transformation group on  $\mathbf{X}$ ,

$$(27) \quad G = \{G : G \ll \mu, GT^{-1} = G \text{ for all } T \in T\},$$

and  $\mathcal{B}_T$  is the corresponding invariant  $\sigma$ -field. It is often convenient to strengthen (27) by requiring tail conditions. For example,

$$(28) \quad G = \{G : G \ll \mu, GT^{-1} = G \text{ for all } T \in T, \int b dG \leq M\}$$

for some  $b \geq 0, M < \infty$ . In this situation,  $G$  is typically no longer identifiable



but its “shape” is. Sometimes it is even adaptively estimable. We do not define “shape” formally, but illustrate what can happen with three examples.

**Example 7.**  $A = \{\text{location group}\}$ ,  $T = \{\text{identity}\}$ .

One treatment of this example was given in example 5.5.1. We now give a slightly different approach. In this case  $X = R$ , and we take  $G = \{G \ll \mu : \int |x|^{2+\delta} dG(x) \leq M\}$  for some  $\delta > 0$ ,  $M < \infty$ . If  $a(x, \theta) \equiv x + \theta$ , then  $P_{(\theta, G)} = G(\cdot - \theta)$  and  $P_{(\theta, G)} = P_{(0, G(\cdot - \theta))}$  so that  $G$  is not identifiable. The obvious solution is to recognize that  $P$  is (essentially) all distributions and fix  $\theta = 0$ . Then,  $G = P$  and is, of course, identifiable. An alternative classical remedy is to let

$$(29) \quad G_0 = \{G \in G : \mu(G) = 0\},$$

where  $\mu(G)$  is any location functional, for instance,  $\mu(G) = \int x dG(x)$  (see Bickel and Lehmann (1975)). Then

$$(30) \quad P = \{P_{(\theta, G)} : \theta \in R, G \in G_0\}$$

where both parameters are identifiable, namely by  $\theta = \mu(P)$  and  $G(\cdot) = P((-\infty, \cdot + \theta])$ .

By applying corollary 5.5.2.B, it is not hard to show that smooth functions of  $G \in G_0$  are pathwise differentiable and to determine the efficient influence operator. It is easier though to calculate these influence operators or functions directly. We use the fact that  $\dot{P}$  is saturated (i.e.,  $\dot{P} = \dot{M}$ ) and the results of example 5.3.8. For instance, if  $v(P)(f) \equiv \int f dG, f \in l^\infty(R)$ , then

$$(31) \quad v(P)(f) = \int f(x) dP(x + \mu(P)).$$

For  $\mu(P) = \int x dP(x)$ , a chain rule argument using (5.3.28) gives the efficient influence function  $\tilde{I}$  by

$$(32) \quad \begin{aligned} \pi_f \tilde{I}(X) &= \tilde{I}_v(\pi_f)(X) \\ &= f(X - \theta) - v(P)(f) + (X - \theta) \int f(t) g'(t) dt \in \dot{P}, \end{aligned}$$

which for  $f = 1_{(-\infty, t]}$  reduces to the efficient influence function given in example 5.5.1. Note that the influence function depends on which  $\mu(G)$  is used to define  $G_0$ . For instance, if  $G_0$  is defined by  $\mu(G) \equiv \int x dG(x)$  in (29), then  $\mu(G)$  may be estimated perfectly, while

$$\mu_1(G) \equiv G^{-1}\left(\frac{1}{2}\right) = F^{-1}\left(\frac{1}{2}\right) - \int x dF(x),$$

where  $F$  is the distribution function of  $P$ , cannot be so estimated. Of course, the roles of  $\mu(G)$  and  $\mu_1(G)$  can be inverted. □

**Example 8. Symmetric location model.**

In this case also  $X = R$ . We take  $T = \{x \rightarrow x, x \rightarrow -x\}$ , and by adding tail and finite Fisher information conditions end up with

(33)  $G = \{G \ll \mu: G \text{ symmetric about } 0,$

$$\int |x|^{4+\delta} dG(x) \leq M, \int (1+x^2) \frac{[g']^2}{g}(x) dx < \infty\}$$

for some  $\delta > 0$ . If  $A$  is the location group and  $a(x, \theta) = x + \theta$ , then  $G$  is identifiable. Let  $P \equiv P_{(\theta, G)}$ . We have seen in example 3.4.1 that

$$\dot{P}_1 = \left[ -\frac{g'}{g}(\cdot - \theta) \right],$$

$$\dot{P}_2 = \{h \in L_2^0(P): h(x) = h(2\theta - x) \text{ for all } x\},$$

and

$$\dot{P}_1 \perp \dot{P}_2.$$

Hence it should be possible to estimate not only  $\theta$  but also  $G$  adaptively. The efficient influence function  $\tilde{I}$  of  $v(P) \equiv G$  is then, by example 5.3.3, given by

$$(34) \quad \pi_t \tilde{I}(X) = \tilde{I}_v(\pi_t)(X) = \frac{1}{2} (1_{[X \leq t+\theta]} + 1_{[X \geq -t+\theta]}) - F(t+\theta),$$

where  $F$  is the distribution function of  $P$ .

If we now let  $G$  be as above, but take

$$A = \{x \rightarrow \theta_1 + \theta_2 x, \theta_1 \in R, \theta_2 > 0\},$$

the location-scale group,  $G$  and  $\theta_2$  are not identifiable. Again we can change  $\theta$  to  $\theta_1$ , or equivalently change  $\Theta$  to  $\Theta_0 = \{(\theta_1, \theta_2): \theta_2 = 1\}$ . Alternatively, if we take  $G_0 = \{G \in G: \int x^2 dG(x) = 1\}$  and restrict  $G$  to  $G_0$ , we can identify  $G$  and  $\theta_2$ . To obtain the influence functions of smooth functions of  $G$  we note, analogously to (31), that, for  $G \in G_0$ ,

$$(35) \quad G(t) = F(\theta_1 + \theta_2 t),$$

where  $F$  is as above and

$$\theta_1(P) \equiv \int x dP(x),$$

$$\theta_2^2(P) = \int (x - \theta_1(P))^2 dP(x).$$

The efficient influence function  $\tilde{I}$  of the shape  $G$  viewed as the abstract parameter given by

$$\begin{aligned} v(P)(f) &\equiv \int f(t) dP(\theta_1 + \theta_2 t) \\ &= \int f(t) dG(t) = \frac{1}{2} \int \{f(t) + f(-t)\} dG(t) \end{aligned}$$

is, by the same calculations as in example 7, given by

$$\begin{aligned} \pi_f \tilde{I}(X) &= \tilde{I}_v(\pi_f)(X) \\ &= \frac{1}{2} \left\{ f\left(\frac{X - \theta_1}{\theta_2}\right) + f\left(-\frac{X - \theta_1}{\theta_2}\right) \right\} - v(P)(f) \\ &\quad + \frac{1}{2} (X - \theta_1) \int [f(z) + f(-z)] g'(z) dz \end{aligned}$$

$$\begin{aligned}
 & + \frac{1}{4\theta_2} \{(X - \theta_1)^2 - \theta_2^2\} \int [f(z) + f(-z)] [zg'(z) + g(z)] dz \\
 & \in \dot{P},
 \end{aligned}$$

which, for  $f = 1_{(-\infty, t]}$ , reduces to

$$\begin{aligned}
 (36) \quad \tilde{I}_v(\pi_f)(X) &= \frac{1}{2} \{1_{[X \leq \theta_1 + \theta_2 t]} + 1_{[X \geq \theta_1 - \theta_2 t]}\} - G(t) \\
 & + \frac{\theta_2}{2} t g(t) \left\{ \frac{(X - \theta_1)^2}{\theta_2^2} - 1 \right\}.
 \end{aligned}$$

By the usual calculation of the tangent space for a constrained model (see, e.g., (6.2.3))

$$\begin{aligned}
 \dot{P}_3 &= \dot{P}_3(P_{(0,1,G)}) \\
 &= \{h \in L_2^0(G) : h \text{ symmetric about } 0, \int h(x)(x^2 - 1) dG(x) = 0\}.
 \end{aligned}$$

Furthermore,

$$\begin{aligned}
 \dot{P}_1 &= \dot{P}_1(P_{(0,1,G)}) = \left[ \frac{g'}{g}(X) \right], \\
 \dot{P}_2 &= \dot{P}_2(P_{(0,1,G)}) = \left[ X \frac{g'}{g}(X) + 1 \right], \\
 \dot{P}_3(P_{(\theta,G)}) &= \{h(X - \theta) : h \in \dot{P}_3(P_{(0,1,G)})\},
 \end{aligned}$$

etc. Because  $g'/g$  is an odd function,  $\dot{P}_1 \perp \dot{P}_2, \dot{P}_3$ . So the necessary orthogonality condition for adaptive estimation of  $G$  in the presence of the location parameter  $\theta_1$  holds.

Since  $\dot{P}_3$  is the orthocomplement of  $x \rightarrow x^2 - 1$  within the space of symmetric functions,  $\dot{P}_2 \perp \dot{P}_3$  if and only if

$$x \frac{g'}{g}(x) + 1 = c(x^2 - 1), \quad x \in R,$$

for some  $c$ . A little computation shows that the necessary orthogonality condition for estimating the shape  $G$  adaptively in the presence of both  $\theta_1$  and  $\theta_2$  is valid in the model

$$P = \{P_{(\theta_1, \theta_2, G)} : \theta_1 \in R, \theta_2 > 0, G \in \mathbf{G}_0\}$$

at all  $G$  such that if  $X$  has distribution  $G$ , then  $X^2$  has a gamma distribution with mean 1. Note that the standard normal distribution has this property.  $\square$

Here is our final example where shape can be estimated adaptively.

**Example 9. Elliptic distributions.**

Here  $\mathbf{X} = R^d$ ,  $\mathbf{T} = \{\text{orthogonal transformations}\}$ ,  $\mathbf{A}$  is the affine group,  $\Theta = \{S_{d \times d} \text{ nonsingular}, \Delta \in R^d\}$  with the identification (26). We restrict by tail conditions and Fisher information and let

$$(37) \quad \mathbf{G} = \{G \text{ spherically symmetric, } \int |x|^{4+\delta} dG(x) < \infty \\ I_0(G), I_1(G) \text{ given by (4.2.23) finite}\}.$$

Of course, neither  $G$  nor  $S$  are identifiable. As we have noted in example 4.2.3, we have identifiability of  $q(\theta) \equiv (SS^T / \text{trace}(SS^T), \Delta)$ . In fact we can estimate  $q$  adaptively in  $\mathbf{P}$ .

Let  $\mathbf{G}_0 \equiv \{G \in \mathbf{G} : \int |x|^2 dG(x) = d\}$ . Identifiability of  $S$  is impossible, but if  $G$  ranges over  $\mathbf{G}_0$  we can identify  $SS^T$  and  $G$ . The influence function of  $v(P)$  defined by  $v(P)(f) = \int f dG$  may be obtained as in example 8. Reparametrize

$$\tilde{\Theta} = \{(\Delta, S) : \Delta \in R^d, S \text{ positive definite symmetric, } \text{trace}(S^2) = d\},$$

and write

$$(38) \quad \mathbf{P} = \{P_{(\theta, \sigma, G)} : \theta \in \tilde{\Theta}, \sigma > 0, G \in \mathbf{G}_0\},$$

$$\sigma^2 = \sigma^2(P) \equiv \int |x|^2 dP(x). \text{ Note that}$$

$$\dot{\mathbf{P}}_3(P_{(\theta, \sigma, G)}) = \{h(S^{-1}\sigma^{-1}(X - \Delta)) : h \in \dot{\mathbf{P}}_3(P_{(e, 1, G)})\},$$

where  $e \equiv (0, J)$  and  $J$  is the identity. By the usual computation for constrained models,

$$\dot{\mathbf{P}}_3(P_{(e, 1, G)}) = \{h \in L_2^0(G) : h \text{ is almost invariant under } \mathbf{T} \\ \int h(x)(|x|^2 - d) dG(x) = 0\}.$$

Therefore, as some computation shows,

$$(39) \quad \dot{\mathbf{P}}_3 \perp \dot{\mathbf{P}}_2.$$

provided  $G$  is such that if  $X$  has distribution  $G$ , then  $|X|^2$  has a gamma distribution with mean  $d$ . Write  $\tilde{I}_t$  for  $\tilde{I}_v(\pi_t)$ . This influence function may be computed as in the previous examples but its form is uninformative. However, we can easily prove:

**Proposition 1.** The necessary orthogonality condition for estimating the shape  $G$  adaptively is valid in the model  $\mathbf{P}$  given in (38) at all  $G$  such that if  $X$  has distribution  $G$ , then  $|X|^2$  has a gamma distribution with mean  $d$ .

**Proof.** Note that (39) implies

$$(a) \quad \tilde{I}_t(\cdot | \mathbf{P}_{23}) = \tilde{I}_t(\cdot | \mathbf{P}_3).$$

By example 4.2.3

$$(b) \quad \tilde{I}(\cdot | \theta, \mathbf{P}) \perp \dot{\mathbf{P}}_{23}.$$

Furthermore,

$$(c) \quad \dot{\mathbf{P}}_{12} = [\tilde{I}(\cdot | \theta, \mathbf{P})] + \dot{\mathbf{P}}_2.$$

From (a),  $\tilde{I}_r(\cdot | P_3) \perp \dot{P}_2$ , and from (b),  $\tilde{I}_r(\cdot | P_3) \perp [\tilde{I}(\cdot | \theta, P)]$ , and the proposition follows from (c).  $\square$

Note that the  $N(0, J)$  distribution satisfies the condition in proposition 1.  $\square$

We now consider the regression models of section 4.3. Let  $X = (Z, Y)$ ,  $Y \in R^k$ ,  $Z \in Z$ ,  $R = \{r(\cdot, v) : v \in V\}$ , where  $r(\cdot, v) : Z \rightarrow \Gamma$ ,  $\Gamma$  open  $\subset R^s$ ,  $V$  is a subset of a Banach space  $B$  and  $A$  is as in (26). Then, let

$$(40) \quad P \equiv \{P \ll \mu \times m : Z \sim H, L(Y | Z) = Ga^{-1}(\cdot, r(Z, v)), \\ G \in G, v \in V, H \in H\},$$

where  $\mu$  is Lebesgue measure,  $G$  is a nonparametric group model as in (27),  $V$  and  $H$  are arbitrary. That is,  $Z \sim H$ ,

$$Y = a(\varepsilon, r(Z, v)),$$

where  $\varepsilon \sim G$  and  $Z$  and  $\varepsilon$  are independent.

Let  $Q$  be the semiparametric group model generating  $P$ . That is  $X \sim Q \in Q$ , where

$$(41) \quad Q = \{Ga^{-1}(\cdot, \gamma) : \gamma \in \Gamma, G \in G\}.$$

Let  $G_0 \subset G$  as in examples 7-9 be such that  $P$  and  $Q$  are unchanged if we replace  $G$  by  $G_0$ .

We now argue formally that if, as in examples 8 and 9,  $G$  may be estimated adaptively in  $Q$  when it is restricted to  $G_0$  (possibly only at some particular points  $P$  as in examples 8 and 9), then  $G$  may similarly be adaptively estimated in  $P$ . Note that adaptation in  $Q$  simply means that if  $G \in G_0$  is viewed as the second variable in  $Q_{(\gamma, G)}$  then

$$(42) \quad \dot{Q}_2 \perp \dot{Q}_1.$$

Further, (4.2.5) yields formally

$$(43) \quad \dot{Q}_1(\gamma, G) = \dot{\omega}^T(\gamma) [\dot{I}_1(\varepsilon)]$$

where  $\varepsilon = a^{-1}(Y, \gamma)$  and where  $\dot{\omega}$  and  $\dot{I}_1(\varepsilon)$  are defined in (i) and (ii) of section 4.2. Formally

$$(44) \quad \dot{P}_1(v, G, H) = [h(Z) \dot{\omega}^T(r(Z, v)) \dot{I}_1(\varepsilon) : h \in \dot{R}],$$

where  $\varepsilon = a^{-1}(Y, r(Z, v))$ . We now state:

**Theorem 1.** Suppose  $\dot{P} = \dot{P}_1 + \dot{P}_2 + \dot{P}_3$ ,  $\dot{Q} = \dot{Q}_1 + \dot{Q}_2$ . If  $\dot{Q}_1, \dot{P}_1$  are given by (43) and (44) respectively, and (42) holds, then

$$\dot{P}_2 \perp \dot{P}_1 + \dot{P}_3.$$

That is, for estimation of  $G$ , adaptation to both  $v$  and  $H$  should be possible.

**Proof.** Note that  $\dot{P}_2$  consists of functions of  $\varepsilon$  only, while  $\dot{P}_3$  consists of functions of  $Z$  only. Hence  $\dot{P}_2 \perp \dot{P}_3$  by the assumed independence of  $\varepsilon, Z$ . On the other hand, if  $b(\varepsilon) \in \dot{Q}_2$ , then (42) implies that

$$(a) \quad E_Q(\dot{\omega}^T(\gamma) b(\varepsilon) \dot{I}_1(\varepsilon)) = 0,$$

and hence, by the independence of  $\varepsilon$  and  $Z$ ,

$$(b) \quad E_P(h(Z) \dot{\omega}^T(r(Z, v)) b(\varepsilon) \dot{I}_1(\varepsilon)) = 0$$

for  $h \in \dot{\mathbf{R}}$ . The theorem follows from (44) and (b).  $\square$

We conclude with:

**Example 10. Partial splines, projection pursuit, and periodic regression, generalized.**

In all of the models of examples 4.3.3–4.3.5 we have assumed that the errors are Gaussian. If instead we simply assume the errors are symmetrically distributed according to  $G$  with  $G$  restricted as in (33), say, then theorem 1 and example 8 enable us to conclude that  $G$  may be adaptively estimated at some  $G$ . In fact, the last part of example 8 gives more. If we replace the error in these models by  $\sigma(Z, v)\varepsilon$  where  $\varepsilon \sim G$  as above, then the shape of  $G$  can be estimated adaptively at the same  $G$ . If we drop the restriction that  $G$  is symmetric about 0, we can conclude from example 7 that the shape of  $G$  is, in general, estimable but not adaptively.  $\square$

We can similarly use example 9 to obtain results about multivariate regression models with (possibly conditional on  $Z$ ) elliptic errors. We do not pursue this here.

## 6.4 BIASED SAMPLING MODELS

Now we return to the biased sampling models introduced in section 4.4, but with emphasis on estimation of parameters related to the infinite-dimensional components of the models. As in section 4.4, biased sampling models have three ingredients: an underlying model  $\mathbf{Q}$ , a set of “stratum weight functions”  $w_i : \mathbf{X} \rightarrow R^+$ ,  $i = 1, \dots, s$ , and selection probabilities  $\lambda_1, \dots, \lambda_s \geq 0$  with  $\sum_{i=1}^s \lambda_i = 1$ . In our examples here,  $\mathbf{Q}$  is semiparametric or nonparametric. This is the general i.i.d. biased sampling model with  $s$  “strata;” if  $Q \in \mathbf{Q}$  has density  $q$ , then  $(I, X)$  from  $P \in \mathbf{P}$  has density

$$p(i, x; q) = \lambda_i \frac{w_i(x)}{W_i(Q)} q(x) \quad \text{for } x \in \mathbf{X},$$

with  $W_i(Q) \equiv \int w_i dQ$ ,  $i = 1, \dots, s$ .

Often  $s = 1$ ,  $w$  is the indicator of a proper subset of  $\mathbf{X}$ , and  $\mathbf{Q}$  has sufficient structure to make parameters of interest identifiable.

We will focus here on example 4.4.3, and we will begin with an important special case of it, the random left truncation model. It deserves close comparison with the random censoring model studied in section 6.6. This will lead us naturally to a deeper study of example 4.4.3, the truncated regression models of Bhattacharya, Chernoff, and Yang (1983) and Jewell (1985).

**Example 1. Random truncation model; nonparametric view.**

Suppose that (under  $Q \in \mathbf{Q}$ )  $U, V$  are independent rv's with distributions  $F, G$  on  $R^+ = [0, \infty)$ . Here  $Q$  is the product measure  $F \times G$  on  $R^{+2}$ . Suppose that the stratum weight function is  $w(u, v) = 1_{[u < v]}$ . Thus we observe  $(U, V)$

only if  $U < V$ . To make both  $F$  and  $G$  identifiable (cf. (8) and (24) below), we assume throughout that

$$(1) \quad F^{-1}(0+) \leq G^{-1}(0+), \quad F^{-1}(1) \leq G^{-1}(1).$$

Then if  $F$  and  $G$  have densities  $f$  and  $g$  respectively,  $\alpha \equiv \int F_- dG > 0$  and the resulting biased sampling density is

$$(2) \quad p(u, v; F, G) = \frac{f(u)g(v)}{\int F_- dG} 1_{[u < v]} \equiv \frac{1}{\alpha} f(u)g(v) 1_{[u < v]}.$$

Here we introduce the notation  $F_-(u) = F(u-)$ . We can regard  $U$  as a random variable which truncates observation of  $V$  on the left (we don't see  $V$  or  $U$  if  $V \leq U$ ), and therefore this model is sometimes called the "random left truncation model." With this viewpoint, the df  $G$  of  $V$  is the parameter of primary interest.

We let  $F^*$  and  $G^*$  denote the marginal distributions of  $U$  and  $V$  respectively calculated under  $P$ ; thus

$$(3) \quad \begin{aligned} F^*(u) &\equiv P(U \leq u) = \frac{1}{\alpha} \int_0^u \bar{G} dF \\ &= \frac{1}{\alpha} \{ \bar{G}(u)F(u) + \int_0^u F_- dG \}, \end{aligned}$$

$$(4) \quad G^*(v) \equiv P(V \leq v) = \frac{1}{\alpha} \int_0^v F_- dG,$$

and

$$(5) \quad \begin{aligned} F^*(u) - G^*(u) &= \frac{1}{\alpha} \bar{G}(u)F(u) \equiv M^*(u) \\ &= \frac{1}{\alpha} P(U \leq u < V). \end{aligned}$$

Now consider estimation of the parameters  $\nu$  and  $\kappa : \mathbf{P} \rightarrow \mathbf{B} \equiv I^\infty([0, \tau])$  with  $0 < \tau < \tau_G \equiv G^{-1}(1) = \inf\{s : G(s) = 1\}$  defined by

$$(6) \quad \nu(P_{(F,G)})(t) = \Lambda_G(t) = \int_0^t \frac{dG}{1 - G_-}, \quad 0 \leq t \leq \tau,$$

and

$$(7) \quad \begin{aligned} \kappa(P_{(F,G)})(t) &= 1 - G(t) \\ &= \prod_{s \leq t} (1 - d\Lambda_G(s)), \quad 0 \leq t \leq \tau, \\ &\equiv \exp(-\Lambda_G^c(t)) \prod_{s \leq t} (1 - \Delta\Lambda_G(s)) \\ &= \exp(-\Lambda_G(t)) \quad \text{when } \Lambda_G \text{ is continuous.} \end{aligned}$$

Here  $\Lambda^c$  is the continuous part of  $\Lambda$  and  $\Delta\Lambda$  denotes the jump height. It follows from (1), (4), and (5) that

$$(8) \quad v(P_{(F,G)})(t) = \Lambda_G(t) = \int_0^t \frac{dG^*}{M_-^*};$$

thus  $\Lambda_G$  is expressed directly and explicitly by (8) as a function of  $P$  (and its marginals  $F^*$  and  $G^*$ ), and defines an extension  $v_e$  of  $v$  from  $\mathbf{P}$  to  $\mathbf{M}$ , the collection of all measures on the sample space  $\mathbf{X} = \{(u, v) : u < v\}$ . Note that (8) shows identifiability of  $G$ . In (6) and (8), as in example 5.3.5, we have indicated the left limits in the denominators that are appropriate for discontinuous cases since these help in forming estimators. However, from here on we will not indicate these left limits anymore.

We first take a completely nonparametric approach to bounds for estimation of  $v = \Lambda_G$  and  $\kappa = 1 - G$  in the spirit of sections 3.3 and 5.2. In example 2 we will reconsider the same model using score operators (and some martingale theory) as in sections 3.4 and 5.4–5.5.

We will assume throughout that (1) is strengthened to

$$(9) \quad \int \frac{1}{F} dG < \infty \quad \text{and} \quad \int \frac{1}{G} dF < \infty.$$

Under the assumption (9), it is straightforward to show that  $v_e$  defined by the right side of (8) is pathwise differentiable with derivative  $\dot{v}_e$  given by

$$(10) \quad \dot{v}_e(h)(t) = \iint \left\{ \frac{1}{M^*(v)} 1_{[v \leq t]} - \int_0^t \frac{1_{[u \leq s < v]}}{M^*(s)^2} dG^*(s) \right\} h(u, v) dP(u, v)$$

for  $0 \leq t \leq \tau$  and  $h \in \dot{\mathbf{M}}$ . Hence in view of (5.2.26),  $\dot{v}_{et} = \dot{v}_e^T(\pi_t)$  is given by

$$(11) \quad \begin{aligned} \dot{v}_{et}(u, v) &= \frac{1}{M^*(v)} 1_{[v \leq t]} - \int_0^t \frac{1_{[u \leq s < v]}}{M^*(s)^2} dG^*(s) \\ &= \frac{1}{M^*(v)} 1_{[v \leq t]} - \{C_v(t \wedge v) - C_v(t \wedge u)\}, \end{aligned}$$

where

$$(12) \quad C_v(t) \equiv \int_0^t \frac{1}{M^*{}^2} dG^* = \alpha \int_0^t \frac{1}{G^2 F} dG = \int_0^t \frac{1}{M^*} d\Lambda_G.$$

To use (10) or (11) to find  $\dot{v}_t = \tilde{\mathbf{I}}_v(\pi_t)$  of (5.2.25), (5.2.8), and (5.2.7), we first need to know more about  $\dot{\mathbf{P}}$ . It follows immediately from theorem 4.4.3.A that the left truncation model is differentiable with score operator  $\dot{\mathbf{i}} : \dot{\mathbf{G}}_1 \times \dot{\mathbf{G}}_2 \rightarrow L_2^0(P)$  given by

$$\dot{\mathbf{i}}(a, b) = \dot{\mathbf{i}}_1 a + \dot{\mathbf{i}}_2 b, \quad a \in \dot{\mathbf{G}}_1 = L_2^0(F), \quad b \in \dot{\mathbf{G}}_2 = L_2^0(G),$$

where

$$(13) \quad \begin{aligned} \dot{\mathbf{i}}_1 a(u, v) &= a(u) - Ea(U) = a(u) - \int a dF^*, \\ \dot{\mathbf{i}}_2 b(u, v) &= b(v) - Eb(V) = b(v) - \int b dG^*. \end{aligned}$$



Thus the tangent set for the model  $\mathbf{P}$  satisfies

$$(14) \quad \begin{aligned} \dot{\mathbf{P}}^0 &\supset \mathbf{R}(\dot{\mathbf{i}}) \\ &= \left\{ \int a(u) + b(v) : a \in L_2(F), \int a dF^* = 0, b \in L_2(G), \int b dG^* = 0 \right\} \\ &\subset L_2^0(F^*) + L_2^0(G^*), \end{aligned}$$

where, as usual,  $L_2^0(F^*) = \{a \in L_2(F^*) : \int a dF^* = 0\}$  and similarly for  $L_2^0(G^*)$ . The last inclusion in (14) follows by (3) and (4):

$$\int a^2 dF^* = \frac{1}{\alpha} \int a^2 \bar{G} dF \leq \frac{1}{\alpha} \int a^2 dF < \infty,$$

and

$$\int b^2 dG^* = \frac{1}{\alpha} \int b^2 F dG \leq \frac{1}{\alpha} \int b^2 dG < \infty.$$

But since any function of the form  $a(u) + b(v)$  with  $\int a dF^* = 0 = \int b dG^*$  and  $a$  and  $b$  bounded, is contained in both  $\mathbf{R}(\dot{\mathbf{i}})$  and  $L_2^0(F^*) + L_2^0(G^*)$  of (14),

$$(15) \quad \dot{\mathbf{P}} \supset \mathbf{R}(\dot{\mathbf{i}}) = \overline{L_2^0(F^*) + L_2^0(G^*)} \supset L_2^0(F^*) + L_2^0(G^*).$$

In fact, as we will show at the end of this section, (15) can be sharpened under (9) to

$$(16) \quad \dot{\mathbf{P}} \supset \mathbf{R}(\dot{\mathbf{i}}) = L_2^0(F^*) + L_2^0(G^*) = \dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2.$$

But (15) implies that under the assumption (9),  $\dot{v}_{et}$  given in (11) is (more by accident than by design) already in  $\dot{\mathbf{P}}$ , and hence no projection is necessary to calculate  $\dot{v}_t$  in spite of the fact that  $\dot{\mathbf{P}}$  is a proper subset of  $\dot{\mathbf{M}} = L_2^0(P)$ :

$$\dot{v}_{et}(u, v) = C_v(t \wedge u) + \left\{ \frac{1}{M^*(v)} 1_{[v \leq t]} - C_v(t \wedge v) \right\}$$

$$\equiv a_t(u) + b_t(v),$$

where  $a_t \in L_2(F^*)$ ,  $b_t \in L_2(G^*)$ , since  $t \leq \tau < \tau_G$ . It follows that the efficient influence function  $\dot{\mathbf{l}}$  for  $v$  is given by  $\pi_t \dot{\mathbf{l}} = \dot{\mathbf{l}}_v(\pi_t) = \dot{v}_t = \dot{v}_{et}$ , and hence we obtain from (5.2.23), after two integrations by parts,

$$(17) \quad \Gamma_v^{-1}(s, t) = E(\dot{v}_s(U, V) \dot{v}_t(U, V)) = C_v(s \wedge t)$$

where  $C_v$  is given by (12).

Thus the Gaussian process  $\mathbf{Z}_0$  of theorem 5.2.1 is

$$(18) \quad \mathbf{Z}_0(t) = B(C_v(t)), \quad 0 \leq t \leq \tau < \tau_G,$$

where  $B$  is a standard Brownian motion.

Of course, this carries over to give a bound for estimation of  $\kappa \equiv 1 - G$  on  $[0, \tau]$ : since

$$\kappa(P_{(F,G)})(t) \equiv 1 - G(t) = \exp(-\Lambda_G(t)) \equiv \phi(v(P_{(F,G)})(t))$$

for continuous  $\Lambda_G$  where  $\phi(b) \equiv \exp(-b) \in l^\infty([0, \tau])$  for  $b \in l^\infty([0, \tau])$ , and since  $\dot{\phi} : l^\infty \rightarrow l^\infty$  is the multiplication operator

$$\dot{\phi}(b)(h) = -\exp(-b)h = -\phi h,$$

the derivative  $\dot{\kappa}$  is simply related to  $\dot{v}$  by the chain rule for pathwise derivatives. We compute

$$(19) \quad \dot{\kappa}_t(u, v) = \dot{\phi}(\Lambda_G) \dot{v}_t(u, v) = -\bar{G}(t) \dot{v}_t(u, v).$$

Therefore the inverse information covariance function (5.2.23) for estimation of  $\kappa = \bar{G}$  is, via (17), with  $1 - K = \bar{K} \equiv 1/(1 + C_v)$ ,

$$(20) \quad I_{\kappa}^{-1}(s, t) = \bar{G}(s) \bar{G}(t) C_v(s \wedge t) \\ = \frac{\bar{G}}{\bar{K}}(s) \frac{\bar{G}}{\bar{K}}(t) (K(s \wedge t) - K(s)K(t)).$$

Thus for estimation of  $\kappa$  the process  $ZZ_0$  of theorem 5.2.1 is

$$(21) \quad ZZ_0(t) = -\bar{G}(t)B(C_v(t)) = -\frac{\bar{G}(t)}{\bar{K}(t)}B_0(K(t)),$$

where  $B$  is standard Brownian motion and  $B_0(t) \equiv (1 - t)B(t/(1 - t))$  is the Brownian bridge.

What about bounds for estimation of  $F$ ? These follow in a symmetric fashion via estimation of the "reverse" or "retro" cumulative hazard function for  $F$  defined as follows:

$$(22) \quad \gamma(P_{(F,G)})(t) \equiv \bar{\Lambda}_F(t) = \int_t^\infty \frac{dF}{F}, \quad \bar{\tau}_F < \tau \leq t < \infty,$$

where  $\bar{\tau}_F \equiv \sup\{s : F(s) = 0\}$ . Then

$$(23) \quad \xi(P_{(F,G)})(t) \equiv F(t) = \prod_{s \geq t} (1 - d\bar{\Lambda}_F(s)) \\ \equiv \exp(-\bar{\Lambda}_F^c(t)) \prod_{s \geq t} (1 - \Delta\bar{\Lambda}_F(s)) \\ = \exp(-\bar{\Lambda}_F(t)) \text{ if } \bar{\Lambda}_F \text{ is continuous.}$$

As in (8),

$$(24) \quad \gamma_e(P_{(F,G)})(t) \equiv \int_t^\infty \frac{1}{M^*} dF^*$$

defines an extension  $\gamma_e$  of  $\gamma$  which is pathwise differentiable, with

$$\dot{\gamma}_{et}(u, v) = \frac{1}{M^*(u)} 1_{[t \leq u]} - \int_t^\infty \frac{1_{[u \leq s < v]}}{M^*(s)^2} dF^*(s)$$

$$= \frac{1}{M^*(u)} 1_{[t \leq u]} - \{C_\gamma(t \vee u) - C_\gamma(t \vee v)\},$$

where

$$C_\gamma(t) \equiv \int_t^\infty \frac{1}{M^{*2}} dF^* = \alpha \int_t^\infty \frac{1}{GF^2} dF = \int_t^\infty \frac{1}{M^*} d\bar{\Lambda}_F$$

is a decreasing function. Note that (24) implies identifiability of  $F$ . As before,  $\dot{\gamma}_{et} \in \dot{P}$  if (9) holds, and hence  $\dot{\Gamma}_\gamma(\pi_t) = \dot{\gamma}_t = \dot{\gamma}_{et}$ , and the inverse information covariance function (5.2.23) is

$$(25) \quad I_\gamma^{-1}(s, t) = E(\dot{\gamma}_s(U, V) \dot{\gamma}_t(U, V)) = C_\gamma(s \vee t)$$

after two integrations by parts. So for estimation of  $\gamma = \bar{\Lambda}_F$  the process  $\mathbf{Z}_0$  of theorem 5.2.1 is

$$\mathbf{Z}_0(t) = B(C_\gamma(t)), \quad \bar{\tau}_F < \tau \leq t < \infty.$$

In parallel to (19)–(21),

$$(26) \quad \dot{\xi}_t(u, v) = -F(t) \dot{\gamma}_t(u, v)$$

and

$$(27) \quad I_\xi^{-1}(s, t) = F(s)F(t)C_\gamma(s \vee t),$$

so for estimation of  $\xi \equiv F$ , the process  $\mathbf{Z}_0$  of theorem 5.2.1 is

$$(28) \quad \mathbf{Z}_0(t) = -F(t)B(C_\gamma(t)) = -\frac{F(t)}{\bar{J}(t)} B_0(J(t)), \quad \bar{\tau}_F < \tau \leq t < \infty,$$

where  $B$  is standard Brownian motion,  $\bar{J} \equiv 1/(1 + C_\gamma)$ , and  $B_0$  is Brownian bridge.

The information bound calculations as carried out above are straightforward and elementary, but somewhat tedious. A more sophisticated approach, via martingale theory, is less elementary, but avoids the tedious calculations. We sketch this approach. Consider the counting processes

$$N_1(t) \equiv 1_{[U \leq t]}, \quad N_2(t) \equiv 1_{[U < V \leq t]}, \quad t \geq 0,$$

with their natural history (or self-exciting) filtration  $\{\mathcal{F}_t\}_{t \geq 0}$  given by

$$\mathcal{F}_t \equiv \sigma\{1_{[U \leq s]}, 1_{[U < V \leq s]} : 0 \leq s \leq t\} \quad \text{for } t \geq 0.$$

Then, under  $P$ ,  $N_1$  and  $N_2$  have (predictable) compensators  $A_1$  and  $A_2$  given by

$$A_1(t) \equiv \int_0^t 1_{[U \geq s]} d\Lambda_F(s)$$

and

$$A_2(t) \equiv \frac{1}{\alpha} \int_0^t 1_{[V \geq s > U]} d\Lambda_G(s),$$

so

$$M_i \equiv N_i - A_i, \quad i = 1, 2,$$

are square integrable (orthogonal) martingales with predictable variation processes  $\langle M_i \rangle = A_i, i = 1, 2,$  and  $\langle M_1, M_2 \rangle = 0.$  Then it is easy to write  $\dot{v}_t$  and  $\dot{k}_t$  in terms of  $M_2$  as

$$\dot{v}_t(U, V) = \int_0^t \frac{1}{M^*} dM_2;$$

i.e., the process  $\{\dot{v}_t(U, V) : 0 \leq t \leq \tau\}$  is a martingale transform of  $M_2.$  Hence by martingale calculus

$$\begin{aligned} \Gamma_V^{-1}(s, t) &= E \dot{v}_s(U, V) \dot{v}_t(U, V) = E \int_0^{s \wedge t} \frac{1}{M^{*2}} d\langle M_2 \rangle \\ &= \int_0^{s \wedge t} \frac{1}{M^*} d\Lambda_G = C_v(s \wedge t) \end{aligned}$$

as in (17), and these calculations extend immediately to yield (20) as well.

The calculations for  $\gamma \equiv \bar{\Lambda}_F$  and  $\xi \equiv F$  can be done in an analogous way using the “reverse time” or “retro” counting processes

$$\bar{N}_1(t) \equiv 1_{[t \leq U < V]}, \quad \bar{N}_2(t) \equiv 1_{[t \leq V]}, \quad t \geq 0,$$

with natural history filtration  $\{\bar{\mathcal{F}}_t\}_{t \geq 0}$  and compensators

$$\bar{A}_1(t) \equiv \frac{1}{\alpha} \int_t^\infty 1_{[U \leq s < V]} d\bar{\Lambda}_F(s)$$

and

$$\bar{A}_2(t) \equiv \int_t^\infty 1_{[V \leq s]} d\bar{\Lambda}_G(s).$$

Note that under (9), we have both

$$\begin{aligned} \sup_{0 \leq t < \infty} \Gamma_{\bar{\kappa}}^{-1}(t, t) &= \sup_{0 \leq t < \infty} \alpha \bar{G}^2(t) \int_0^t \frac{1}{\bar{G}^2 F} dG \\ &\leq \alpha \int_0^\infty \frac{1}{F} dG < \infty \end{aligned}$$

and

$$\begin{aligned} \sup_{0 \leq t < \infty} \Gamma_{\bar{\xi}}^{-1}(t, t) &= \sup_{0 \leq t < \infty} \alpha F^2(t) \int_t^\infty \frac{1}{\bar{G} F^2} dF \\ &\leq \alpha \int_0^\infty \frac{1}{\bar{G}} dF < \infty. \end{aligned}$$

This suggests that when (9) holds the maps  $\kappa(P_{(F,G)}) = \bar{G}$  and  $\xi(P_{(F,G)}) = F$  from  $\mathbf{P}$  to  $l^\infty(R^+)$  are both pathwise differentiable; i.e. the restrictions to  $t \leq \tau < \tau_G$  and  $\bar{\tau}_F < \tau \leq t$  are unnecessary for  $\kappa$  and  $\xi$  respectively. We will verify this in example 2 by taking a different approach.

In view of (8) and (24), it is natural to estimate  $v = \Lambda_G$  and  $\gamma = \bar{\Lambda}_F$  by the Nelson-Aalen estimators

$$(29) \quad \hat{v}_n(t) \equiv \int_0^t \frac{1_{[M_n^*(s-) > 0]}}{M_n^*(s-)} dG_n^*(s)$$

and

$$(30) \quad \hat{\gamma}_n(t) \equiv \int_t^\infty \frac{1_{[M_n^*(s) > 0]}}{M_n^*(s)} dF_n^*(s),$$

where

$$F_n^*(t) \equiv n^{-1} \sum_{i=1}^n 1_{[U_i \leq t]}, \quad G_n^*(t) \equiv n^{-1} \sum_{i=1}^n 1_{[V_i \leq t]},$$

and

$$M_n^*(t) \equiv (F_n^* - G_n^*)(t) = n^{-1} \sum_{i=1}^n 1_{[U_i \leq t < V_i]}.$$

As shown by Woodroffe (1985), Wang, Jewell, and Tsai (1986), Keiding and Gill (1990), and Van der Vaart (1991), under the assumption (9) these estimators achieve the bounds (17) and (25).

Of course, the corresponding (product limit) estimators of  $\bar{G}$  and  $F$  are just

$$(31) \quad \bar{G}_n(t) \equiv \prod_{s \leq t} (1 - \Delta \hat{v}_n(s)), \quad 0 \leq t < \infty,$$

and

$$(32) \quad F_n(t) \equiv \prod_{s \geq t} (1 - \Delta \hat{\gamma}_n(s)), \quad 0 < t \leq \infty,$$

and these estimators are both efficient if (9) holds.

What if (9) fails? Then,  $v(P_{(F,G)}) = \Lambda_G$ ,  $\kappa = \bar{G}$ ,  $\gamma(P_{(F,G)}) = \bar{\Lambda}_F$ , and  $\xi(P_{(F,G)}) = F$  are, in view of proposition 2 given in example 2 below, *not* pathwise differentiable and theorem 5.2.1 does *not* apply. In this case we can only prove differentiability of the increments of  $\Lambda_G$  starting at some fixed  $\varepsilon > 0$ . That is, for  $\varepsilon > 0$  we define

$$(33) \quad v_\varepsilon(P)(t) \equiv \int_\varepsilon^t \frac{1}{1 - G_-} dG = \Lambda_G(t) - \Lambda_G(\varepsilon), \quad \varepsilon \leq t \leq \tau < \tau_G.$$

The corresponding parameter  $\kappa_\varepsilon$  is the conditional distribution of  $V$  given  $V > \varepsilon$ :

$$(34) \quad \kappa_\varepsilon(P)(t) \equiv \frac{1 - G(t)}{1 - G(\varepsilon)}, \quad \varepsilon \leq t < \infty. \quad \square$$

### Example 2. Random truncation model; score operator approach.

We reconsider bounds for estimation of  $\kappa = \bar{G}$  and  $\xi = F$  using score and information operator methods of sections 5.4 and 5.5. We also use the  $R$  and  $L$  operators introduced in appendices A.1 and A.3.

As seen in example 1, the random truncation model is differentiable with score operators as given in (13). We will continue to assume that (9) holds, so that by (16)

$$(35) \quad \dot{\mathbf{P}} = L_2^0(F^*) + L_2^0(G^*) \equiv \dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2.$$

Thus, after straightforward calculation via Fubini's theorem, for  $h(u, v) = a(u) + b(v) \in \dot{\mathbf{P}}$ , the adjoints  $\dot{\mathbf{I}}_1^T$  and  $\dot{\mathbf{I}}_2^T$  of the score operators  $\dot{\mathbf{I}}_1$  and  $\dot{\mathbf{I}}_2$  given in (13) are

$$(36) \quad \begin{aligned} \dot{\mathbf{I}}_1^T h(u) &= \frac{1}{\alpha} \int_u^\infty h(u, v) dG(v) \\ &= \frac{1}{\alpha} \bar{G}(u) a(u) + \frac{1}{\alpha} \int_u^\infty b dG \end{aligned}$$

and

$$(37) \quad \begin{aligned} \dot{\mathbf{I}}_2^T h(v) &= \frac{1}{\alpha} \int_0^v h(u, v) dF(u) \\ &= \frac{1}{\alpha} \int_0^v a dF + \frac{1}{\alpha} F(v) b(v). \end{aligned}$$

Of course,  $\dot{\mathbf{P}}_1 = L_2^0(F^*)$  is *not* orthogonal to  $\dot{\mathbf{P}}_2 = L_2^0(G^*)$ , and

$$\dot{\mathbf{I}}_2^T \dot{\mathbf{I}}_1 a(v) = \frac{1}{\alpha} \left\{ \int_0^v a dF - F(v) E a(U) \right\}$$

does not vanish, typically. Now we use (36) and (37) to strengthen and extend the differentiability of  $\kappa \equiv \bar{G}$  and  $\xi = F$  given in example 1. The following proposition proved at the end of this section and due to Van der Vaart (1991) shows that these functions are pathwise differentiable as maps from  $\mathbf{P}$  to  $\mathbf{B} = l^\infty(R^+)$  if and only if (9) holds.

**Proposition 1.** In the random left truncation model the following are equivalent:

- A.  $\int_0^\infty (1/F) dG < \infty$  and  $\int_0^\infty (1/\bar{G}) dF < \infty$  (i.e., (9) holds).
- B.  $\kappa(P_{(F,G)}) = \bar{G}$  and  $\xi(P_{(F,G)}) = F$  are pathwise differentiable maps from  $\mathbf{P}_2$  and  $\mathbf{P}_1$  respectively to  $l^\infty(R^+)$ .
- C. For some  $t \in R^+$ ,  $\pi_t \kappa(P_{(F,G)}) = \bar{G}(t)$  and  $\pi_t \xi(P_{(F,G)}) = F(t)$  are pathwise differentiable maps from  $\mathbf{P}_2$  and  $\mathbf{P}_1$  respectively to  $R$ .

To calculate the information bounds for estimation of  $\kappa(P_{(F,G)}) = 1 - G$  via the semiparametric, or score operator, approach, we need to use the results for "models with composite  $g$ " in section 5.5 and calculate

$$(38) \quad \dot{\mathbf{I}}_2^* \equiv (I - P_1) \dot{\mathbf{I}}_2,$$

as defined in (5.5.14), where  $P_1 = \Pi_0(\cdot | \dot{\mathbf{P}}_1)$  is the projection onto  $\dot{\mathbf{P}}_1$ . But

$$(39) \quad \dot{\mathbf{P}}_1 = \overline{\mathbf{R}(\dot{\mathbf{I}}_1)} = L_2^0(F^*),$$

and hence

$$(40) \quad P_1 h = E(h | U) - Eh \quad \text{for } h \in L_2(P).$$

Thus

$$P_1 \dot{I}_2 b(u) = \frac{\int_u^\infty (b(v) - Eb(V)) dG(v)}{\bar{G}(u)},$$

and

$$(41) \quad I_2^* b(u, v) = b(v) - \frac{\int_u^\infty b dG}{\bar{G}(u)}.$$

We can now apply corollary 5.5.1.A, by first calculating  $I_2^{*T}$  and then verifying (5.5.16).

By straightforward calculation we find, for  $h(u, v) = a(u) + b(v) \in \dot{P}$ ,

$$(42) \quad I_2^{*T} h(v) = \frac{1}{\alpha} \{ F(v)b(v) - \int b(t) \int_0^{t \wedge v} \frac{dF}{\bar{G}} dG(t) \}$$

$$(43) \quad = L_G D R_G b(v),$$

where

$$L_G b(v) = b(v) - \int_0^v \frac{b}{\bar{G}} dG,$$

$$R_G b(v) = b(v) - \frac{\int_v^\infty b dG}{\bar{G}(v)}, \quad b \in L_2(G),$$

are the  $L$  and  $R$  operators for the df  $G$  as defined in (A.1.19) and (A.1.20), and  $D$  is the diagonal (or multiplication) operator

$$(44) \quad Db(v) = \frac{F(v)}{\alpha} b(v).$$

(This should be compared with formulas (6.6.27) and (6.6.34) for  $\dot{I}_1^T \dot{I}_1$  in the case of the random censoring model, examples 2 and 1.C.) Since  $\chi(g) \equiv \bar{G}$  has derivative  $\dot{\chi}$  with  $\dot{\chi}^T(\pi_t) = \dot{\chi}_t$  given by (cf. (5.3.2))

$$(45) \quad \dot{\chi}_t(v) = - (1_{[0,t]}(v) - G(t)),$$

to find  $\tilde{I}_\kappa(\pi_t)$  for  $\kappa(P_{(F,G)}) \equiv \bar{G}$  we want to solve, as in (5.5.16),

$$(46) \quad \dot{\chi}_t = I_2^{*T} \tilde{I}_\kappa(\pi_t)$$

for  $\tilde{I}_\kappa(\pi_t) = \tilde{I}_\kappa^{(1)}(\pi_t) + \tilde{I}_\kappa^{(2)}(\pi_t)$  with  $\tilde{I}_\kappa^{(i)}(\pi_t) \in \dot{P}_i$ . Of course,  $I_2^{*T}(\dot{P}_1) = \{0\}$ , so

$$(47) \quad I_2^{*T} \tilde{I}_\kappa(\pi_t) = I_2^{*T} \tilde{I}_\kappa^{(2)}(\pi_t),$$

so we can at least solve for  $\tilde{I}_\kappa^{(2)}(\pi_t) \in \dot{P}_2$ , and then proceed to find  $\tilde{I}_\kappa(\pi_t)$ .

Now

$$(48) \quad R_G \dot{\chi}_t = - \frac{\bar{G}(t) 1_{[0,t]}}{\bar{G}},$$

and hence (43), (46), (47), and (48) yield

$$(49) \quad \begin{aligned} \tilde{\mathbf{I}}_{\kappa}^{(2)}(\pi_t) - E_G \tilde{\mathbf{I}}_{\kappa}^{(2)}(\pi_t) &= L_G D^{-1} R_G (\dot{\chi}_t) \\ &= - \bar{G}(t) L_G \left( \frac{\alpha 1_{[0,t]}}{F \bar{G}} \right) = - \bar{G}(t) L_G \left( \frac{1_{[0,t]}}{M^*} \right), \end{aligned}$$

using  $R_G \circ L_G b = b$  and  $L_G \circ R_G b = b - E_G b$  from proposition A.1.8.B and (5); or by (12), for some  $c$ ,

$$(50) \quad \begin{aligned} \tilde{\mathbf{I}}_{\kappa}^{(2)}(\pi_t)(v) &= - \bar{G}(t) \left\{ \frac{1}{M^*(v)} 1_{[0,t]}(v) - \int_0^v \frac{1_{[0,t]}}{M^*} d\Lambda_G \right\} + c \\ &= - \bar{G}(t) \left\{ \frac{1}{M^*(v)} 1_{[0,t]}(v) - C_v(t \wedge v) \right\} + c. \end{aligned}$$

Furthermore, note that (cf. (36))  $\dot{\mathbf{I}}_1^T : \dot{\mathbf{P}} \rightarrow L_2^0(F)$  is given by

$$\dot{\mathbf{I}}_1^T h(u) = \frac{\bar{G}(u)}{\alpha} E(h \mid U = u).$$

Consequently, the second equation in (5.5.16) becomes

$$0 = E(\tilde{\mathbf{I}}_{\kappa}(\pi_t) \mid U) = \tilde{\mathbf{I}}_{\kappa}^{(1)}(\pi_t) + E(\tilde{\mathbf{I}}_{\kappa}^{(2)}(\pi_t) \mid U),$$

and hence

$$(51) \quad \tilde{\mathbf{I}}_{\kappa}(\pi_t) = \tilde{\mathbf{I}}_{\kappa}^{(2)}(\pi_t) - E(\tilde{\mathbf{I}}_{\kappa}^{(2)}(\pi_t) \mid U).$$

But

$$E(\tilde{\mathbf{I}}_{\kappa}^{(2)}(\pi_t) \mid U = u) = \bar{G}(t) C_v(t \wedge u) + c$$

by direct calculation, so combining this with (50) and (51) gives

$$(52) \quad \begin{aligned} \dot{\kappa}_t(u, v) &= \tilde{\mathbf{I}}_{\kappa}(\pi_t)(u, v) \\ &= - \bar{G}(t) \left\{ \frac{1}{M^*(v)} 1_{[0,t]}(v) - (C_v(t \wedge v) - C_v(t \wedge u)) \right\} \end{aligned}$$

in agreement with (19) in example 1.

Note that we have solved (46) and hence (5.5.16) for  $b^* = \pi_t$ . Consequently  $\dot{\chi}_t^T(\pi_t) \in \mathbf{R}(I_2^{*T})$  and corollary 5.5.1.A yields the pathwise differentiability of  $\kappa(P_{(F,G)})(t) = \bar{G}(t)$  once more.

Now suppose that the df  $F$  is known, and consider information bounds for estimation of  $v \equiv \Lambda_G$ . Thus  $\mathbf{P} = \mathbf{P}_2$  and  $\dot{\mathbf{P}} = \dot{\mathbf{P}}_2$ . In the nonparametric approach we would project  $\dot{v}_t = \dot{v}_{et}$  of (11) onto  $\dot{\mathbf{P}}_2$  to get the new efficient influence function. Since  $P_2 = E(\cdot \mid V)$ , this is easy to compute, and we find that the efficient influence function  $\mathbf{I}(\cdot \mid \mathbf{P}_2)$  satisfies



$$\begin{aligned}
 \pi_t \tilde{I}(u, v | P_2) &= \tilde{I}_v(\pi_t)(u, v | P_2) \\
 &= \frac{1_{[v \leq t]}}{M^*(v)} - \frac{1}{F(v)} \int_0^{t \wedge v} \frac{F}{M^*} d\Lambda_G \\
 (53) \qquad \qquad \qquad &= \frac{\alpha(1_{[v \leq t]} - G(t))}{F(v)\bar{G}(t)}
 \end{aligned}$$

by direct calculation using  $F/M^* = \alpha/\bar{G}$  by (5).

Taking the score operator approach, we need to solve, by (5.4.7) of theorem 5.4.1,

$$(54) \quad \dot{\chi}_t = \dot{I}_2^T \tilde{I}_v(\pi_t)$$

for  $\tilde{I}_v(\pi_t)(u, v | P_2) = \tilde{I}_v(\pi_t)(v | P_2) \in \dot{P}_2$ , where by (45) and (19),

$$(55) \quad \dot{\chi}_t(v) = \frac{1_{[v \leq t]} - G(t)}{\bar{G}(t)}$$

holds. But for  $h(u, v) = b(v) \in \dot{P}_2$ , it follows from (37) that

$$(56) \quad \dot{I}_2^T b(v) = \frac{1}{\alpha} F(v) b(v),$$

so we want to solve

$$(57) \quad \frac{1}{\alpha} F(v) \tilde{I}_v(\pi_t)(v | P_2) = \dot{\chi}_t(v).$$

Hence  $\tilde{I}_v(\pi_t)(v | P_2)$  is just as in (53). In any case,

$$\begin{aligned}
 \Gamma_v^{-1}(s, t | P_2) &= E \tilde{I}_v(\pi_s) \tilde{I}_v(\pi_t) \\
 (58) \qquad \qquad \qquad &= \frac{\alpha^2}{\bar{G}(s)\bar{G}(t)} E \frac{(1_{[v \leq s]} - G(s))(1_{[v \leq t]} - G(t))}{F^2(v)} \\
 &= \frac{\alpha B(\infty)}{\bar{G}(s)\bar{G}(t)} \{H(s \wedge t) - G(s)H(t) \\
 &\qquad \qquad \qquad - G(t)H(s) + G(s)G(t)\},
 \end{aligned}$$

where

$$B(t) \equiv \int_0^t \frac{dG}{F} \quad \text{and} \quad H(t) \equiv \frac{B(t)}{B(\infty)}.$$

The corresponding bound for estimation of  $\kappa \equiv \bar{G}$  is given by (cf. (19))

$$\begin{aligned}
 (59) \quad \Gamma_\kappa^{-1}(s, t | P_2) \\
 = \alpha B(\infty) \{H(s \wedge t) - G(s)H(t) - G(t)H(s) + G(s)G(t)\}.
 \end{aligned}$$

Construction of an estimator achieving this bound is straightforward. When  $F$  is known, the conditional density of  $U$  given  $V = v$  is  $p(u|v) = f(u)/F(v)$  for  $u < v$ , which does not depend on  $g$ , so  $V$  is sufficient for  $g$ , and therefore it suffices to consider the (marginal) distribution of  $V$ . From (4) this is just the general biased sampling model with new stratum weight function  $w \equiv F$  on the sample space  $R^+$ , and the results of Vardi (1985) and Gill, Vardi, and Wellner (1988) apply. In this simple case, the nonparametric maximum likelihood estimator  $\hat{G}_n$  of  $G$  is obtained by solving (4) for  $G$  and estimating  $G^*$  by the empirical df  $G_n^*$  of the  $V$ 's: inversion of (4) yields

$$(60) \quad G(v) = \frac{\int_0^v \frac{1}{F} dG^*}{\int_0^\infty \frac{1}{F} dG^*},$$

and hence the "reduced sample estimator"

$$(61) \quad \hat{G}_n(v) = \frac{\int_0^v \frac{1}{F} dG_n^*}{\int_0^\infty \frac{1}{F} dG_n^*},$$

of  $G$ . This estimator achieves the bound given by (59) as shown by Gill, Vardi, and Wellner (1988, example 4.2). (When  $w = F$ , equation (4.8) in Gill, Vardi, and Wellner (1988) reduces to (59).)  $\square$

Table 1 gives a summary and comparison of the random truncation model and the random censoring model which will be studied in section 6.6.

Random Censoring	Random Truncation
$\dot{P} = \dot{M}$	$\dot{P} \subset \dot{M}$
saturated	unsaturated
$\dot{P} = \dot{P}_1 + \dot{P}_2$	$\dot{P} = \dot{P}_1 + \dot{P}_2$
$\dot{P}_1 \perp \dot{P}_2$	$\dot{P}_1$ and $\dot{P}_2$ are not orthogonal
Bound for $F$ if $G$ is known is unchanged	Bound for $F$ if $G$ is known changes
Reduced sample estimator inefficient	Reduced sample estimator efficient

Table 1. Comparison of Random Censoring and Random Truncation.

### A Formal Principle

Before proceeding with example 2, we briefly discuss a useful principle. Suppose that

$$P = \{P_{(\theta, G)} : \theta \in \Theta \subset R^k, G \in \mathcal{G}\}$$

is a semiparametric model. Suppose that we know the efficient influence functions for estimation of  $P$  in the submodel

$$(62) \quad \mathbf{P}_2 \equiv \{P_{(\theta_0, G)} : G \in \mathbf{G}\}$$

for fixed  $\theta_0 \in \Theta$ . Then we can calculate all projections of interest in the full model  $\mathbf{P}$ , at least in principle. First note that if we can calculate projections on  $\dot{\mathbf{P}}_2$ , then we can project  $\dot{\mathbf{I}}_1$  on  $\dot{\mathbf{P}}_2$ , determine the efficient score function  $\mathbf{I}_1^*$  for  $\theta$ , and hence calculate projections on  $\dot{\mathbf{P}} = [\mathbf{I}_1^*] + \dot{\mathbf{P}}_2$  in a straightforward manner since  $\mathbf{I}_1^* \perp \dot{\mathbf{P}}_2$ . Here is an outline of why it should be possible to compute projections on  $\dot{\mathbf{P}}_2$ . Regard

$$(63) \quad v \equiv P_{(\theta_0, G)}(dx)$$

as a real parameter with efficient influence function

$$(64) \quad \tilde{\mathbf{I}}_v(X) = \gamma(dx, X)$$

at  $P_0 \equiv P_{(\theta_0, G_0)}$ . That is, if  $\mathbf{X} = R$  and  $\kappa \equiv P_{(\theta_0, G)}((-\infty, x])$  has efficient influence function  $\tilde{\mathbf{I}}_\kappa(X) = v(x, X)$ , then

$$(65) \quad \tilde{\mathbf{I}}_v(X) = \gamma(dx, X) = \frac{\partial}{\partial x} v(x, X) dx.$$

Consider estimation of  $\int h dP_{(\theta, G)}$  for a bounded function  $h$ . The natural non-parametric estimator of this parameter, namely  $\int h dP_n$ , has influence function  $h - E_0 h$  at  $P_0$ , and hence, by (3.3.28), the efficient influence function for estimation of  $\int h dP$  in  $\mathbf{P}_2(P_0)$  is

$$(66) \quad \Pi_0(h | \dot{\mathbf{P}}_2(P_0))(X).$$

On the other hand, it is also

$$(67) \quad \int h(x) \gamma(dx, X),$$

and, for  $h \in L_2^0(P_0)$  we expect

$$(68) \quad \Pi_0(h | \dot{\mathbf{P}}_2(P_0))(X) = \int h(x) \gamma(dx, X).$$

This is essentially what is going on in the following example: if the parametric part of the model is known, the truncated regression model 2 below reduces to the left truncation model studied in example 1 where projections are known explicitly. Hence we also expect to be able to calculate explicitly for the larger model. However, we will carry out the calculations only to determine  $\mathbf{I}_1^*$ .

**Example 3. Truncated regression.**

Now consider the truncated regression model of example 4.4.3. Here we relabel slightly, and write the density of an observation  $X = (Y, Z) \in R^2$  as

$$(69) \quad p(y, z, v, F, H) = 1_{[y \leq y_0]} \frac{f(y - vz) h(z)}{W(v, F, H)},$$

where  $y_0$  is a fixed number ( the threshold or truncation parameter) and

$$(70) \quad W(v, F, H) \equiv Q(Y \leq y_0) = \iint_{vz + \varepsilon \leq y_0} dF(\varepsilon) dH(z)$$

$$= \int F(y_0 - vZ) dH(z);$$

here  $Q$  denotes the product measure  $F \times H$  for  $(\epsilon, Z)$ . As noted in section 4.4,  $\sqrt{n}$ -consistent estimators for  $v$  based on ranks were studied by Bhattacharya, Chernoff, and Yang (1983), who were motivated by problems in astronomy.

Let  $v$ ,  $f$  (or  $F$ ), and  $h$  (or  $H$ ) be labeled by 1, 2, and 3 respectively. We want to consider estimation of the backward cumulative hazard function corresponding to  $F$ , defined by

$$(71) \quad \gamma(P_{(v,F,H)})(t) \equiv \bar{\Lambda}_F(t) = \int_t^\infty \frac{dF}{F}$$

as in (22), and

$$(72) \quad \xi(P_{(v,F,H)})(t) \equiv F(t) = \prod_{s \geq t} (1 - d\bar{\Lambda}_F(s))$$

as in (23).

Suppose first that  $v \neq 0$  is known (so that  $P \in P_{23}$ ); we suppose without loss of generality that  $v > 0$ . Then the model is the same as in examples 1 and 2: we can rewrite and reparametrize the density (69) with  $v$  known as

$$(73) \quad p(y, z, F, H) = 1_{[e \leq y_0 - vz]} \frac{f(\epsilon)h(z)}{W(F, H)} = 1_{[e \leq v]} \frac{f(\epsilon)g(v)}{W(F, G)},$$

where  $V \equiv y_0 - vZ$  and  $g(v) = v^{-1}h((y_0 - v)/v)$  is the density of  $V$ ; or equivalently,  $(U, V) \equiv (Y - vZ, y_0 - vZ) = (\epsilon, y_0 - vZ)$  has density given by (2). Thus the results of examples 1 and 2 apply to estimation of  $\gamma = \bar{\Lambda}_F$ : by (25), assuming that (9) holds,

$$(74) \quad I_\gamma^{-1}(s, t) = C_\gamma(s \vee t)$$

with

$$C_\gamma(t) \equiv \int_t^\infty \frac{1}{M^*} d\bar{\Lambda}_F = W \int_t^\infty \frac{1}{F\bar{G}} d\bar{\Lambda}_F$$

and where

$$\bar{G}(v) = P(y_0 - vZ \geq v) = P(Z \leq \frac{y_0 - v}{v}) = H(\frac{y_0 - v}{v}).$$

Thus for estimation of  $\gamma$  the process  $Z_0$  of theorem 5.2.1 is

$$Z_0(t) \equiv B(C_\gamma(t)), \quad 0 < \tau \leq t < \infty.$$

Similarly, the results for  $\xi = F$  given by (27) and (28) apply as well.

Now suppose that  $v$  is unknown, so  $P \in P \equiv P_{123}$ . Then we must also deal with the score for the regression parameter  $v$ : with  $\epsilon \equiv Y - vZ$  and  $\psi \equiv -f'/f$ ,

$$(75) \quad \dot{I}_1(Y, Z) = Z\psi(\epsilon) - EZ\psi(\epsilon),$$

and, as in example 1,

$$(76) \quad \dot{I}_2 a(Y, Z) = a(\varepsilon) - Ea(\varepsilon), \quad a \in L_2^0(F),$$

$$(77) \quad \dot{I}_3 b(Y, Z) = b(Z) - Eb(Z), \quad b \in L_2^0(H).$$

Thus by (16), under (9),

$$(78) \quad \begin{aligned} \dot{P} &\supset [\dot{I}_1] + L_2^0(F^*) + L_2^0(H^*) \\ &= [\dot{I}_1] + \dot{P}_2 + \dot{P}_3 \equiv [\dot{I}_1] + \dot{P}_{23}, \end{aligned}$$

where  $F^*$  and  $H^*$  denote the marginal distributions of  $\varepsilon$  and  $Z$  respectively. To find  $I_1^*$ , we need to compute

$$(79) \quad I_1^* = (I - P_{23})\dot{I}_1,$$

where  $P_{23} \equiv \Pi_0(\cdot | \dot{P}_{23})$ . As shown in example 4.4.3, because  $\dot{P}_{23} = \dot{P}_2 + \dot{P}_3$  is a sumspace,  $P_{23}$  is the ACE projection studied in section A.4; i.e.,

$$(80) \quad I_1^* = \dot{I}_1 - ACE(\dot{I}_1 | \varepsilon, Z).$$

This projection is not generally computable in closed form. In the present case, however, the ACE equations (A.4.12) and (A.4.13), or (A.4.14) and (A.4.15), can be solved explicitly because of the special structure of  $P$  and its conditional expectation operators. Note that

$$(81) \quad E(a(U) | V = v) = \frac{1}{F(v)} \int_{-\infty}^v a dF = (I - \bar{R}_F)a(v)$$

(or  $E_P(a(U) | V = v) = E_Q(a(U) | U < V = v)$ ) and

$$(82) \quad E(b(V) | U = u) = \frac{1}{G(u)} \int_u^\infty b dG = (I - R_G)b(u)$$

(or  $E_P(b(V) | U = u) = E_Q(b(V) | V > U = u)$ ), where

$$\bar{R}_F a(u) \equiv a(u) - \frac{\int_{-\infty}^u a dF}{F(u)}$$

and

$$R_G b(v) \equiv b(v) - \frac{\int_v^\infty b dG}{G(v)}.$$

$R_G$  is the  $R$  operator corresponding to the df  $G$ , while  $\bar{R}_F$  is the "reversed" or "retro"  $R$  operator corresponding to the df  $F$ ; the "reversed  $L$  operator"  $\bar{L}_F$  is just

$$\bar{L}_F a(u) = a(u) - \int_u^\infty a d\bar{\Lambda}_F,$$

where  $\bar{\Lambda}_F$  is the reverse cumulative hazard function for  $F$  of (22); recall the definitions given in appendix A.1.

To compute the projection in (79) or (80), we first rewrite  $\dot{I}_1$  as follows:

$$\begin{aligned} \dot{I}_1(Y, Z) &= Z\psi(\varepsilon) - EZ\psi(\varepsilon) \\ &= -\frac{1}{v}(V\psi(U) - E(V\psi(U))) + \frac{y_0}{v}(\psi(U) - E\psi(U)), \end{aligned}$$

where  $V = y_0 - vZ$ . Since the second term is in  $\dot{P}_{23} = L_2^0(F^*) + L_2^0(H^*)$ , it suffices to project the first term, or just  $V\psi(U)$ .

Let  $h(U, V) = V\psi(U)$ . By the discussion in section A.4, to find

$$(83) \quad \Pi_0(V\psi(U) | \dot{P}_2 + \dot{P}_3) \equiv a^*(U) + b^*(V),$$

it suffices to solve one of (A.4.14) or (A.4.15) for  $a^*(u) \equiv h_1^*(u)$  or  $b^*(v) \equiv h_2^*(v)$ . These two equations now take the form

$$(84) \quad a(u) = P_2 Q_3 h(u) + E(E(a(U) | V) | U = u)$$

and

$$(85) \quad b(v) = P_3 Q_2 h(v) + E(E(b(V) | U) | V = v),$$

where  $P_2 \equiv \Pi_0(\cdot | \dot{P}_2) = E(\cdot | U)$  and  $P_3 \equiv \Pi_0(\cdot | \dot{P}_3) = E(\cdot | V)$ . But the second terms on the right side can be written, by (81) and (82), as

$$\begin{aligned} (86) \quad E(E(a(U) | V) | U = u) &= \frac{1}{\bar{G}(u)} \int_u^\infty \frac{\int_{-\infty}^v a \, dF}{F(v)} \, dG(v) \\ &= \frac{1}{\bar{G}(u)} \int a(t) \int_{u \vee t}^\infty \frac{dG}{F} \, dF(t), \end{aligned}$$

and

$$\begin{aligned} (87) \quad E(E(b(V) | U) | V = v) &= \frac{1}{F(v)} \int_{-\infty}^v \frac{\int_u^\infty b \, dG}{\bar{G}(u)} \, dF(u) \\ &= \frac{1}{F(v)} \int b(t) \int_{-\infty}^{v \wedge t} \frac{dF}{\bar{G}} \, dG(t). \end{aligned}$$

Set

$$(88) \quad \begin{aligned} c(u) \equiv P_2 Q_3 h(u) &= \psi(u) E(V | U = u) \\ &\quad - E(VE(\psi(U) | V) | U = u) \end{aligned}$$

and

$$(89) \quad \begin{aligned} d(v) \equiv P_3 Q_2 h(v) &= v E(\psi(U) | V = v) \\ &\quad - E(\psi(U) E(V | U) | V = v). \end{aligned}$$

Substituting (86)–(89) into (84) and (85), multiplying across (84) by  $\bar{G}(u)$ , and across (85) by  $F(v)$ , and then rearranging yields

$$(90) \quad \bar{G}(u)c(u) = \bar{G}(u)a(u) - \int a(t) \int_{u \vee t}^\infty \frac{dG}{F} \, dF(t),$$

and

$$(91) \quad F(v)d(v) = F(v)b(v) - \int_{-\infty}^{v \wedge t} b(t) \frac{dF}{G} dG(t).$$

But the right side of (91) is similar to the one of (42) of example 2; and the right side of (90) is the natural analogue to it. Hence, just as in (42) and (43),

$$(92) \quad \bar{G}(u)c(u) = \bar{L}_F D_G \bar{R}_F a(u)$$

and

$$(93) \quad F(v)d(v) = L_G D_F R_G b(v),$$

where  $D_G a(u) \equiv \bar{G}(u)a(u)$  and  $D_F b(v) \equiv F(v)b(v)$ . Again by proposition A.1.8.B, it follows that

$$(94) \quad a^*(u) = \bar{L}_F D_G^{-1} \bar{R}_G(\bar{G}c)(u)$$

and

$$(95) \quad b^*(v) = L_G D_F^{-1} R_G(Fd)(v),$$

and it remains only to compute the right sides.

To do this we use the fact that

$$(96) \quad \bar{R}_F \psi = \psi + \frac{f}{F} \equiv \psi + \bar{\lambda}_F = -\frac{\bar{\lambda}'_F}{\bar{\lambda}_F}$$

(which is the "reversed" or "retro" analogue of (4.6.23)) and we write

$$M_G(u) \equiv E(V | U = u) = \frac{\int_u^\infty v dG(v)}{\bar{G}(u)} = (I - R_G)(V)(u),$$

and

$$(97) \quad e_G(u) \equiv M_G(u) - u = E(V - u | U = u) = -R_G(V)(u).$$

Then straightforward calculation yields

$$(98) \quad \frac{\bar{R}_F(\bar{G}c)}{\bar{G}}(u) = \{\bar{R}_F \psi(u)\} M_G(u)$$

and

$$(99) \quad \frac{R_G(Fd)}{F}(v) = \bar{\lambda}_F(v) e_G(v),$$

and hence

$$(100) \quad \begin{aligned} a^*(u) &= \bar{L}_F(\{\bar{R}_F \psi\} M_G)(u) \\ &= \{\bar{R}_F \psi(u)\} M_G(u) - \int_u^\infty \{\bar{R}_F \psi\} M_G d\bar{\Lambda}_F \\ &= \{\bar{R}_F \psi(u)\} M_G(u) - \int_u^\infty M_G d\bar{\Lambda}_F \end{aligned}$$

$$\begin{aligned}
& \text{since } \int_u^\infty \bar{R}_F \psi d\bar{\Lambda}_F = - \int_u^\infty \bar{\lambda}'_F(t) dt = \bar{\lambda}_F(u) \\
& = \{\bar{R}_F \psi(u)\} M_G(u) - M_G(u) \bar{\lambda}_F(u) + \int_u^\infty \bar{\lambda}_F dM_G \\
& = \psi(u) M_G(u) + \int_u^\infty \bar{\lambda}_F dM_G \quad \text{by (96),}
\end{aligned}$$

while

$$\begin{aligned}
(101) \quad b^*(v) &= L_G(\bar{\lambda}_F e_G)(v) = \bar{\lambda}_F(v) e_G(v) - \int_{-\infty}^v \bar{\lambda}_F e_G d\Lambda_G \\
&= \bar{\lambda}_F(v) e_G(v) - \int_{-\infty}^v \bar{\lambda}_F dM_G
\end{aligned}$$

since by (97)

$$\begin{aligned}
\int_{-\infty}^u e_G d\Lambda_G &= - \int_{-\infty}^u (R_G V) d\Lambda_G \\
&= (L_G - I) R_G(V)(u) \\
&= (I - R_G)(V)(u) - \int v dG(v) \\
&= M_G(u) - \int v dG(v).
\end{aligned}$$

The upshot is that by (83)

$$\begin{aligned}
(102) \quad I_1^*(y, z) &= - \frac{1}{v} \{ v \psi(u) - EV \psi(U) \\
&\quad - (a^*(u) - Ea^*(U)) - (b^*(v) - Eb^*(V)) \}
\end{aligned}$$

with  $a^*$  and  $b^*$  given by (100) and (101). And now we can calculate the information for estimation of  $v$  explicitly.  $\square$

### Proofs

**Proof of (16).** To prove that the last inclusion of (15) is an equality, we use parts B and C of proposition A.4.2 and the calculations just preceding it: the sum space  $L_2^0(F^*) + L_2^0(G^*)$  is closed if the projection operator (conditional expectation) from  $L_2^0(G^*) \equiv \mathbf{H}_2$  to  $L_2^0(F^*) \equiv \mathbf{H}_1$  is Hilbert-Schmidt and hence compact, and a sufficient condition for this is:

$$\iint K^2(u, v) dF^*(u) dG^*(v) = \iint \frac{p^2(u, v)}{f^*(u) g^*(v)} du dv < \infty,$$

where  $K(u, v) \equiv p(u, v)/(f^*(u) g^*(v))$  is the kernel of the conditional expectation operators. But from (2)–(4),

$$\begin{aligned}
\iint \frac{p^2(u, v)}{f^*(u) g^*(v)} du dv &= \iint_{u < v} \frac{f^2(u) g^2(v)}{\bar{G}(u) f(u) F(v) g(v)} du dv \\
&\leq \int \frac{1}{\bar{G}} dF \int \frac{1}{F} dG < \infty \quad \text{by (9).}
\end{aligned}$$

Consequently,  $L_2^0(F^*) + L_2^0(G^*)$  is closed.



Furthermore, if  $\{P_{(F,\eta,G)}\}$  is a curve with tangent  $h \in \dot{P}_1^0$ , then by proposition A.5.5.A with  $T(U,V) = U$ ,  $\{F_\eta^*\}$  is a curve with tangent  $E(h(U,V) | U) \in L_2^0(F^*)$ . In fact, with  $f_\eta$  and  $g$  the densities of  $F_\eta$  and  $G$  with  $\alpha(\eta) = \int F_\eta dG$ , by Cauchy-Schwarz,

$$\begin{aligned} o(\eta^2) &= \iint 1_{|u \leq v|} \{ \alpha^{-1/2}(\eta) f_\eta^{1/2}(u) g^{1/2}(v) - \alpha^{-1/2}(0) f_0^{1/2}(u) g^{1/2}(v) \\ &\quad \cdot [1 + \frac{1}{2} \eta h(u,v)] \}^2 dv du \\ &\geq \int \frac{1}{\bar{G}(u)} \left\{ \int 1_{|u < v|} \{ \alpha^{-1/2}(\eta) f_\eta^{1/2}(u) g^{1/2}(v) \right. \\ &\quad \left. - \alpha^{-1/2}(0) f_0^{1/2}(u) g^{1/2}(v) [1 + \frac{1}{2} \eta h(u,v)] \} \right. \\ &\quad \left. \cdot g^{1/2}(v) dv \right\}^2 du \end{aligned}$$

$$\begin{aligned} (a) \quad &= \int \frac{1}{\bar{G}(u)} \left\{ \alpha^{-1/2}(\eta) f_\eta^{1/2}(u) - \alpha^{-1/2}(0) f_0^{1/2}(u) \right. \\ &\quad \left. \cdot [1 + \frac{1}{2} \eta E(h(U,V) | U = u)] \right\}^2 du \\ &= \iint 1_{|u \leq v|} \left\{ \alpha^{-1/2}(\eta) f_\eta^{1/2}(u) g^{1/2}(v) - \alpha^{-1/2}(0) f_0^{1/2}(u) g^{1/2}(v) \right. \\ &\quad \left. \cdot [1 + \frac{1}{2} \eta E(h(U,V) | U = u)] \right\}^2 dv du . \end{aligned}$$

Consequently,  $h(U,V) = E(h(U,V) | U)$  a.s. and  $\dot{P}_1^0 \subset L_2^0(F^*)$ . A similar argument shows that  $\dot{P}_2^0 \subset L_2^0(G^*)$ . □

**Proof of proposition 1.** To see that A implies B, first note that for any path  $\{P_{(F,G,\eta)}\}$  in  $P_2$  with tangent  $h(u,v) = a^*(u) + b^*(v) \in \dot{P}_2^0$ , it follows from (16) and its proof that  $a^*(u) = 0$  and  $b^*(v) \in L_2^0(G^*)$ . Restricting  $b^*$  to  $L_2^0(G^*) \cap L_2(G)$  we see that

$$\begin{aligned} \eta^{-1}(\kappa(P_{(F,G,\eta)}) - \kappa(P_{(F,G)})) &= \eta^{-1}(\bar{G}_\eta - \bar{G}) \\ &\rightarrow \int_0^\infty (b^* - \int b^* dG) dG \\ &= - \int_0^\infty (b^* - \int b^* dG) dG, \end{aligned}$$

and similarly for  $a^* \in L_2^0(F^*) \cap L_2(F)$ ,

$$\eta^{-1}(\xi(P_{(F,\eta,G)}) - \xi(P_{(F,G)})) = \eta^{-1}(F_\eta - F) \rightarrow \int_0^\infty (a^* - \int a^* dF) dF .$$

Thus it suffices to show for  $a \in L_2^0(F)$  and  $b \in L_2^0(G)$  that  $\kappa$  and  $\xi: \dot{P} \rightarrow l^\infty(R^+)$  defined by

$$(a) \quad \dot{\kappa}(\dot{\mathbf{i}}(a, b)) = - \int_0^\cdot b \, dG$$

and

$$(b) \quad \dot{\xi}(\dot{\mathbf{i}}(a, b)) = \int_0^\cdot a \, dF$$

are continuous. Consider  $\dot{\kappa}$ . By the Cauchy-Schwarz inequality,

$$\begin{aligned} (c) \quad \|\dot{\kappa}(\dot{\mathbf{i}}(a, b))\|_\infty^2 &= \sup_{0 \leq t < \infty} \left| \int_0^t b F^{1/2} \frac{dG}{F^{1/2}} \right|^2 \\ &\leq \int_0^\infty b^2 F \, dG \int_0^\infty \frac{dG}{F} \\ &\leq \text{Var}_G(b(V)) \int_0^\infty \frac{dG}{F} \equiv \|b\|_{L_2^0(G)}^2 \int_0^\infty \frac{dG}{F}. \end{aligned}$$

The map  $\dot{\mathbf{i}} : L_2^0(F) \times L_2^0(G) \rightarrow L_2^0(F^*) + L_2^0(G^*)$  defined by (13) is one-to-one, linear, and continuous from a Banach space onto a Banach space. But since A implies that (16) holds and, in particular, that the sum space  $L_2^0(F^*) + L_2^0(G^*)$  is closed, Banach's theorem proposition A.1.7.B, shows that  $\dot{\mathbf{i}}$  has a bounded inverse  $\dot{\mathbf{i}}^{-1}$ , and

$$(d) \quad \|(a, b)\|_{L_2^0(F) \times L_2^0(G)} \leq \|\dot{\mathbf{i}}^{-1}\| \|\dot{\mathbf{i}}(a, b)\|_0,$$

where  $\|(a, b)\|_{L_2^0(F) \times L_2^0(G)}^2 = \|a\|_{L_2^0(F)}^2 + \|b\|_{L_2^0(G)}^2$ . Combining (c) and (d) yields

$$(e) \quad \|\dot{\kappa}(\dot{\mathbf{i}}(a, b))\|_\infty^2 \leq \int_0^\infty \frac{dG}{F} \|\dot{\mathbf{i}}^{-1}\| \|\dot{\mathbf{i}}(a, b)\|_0^2;$$

i.e.,  $\dot{\kappa} : \dot{\mathbf{P}} \rightarrow \mathbf{B} = l^\infty(\mathbb{R}^+)$  is bounded, and hence continuous.

The argument for  $\dot{\xi}$  is similar, and hence B holds.

That B implies C is trivial.

To see that C implies A, we first note that, if  $\pi_t \kappa(P_{(F, G)}) = \bar{G}(t) \equiv \psi(g)$  is differentiable at  $P$  in  $\mathbf{P}_2(F) \equiv \{P_{(F, G)} : F \text{ known}\}$ , then, by theorem 5.4.1.A and (5.3.2) of example 5.3.1,

$$(f) \quad -(1_{[0, t]} - G(t)) \in \mathbf{R}(\dot{\mathbf{i}}_2^T),$$

with  $\dot{\mathbf{i}}_2^T : L_2^0(G^*) \rightarrow L_2^0(G)$  here. But by (37), for  $h(u, v) = b(v)$ ,

$$\dot{\mathbf{i}}_2^T h(v) = \frac{1}{\alpha} F(v) b(v),$$

so (f) becomes

$$-(1_{[0, t]}(v) - G(t)) = \frac{1}{\alpha} F(v) b(v)$$

for some  $b \in L_2^0(G^*)$ . Thus we must have

$$(g) \quad \int_0^\infty \left\{ \frac{1_{[0, t]}(v) - G(t)}{F(v)} \right\}^2 F(v) \, dG(v) < \infty,$$

and (g) is equivalent to  $\int_0^\infty (1/F) dG < \infty$ . Similarly, pathwise differentiability of  $\pi_{t\xi}(P_{(F,G)}) = F(t)$  implies  $\int_0^\infty (1/\bar{G}) dF < \infty$ , by using theorem 5.4.1.A and (36), so A holds.  $\square$

### 6.5 MIXTURE MODELS AND MODELS WITH MONOTONICITY CONSTRAINTS

In this section we consider estimation of infinite-dimensional parameters for several interrelated models connected with the mixture models of section 4.5, the Has'minskii-Ibragimov models (also in section 4.5), and models obtained by monotonicity or inequality constraints.

#### *Mixture Models: Connections with Monotonicity Constraints*

To begin, let  $\{f(\cdot, \eta) : \eta \in H \subset R^d\}$  be a parametric family of densities with respect to a dominating measure  $\mu$ , and consider the family

$$(1) \quad \mathbf{P} = \{P_G : G \in \mathbf{G}\}$$

with densities

$$(2) \quad p(x, G) = \int f(x, \eta) dG(\eta)$$

with respect to  $\mu$ . Many models of this type are intimately related to models with monotonicity constraints. Here are two examples:

**Example 1. Mixtures of uniforms: distributions with monotone density.**

Suppose that  $f(x, \eta) = \eta^{-1} 1_{[0, \eta]}(x)$ , so that

$$(3) \quad p(x, G) = \int_{[x \leq \eta]} \eta^{-1} dG(\eta), \quad x > 0,$$

is decreasing (i.e., nonincreasing). Furthermore, if  $p$  is decreasing with df  $P(x) \equiv P(X \leq x)$ , then

$$(4) \quad P(x) - xp(x) \equiv G(x), \quad x \geq 0,$$

is increasing from 0 to 1 and hence defines a df  $G$ , for which

$$\int_x^\infty \frac{1}{\eta} dG(\eta) = -\frac{G(x)}{x} + \int_x^\infty \frac{G(\eta)}{\eta^2} d\eta = p(x),$$

i.e., for which (3) holds. Thus the class of mixtures of  $Uniform(0, \eta)$  and the class of distributions with monotone density are the same.  $\square$

**Example 2. Scale mixtures of exponentials: distributions with completely monotone density.**

Suppose that  $f(x, \eta) = \eta \exp(-\eta x)$ ,  $x \geq 0$ ,  $\eta \geq 0$ ; the resulting family  $\mathbf{P}$  given by (2) has density functions

$$(5) \quad p(x, G) = \int_0^\infty \eta e^{-\eta x} dG(\eta), \quad x \geq 0.$$

Recall that a function  $\phi$  is *completely monotone* if it has derivatives  $\phi^{(k)}$  of all orders and  $(-1)^k \phi^{(k)}(x) \geq 0$ ,  $x \geq 0$ . Furthermore (see, e.g., Feller (1971, Vol. II, theorem XIII.4.1a, page 439, and problem XIII.11.11, page 464)), a density function  $p$  is completely monotone if and only if it is a mixture of exponential densities. Again the mixture class corresponds exactly to the monotonicity class.  $\square$

Note that  $\mathbf{P}_{\text{example1}} \supset \mathbf{P}_{\text{example2}}$ .

Of course there are many classes  $\mathbf{P}$  defined by monotonicity or order constraints that are not related to some class of mixtures; e.g., the classes of distributions on  $R^+$  with increasing failure rate (IFR) or increasing failure rate average (IFRA); see examples 6 and 7 below. To broaden the scope of the examples somewhat and to serve as illustration, we consider three more simple examples of the mixture models given by (1) and (2).

**Example 3. Poisson mixture model.**

Suppose that  $f(x, \eta) = \eta^x e^{-\eta}/x!$ ,  $x = 0, 1, \dots$ ,  $\eta > 0$ . The resulting Poisson mixture family  $\mathbf{P}$  has probability mass function

$$p(x, G) = \int_0^\infty \frac{\eta^x e^{-\eta}}{x!} dG(\eta), \quad x = 0, 1, \dots \quad \square$$

**Example 4. Location mixtures of Gaussians.**

In this example  $f(x, \eta) = \phi(x - \eta)$  for  $x, \eta \in R$  where  $\phi$  is the standard normal density, so that

$$(6) \quad p(x, G) = \int \phi(x - \eta) dG(\eta), \quad -\infty < x < \infty. \quad \square$$

**Example 5. Scale mixtures of centered Gaussians.**

In this example  $f(x, \eta) = \eta^{-1} \phi(x/\eta)$  for  $x \in R$  and  $\eta > 0$  where  $\phi$  is the standard normal density, so that

$$(7) \quad p(x, G) = \int_0^\infty \frac{1}{\eta} \phi\left(\frac{x}{\eta}\right) dG(\eta), \quad -\infty < x < \infty. \quad \square$$

To get started, here is a basic proposition giving scores and information operators for a class of mixture models of the type given by (1) and (2).

**Proposition 1.** Suppose that  $\{f(\cdot, \eta) : \eta \in H\}$  is any fixed family of (conditional, given  $\eta$ ) densities. Then, with  $U \sim G$  and  $L(X | U = u)$  having density  $f(\cdot, u)$ , the score operator  $\dot{\mathbf{I}} : L_2^0(G) \rightarrow \dot{\mathbf{P}}$  satisfies

$$A. \quad \dot{\mathbf{I}}a(x) = E(a(U) | X = x) = \frac{\int a(u) f(x, u) dG(u)}{p(x, G)}, \quad a \in L_2^0(G).$$

$$B. \quad \dot{\mathbf{I}}^T h(u) = E(h(X) | U = u) = \int h(x) f(x, u) d\mu(x), \quad h \in L_2^0(P).$$

$$C. \quad \dot{\mathbf{I}}^T \dot{\mathbf{I}}a(u) = E(E(a(U) | X) | U = u) = \int a(u') K(u, u') dG(u'),$$

where

$$K(u, u') \equiv \int \frac{f(x, u)f(x, u')}{p(x, G)} d\mu(x).$$

**Proof.** A follows immediately from proposition A.5.5.A and (3.2.5). B and C follow from A by direct calculation.  $\square$

The main points we want to illustrate are:

1. Frequently, when  $\mathbf{P}$  satisfies monotonicity or order restrictions, or is a mixture model, and  $P_0$  is on the "boundary" of  $\mathbf{P}$ , the set of all "one-sided" derivatives is a cone strictly larger than  $\dot{\mathbf{P}}^0$ . That is, there are paths  $\{P_\varepsilon : \varepsilon > 0\}$  such that  $\|s_\varepsilon - s_0 - 2^{-1}\varepsilon h s_0\| = o(\varepsilon)$  as  $\varepsilon \downarrow 0$  where  $s_\varepsilon^2 = dP_\varepsilon/d\mu$  and  $h \in L_2^0(P_0)$ . (Recall from section 3.2 that  $h$  would belong to  $\dot{\mathbf{P}}^0$  only if the path is *two-sided*; i.e. defined for  $|\varepsilon| < 1$ .) Typically this happens when  $P_0$  is on the boundary of  $\mathbf{P}$  in the sense of not being strictly monotone, or if the mixing distribution has a discrete support. In other cases, although  $\mathbf{P}$  is strictly less than the family of all continuous distributions, and  $\mathbf{R}(\dot{\mathbf{I}})$  is not closed, it or some other set of tangents is dense in  $L_2^0(P)$  and hence  $\mathbf{R}(\dot{\mathbf{I}}) = L_2^0(P)$ . Then the asymptotic information bounds for regular estimation are the same as if  $\mathbf{P} = \mathbf{M}$ ; i.e., there is no gain (at least asymptotically) from monotonicity or the fact that the model is a mixture in estimation of pathwise differentiable functions  $v$  of  $P$ . One formulation of this type of result was given by Millar (1979).
2. In pure mixture models given by (1) and (2)—and most of the mixture models studied in section 4.5—even when the mixing distribution function  $G$  is identifiable, it is *not* a pathwise differentiable function of  $P = P_G$ . This means that it is *not* estimable at rate  $n^{-1/2}$ .
3. If additional observations from the mixing distribution  $G$  are available (this is the Has'minskii-Ibragimov model introduced in section 4.5), then the mixing distribution function  $G$  is a pathwise differentiable function of  $P = P_G$ , and it is potentially estimable at rate  $n^{-1/2}$ .

Note that the information operator  $\dot{\mathbf{I}}^T \dot{\mathbf{I}}$  for  $g$  has the form of an integral operator with kernel  $K$ . Typically this type of operator is not boundedly invertible even if it is one-to-one. By Banach's theorem, Proposition A.1.7.B, invertibility can happen then only if the range of  $\dot{\mathbf{I}}^T \dot{\mathbf{I}}$  is closed in  $L_2^0(G)$ , and this typically fails for mixture models. We now apply proposition 5.4.1, primarily part B, to examples 1–5. The key question in each case is one-to-oneness of  $\dot{\mathbf{I}}^T$ ; and this is, as in section 4.5, a kind of  $L_2$ -completeness question.

**Example 1. Mixtures of uniforms, continued.**

In this case

$$\dot{\mathbf{I}}^T h(u) = \int h(x) \frac{1}{u} 1_{[0, u]}(x) dx = \frac{1}{u} \int_0^u h(x) dx, \quad u > 0,$$

so  $\dot{\mathbf{I}}^T h(u) = 0$  a.e.  $G$  implies  $h(u) = 0$  a.e.  $G$  if the points of increase of  $G$  are

dense in the support of  $G$ . For these  $G$ 's we have  $N(\dot{\mathbf{I}}^T) = \{0\}$  and  $\overline{\mathbf{R}(\dot{\mathbf{I}})} = N(\dot{\mathbf{I}}^T)^\perp = L_2^0(P) = \dot{\mathbf{M}}$ . Alternatively, this also follows from theorem 4.5.1 with  $T(X, \theta) \equiv X$ . It follows that the bounds for estimation of the df  $\nu(P) \equiv P$  are just the same as if  $\mathbf{P} = \mathbf{M}$  as in example 5.3.1. This example was treated from the perspective of the family of measures on  $R^+$  with decreasing densities by Millar (1979).  $\square$

**Example 2. Scale mixtures of exponentials, continued.**

Here

$$\dot{\mathbf{I}}^T h(u) = \int_0^\infty h(x) u e^{-ux} dx, \quad u \geq 0,$$

so  $\dot{\mathbf{I}}^T h = 0$  a.e.  $G$  implies  $h(x) = 0$  for all  $x$  if the support of  $G$  has nonempty interior by standard exponential family theory (see, e.g. Lehmann (1986, theorem 4.3.1, page 142)). Thus if the support of  $G$  has nonempty interior,  $N(\dot{\mathbf{I}}^T) = \{0\}$  and  $\overline{\mathbf{R}(\dot{\mathbf{I}})} = L_2^0(P) = \dot{\mathbf{M}}$ , so bounds for estimation of the df  $\nu(P) \equiv P$  are just the same as if  $\mathbf{P} = \mathbf{M}$  (as in example 5.3.1).  $\square$

**Example 3. Poisson mixture model, continued.**

In this example

$$\dot{\mathbf{I}}^T h(u) = \sum_{x=0}^\infty h(x) \frac{u^x e^{-u}}{x!} \quad \text{for } h \in L_2^0(P).$$

Thus  $\dot{\mathbf{I}}^T h = 0$  a.e.  $G$  implies  $h(x) = 0$  for all  $x = 0, 1, \dots$  if the support of  $G$  has nonempty interior by exponential family theory. Under this condition we again have  $N(\dot{\mathbf{I}}^T) = \{0\}$  and  $\overline{\mathbf{R}(\dot{\mathbf{I}})} = L_2^0(P) = \dot{\mathbf{M}}$ , so bounds for estimation of the df  $\nu(P) = P$  are just the same as if  $\mathbf{P} = \mathbf{M}$ . This was also shown by Tierney and Lambert (1984).  $\square$

**Example 4. Location mixtures of Gaussians, continued.**

In this case

$$\dot{\mathbf{I}}^T h(u) = \int_{-\infty}^\infty h(x) \phi(x-u) dx, \quad u \in R,$$

so  $\dot{\mathbf{I}}^T h = 0$  a.e.  $G$  implies  $h(x) = 0$  for all  $x$  if the support of  $G$  has nonempty interior, again by standard exponential family theory. Thus, under a mild condition on  $G$ ,  $\dot{\mathbf{I}}^T$  is one-to-one and  $\overline{\mathbf{R}(\dot{\mathbf{I}}^T)} = L_2^0(P) = \dot{\mathbf{M}}$ .  $\square$

**Example 5. Scale mixtures of centered Gaussians, continued.**

In this example

$$\dot{\mathbf{I}}^T h(u) = \int_{-\infty}^\infty h(x) \frac{1}{u} \phi\left(\frac{x}{u}\right) dx, \quad u \in R^+.$$

Now  $\dot{\mathbf{I}}^T h = 0$  for all  $h \in L_2^0(P) \cap \{\text{odd functions}\}$ , so  $\dot{\mathbf{I}}^T$  is *not* one-to-one,  $N(\dot{\mathbf{I}}^T) \neq \{0\}$ , and  $\overline{\mathbf{R}(\dot{\mathbf{I}})} \neq L_2^0(P)$ . Again, if  $G$  has support with nonempty interior, then

$$N(\dot{\mathbf{I}}^T) = L_2^0(P) \cap \{\text{odd functions}\},$$

and

$$\overline{\mathbf{R}(\dot{\mathbf{I}})} = \mathbf{N}(\dot{\mathbf{I}}^T)^\perp = L_2^0(P) \cap \{\text{even functions}\}.$$

Thus in example 5, under mild conditions on  $G$  the information bounds for estimation of  $v(P) \equiv P$  are the same as in example 5.3.3 where we discussed estimation of a symmetric df. □

What if  $G \equiv G_0$  is discrete or has finite support? Let  $P_0 \equiv P_{G_0}$ . Then there are many tangents in  $\dot{\mathbf{P}}^0(P_0)$  that are *not* in  $\mathbf{R}(\dot{\mathbf{I}})$ . If  $G_0$  has finite support  $\{\eta_1, \dots, \eta_m\}$ , say, then  $\dim[\mathbf{R}(\dot{\mathbf{I}})] \leq m - 1$  since there are only  $m - 1$  linearly independent  $L_2^0(G)$  functions  $a$  (view  $a$  as an  $m$ -vector subject to the one constraint  $\sum_{i=1}^m a_i G(\{\eta_i\}) = 0$ ). But the collection of tangents to curves  $\{P_{G_\varepsilon}\}$  such that  $\eta_{\varepsilon i} \rightarrow \eta_i$ ,  $G_\varepsilon(\{\eta_{\varepsilon i}\}) \rightarrow G(\{\eta_i\})$ ,  $i = 1, \dots, m$ , and  $\sum_{i=1}^m G_\varepsilon(\{\eta_{\varepsilon i}\}) = 1$ , has  $2m - 1$  linearly independent tangents (if  $H = R$ ), so this gives an example in which the hypothesis  $\dot{\mathbf{P}} = \overline{\mathbf{R}(\dot{\mathbf{I}})}$  in theorem 5.4.1.B fails. Moreover,  $\dot{\mathbf{P}}^0$  as we have defined it (in terms of two-sided curves in  $\mathbf{P}$ ) does not exhaust all the directions of approach to  $P_0$ . For example, consider the family of *one-sided* curves

$$\{p_\varepsilon \equiv p(\cdot, (1 - \varepsilon)G_0 + \varepsilon G) : \varepsilon \in [0, 1]\}$$

for fixed  $G \in \mathbf{G}$ ,  $G \neq G_0$ . Then

$$\frac{p_\varepsilon - p_0}{p_0} = \varepsilon \frac{p(\cdot, G) - p(\cdot, G_0)}{p(\cdot, G_0)} \equiv \varepsilon a$$

where  $a \in L_2^0(P_0)$  if

$$\int \frac{p^2(x, G)}{p(x, G_0)} d\mu(x) < \infty.$$

This collection of tangents  $\mathbf{C}$ , namely

$$\mathbf{C} \equiv \left\{ \alpha \frac{p(\cdot, G) - p(\cdot, G_0)}{p(\cdot, G_0)} : \int \frac{p^2(x, G)}{p(x, G_0)} d\mu(x) < \infty, \alpha > 0 \right\},$$

is a convex cone and it typically spans all of  $L_2^0(P_0)$ .

*Non-Differentiability of the Function  $v(P_G) = G$*

Now suppose that  $G$  is a distribution on  $H = R$ , let  $t \in R$ , and let  $v : \mathbf{P} \rightarrow [0, 1]$  be the real-valued parameter defined by

$$v(P_G) = G(t) \equiv \psi(g)$$

in examples 1–4, where  $\dot{\mathbf{P}} = \overline{\mathbf{R}(\dot{\mathbf{I}})}$ . By example 5.3.1 we know that

$$(8) \quad \dot{\psi}(a) = \int (1_{[0,t]}(u) - G(t)) a(u) dG(u) = \langle \dot{\psi}_t, a \rangle_G$$

for  $a \in L_2^0(G)$  with  $\dot{\psi}_t = 1_{[0,t]} - G(t)$ . Note that  $\dot{\psi}^T : R \rightarrow L_2^0(G)$  is given by  $\dot{\psi}^T(r) = r \dot{\psi}_t$ . By theorem 5.4.1,  $v(P)$  is differentiable if and only if

$$(9) \quad \dot{\psi}_t \in \mathbf{R}(\dot{\mathbf{I}}^T).$$

But (9) fails in all of examples 1–5: note that in each of these examples all the functions in  $\mathbf{R}(\dot{\mathbf{I}}^T)$  are *continuous* functions of  $u$ , while  $\dot{\psi}_t$  is *discontinuous*. Hence  $v(P_G) = G(t)$  is *not* differentiable in any of these cases, and the “information for  $v(P_G) = G(t)$ , and hence also for  $v(P_G) = G$ , is zero.” Since  $\mathbf{N}(\dot{\mathbf{I}}) = \{0\}$  and  $\mathbf{R}(\dot{\mathbf{I}}^T \dot{\mathbf{I}}) \subset \mathbf{R}(\dot{\mathbf{I}}^T) \neq L_2^0(G)$  in these examples, it follows from proposition 5.4.1.C that  $(\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{-1}$  is *not* bounded.

What linear functionals  $v(P_G) = \int b \, dG$  are differentiable? If  $\psi(g) = \int b \, dG$ , then  $\dot{\psi} = b - \int b \, dG$ , and for this to be in  $\mathbf{R}(\dot{\mathbf{I}}^T)$  we must have, by the formula for  $\dot{\mathbf{I}}^T$  in proposition 1,

$$\begin{aligned} b(u) - \int b \, dG &= \int h(x) f(x, u) \, d\mu(x) \\ &= \int \tilde{h}(x) f(x, u) \, d\mu(x) - E_0 \tilde{h}(X) \end{aligned}$$

for some  $h \in L_2^0(P)$  or  $\tilde{h} \in L_2(P)$ . Hence

$$b(u) = \int \tilde{h}(x) f(x, u) \, d\mu(x)$$

for some  $\tilde{h} \in L_2(P)$ . Then

$$v(P_G) \equiv \int b \, dG = \int \tilde{h} \, dP_G,$$

so  $v$  is just the fixed linear functional of  $P = P_G$  corresponding to  $\tilde{h}$ , and is obviously estimable (at rate  $n^{-1/2}$ ) by  $v(IP_n) = \int \tilde{h} \, dIP_n$ . All the calculations in this subsection are from Van der Vaart (1991).

What is known about estimation of  $G$  in these examples? While we have shown that  $G$  or  $G(t)$  is not a differentiable function of  $P$  and hence not regularly estimable, in examples 1–5, it would probably be worthwhile to briefly review what is known about estimation of  $G$  in these examples.

**Example 1. Mixtures of uniforms, continued.**

In this model the nonparametric maximum likelihood estimator (NPMLE) of the density  $p$  is the Grenander estimator given by the density  $\hat{p}_n$  corresponding to the least concave majorant  $\hat{IP}_n$  of the empirical df  $IP_n$  of the sample; see, e.g., Grenander (1956), Groeneboom (1985), and Birgé (1989). Hence the NPMLE  $\hat{G}_n$  of  $G$  is easily obtained from (4) as

$$(10) \quad \hat{G}_n(u) = \hat{IP}_n(u) - u \hat{p}_n(u).$$

It is known that  $\hat{p}_n$  converges to  $p$  at rate  $n^{-1/3}$  pointwise at each point  $u$  at which  $p'(u)$  exists and is positive, and in other senses as well; see, e.g., Prakasa Rao (1983) and Groeneboom (1985). So this immediately carries over to  $\hat{G}_n$  too via (10).  $n^{-1/3}$  is also known to be the best possible rate of estimation in this model; see Birgé (1987a, 1987b), (1989). □

**Examples 2–5 continued.**

From Jewell (1982), and Lindsay (1983a, 1983b), the NPMLE  $\hat{G}_n$  of  $G$  is known to exist; by Kiefer and Wolfowitz (1956), Jewell (1982), and Pfanzagl



(1988),  $\hat{G}_n$  is consistent. Ritov (1986), Fan (1991), and Millar (1989) construct procedures that converge to  $G$  at rate  $(\log n)^{-1}$  in example 2, and  $(\log n)^{-1/2}$  in example 4, and show that these are the best rates possible. Tierney and Lambert (1984) have studied the NPMLE  $P_{\hat{G}_n}$  of  $P_G$  in the Poisson mixture model of example 3. For examples 2, 4, and 5 apparently little or nothing beyond consistency (which follows from Kiefer and Wolfowitz (1956)) is known about the NPMLE  $P_{\hat{G}_n}$  of  $P_G$ . □

*The Has'minskii-Ibragimov Model*

Now consider the Has'minskii-Ibragimov model as in (4.5.47) of section 4.5, but without the dependence on the finite-dimensional parameter  $\theta$  for simplicity. Thus, assuming that  $G \ll \lambda$  for all  $G \in \mathbf{G}$ ,  $P = P_G \in \mathbf{P}$  has density

$$(11) \quad p(x, G) = g(u) \int f(y, \eta) g(\eta) d\eta \quad \text{for } x = (u, y) \in \mathbf{U} \times \mathbf{Y}.$$

**Proposition 2.** Suppose that  $\{f(\cdot, \eta) : \eta \in H\}$  is any fixed family of (conditional, given  $\eta$ ) densities. Then, with  $X = (U, Y) \sim P$  with density given by  $p$  in (11) with respect to  $\lambda \times \mu$ , and  $V \sim G$  unobservable with  $L(Y | V = v)$  having density  $f(\cdot, v)$ , the score function  $\dot{\mathbf{I}}$  satisfies:

A.  $\dot{\mathbf{I}}a(x) = a(u) + E(a(V) | Y = y)$  for  $a \in L_2^0(G)$ .

B.  $\dot{\mathbf{I}}^T h(u) = E(h(U, Y) | U = u) + E(h(U, Y) | V = u)$ .

C.  $\dot{\mathbf{I}}^T \dot{\mathbf{I}}a(u) = a(u) + E(E(a(V) | Y) | V = u)$   
 $\equiv (I + E^V E^Y)a(u) \equiv (I + K)a(u)$

where  $K$  is the information operator from proposition 1.C.

D.  $(\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{-1} : L_2^0(G) \rightarrow L_2^0(G)$  is bounded.

E.  $\dot{\mathbf{P}} = \overline{\mathbf{R}(\dot{\mathbf{I}})}$  and

$$\begin{aligned} \Pi_0(h | \dot{\mathbf{P}}) &= \dot{\mathbf{I}}(\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{-1} \dot{\mathbf{I}}^T h \\ &= (I + E^Y)(I + E^V E^Y)^{-1} (E^U + E^V)h \end{aligned}$$

for  $h \in L_2^0(P)$ . In particular, if  $h(x) = h(u)$ ,  $(E^U + E^V)h = h$  and

$$\Pi_0(h | \dot{\mathbf{P}}) = (I + E^Y)(I + E^V E^Y)^{-1} h.$$

**Proof.** A follows from proposition A.5.5.A. B then follows from A by conditional expectation calculations:

$$\begin{aligned} \langle \dot{\mathbf{I}}a, h \rangle_0 &= E(a(U)h(U, Y)) + E(E(a(V) | Y)h(U, Y)) \\ &= E(a(U)E(h(U, Y) | U)) + E(a(V)h(U, Y)) \end{aligned}$$

$$\begin{aligned}
 &= E(a(U)E(h(U,Y) | U)) + E(a(V)E(h(U,Y) | V)) \\
 &= \langle a, \dot{\mathbf{I}}^T h \rangle_G
 \end{aligned}$$

with  $\dot{\mathbf{I}}^T$  as claimed.

The formula for  $\dot{\mathbf{I}}\dot{\mathbf{I}}$  in C follows by straightforward calculation from A and B. That  $K$  is nonnegative definite is immediate since it is of the form  $A^T A$  for a bounded linear operator  $A$ . Boundedness of  $(\dot{\mathbf{I}}\dot{\mathbf{I}})^{-1} = (I + E^V E^Y)^{-1}$  follows, as in the proof of proposition 4.5.1 by an application of corollary A.1.2, from the fact that  $I + K$  is a self-adjoint operator with  $\|(I + K)a\| \geq \|a\|$ . To see the first part of E, namely that  $\dot{\mathbf{P}} = \mathbf{R}(\dot{\mathbf{I}})$ , note that if  $\{P_\eta = P_{G_\eta}\}$  is a curve through  $P_0$  with tangent  $h$ , then its induced marginal distributions are differentiable by proposition A.5.5.A. In particular,  $\{G_\eta\}$  is a curve in  $\mathbf{G}$  with a tangent  $a$ , say. Then  $h$  has to be equal to  $\dot{I}a \in \mathbf{R}(\dot{\mathbf{I}})$ , so  $\dot{\mathbf{P}}^0 = \mathbf{R}(\dot{\mathbf{I}})$ . Finally, the projection formula given in E is a consequence of D and (A.2.13) of theorem A.2.2.  $\square$

Now we specialize to  $U = R$ , and let  $G$  also denote the distribution function corresponding to the mixing distribution (measure)  $G$ . It follows immediately from proposition 3.E and theorem 5.4.1.B that the function  $v : \mathbf{P} \rightarrow \mathbf{B} = l^\infty(R)$  defined by  $v(P_G) = G \equiv \psi(g)$  is pathwise differentiable in this model. Moreover, by (5.4.16), the efficient influence operator  $\tilde{I}_v$  is given by

$$\begin{aligned}
 (12) \quad \tilde{I}_v b^* &= \dot{\mathbf{I}}(\dot{\mathbf{I}}\dot{\mathbf{I}})^{-1} \dot{\psi}^T b^* \\
 &= (I + E^Y)(I + E^V E^Y)^{-1} \dot{\psi}^T b^*,
 \end{aligned}$$

where  $\dot{\psi}$  is given in (8). In particular, for  $t \in H = R$ ,

$$(13) \quad \tilde{I}_v \pi_t = (I + E^Y)(I + E^V E^Y)^{-1} \dot{\psi}_t,$$

and the inverse information covariance function (5.2.23) for estimation of  $v = G$  is

$$\begin{aligned}
 (14) \quad I_v^{-1}(s, t) &= E_0(\tilde{I}_v \pi_s, \tilde{I}_v \pi_t) \\
 &= \langle \dot{\mathbf{I}}(\dot{\mathbf{I}}\dot{\mathbf{I}})^{-1} \dot{\psi}_s, \dot{\mathbf{I}}(\dot{\mathbf{I}}\dot{\mathbf{I}})^{-1} \dot{\psi}_t \rangle_0 \\
 &= \langle \dot{\psi}_s, (\dot{\mathbf{I}}\dot{\mathbf{I}})^{-1} \dot{\psi}_t \rangle_G \\
 &= \langle \dot{\psi}_s, (I + E^V E^Y)^{-1} \dot{\psi}_t \rangle_G.
 \end{aligned}$$

For an estimator achieving this bound, see example 7.6.6.

Of course these calculations can be generalized considerably. For example, it is easy to do the computations for mixing distributions  $G$  on a Euclidean space  $U$  with  $v : \mathbf{P} \rightarrow l^\infty(\mathbf{C})$  defined by  $v(P_G)(C) = G(C)$  for  $C \in \mathbf{C}$ , a nice class of subsets of  $U$ . They also extend straightforwardly to the two-sample version of the model, with a sample of  $m$   $U$ 's i.i.d.  $G$  and  $n$   $Y$ 's i.i.d. according to the mixed density  $\int f(\cdot, \eta) dG(\eta)$ , and more generally to missing value models.

*More on Models with Monotonicity Constraints*

Now we return to models defined by monotonicity constraints. Here are several examples of models of this type which are not equivalent to any mixture family.

**Example 6. Increasing failure rate (IFR) distributions.**

Suppose that  $X = R^+$  and  $P$  is the family of all absolutely continuous distributions with increasing failure rate (IFR): if  $P(x) \equiv P(X \leq x)$ ,  $\bar{P}(x) = 1 - P(x)$ ,  $p(x)$  the density, and  $\lambda(x) \equiv p(x)/\bar{P}(x)$  the hazard or failure rate for  $x \geq 0$ , then  $\lambda(x) \leq \lambda(y)$  for all  $x \leq y$ . Alternatively, the cumulative hazard function  $\Lambda_P$ , given by

$$\Lambda_P(x) \equiv \int_{[0,x]} \frac{1}{\bar{P}(s)} dP(s),$$

is convex on the support of  $P$ . Note however, that  $P$  is *not* a mixture model since it is not a convex family:  $U(0,1)$  and  $U(1,2)$  both have IFR, but  $pU(0,1) + (1-p)U(1,2)$  does not have IFR if  $1/2 < p < 1$ .  $\square$

**Example 7. Increasing failure rate average (IFRA) distributions.**

If  $X = R^+$  and  $P$  is the family of all distributions such that the functions  $A$  defined by  $A(x) \equiv x^{-1} \Lambda_P(x)$ ,  $x > 0$ , are nondecreasing (i.e.,  $\Lambda_P$  is star-shaped), then  $P$  is called the family of distributions with increasing failure rate average. This class has importance in reliability theory because it is the smallest class containing the exponential distributions which is closed under the formation of coherent systems and limits in distribution; see Birnbaum, Esary, and Marshall (1966). It is known that the nonparametric maximum likelihood estimator (NPMLE) of the df in this class is *inconsistent*, as is the NPMLE of a star-shaped df; see Boyles, Marshall, and Proschan (1985), and Barlow, Bartholomew, Bremner, and Brunk (1972, pages 257-260).

There are many other possible examples in one dimension; for more classes like  $P_{IFR}$  and  $P_{IFRA}$  and relationships between them, see, e.g., Klefsjö (1983).  $\square$

Here is one further example in  $R^d$ .

**Example 8. Schur-concave distributions.**

Suppose that  $X = R^d$ , and that  $P$  is the collection of all measures on  $X$  with Schur-concave densities  $p$  (with respect to Lebesgue measure): i.e.,

$$(15) \quad x <_m y \quad \text{implies} \quad p(x) \geq p(y).$$

Here  $x <_m y$  means that

$$\sum_{i=1}^j x_{[i]} \leq \sum_{i=1}^j y_{[i]} \quad \text{for } j = 1, \dots, d-1$$

and

$$\sum_{i=1}^d x_{[i]} = \sum_{i=1}^d y_{[i]},$$

where  $x_{[i]}$  denotes the  $i$ th largest component of  $x = (x_1, \dots, x_d)$ ; thus  $x_{[1]} \geq x_{[2]} \geq \dots \geq x_{[d]}$ . See Marshall and Olkin (1979, pages 7 and 54).

This implies that  $p$  is exchangeable or permutation symmetric:  $p(x) = p(\pi x)$  for all  $x \in R^d$  and for all permutations  $\pi \in \Pi$  (here  $\pi x = (x_{\pi(1)}, \dots, x_{\pi(d)})$ ).  $\square$

A set of basic tools for establishing information bounds for these types of examples is provided by the work of Van der Vaart (1988a), (1989). He shows that the convolution theorems 2.3.1, 3.3.2, and 5.2.1 continue to hold if the set of one-sided tangents is *not* necessarily linear, but just contains a convex cone  $\mathbf{T}$  with closed linear span. The basic idea is that averages of limiting distributions of regular estimators with respect to Gaussian distributions can always be represented in the form of a convolution; see Van der Vaart's (1989) theorem 2.2, page 1489, and the discussion on pages 1491–1492.

When  $P$  is not an interior point of  $\mathbf{P}$ , then  $\dot{\mathbf{P}}^0$  is a proper subspace of  $[\mathbf{T}]$ , and Van der Vaart's (1989) generalization of theorem 5.2.2.A is needed to obtain a convolution result for regular estimators.

On the other hand, it is not clear that the restriction to regular estimators is acceptable in these situations. This is similar to what happens in the simple parametric model  $\{\phi(\cdot - \theta) : \theta \in \Theta\}$  in which the restriction to invariant estimators makes sense if  $\Theta = R$ , but is not sensible if  $\Theta = R^+ \equiv [0, \infty)$ .

These problems are avoided if we assume that  $P$  is an interior point of  $\mathbf{P}$  in the sense that  $[\mathbf{T}] = \dot{\mathbf{P}}$ . This was done in examples 1–4, and will be equivalent to the restriction to strict monotonicity in examples 5–7.

**Example 6. IFR distributions, continued.**

Fix  $P \in \mathbf{P}_{IFR}$  with density  $p$ , df  $P$ , and strictly monotone hazard rate  $\lambda = p/\bar{P}$ . Let

$$\tilde{\mathbf{H}} = \{h : \|\frac{\partial h}{\partial \lambda}\|_{\infty} < \infty, h(0) = 0\}.$$

For  $h \in \tilde{\mathbf{H}}$  and  $\eta \in R$  define

$$\lambda_{\eta}(t) \equiv \lambda(t) + \eta h(t), \quad t \geq 0.$$

Note that for  $\eta$  with  $|\eta| \|\partial h / \partial \lambda\|_{\infty} \leq 1$ , the function  $\lambda_{\eta}$  is nondecreasing and hence nonnegative and a hazard rate. Thus  $P_{\eta}$  with density

$$p_{\eta}(x) = \lambda_{\eta}(x) \exp(-\int_0^x \lambda_{\eta}(t) dt)$$

has increasing failure rate. Furthermore, the curve  $\{p_{\eta}\}$  has tangent (at  $\eta = 0$ )

$$L_P\left(\frac{h}{\lambda}\right)(x) = \frac{h}{\lambda}(x) - \int_0^x \frac{h}{\lambda}(t) \frac{1}{\bar{P}(t)} dP(t) \in \dot{\mathbf{P}}^0,$$

where  $L_P$  is the  $L$ -operator corresponding to the df  $P$ ; see the paragraph following (A.1.22). Now

$$\mathbf{H} \equiv \left\{L_P\left(\frac{h}{\lambda}\right) : h \in \tilde{\mathbf{H}}\right\} \subset \dot{\mathbf{P}}^0 \subset L_2^0(P).$$

Since  $\mathbf{N}(L_P) = \{0\}$  and  $\mathbf{R}(L_P) = L_2^0(P)$  by proposition A.1.8.E and  $\tilde{\mathbf{H}}/\lambda$  is dense in  $L_2(P)$  whenever  $\lambda$  is strictly increasing, it follows that  $[\mathbf{H}] = L_2^0(P) = \dot{\mathbf{P}} = \dot{\mathbf{M}}$  for every  $P \in \mathbf{P}_{IFR}$  with strictly increasing failure

rate. Once again, bounds for estimation of the df  $P$ , or any other pathwise differentiable function  $v$  of  $P$  are the same as for  $\mathbf{P} = \mathbf{M}$ .

This is a slight extension of the calculations of Millar (1979, page 243). Millar was able to fix  $P = \text{Uniform}(0,1)$  because of his focus on global asymptotic minimax theorems.

Since  $\mathbf{P}_{IFR} \subset \mathbf{P}_{IFRA}$ , it is clear that the same calculations apply to example 7, too. See Wang (1986), (1987) for construction of efficient estimators within the IFR and IFRA families specified by examples 6 and 7.  $\square$

**Example 8. Schur-concave distributions, continued.**

Suppose that  $P \in \mathbf{P}$  has density  $p$  satisfying  $p(x) > p(y)$  for any  $y \neq x$  with  $x <_m y$ . Let

$$\mathbf{H} = \left\{ h \in L_2^0(P) : \|h\|_\infty < \infty, \sup_{x \neq y, x <_m y} \left| \frac{h(x) - h(y)}{p(x) - p(y)} \right| < \infty, \right. \\ \left. h(\pi \cdot) = h(\cdot) \text{ for all } \pi \in \Pi \right\}.$$

For  $h \in \mathbf{H}$  and  $\eta \in R$  let  $\tilde{p}_\eta(x) = p(x) + \eta h(x)$ . Then

$$\tilde{p}_\eta(x) - \tilde{p}_\eta(y) = (p(x) - p(y)) \left\{ 1 + \eta \frac{h(x) - h(y)}{p(x) - p(y)} \right\}, \quad x \neq y,$$

and for  $|\eta|$  sufficiently small  $\tilde{p}_\eta$  is Schur-concave. Consequently,

$$p_\eta(x) = \frac{\tilde{p}_\eta(x) \vee 0}{\int \tilde{p}_\eta(y) \vee 0 dy}$$

is a Schur-concave density on  $R^d$  for  $|\eta|$  sufficiently small,  $\{p_\eta\}$  has tangent  $h$ , and hence

$$\dot{\mathbf{P}} \supset [\mathbf{H}] = \{ h \in L_2^0(P) : h(x) = h(\pi x) \text{ a.s. for all } \pi \in \Pi \}.$$

But the latter is just the tangent space for the collection  $\mathbf{P}_\Pi$  of all permutation symmetric distributions on  $R^d$ , example 6.2.3. Hence  $\dot{\mathbf{P}} = \dot{\mathbf{P}}_\Pi$ , and bounds for estimation of pathwise differentiable functions of  $P$  in the collection of Schur-concave distributions are just the same as in the family of all permutation symmetric distributions; see example 6.2.3.  $\square$

There are undoubtedly many other examples of this same type involving other orderings.

## 6.6 MISSING DATA AND CENSORING

Here we pick up the thread of section 4.6, but with emphasis now on estimation of parameters related to the infinite-dimensional components of the models introduced there. In fact, we will concentrate on random censoring models: most of the models we study here can be viewed as either special cases or generalizations of the censored regression model example 4.6.4.

As in section 4.6, we suppose that  $X^0 \sim Q \in \mathbf{Q}$  is unobservable, but there is a (fixed, known) function  $T$  such that  $X = T(X^0)$  is observable, and the

model induced by  $T$  is called  $\mathbf{P} \equiv \mathbf{Q}T^{-1} = \{QT^{-1} : Q \in \mathbf{Q}\}$ . As in (4.6.1) and (4.6.2), proposition A.5.5.A yields, with  $P = QT^{-1}$ ,

$$\dot{\mathbf{P}} \supset \{E(a(X^0) | X) : a \in \dot{\mathbf{Q}}\},$$

and we expect that equality holds. Thus the score operator  $\dot{\mathbf{i}} = \dot{\mathbf{i}}_Q : \dot{\mathbf{Q}} \rightarrow L_2^0(P)$  is given by

$$\dot{\mathbf{i}}a(X) = E(a(X^0) | X) \quad \text{for } a \in \dot{\mathbf{Q}}.$$

When  $\mathbf{Q}$  is nonparametric in the sense that  $\dot{\mathbf{Q}} = L_2^0(Q)$ , straightforward calculation shows that  $\dot{\mathbf{i}}^T : L_2^0(P) \rightarrow \dot{\mathbf{Q}}$  is given by

$$\dot{\mathbf{i}}^T b(X^0) = E(b(X) | X^0) \quad \text{for } b \in L_2^0(P).$$

Then the information operator  $\dot{\mathbf{i}}^T \dot{\mathbf{i}}$  is given by

$$\dot{\mathbf{i}}^T \dot{\mathbf{i}}a(X^0) = E(E(a(X^0) | X) | X^0) = E^{X^0} E^X a(X^0),$$

where

$E^X$  is the conditional expectation given  $X = T(X^0)$ , and  
 $E^{X^0}$  is the conditional expectation given  $X^0$ .

All of our examples in this section will have information operators with this structure.

An important question in these types of models is: how big is  $\dot{\mathbf{P}}$ ? If  $\dot{\mathbf{P}} = \dot{\mathbf{M}} \equiv L_2^0(P)$  for  $P \in \mathbf{P}$ , then we say that the tangent space is *saturated* or *full*. This is the case when  $\mathbf{P}$  is  $\mathbf{M}$  or is a substantial subset of  $\mathbf{M}$ , and then, in view of proposition 3.3.1 and corollary 5.2.1, there is at most one regular, asymptotically linear estimator (up to asymptotic equivalence and symmetry in the observations), and it is efficient. Thus if we know  $\dot{\mathbf{P}} = \dot{\mathbf{M}}$  and we know an estimator which is regular and asymptotically linear, it must be efficient. Since  $N(\dot{\mathbf{i}}^T)^\perp = \overline{R(\dot{\mathbf{i}})}$  and  $R(\dot{\mathbf{i}}) \subset \dot{\mathbf{P}} \subset \dot{\mathbf{M}}$ , it follows that  $\dot{\mathbf{P}} = \dot{\mathbf{M}} = L_2^0(P)$  whenever  $N(\dot{\mathbf{i}}^T) = \{0\}$ .

When  $\dot{\mathbf{P}} \subset \dot{\mathbf{M}}$  with strict inclusion, the situation is more complicated: there may exist many regular, asymptotically linear estimators which are inefficient. Finding an efficient estimator requires us to find an asymptotically linear estimator with influence function in the tangent space  $\dot{\mathbf{P}}$ .

The phenomenon of unsaturated tangent spaces  $\dot{\mathbf{P}} (\subset \dot{\mathbf{M}}$  with strict inclusion) will occur repeatedly in the following examples. Sometimes this happens because the parameter space is not as rich as possible (as in the random censoring model with  $G$  known). In other cases, as in the bivariate censoring model to be discussed in example 6 below, it occurs because the map from  $\mathbf{Q}$  to  $\mathbf{P}$  induced by  $T$  has range  $\mathbf{P}$  which is substantially smaller than the collection of all measures  $\mathbf{M}(X)$  on the sample space  $X = T(X^0)$  of  $X$ .

#### Example 1. Random censoring; nonparametric view.

Suppose that  $Y, C$  are independent rv's with distribution functions  $F$  and  $G$  on  $R^+ = [0, \infty)$ . At most one of these may be substochastic, i.e., may put mass

at  $\infty$ . Think of  $Y$  as a “survival time” and  $C$  as a “censoring time;” in this model with  $X^0 = (Y, C)$  we observe only  $X \equiv (T, \Delta)$  where

$$T \equiv Y \wedge C, \quad \Delta \equiv 1_{[Y \leq C]}.$$

If  $F$  and  $G$  have densities  $f$  and  $g$  with respect to Lebesgue measure, then the joint density  $p = p(\cdot; F, G)$  of  $X = (T, \Delta)$  on  $\mathbf{X} = R^+ \times \{0, 1\}$  (with respect to  $\mu = \text{Lebesgue} \times \text{counting measure}$ ) is

$$(1) \quad p(t, \delta) = \{(1 - G(t))f(t)\}^\delta \{(1 - F(t))g(t)\}^{1-\delta}$$

for  $t \in R^+, \delta \in \{0, 1\}$ . Furthermore, we define

$$\begin{aligned} H(t) &\equiv P(T \leq t) = H_0(t) + H_1(t) \\ &\equiv \int_{[0, t]} p(s, 0) ds + \int_{[0, t]} p(s, 1) ds = \int_0^t \bar{F} dG + \int_0^t \bar{G} dF \end{aligned}$$

and

$$\begin{aligned} \bar{H}(t) &\equiv 1 - H(t) = P(T > t) \\ &= P(Y > t, C > t) = \bar{F}(t)\bar{G}(t), \end{aligned}$$

and hence, by (1),

$$(2) \quad \begin{aligned} \frac{f(t)}{\bar{F}(t)} &= \frac{p(t, 1)}{\bar{H}(t)}, & \bar{F}(t) &= \exp\left(-\int_0^t \frac{p(s, 1)}{\bar{H}(s)} ds\right), \\ \frac{g(t)}{\bar{G}(t)} &= \frac{p(t, 0)}{\bar{H}(t)}, & \bar{G}(t) &= \exp\left(-\int_0^t \frac{p(s, 0)}{\bar{H}(s)} ds\right). \end{aligned}$$

This shows that  $F$  and  $G$  are identifiable on  $[0, \tau_H]$  with  $\tau_H = H^{-1}(1) \equiv \inf\{s : H(s) = 1\} = F^{-1}(1) \wedge G^{-1}(1)$ . In fact

$$(3) \quad \mathbf{P} = \{P \text{ on } R^+ \times \{0, 1\} \text{ with } p = \frac{dP}{d\mu} \text{ of the form (1)}\} = \mathbf{M}_\mu,$$

where  $\mathbf{M}_\mu$  consists of all proper probability distributions on  $\mathbf{X}$  dominated by  $\mu$ . To see this we define, for  $P \in \mathbf{M}_\mu$  with density  $p$ ,

$$\bar{F}(t) = \exp\left(-\int_0^t \frac{p(s, 1)}{\bar{H}(s)} ds\right), \quad \bar{G}(t) = \exp\left(-\int_0^t \frac{p(s, 0)}{\bar{H}(s)} ds\right),$$

where

$$\bar{H}(t) = 1 - H(t) = \int_t^\infty (p(s, 0) + p(s, 1)) ds.$$

Indeed, since  $p$  is a proper density,

$$\begin{aligned} -\log(\bar{F}(t)\bar{G}(t)) &= \int_0^t \frac{p(s, 0) + p(s, 1)}{\bar{H}(s)} ds \\ &= \int_0^t \frac{1}{\bar{H}(s)} dH(s) \\ &= -\log \bar{H}(t) \rightarrow \infty \end{aligned}$$

as  $t \rightarrow \infty$ , and hence at least one of  $F$  and  $G$  is a proper distribution function on  $R^+$ . This proves (3).

Therefore, by example 3.2.1, for any fixed  $P \in \mathbf{P}$ ,

$$(4) \quad \dot{\mathbf{P}} = \{h \in L_2(P) : Eh = 0\} = L_2^0(P).$$

Of course, this depends crucially on the fact that both  $F$  and  $G$  are unknown.

Now consider the parameters  $v(P) \equiv \Lambda$  and  $\kappa(P) \equiv 1 - F$  given by

$$(5) \quad v(P)(t) \equiv \Lambda(t) = \int_0^t \frac{dF}{1 - F_-} = \int_0^t \frac{dH_1}{\bar{H}_-}, \quad 0 \leq t \leq \tau,$$

and

$$(6) \quad \begin{aligned} \kappa(P)(t) &\equiv 1 - F(t) = \prod_{0 \leq s \leq t} (1 - d\Lambda(s)) \\ &\equiv \exp(-\Lambda^c(t)) \prod_{s \leq t} (1 - \Delta\Lambda(s)), \quad 0 \leq t \leq \tau, \end{aligned}$$

with  $\bar{H}(\tau) > 0$  and  $\Delta\Lambda \equiv \Lambda - \Lambda_-$ ,  $\Lambda^c(t) \equiv \Lambda(t) - \sum_{s \leq t} \Delta\Lambda(s)$ . In (5), as in examples 5.3.5 and 6.4.1, we have indicated the left limits in the denominators which are correct in the general discontinuous case since these are useful when forming estimators.

We first take a completely nonparametric approach to bounds for estimation of  $v$  and  $\kappa$  in the spirit of sections 3.3 and 5.2. (In examples 2 and 3 we will reconsider the same model and the functions  $v$  and  $\kappa$  from a more semi-parametric perspective, using score operators as in sections 3.4 and 5.4.)

It is straightforward to show that  $v$  is pathwise differentiable with pathwise derivative  $\dot{v}$  given by

$$(7) \quad \begin{aligned} \dot{v}(h)(t) &= \int_0^\infty \left\{ \frac{1_{[0,t]}(s)}{\bar{H}(s)} - C(s \wedge t) \right\} h(s, 1) p(s, 1) ds \\ &\quad - \int_0^\infty C(s \wedge t) h(s, 0) p(s, 0) ds, \end{aligned}$$

where  $h \in \dot{\mathbf{P}}$  and

$$(8) \quad C(t) \equiv \int_0^t \frac{1}{\bar{H}^2} dH_1 = \int_0^t \frac{1}{\bar{H}} d\Lambda.$$

Hence  $\dot{v}_t$  of (5.2.25) is given by

$$(9) \quad \dot{v}_t(s, \delta) = \delta \frac{1_{[0,t]}(s)}{\bar{H}(s)} - C(s \wedge t)$$

for  $(s, \delta) \in R^+ \times \{0, 1\}$  and  $0 \leq t \leq \tau$ , with  $\bar{H}(\tau) > 0$ . Thus from (5.2.8) and (5.2.25) it follows that  $\tilde{I}_t \equiv \tilde{I}_v(\pi_t) = \dot{v}_t^T(\pi_t)$  is just  $\dot{v}_t(s, \delta)$  given by (9), and the inverse information covariance function  $I_v^{-1}$  of (5.2.23) is

$$(10) \quad I_v^{-1}(s, t) = E \dot{v}_s(T, \Delta) \dot{v}_t(T, \Delta) = C(s \wedge t)$$



after two integrations by parts:

$$\begin{aligned} \int_0^t C^2 dH &= - \int_0^t C^2 d\bar{H} = \int_0^t \bar{H} d(C^2) - C^2(t)\bar{H}(t) \\ &= 2 \int_0^t \frac{C}{\bar{H}} dH_1 - C^2(t)\bar{H}(t), \end{aligned}$$

and, for  $s \leq t$ ,

$$\int_s^t C dH = - \int_s^t C d\bar{H} = \int_s^t \frac{1}{\bar{H}} dH_1 - \bar{H}(t)C(t) + \bar{H}(s)C(s).$$

Thus the Gaussian process  $\mathbf{Z}_0$  of theorem 5.2.1 is

$$(11) \quad \mathbf{Z}_0(t) = B(C(t)), \quad 0 \leq t \leq \tau,$$

where  $B$  is standard Brownian motion. (Done this way the calculation is straightforward but tedious; we will do it again via score operators and martingale methods in example 2.)

Now consider estimation of

$$\begin{aligned} (12) \quad \kappa(P)(t) &\equiv 1 - F(t) = \prod_{0 \leq s \leq t} (1 - d\Lambda(s)) \\ &= \exp(-\Lambda(t)) \quad \text{since } F, \text{ and hence } \Lambda, \text{ is continuous} \\ &= \alpha(v(P))(t) \end{aligned}$$

for  $0 \leq t \leq \tau$ , where  $\phi(b) \equiv \exp(-b) \in l^\infty([0, \tau])$  for  $b \in l^\infty([0, \tau]) \equiv l^\infty$ . Then  $\dot{\phi}: l^\infty \rightarrow l^\infty$  is simply the multiplication operator

$$(13) \quad \dot{\phi}(b)h = -\exp(-b)h = -\phi h,$$

and hence for  $0 \leq t \leq \tau$  and  $(s, \delta) \in R^+ \times \{0, 1\}$  (cf. (6.4.19))

$$\begin{aligned} (14) \quad \dot{\kappa}_t(s, \delta) &= \dot{\phi}(\Lambda) \dot{v}_t(s, \delta) = -\exp(-\Lambda(t)) \dot{v}_t(s, \delta) \\ &= -\bar{F}(t) \dot{v}_t(s, \delta). \end{aligned}$$

Therefore, as in (6.4.20), the inverse information covariance function (5.2.23) for estimation of  $\kappa \equiv \bar{F}$  is

$$\begin{aligned} (15) \quad I_\kappa^{-1}(s, t) &= \bar{F}(s)\bar{F}(t)C(s \wedge t) \\ &= \frac{\bar{F}}{K}(s) \frac{\bar{F}}{K}(t) \{K(s \wedge t) - K(s)K(t)\} \end{aligned}$$

where

$$(16) \quad K \equiv \frac{C}{1+C}, \quad \bar{K} \equiv 1 - K = \frac{1}{1+C}.$$

Thus for estimation of  $\kappa$  the process  $\mathbf{Z}_0$  of theorem 5.2.1 has

$$(17) \quad \mathbf{L}(\mathbf{Z}_0) = \mathbf{L}(\bar{F}B(C)) \quad \text{where } B \text{ is standard Brownian motion} \\ = \mathbf{L}\left(\frac{\bar{F}}{\bar{K}}B_0(K)\right) \quad \text{where } B_0 \text{ is the Brownian bridge .}$$

As is well known, the bound for estimation of  $v = \Lambda$  given by (10) and (11) is achieved by the nonparametric maximum likelihood, or Nelson-Aalen, estimator  $\hat{v}_n \equiv \hat{\Lambda}_n = v(P_n)$  obtained from (5):

$$(18) \quad \hat{v}_n(t) = \hat{\Lambda}_n(t) = \int_0^t \frac{dH_{n1}}{\bar{H}_{n-}}, \quad 0 \leq t \leq \tau,$$

where  $H_{n1}$  and  $H_n$  are the natural empirical estimators of  $H_1$  and  $H$  respectively, and  $\bar{H}_{n-}(t) = 1 - H_n(t-)$ . Similarly, the nonparametric maximum likelihood (or product-limit or Kaplan-Meier) estimator obtained from (18) and (6),

$$(19) \quad \hat{\kappa}_n(t) \equiv 1 - \hat{F}_n(t) = \prod_{0 \leq s \leq t} (1 - d\hat{\Lambda}_n(s)),$$

achieves the bound given by (15) and (17). See e.g. Breslow and Crowley (1974), Gill (1980), (1983), or Shorack and Wellner (1986, chapter 7). Identification of  $\mathbf{Z}_0$  in terms of a Brownian bridge process  $B_0$  is from Hall and Wellner (1980b). The preceding calculations are essentially the same as those in Wellner (1982).  $\square$

### Example 2. Random censoring; score operators and martingale theory.

We now reconsider estimation of both  $v$  and  $\kappa$  from a more semiparametric perspective using the tools and methods of sections 5.4 and 5.5. We also use the  $R$  and  $L$  operators introduced in appendices A.1 and A.3.

It follows from proposition A.5.5.A that the model  $\mathbf{P}$  is differentiable and, for  $a \in \dot{\mathbf{G}}_1 \equiv L_2^0(F)$ ,  $b \in \dot{\mathbf{G}}_2 \equiv L_2^0(G)$ ,

$$\dot{\mathbf{i}}(a, b) = \dot{\mathbf{i}}_1 a + \dot{\mathbf{i}}_2 b,$$

where  $\dot{\mathbf{i}}_i : \dot{\mathbf{G}}_i \rightarrow \dot{\mathbf{P}}$ ,  $i = 1, 2$ , are given, much as in proposition 4.6.1, by

$$(20) \quad \dot{\mathbf{i}}_1 a = \Delta a + (1 - \Delta) \frac{\int_0^\infty a dF}{1 - F} = \Delta a + (1 - \Delta)E(a(Y) | Y > \cdot)$$

or

$$(20') \quad \dot{\mathbf{i}}_1 a(T, \Delta) = \int R_F a dM_{uc};$$

and

$$(21) \quad \dot{\mathbf{i}}_2 b = (1 - \Delta)b + \Delta \frac{\int_0^\infty b dG}{1 - G} = (1 - \Delta)b + \Delta E(b(C) | C > \cdot),$$

or

$$(21') \quad \dot{\mathbf{i}}_2 b(T, \Delta) = \int R_G b dM_c.$$

Here

$$M_{uc}(t) \equiv 1_{[T \leq t, \Delta = 1]} - \int_0^t 1_{[T \geq s]} d\Lambda(s)$$

and

$$M_c(t) \equiv 1_{[T \leq t, \Delta = 0]} - \int_0^t 1_{[T \geq s]} \frac{dG(s)}{1 - G(s)}$$

and  $R_F$  and  $R_G$  are the  $R$  operators corresponding to  $F$  and  $G$ ; see (A.1.19). Now by a martingale calculation

$$(22) \quad \langle \dot{I}_1 a, \dot{I}_2 b \rangle_0 = E \int_0^\infty R_F a R_G b d \langle M_{uc}, M_c \rangle = 0$$

since  $\langle M_{uc}, M_c \rangle = 0$  a.s. Thus  $\mathbf{R}(\dot{I}_1) \perp \mathbf{R}(\dot{I}_2)$  and

$$\dot{I}_1^T \dot{I}_2 = 0, \quad \dot{I}_2^T \dot{I}_1 = 0.$$

Now we extend this orthogonality to  $\dot{P}_1 \perp \dot{P}_2$  with  $\dot{P}_1 \supset \overline{\mathbf{R}(\dot{I}_1)}$  and  $\dot{P}_2 \supset \mathbf{R}(\dot{I}_2)$ . In the process we will improve our understanding of  $\dot{P}_i, i = 1, 2$  and  $\dot{P}$ . Let  $P_\eta \in \mathbf{P}_1$  with density  $p_\eta$  corresponding to  $f_\eta$  and  $g$  be such that  $\{p_\eta\}$  is a curve with tangent  $h \in \dot{P}_1^0$ . Then by Cauchy-Schwarz or lemma A.5.1 with  $Z \equiv 0$ ,

$$\begin{aligned} o(\eta^2) &= \sum_{\delta=0}^1 \int_0^\infty \{p_\eta^{1/2}(t, \delta) - p_0^{1/2}(t, \delta) - \frac{1}{2}\eta h(t, \delta)p_0^{1/2}(t, \delta)\}^2 dt \\ &\geq \int_0^s \{[1 + \frac{1}{\bar{F}(t)} \int_t^\infty (f_\eta - f_0)]^{1/2} - [1 + \frac{1}{2}\eta h(t, 0)]\}^2 \bar{F}(t) dG(t) \\ &\quad + \int_0^s \{f_\eta^{1/2}(t) - f_0^{1/2}(t)[1 + \frac{1}{2}\eta h(t, 1)]\}^2 \bar{G}(t) dt \\ &\geq \int_0^s \bar{F} dG \left\{ [1 + \frac{\int_0^s \int_0^t (f_0 - f_\eta)(u) du dG(t)}{\int_0^s \bar{F} dG}]^{1/2} \right. \\ &\quad \left. - [1 + \eta \frac{\int_0^s h(t, 0) \bar{F}(t) dG(t)}{\int_0^s \bar{F} dG} \right. \\ &\quad \left. + \frac{1}{4}\eta^2 \frac{\int_0^s h^2(t, 0) \bar{F}(t) dG(t)}{\int_0^s \bar{F} dG}]^{1/2} \right\}^2 \\ &\quad + s\bar{G}(s) \left\{ [\frac{1}{s} \int_0^s f_\eta]^{1/2} - [\frac{1}{s} \int_0^s [1 + \frac{1}{2}\eta h(t, 1)]^2 f_0(t) dt]^{1/2} \right\}^2 \end{aligned}$$

uniformly in  $s$ . Consequently

$$\int_0^s \int_0^t (f_0 - f_\eta)(u) du dG(t) = \eta \int_0^s h(t, 0) \bar{F}(t) dG(t) + o(|\eta|),$$

$$\int_0^s f_\eta = \int_0^s [1 + \frac{1}{2}\eta h(t, 1)]^2 f_0(t) dt + [\bar{G}(s)]^{-1/2} o(|\eta|),$$

and hence

$$\begin{aligned} & - \int_0^s [\eta \int_0^t h(u, 1) f_0(u) du \\ & + \frac{1}{4}\eta^2 \int_0^t h^2(u, 1) f_0(u) du + [\bar{G}(t)]^{-1/2} o(|\eta|)] dG(t) \\ & = \eta \int_0^s h(t, 0) \bar{F}(t) dG(t) + o(|\eta|) \end{aligned}$$

uniformly in  $s$ . It follows that

$$\eta \int_0^s [h(t, 0) + \frac{1}{\bar{F}(t)} \int_0^t h(u, 1) dF(u)] \bar{F}(t) dG(t) = o(|\eta|),$$

and hence

$$(23) \quad h(t, 0) = - \frac{1}{\bar{F}(t)} \int_0^t h(u, 1) dF(u), \quad 0 \leq t \leq \tau_H.$$

Let  $F^*(t) \equiv \int_0^t \bar{G} dF / \int_0^\infty \bar{G} dF$ . For  $a \in L_2(F^*)$  we define

$$\dot{I}_1 a = \Delta a - \frac{1-\Delta}{\bar{F}} \int_0^\cdot a dF.$$

Note that this agrees with (20) for  $a \in L_2^0(F)$  and that (20') can be generalized in the same vein. With this notation (23) implies

$$(24) \quad \overline{\mathbf{R}(\dot{I}_1)} = [\dot{I}_1 a : a \in L_2^0(F)] \subset \dot{\mathbf{P}}_1 \subset [\dot{I}_1 a : a \in L_2^0(F^*)]$$

in  $L_2^0(P)$ . A similar string of inequalities holds for  $\dot{I}_2$ , and (22) is also valid with  $a \in L_2(F^*)$  and  $b \in L_2(G^*)$ . Consequently,

$$(25) \quad \dot{\mathbf{P}}_1 \perp \dot{\mathbf{P}}_2,$$

which is an infinite-dimensional example of the orthogonality condition for adaptation given in proposition 3.4.3. In fact, equalities hold in (24) if  $\int (1/\bar{F}) dG < \infty$ . To prove this, it suffices to consider  $a \in L_2(F^*)$  and to construct  $a_M \in L_2^0(F)$  with  $\dot{I}_1 a_M$  arbitrarily close to  $\dot{I}_1 a$ . Since (the measurable subset of)  $l^\infty([0, \infty))$  is a dense subset of  $L_2(F^*)$ , we assume without loss of generality that  $a$  is bounded. Define

$$a_M(t) = a(t) 1_{[t \leq M]} - \frac{1}{\bar{F}(M)} \int_0^M a dF 1_{[t > M]}$$

with  $0 < M < \tau_H = \tau_G$ . Indeed,  $a_M \in L_2^0(F)$  and by the boundedness of  $a$

$$\begin{aligned} & E[\dot{I}_1(a - a_M)(T, \Delta)]^2 \\ & = \int_M^\infty [a + \frac{1}{\bar{F}(M)} \int_0^M a dF]^2 \bar{G} dF \end{aligned}$$

$$\begin{aligned}
 & + \int_M^\infty \frac{1}{\bar{F}(t)} \left\{ \int_M^t \left[ a + \frac{1}{\bar{F}(M)} \int_0^M a dF \right] dF \right\}^2 dG(t) \\
 & = o(1) + O\left( \frac{1}{\bar{F}^2(M)} \int_M^\infty \bar{G} dF \right) + O\left( \int_M^\infty \frac{1}{\bar{F}} dG \right) \\
 & = o(1)
 \end{aligned}$$

as  $M \uparrow \tau_H$  since

$$\frac{1}{\bar{F}^2(M)} \int_M^\infty \bar{G} dF \leq \frac{\bar{G}(M)}{\bar{F}(M)} \leq \int_M^\infty \frac{1}{\bar{F}} dG .$$

Hence the extreme sides of (24) are equal and

$$(26) \quad \overline{\mathbf{R}(\dot{\mathbf{I}}_1)} = \dot{\mathbf{P}}_1 .$$

The information operator for  $f, \dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1$ , is given by

$$\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1 a(s) = R_F^T (L_{uc}^T L_{uc}) R_F a(s)$$

where  $L_{uc} : L_2(F) \rightarrow L_2(P)$  is defined by

$$(27) \quad L_{uc} a(T, \Delta) = \int a dM_{uc} .$$

Since  $R_F^T = L_F$  and  $R_F$  are inverses of each other on  $L_2^0(F)$  by proposition A.1.8, to calculate (or invert)  $\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1$  it suffices to calculate (or invert)  $L_{uc}^T L_{uc}$ . To calculate  $L_{uc}^T L_{uc}$ , let  $a, b \in L_2(F)$ . Then

$$\begin{aligned}
 \langle L_{uc}^T L_{uc} a, b \rangle_{L_2(F)} & = \langle L_{uc} a, L_{uc} b \rangle_0 \quad \text{by definition of } L_{uc}^T \\
 & = E \int_0^\infty a(s) b(s) 1_{[T \geq s]} d\Lambda(s) \\
 & \quad \text{by (27) and martingale calculus} \\
 & = E \int_0^\infty a(s) b(s) 1_{[Y \geq s]} 1_{[C \geq s]} d\Lambda(s) \\
 & = \int_0^\infty a(s) b(s) (1 - G(s)) dF(s) \\
 & = \langle \bar{G} a, b \rangle_{L_2(F)} ,
 \end{aligned}$$

where  $\bar{G} \equiv 1 - G$ . Since this holds for all  $b \in L_2(F)$ , it follows that  $L_{uc}^T L_{uc} a$  is the diagonal operator

$$L_{uc}^T L_{uc} a(s) = \bar{G}(s) a(s) \equiv D a(s) .$$

Thus

$$(28) \quad \dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1 = R_F^{-1} D R_F ,$$

with  $D$  diagonal, which is a form of the spectral decomposition of the positive self-adjoint operator  $\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1$ ; cf. Reed and Simon (1972, page 221). Note that  $\bar{G}(t) \downarrow 0$  as  $t \uparrow \tau_G$ . It follows that whenever  $\tau_G = \tau_H \equiv \tau_F \wedge \tau_G$  (which is usually the case of interest in practice),  $L_{uc}^T L_{uc} = D$  does *not* have closed range; see, e.g., examples A.1.10 and A.1.11. Now  $a 1_{[0, \tau_G - \varepsilon]} \in \mathbf{R}(D)$  for every  $a \in L_2(F)$  and  $\varepsilon > 0$ , and hence

$$\overline{\mathbf{R}(D)} = \{a \in L_2(F) : a(x) = 0 \text{ for } F\text{-a.e. } x > \tau_G\}.$$

Assume  $\mathbf{R}(\dot{\mathbf{I}}_1)$  is closed. Then  $\mathbf{R}(\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1) = \mathbf{R}(\dot{\mathbf{I}}_1^T)$  is closed by corollary 5.4.2.C and proposition A.1.7.D. Since  $R_F$  is an isometry,  $\mathbf{R}(\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1)$  is closed if and only if  $\mathbf{R}(D)$  is closed. However, in the practically important case  $\tau_G \leq \tau_F$ , as we just saw,  $\mathbf{R}(D)$  is *not* closed, and hence  $\mathbf{R}(\dot{\mathbf{I}}_1)$  fails to be closed.

Nevertheless, we may still apply theorem 5.4.1 and its corollaries. Consider estimation of  $\kappa(P_{(F,G)})(t) \equiv 1 - F(t)$  for  $0 \leq t \leq \tau < \tau_H$ . As in example 5.3.1, the map  $\psi(f) \equiv 1 - F$  has gradient

$$\dot{\psi}_t \equiv \dot{\psi}^T(\pi_t) = - (1_{[0,t]} - F(t)), \quad 0 \leq t \leq \tau.$$

To verify differentiability of  $\pi_t \kappa(P_{(F,G)}) = 1 - F(t)$  via corollary 5.4.2.C and theorem 5.4.1, we need (26) and we want to check that

$$(29) \quad \dot{\psi}_t = \dot{\psi}^T(\pi_t) \in \mathbf{R}(\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1) = \mathbf{R}(R_F^{-1} D R_F).$$

But

$$(30) \quad R_F \dot{\psi}_t = - 1_{[0,t]} \frac{\bar{F}(t)}{\bar{F}},$$

and hence, with  $a^* \equiv R_F^{-1}(-1_{[0,t]} \bar{F}(t)/\bar{H}) \in L_2(F)$ ,

$$\begin{aligned} \dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1 a^* &= \dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1 (R_F^{-1}(-1_{[0,t]} \frac{\bar{F}(t)}{\bar{H}})) \\ &= R_F^{-1} D (-1_{[0,t]} \frac{\bar{F}(t)}{\bar{H}}) \\ &= R_F^{-1}(-1_{[0,t]} \frac{\bar{F}(t)}{\bar{F}}) = \dot{\psi}_t. \end{aligned}$$

Thus (29) holds, and, by corollary 5.4.2.C and (5.4.16), the efficient influence function for  $\pi_t \kappa$  is given by

$$\begin{aligned} \tilde{\mathbf{I}}_\kappa(\pi_t)(T, \Delta) &= \dot{\mathbf{I}}_1 a^* \\ &= -\bar{F}(t) \int_0^\infty \frac{1_{[0,t]}}{\bar{H}} dM_{uc} \quad \text{by (20')} \\ (31) \quad &= -\bar{F}(t) \left\{ \frac{\Delta 1_{[0,t]}(T)}{\bar{H}(T)} - C(T \wedge t) \right\}, \end{aligned}$$

in agreement with (14) and (9). By the martingale representation in (31)

$$\begin{aligned}
 I_{\kappa}^{-1}(s, t) &= \bar{F}(s)\bar{F}(t) \int_0^{s \wedge t} \frac{1}{\bar{H}^2(y)} E 1_{[T \geq y]} d\Lambda(y) \\
 &= \bar{F}(s)\bar{F}(t) \int_0^{s \wedge t} \frac{1}{\bar{H}(y)} d\Lambda(y) \\
 (32) \qquad &= \bar{F}(s)\bar{F}(t) C(s \wedge t)
 \end{aligned}$$

just as in (15).

To discuss differentiability of the whole function  $\kappa(P_{(F,G)}) \equiv 1 - F \in I^\infty([0, \infty)) \equiv \mathbf{B}$  on  $[0, \tau_H]$  (here we will assume, for simplicity, that  $\tau_H = \tau_F \wedge \tau_G < \infty$ ), we need to verify the hypothesis (5.4.6) of theorem 5.4.1. To do this, we first compute  $\dot{\psi}$  and  $\dot{\psi}^T$ : from example 5.3.1,  $\psi$  has derivative  $\dot{\psi}$  given by

$$\dot{\psi}(a)(t) = \int \dot{\psi}_t a dF = \langle \dot{\psi}_t, a \rangle_{L_2(F)}$$

for  $a \in L_2^0(F)$  and  $t \in [0, \tau_H]$ . Thus for  $F$  continuous,  $\dot{\psi}: L_2^0(F) \rightarrow \mathbf{B}_0 \equiv C[0, \tau_H] \subset I^\infty(R^+) = \mathbf{B}$ . Thus  $\mathbf{R}(\dot{\psi}^T) = \mathbf{R}(\dot{\psi}_0^T)$ , where  $\dot{\psi}_0^T: \mathbf{B}_0^* \rightarrow L_2^0(F)$  is given by

$$\dot{\psi}_0^T m(x) = \int \dot{\psi}_t(x) dm(t), \quad x \in [0, \tau_H],$$

for  $m \in \mathbf{B}_0^* = \{\text{all finite signed measures on } [0, \tau_H]\}$ ; see example A.1.9. Note that  $\dot{\psi}_0^T m$  is (uniformly) bounded as a function of  $x$  for each fixed measure  $m$ . Thus,

$$\mathbf{R}(\dot{\psi}^T) = \{ \dot{\psi}_0^T m : m \in \mathbf{B}_0^* = C([0, \tau_H])^* \}.$$

Now, by the polar decomposition of  $\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1$  (see proposition A.1.6),

$$\mathbf{R}(\dot{\mathbf{I}}_1^T) = \mathbf{R}((\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1)^{1/2}) = \mathbf{R}(R_F^{-1} D^{1/2} R_F),$$

where  $D^{1/2} a \equiv a \bar{G}^{1/2}$ . Thus to verify (5.4.6), we want to show that

$$\dot{\psi}_0^T m = R_F^{-1} D^{1/2} R_F a \quad \text{for some } a \in L_2^0(F).$$

Since  $R_F^{-1} = L_F$  is an isometry, this reduces to showing that

$$(33) \quad \int_0^{\tau_H} \frac{(R_F \dot{\psi}_0^T m)^2}{\bar{G}} dF < \infty.$$

But if  $\int_0^{\tau_H} \bar{G}^{-1} dF < \infty$ , (33) holds by boundedness of  $\dot{\psi}_0^T m$  since  $R_F$  maps bounded functions to bounded functions. Hence (5.4.6) holds and  $\kappa \equiv \bar{F}$  is pathwise differentiable if  $\int_0^{\tau_H} \bar{G}^{-1} dF < \infty$ .

On the other hand, if  $\kappa = 1 - F$  is pathwise differentiable (as a map from  $\mathbf{P}$  to  $\mathbf{B}_0 = C([0, \tau_H])$ ), then so is  $\kappa(P_{(F,G)}) = 1 - F(\tau_H)$  (as a map from  $\mathbf{P}$  to  $R$ ); so by theorem 5.4.1

$$\dot{\Psi}_{\tau_H} \in \mathbf{R}(\dot{\mathbf{I}}^T) = \mathbf{R}((\dot{\mathbf{I}}^T \dot{\mathbf{I}})^{1/2}),$$

and hence

$$\dot{\Psi}_{\tau_H} = R_F^{-1} D^{1/2} R_F a$$

for some  $a \in L_2^0(F)$ . But since  $R_F^{-1} = L_F$  is an isometry, this implies that

$$\int_0^{\tau_H} \frac{(R_F \dot{\Psi}_{\tau_H})^2}{\bar{G}} dF < \infty.$$

But by (30),  $R_F \dot{\Psi}_{\tau_H} = -1_{[0, \tau_H]} \bar{F}(\tau_H) / \bar{F}$ , and this yields

$$\bar{F}^2(\tau_H) \int_0^{\tau_H} \frac{1}{\bar{G}} dF \leq \bar{F}^2(\tau_H) \int_0^{\tau_H} \frac{1}{\bar{F}^2 \bar{G}} dF < \infty,$$

and hence  $\int_0^{\tau_H} \bar{G}^{-1} dF < \infty$  if  $\bar{F}(\tau_H) > 0$  (which of course implies  $\tau_H = \tau_G < \tau_F$ ).

What happens if the censoring distribution  $G$  is known? Since

$$\overline{\mathbf{R}(\dot{\mathbf{I}}_1)} \subset \dot{\mathbf{P}}_1 \subset [\dot{\mathbf{I}}_1 a : a \in L_2(F^*)] \subset \dot{\mathbf{P}} = L_2^0(P),$$

where the third inclusion is strict, the tangent space of the model  $\mathbf{P}_1$  is *not* "saturated" or "full". Because  $\dot{\mathbf{P}}_2$  is orthogonal to  $\dot{\mathbf{P}}_1$  (recall (25)), the information bound for estimation of  $\nu$  or  $\kappa$  in the smaller model  $\mathbf{P}_1$  in which the censoring distribution  $G$  is known, is exactly the same as in the bigger model  $\mathbf{P}$  with  $G$  unknown: for example,

$$\Gamma_{\kappa}^{-1}(s, t | \mathbf{P}_1) = \Gamma_{\kappa}^{-1}(s, t | \mathbf{P}) = \bar{F}(s) \bar{F}(t) C(s \wedge t)$$

as given by (15) and (32). Note that  $\dot{\kappa}_t$  and  $\dot{\nu}_t$  given by (9), (14), and (31) are in  $\mathbf{R}(\dot{\mathbf{I}}_1) \subset \overline{\mathbf{R}(\dot{\mathbf{I}}_1)} \subset \dot{\mathbf{P}}_1$ . But because  $\dot{\mathbf{P}}_1$  is a proper subspace of  $\dot{\mathbf{P}}$ , there are now other estimators which are  $\sqrt{n}$ -consistent but *inefficient*: e.g., the ("reduced sample") estimator using only the  $T$ 's with  $\Delta = 1$ ; since

$$F(t) = \frac{\int_0^t \frac{1}{\bar{G}(s)} dH_1(s)}{\int_0^{\infty} \frac{1}{\bar{G}(s)} dH_1(s)}, \quad 0 \leq t \leq \tau < \tau_H,$$

where  $\bar{G}$  is known and  $H_1$  is defined after (1), an alternative estimator (to the Kaplan-Meier estimator) is

$$\mathbf{F}_n^{(1)}(t) = \frac{\int_0^t \frac{1}{\bar{G}} dH_{n1}}{\int_0^{\infty} \frac{1}{\bar{G}} dH_{n1}}, \quad 0 \leq t \leq \tau.$$



Assuming that  $\int_0^{\tau_n} \bar{G}^{-1} dF < \infty$ , it is easily shown that the influence function of  $F_n^{(1)}(t)$  is  $\Delta(1_{[0,t]}(v) - F(t))/\bar{G}(v) \in L_2^0(P)$ , does not satisfy (23) and hence cannot belong to  $\dot{P}_1$  in view of (24). Consequently,  $F_n^{(1)}$  is inefficient.  $\square$

**Example 3. Random censoring; score operators and integral equation theory.**

In example 2 we used the score operators  $\dot{I}_i$  together with martingale methods to derive the efficient influence operators for estimation of  $v = \Lambda$  and  $\kappa = 1 - F$ . In more complicated censoring models, the martingales introduced there will no longer be available, and it will be useful to consider the information operator  $\dot{I}_1^T \dot{I}_1$  by use of classical integral equation methods. Although this is not necessary in the simple random censoring model, consideration of the more classical methods in the simple problem will provide a useful point of comparison and reference when we consider the more complicated double and bivariate censoring models in examples 5 and 6.

The model is exactly as in examples 1 and 2, and, in particular, the score operators  $\dot{I}_i$ ,  $i = 1, 2$ , are given by (20) and (21). But we now calculate the adjoint  $\dot{I}_1^T : \dot{P} \rightarrow \dot{G}_1$  and information operator  $\dot{I}_1^T \dot{I}_1 : \dot{G}_1 \rightarrow \dot{G}_1 = L_2^0(F)$  directly. By straightforward calculation

$$(34) \quad \dot{I}_1^T h(x) = \bar{G}(x)h(x, 1) + \int_0^x h(y, 0) dG(y)$$

for  $h \in \dot{P} = L_2^0(P)$ . Therefore, by (20),

$$\begin{aligned} \dot{I}_1^T \dot{I}_1 a(x) &= \bar{G}(x)a(x) + \int_0^x \frac{\int_y^\infty a dF}{\bar{F}(y)} dG(y) \\ &= \bar{G}(x)a(x) + \int_0^\infty \left( \int_0^{x \wedge y} \frac{1}{\bar{F}} dG \right) a(y) dF(y) \\ (35) \quad &\equiv \bar{G}(x)a(x) + \int_0^\infty K(x,y)a(y) dF(y). \end{aligned}$$

While we know from example 2 that the information operator  $\dot{I}_1^T \dot{I}_1$  is usually not boundedly invertible (it is one-to-one but not onto), the integral equation  $b = \dot{I}_1^T \dot{I}_1 a$ , or

$$(36) \quad b(x) = \bar{G}(x)a(x) + \int_0^\infty K(x,y)a(y) dF(y)$$

may have (does have, in view of our calculations in example 2) a solution for particular functions  $b$ .

We consider estimation of  $F$  on  $[0, \tau]$  where  $\tau$  satisfies

$$(37) \quad \bar{G}(\tau) \equiv P(C \geq \tau) \equiv \varepsilon > 0.$$

Since we are considering estimation of  $F$  only on  $[0, \tau]$ , we can restrict attention

to functions  $a \in L_2^0(F)$  which vanish off  $[0, \tau]$ . Thus we define the Hilbert space

$$(38) \quad \mathbf{H} \equiv \{a \in L_2^0(F) : a = a 1_{[0, \tau]}\} \subset L_2^0(F).$$

We will show that (37) implies that  $\dot{I}_1^T \dot{I}_1 : \mathbf{H} \rightarrow \mathbf{H}$  is one-to-one and onto, and has a bounded inverse on  $\mathbf{H}$  with  $\|(\dot{I}_1^T \dot{I}_1)^{-1}\| \leq \epsilon^{-1}$ .

Since  $\dot{I}_1^T \dot{I}_1$  is self-adjoint, by corollary A.1.2 we need only verify that

$$\langle a, \dot{I}_1^T \dot{I}_1 a \rangle \geq \epsilon \|a\|^2 \quad \text{for all } a \in \mathbf{H}.$$

But from (35), for  $a \in \mathbf{H}$  and if  $\int_{-\infty}^{\tau} (1/\bar{F}) dG < \infty$  (which holds trivially when  $\tau_F > \tau_G > \tau$ )

$$\begin{aligned} \langle a, \dot{I}_1^T \dot{I}_1 a \rangle &= \int a^2 \bar{G} dF + \int_0^{\infty} \int_0^{\infty} a(x) \int_0^{x \wedge y} \frac{1}{\bar{F}} dG a(y) dF(y) dF(x) \\ &= \int a^2 \bar{G} dF + \int_0^{\tau} \frac{1}{\bar{F}(u)} \left\{ \int_u^{\infty} a dF \right\}^2 dG(u) \\ &\geq \epsilon \|a\|^2 \end{aligned}$$

by (37) and (38). Thus  $\dot{I}_1^T \dot{I}_1$  is boundedly invertible on  $\mathbf{H}$ , and, by (5.4.14),  $\kappa(P) = 1 - F$  (restricted to  $[0, \tau]$ ) is again pathwise differentiable. In this case, we know much more, as already shown in examples 1 and 2, but the type of argument given above will prove useful in more complicated problems.  $\square$

**Example 4. Censored linear regression; example 4.6.4 continued.**

Now consider the censored linear regression model given in example 4.6.4:  $\mathbf{P} = \mathbf{Q}T^{-1}$  where  $\mathbf{Q}$  is the linear regression model with error density  $f$  and  $X = T(X^0) = (Z, Y \wedge C, 1_{[Y \leq C]}) \equiv (Z, V, \Delta)$  in the notation used in this section and section 4.6 and continued here. Consider estimation of

$$(40) \quad \kappa(P_{(v, F, H)})(t) = 1 - F(t), \quad -\infty < t \leq \tau < \tau_F \wedge \tau_G,$$

where  $H$  is the distribution of  $(Z, C)$ . Furthermore, we will use  $G(s) \equiv P(C - v^T Z \leq s) \equiv P(\delta \leq s)$ .

Let  $v$  and  $f$  (or  $F$ ) be labeled by 1 and 2 respectively as in section 4.6. Suppose first that  $v$  is known (so  $P \in \mathbf{P}_2$ ). Then by calculations exactly as in example 2,  $\kappa$  defined in (40) is pathwise differentiable: we obtain

$$(41) \quad \tilde{\Gamma}_{\kappa}(\pi_t)(T | \mathbf{P}_2) = -\bar{F}(t) \int_{-\infty}^t \frac{1}{\bar{F} \bar{G}} dM_{uc},$$

where, as in (4.6.22),

$$M_{uc}(t) \equiv 1_{[e \wedge \delta \leq t, \Delta = 1]} - \int_{-\infty}^t 1_{[e \wedge \delta \geq s]} d\Lambda(s),$$

and hence

$$(42) \quad \Gamma_{\kappa}^{-1}(s, t | \mathbf{P}_2) = \bar{F}(s) \bar{F}(t) \int_{-\infty}^{s \wedge t} \frac{1}{\bar{F} \bar{G}} d\Lambda.$$

This is exactly as in (31) and (32), but with  $\bar{H}$  replaced by  $\bar{F}\bar{G}$  and  $F(t) = P(\epsilon \leq t)$ ,  $G(t) = P(\delta \leq t)$ . An obvious estimator is just the Kaplan-Meier estimator based on  $\{(V_i - v^T Z_i, \Delta_i) : i = 1, \dots, n\}$ .

When  $v$  is unknown, so  $P \in \mathbf{P}$ , then we can compute  $\tilde{I}_K$  and  $I_K^{-1}$  via corollary 5.5.2. From proposition 4.6.1 (with  $R = R_F$ )

$$I_1^*(T | \mathbf{P}) = \int (Z - E(Z | V - v^T Z \geq s)) R\psi(s) dM_{uc}(s),$$

$I(v, \mathbf{P})$  is as given in (4.6.28), and (assuming that  $I(v, \mathbf{P})$  is nonsingular)  $\tilde{I}_1 = I^{-1}(v, \mathbf{P})I_1^*$ . Furthermore from the  $a^*$  of the proof of proposition 4.6.1,

$$\begin{aligned} <\dot{\psi}^T(\pi_t), (\dot{I}_2^T \dot{I}_2)^{-1} \dot{I}_2^T \dot{I}_1 >_{L_2(F)} \\ &= < -(1_{[0,t]} - F(t)), L(E(Z | \delta \geq \cdot) R\psi) >_{L_2(F)} \\ &= < -1_{[0,t]} \frac{\bar{F}(t)}{\bar{F}}, E(Z | \delta \geq \cdot) R\psi >_{L_2(F)} \quad \text{by (30)} \\ &= -\bar{F}(t) \int_{-\infty}^t E(Z | \delta \geq s) R\psi(s) d\Lambda(s), \end{aligned}$$

and hence, from (5.5.27) and (5.5.28),

$$(43) \quad \tilde{I}_K(\pi_t)(T | \mathbf{P}) = -\bar{F}(t) \left\{ \int_{-\infty}^t \frac{1}{\bar{F}\bar{G}} dM_{uc} - \left( \int_{-\infty}^t E(Z^T | \delta \geq \cdot) R\psi d\Lambda \right) \tilde{I}_1(T | \mathbf{P}) \right\},$$

and

$$(44) \quad I_K^{-1}(s, t | \mathbf{P}) = \bar{F}(s)\bar{F}(t) \left\{ \int_{-\infty}^{s \wedge t} \frac{1}{\bar{F}\bar{G}} d\Lambda + \left( \int_{-\infty}^s E(Z^T | \delta \geq \cdot) R\psi d\Lambda \right) \cdot I^{-1}(v | \mathbf{P}) \left( \int_{-\infty}^t E(Z | \delta \geq \cdot) R\psi d\Lambda \right) \right\}.$$

Construction of an estimator achieving this bound can be based on the Kaplan-Meier estimator for the pairs  $\{(V_i - \hat{v}_n^T Z_i, \Delta_i) : i = 1, \dots, n\}$  where  $\hat{v}_n$  is an estimator achieving the information bound for estimation of  $v$  given by proposition 4.6.1. The latter involves estimation of  $R\psi = -\lambda' / \lambda$  using residuals from a  $\sqrt{n}$ -consistent preliminary estimator  $\tilde{v}_n$ . As far as we know, this has not yet been carried out in full detail. But see the work by Ritov (1984), (1990), Tsiatis (1990), and Lai and Ying (1988), (1991a), (1991b).  $\square$

**Example 5. Double censoring.**

In this model, the survival time of interest,  $Y$  with df  $F$ , is censored both above and below by a pair of rv's  $(C, D)$  independent of  $Y$  with joint distribution  $G$  satisfying  $P(C \leq D) = 1$ . Here  $X^0 \equiv (Y, C, D)$  with  $Y \sim F$  independent of  $(C, D) \sim G$ , and using the notation of section 4.6, we observe  $X = T(X^0) \equiv (V, \Delta)$  defined by

$$(45) \quad V \equiv (D \wedge Y) \vee C$$

and

$$(46) \quad \Delta = \begin{cases} 1 & \text{if } C \leq Y \leq D \\ 2 & \text{if } D < Y \\ 3 & \text{if } Y < C. \end{cases}$$

If  $X^0 \sim Q \in \mathbf{Q}$ , then  $X = T(X^0) = (V, \Delta) \sim P \in \mathbf{P}$ . This model has been studied by Turnbull (1974), Tsai and Crowley (1985), Chang and Yang (1987), and Chang (1990). Note that this model reduces to the indicator censoring model of example 5.4.1 when  $P(C = D) = 1$  and the distribution of  $C = D$  is continuous so that  $P(\Delta = 1) = 0$ . Since  $Y$  need not be a survival time, we take  $F$  to be a df on  $R$ .

It is easy to calculate that

$$\begin{aligned} H_1(t) &\equiv P(V \leq t, \Delta = 1) = \int_{-\infty}^t M(s) dF(s) \\ &= \int_{-\infty}^t (G_C(s) - G_D(s)) dF(s), \\ H_2(t) &\equiv P(V \leq t, \Delta = 2) = \int_{-\infty}^t \bar{F}(s) dG_D(s), \\ H_3(t) &\equiv P(V \leq t, \Delta = 3) = \int_{-\infty}^t F(s) dG_C(s), \end{aligned}$$

where  $G_C$ ,  $G_D$  are the marginal distributions of  $C$  and  $D$  respectively and

$$(47) \quad M(r) \equiv \int_r^{\infty} \int_{-\infty}^r dG(s, t) = P(C \leq r \leq D) = G_C(r) - G_D(r).$$

Note that only the marginal distributions of  $C$  and  $D$  enter into the distribution  $P$  of  $(V, \Delta)$  given by the  $H_i$ 's.

The model is again differentiable by proposition A.5.5.A, and

$$(48) \quad \begin{aligned} \dot{1}_1 a(V, \Delta) &= E(a(Y) | V, \Delta) \\ &= 1_{[\Delta=1]} a(V) + 1_{[\Delta=2]} \frac{\int_V^{\infty} a dF}{\bar{F}(V)} + 1_{[\Delta=3]} \frac{\int_{-\infty}^V a dF}{F(V)}, \end{aligned}$$

while

$$(49) \quad \begin{aligned} \dot{1}_2 b(V, \Delta) &= E(b(C, D) | V, \Delta) \\ &= 1_{[\Delta=1]} \frac{\int_V^{\infty} \int_{-\infty}^V b(s, t) dG(s, t)}{M(V)} \\ &\quad + 1_{[\Delta=2]} \int_{-\infty}^V b(s, V) dG_{C|D}(s | D=V) \\ &\quad + 1_{[\Delta=3]} \int_V^{\infty} b(V, t) dG_{D|C}(t | C=V). \end{aligned}$$

The independence of  $Y$  and  $(C, D)$  yields  $\langle \dot{I}_1 a, \dot{I}_2 b \rangle_0 = 0$ . This analogue of (22) and an argument similar to the one leading to (25) should show that  $\dot{P}_1$  is orthogonal to  $\dot{P}_2$ . Again this is an instance of the orthogonality condition for adaptation given in proposition 3.4.3 for the case of Euclidean parameters of interest. Hence the bounds for estimation of  $F$  are determined completely by the operator  $\dot{I}_1$ . Note that in example 1 with  $F$  and  $G$  unknown, the map from the distribution  $Q$  of  $X^0 = (Y, C)$  to the distribution  $P$  of  $X = T(X^0) \equiv (T, \Delta)$  takes two functions  $F$  and  $G$  of one real variable into two functions,  $H_0$  and  $H_1$ , of one real variable. In the present example, this map is somewhat more complicated: a function of one variable,  $F$ , and a function of two variables,  $G$ , are first mapped into three functions of one real variable,  $F$ ,  $G_C$ , and  $G_D$ , and then further into three different functions of one real variable,  $H_1, H_2, H_3$ . Since  $P$  is determined by  $H_1, H_2$ , and  $H_3$ , this shows that the joint distribution  $G$  of  $(C, D)$  is not identifiable on the basis of observations of  $T(X^0) = (V, \Delta) \sim P$ . However, the three marginal distributions  $F, G_C, G_D$  are, in fact identifiable under the hypothesis that  $M(r) > 0$  for all  $r$  as shown by Chang and Yang (1987). This leads us to suspect that, as in example 1, the tangent space  $\dot{P}$  of the double censoring model is also "saturated" or "full":  $\dot{P} = L_2^0(P) = \{h \in L_2(P) : E h = 0\}$ . Thus, up to asymptotic equivalence, there is just one asymptotically linear estimator of a given differentiable functional, and it is efficient; see Chang (1990) for such an estimator of  $F$ .

So far this is exactly as in examples 2 and 3. Now we would like to show that a function, such as  $\kappa(P_{(F,G)}) = F$ , is pathwise differentiable, and compute the information bound for estimates thereof. But in this model we lack the martingale structure which we exploited in example 2; so we follow the approach of example 3, calculating  $\dot{I}_1^T$  and then  $\dot{I}_1^T \dot{I}_1$ .

Calculation of  $\dot{I}_1^T$  is straightforward:  $\dot{I}_1^T : L_2^0(P) \rightarrow L_2^0(F)$  is given by

$$\begin{aligned}
 (50) \quad \dot{I}_1^T b(y) &= E(b(V, \Delta) \mid Y = y) \\
 &= M(y) b(y, 1) + \int_{-\infty}^y b(v, 2) dG_D(v) \\
 &\quad + \int_y^{\infty} b(v, 3) dG_C(v).
 \end{aligned}$$

Thus it follows from (47) and (48) that the information operator for  $f$  is given by

$$\begin{aligned}
 (51) \quad \dot{I}_1^T \dot{I}_1 a(y) &= M(y) a(y) + \int_{-\infty}^y \frac{\int_v^{\infty} a dF}{\bar{F}(v)} dG_D(v) \\
 &\quad + \int_y^{\infty} \frac{\int_{-\infty}^v a dF}{F(v)} dG_C(v) \\
 &= M(y) a(y) + \int_{-\infty}^{\infty} K(y, z) a(z) dF(z),
 \end{aligned}$$

where the kernel  $K$  is given by

$$\begin{aligned} K(y, z) &\equiv \int_{y \vee z}^{\infty} \frac{1}{F} dG_C + \int_{-\infty}^{y \wedge z} \frac{1}{\bar{F}} dG_D \\ (52) \quad &\equiv K_1(y, z) + K_2(y, z). \end{aligned}$$

The information operator (51) is very closely related to the information operators which appeared in (35) of example 3 and (5.4.21) of example 5.4.1. In fact, it reduces to the information operator  $\dot{I}_1^T \dot{I}_1$  calculated in example 5.4.1 when  $P(C = D) = 1$  so that  $M = 0$  and  $G_C = G_D$ . To discuss solution of  $\dot{I}_1^T \dot{I}_1 a = b$  for a fixed function  $b \in L_2^0(F)$ , we take the same approach as in example 3.

We consider estimation of  $F$  on  $[\tau_1, \tau_2]$  where  $\tau_1, \tau_2$  satisfy

$$(53) \quad \inf_{\tau_1 \leq r \leq \tau_2} M(r) \equiv \varepsilon > 0.$$

Since we are considering estimation of  $F$  only on  $[\tau_1, \tau_2]$ , we can restrict attention to functions  $a \in L_2^0(F)$  which vanish off  $[\tau_1, \tau_2]$ . Thus we define the Hilbert space

$$(54) \quad \mathbf{H} \equiv \{a \in L_2^0(F) : a = a 1_{[\tau_1, \tau_2]}\} \subset L_2^0(F).$$

We will show that (53) implies that  $\dot{I}_1^T \dot{I}_1 : \mathbf{H} \rightarrow \mathbf{H}$  is one-to-one and onto, and has a bounded inverse on  $\mathbf{H}$  with  $\|(\dot{I}_1^T \dot{I}_1)^{-1}\| \leq \varepsilon^{-1}$ .

Since  $\dot{I}_1^T \dot{I}_1$  is self-adjoint, by corollary A.1.2 we need only verify that

$$\langle a, \dot{I}_1^T \dot{I}_1 a \rangle \geq \varepsilon \|a\|^2 \quad \text{for all } a \in \mathbf{H}.$$

But from (51), for  $a \in \mathbf{H}$  and if  $\int_{\tau_1}^{\infty} (1/F) dG_C < \infty$ ,  $\int_{-\infty}^{\tau_2} (1/\bar{F}) dG_D < \infty$ ,

$$\begin{aligned} \langle a, \dot{I}_1^T \dot{I}_1 a \rangle &= \int a^2 M dF + \int_{\tau_1}^{\tau_2} \frac{1}{F(u)} \left( \int_{-\infty}^u a dF \right)^2 dG_C(u) \\ &\quad + \int_{\tau_1}^{\tau_2} \frac{1}{\bar{F}(u)} \left( \int_u^{\infty} a dF \right)^2 dG_D(u) \end{aligned}$$

$$\geq \varepsilon \|a\|^2$$

by (53) and (54). Thus  $\dot{I}_1^T \dot{I}_1$  is boundedly invertible on  $\mathbf{H}$ , and, by (5.4.14),  $\kappa(P) = 1 - F$  (restricted to  $[\tau_1, \tau_2]$ ) is again pathwise differentiable.

To use the Fredholm recipe for finding the solution of  $\dot{I}_1^T \dot{I}_1 a = b$  for a given bounded function  $b$  vanishing off  $[\tau_1, \tau_2]$ , we rewrite this equation by (51) as

$$\tilde{a}(y) + \int_{-\infty}^{\infty} \tilde{K}(y, z) \tilde{a}(z) dF(z) = \tilde{b}(y),$$

where  $\tilde{a} = M^{1/2} a$ ,  $\tilde{b} = M^{-1/2} b$ , and  $\tilde{K}(y, z) = M^{-1/2}(y) K(y, z) M^{-1/2}(z)$ . Let  $D(s, t; \lambda)$  be the resolvent kernel for this equation:

$$\begin{aligned} (55) \quad D(s, t; \lambda) &= - \sum_{m=0}^{\infty} \frac{(-\lambda)^m}{m!} \int \dots \int \tilde{K} \begin{pmatrix} s, u_1, \dots, u_m \\ t, u_1, \dots, u_m \end{pmatrix} dF(u_1) \dots dF(u_m), \end{aligned}$$

where

$$(56) \quad \tilde{K} \begin{pmatrix} s, u_1, \dots, u_m \\ t, u_1, \dots, u_m \end{pmatrix} \equiv \begin{vmatrix} \tilde{K}(s, t) & \tilde{K}(s, u_1) & \dots & \tilde{K}(s, u_m) \\ \tilde{K}(u_1, t) & \tilde{K}(u_1, u_1) & \dots & \tilde{K}(u_1, u_m) \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \tilde{K}(u_m, t) & \tilde{K}(u_m, u_1) & \dots & \tilde{K}(u_m, u_m) \end{vmatrix};$$

see, e.g., Tricomi (1957, (2.5.12) and (2.5.4), pages 70, 68). Then

$$(57) \quad \tilde{a}(y) = \tilde{b}(y) + \frac{1}{\tilde{D}(-1)} \int \tilde{D}(y, u; -1) \tilde{b}(u) dF(u),$$

where

$$\tilde{D}(\lambda) = 1 + \sum_{m=1}^{\infty} \frac{(-\lambda)^m}{m!} \int \dots \int \tilde{K} \begin{pmatrix} u_1, \dots, u_m \\ u_1, \dots, u_m \end{pmatrix} dF(u_1) \dots dF(u_m)$$

as given in Tricomi (1957, (2.5.5), page 68). For example, for estimation of  $\kappa(P_{(F,G)}) = F(t) \equiv \psi(f)$ , we take  $b = \dot{\psi}_t$  with  $\dot{\psi}_t(y) = 1_{(-\infty, t]}(y) - F(t)$ . Then  $\tilde{b}(y) = (1_{(-\infty, t]}(y) - F(t))M^{-1/2}(y)$ , and (57) yields  $\tilde{a}$  and hence  $a \equiv M^{-1/2} \tilde{a}$ . Finally, the efficient influence function for  $\kappa$  is

$$(58) \quad \tilde{I}_{\kappa}(\pi_t)(V, \Delta) = \dot{I}_1 a(V, \Delta)$$

and

$$I_{\kappa}^{-1}(s, t) = E \tilde{I}_{\kappa}(\pi_s) \tilde{I}_{\kappa}(\pi_t).$$

□

**Example 6. Bivariate censoring.**

Now suppose that  $Y = (Y_1, Y_2) \sim F$  on  $R^{+2}$  with density  $f$ , and  $C = (C_1, C_2) \sim G$  on  $R^{+2}$  with density  $g$  is independent of  $Y$ . In this model  $X^0 = (Y, C)$ , and we observe  $X = T(X^0) \equiv (V, \Delta)$  defined by

$$(59) \quad V \equiv (V_1, V_2) \equiv (Y_1 \wedge C_1, Y_2 \wedge C_2)$$

and

$$(60) \quad \Delta \equiv (\Delta_1, \Delta_2) = (1_{[Y_1 \leq C_1]}, 1_{[Y_2 \leq C_2]}).$$

The joint density  $p$  of  $(V, \Delta)$  is

$$(61) \quad \begin{aligned} p_{11}(v) &= f(v) \bar{G}(v), \\ p_{00}(v) &= g(v) \bar{F}(v), \\ p_{10}(v) &= \int_{v_2}^{\infty} f(v_1, s) ds \int_{v_1}^{\infty} g(t, v_2) dt \equiv \bar{F}_2(v) \bar{G}_1(v), \\ p_{01}(v) &= \int_{v_1}^{\infty} f(s, v_2) ds \int_{v_2}^{\infty} g(v_1, t) dt \equiv \bar{F}_1(v) \bar{G}_2(v), \end{aligned}$$

where

$$\bar{G}(v) \equiv P(C_1 \geq v_1, C_2 \geq v_2) \equiv P(C \geq v)$$

and

$$\bar{F}(v) \equiv P(Y_1 \geq v_1, Y_2 \geq v_2) \equiv P(Y \geq v).$$

We let  $H_{ij}(v) \equiv \int_0^v p_{ij}(u) du$  for  $v \in R^{+2}$  and  $i, j \in \{0, 1\}$ .

Again the model is differentiable (by proposition A.5.5.A), and the score operator for  $f$  is, for  $a \in L_2^0(F)$ ,

$$(62) \quad \begin{aligned} \dot{I}_1 a(V, \Delta) &= E(a(Y) | V, \Delta) \\ &= 1_{[\Delta = (1,1)]} a(V) + 1_{[\Delta = (0,0)]} \frac{1}{\bar{F}(V)} \iint_{[v, \infty)} a dF \\ &\quad + 1_{[\Delta = (0,1)]} \frac{1}{\bar{F}_1(V)} \int_{V_1}^{\infty} a(t, V_2) dF_1(t, V_2) \\ &\quad + 1_{[\Delta = (1,0)]} \frac{1}{\bar{F}_2(V)} \int_{V_2}^{\infty} a(V_1, t) dF_2(V_1, t), \end{aligned}$$

and the score operator for  $g$ ,  $\dot{I}_2 b$ ,  $b \in L_2^0(G)$ , is given by a completely similar expression. Again it should be possible to show along the lines of (22)–(25) that  $\dot{P}_1$  is orthogonal to  $\dot{P}_2$ , and hence bounds for estimation of  $F$  are completely determined by the score operator  $\dot{I}_1$ . But now, unlike the situation in examples 1 and 5, the tangent space  $\dot{P}$  is a *proper* subset of  $\dot{M} = L_2^0(P)$ . We will not show this directly, but give an indirect proof: because there exist many different estimators of  $F$  which are asymptotically linear but *not* asymptotically equivalent (e.g., the pathwise estimators of Campbell (1981), (1982) with two different paths, or the estimator of Dabrowska (1988), (1989) or the estimator based on solving

$$(63) \quad \bar{F}(x) = 1 - F_1(x_1) - F_2(x_2) + \int_0^{x_1} \int_0^{x_2} \bar{F}(y-) d\Lambda(y)$$

for  $\bar{F}$  given marginals  $F_1, F_2$  and a joint hazard function  $\Lambda$ ), the tangent space cannot be  $\dot{M}$ . (If it were, all these estimators would be asymptotically equivalent!) Another way to understand this is via the map which sends the distribution  $Q$  of  $X^0$  into the distribution  $P$  of  $X = T(X^0)$ , as discussed in example 5. Here we start with two functions of two variables,  $F$  and  $G$ , and end up for the distribution of  $X$ , with *four* functions of two variables,  $H_{11}, H_{10}, H_{01}, H_{00}$ . There are two too many functions governing the distribution of  $X$  (or, looked at another way, two too few functions governing the distribution of  $X^0$ ), and hence there are many relationships between these functions, which leads to  $\dot{P}$  being substantially smaller than all possible distributions of  $X = (V, \Delta)$ , and hence  $\dot{P}$  is strictly contained in  $\dot{M}$ .



Nonetheless, we can still use the approach of examples 3 and 5 calculating  $\dot{\mathbf{I}}_1^T$  and the information operator  $\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1$ . From (62) it follows easily that

$$(64) \quad \begin{aligned} \dot{\mathbf{I}}_1^T h(y) &= E\{h(V, \Delta) \mid Y = y\} \\ &= \bar{G}(y) h(y; 1, 1) + \int_0^{y_1} \int_0^{y_2} h(u; 0, 0) dG(u) \\ &\quad + \int_0^{y_2} h(y_1, t; 1, 0) d\bar{G}_1(y_1, t) \\ &\quad + \int_0^{y_1} h(s, y_2; 0, 1) d\bar{G}_2(s, y_2) \end{aligned}$$

for  $h \in \dot{\mathbf{P}}$ . From (62) and (64) we calculate the information operator:

$$(65) \quad \begin{aligned} \dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1 a(y) &= \bar{G}(y) a(y) + \int_0^\infty \int_0^\infty K_{00}(y, u) a(u) dF(u) \\ &\quad + \int_0^\infty K_{10}(y, u_2) a(y_1, u_2) dF_2(y_1, u_2) \\ &\quad + \int_0^\infty K_{01}(y, u_1) a(u_1, y_2) dF_1(u_1, y_2) \end{aligned}$$

where

$$(66) \quad \begin{aligned} K_{00}(y, u) &\equiv \int_0^{u_1 \wedge y_1} \int_0^{u_2 \wedge y_2} \frac{1}{\bar{F}(s)} dG(s), \\ K_{10}(y, u_2) &\equiv \int_0^{u_2 \wedge y_2} \frac{1}{\bar{F}_2(y_1, t)} d\bar{G}_1(y_1, t), \\ K_{01}(y, u_1) &\equiv \int_0^{u_1 \wedge y_1} \frac{1}{\bar{F}_1(s, y_2)} d\bar{G}_2(s, y_2). \end{aligned}$$

The first and second terms on the right side of (65) are analogous to the information operator we computed in (35) of example 3 in the one-dimensional case, but now the third and fourth terms are different, and apparently "intermediate" between the first and second in character. However, they do not seem to yield compact operators on  $L_2^0(F)$ .

Nonetheless, the information operator is invertible under natural hypotheses. To see this we may use the same argument as in examples 3 and 5, but here, to reduce calculations, we will rewrite the information operator  $\dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1$  as

$$(67) \quad \dot{\mathbf{I}}_1^T \dot{\mathbf{I}}_1 a(Y_1, Y_2) = E^Y E^T(a(Y_1, Y_2)),$$

where

$E^Y$  is conditional expectation given  $Y_1, Y_2$ , and  
 $E^T$  is conditional expectation given  $T = (V, \Delta)$ .

We consider estimation of  $F$  on  $[0, \tau] = [0, \tau_1] \times [0, \tau_2]$  where  $\tau$  is such that

$$(68) \quad \bar{G}(\tau) = P(C_1 \geq \tau_1, C_2 \geq \tau_2) \equiv \varepsilon > 0.$$

Since we are considering estimation of  $F$  only on  $[0, \tau] \subset R^2$ , we can restrict

attention to functions  $a \in L_2^0(F)$  which vanish off  $[0, \tau]$ . Thus we define the Hilbert space

$$(69) \quad \mathbf{H} = \{a \in L_2^0(F) : a = a 1_{[0, \tau]} \} \subset L_2^0(F).$$

We claim that (68) implies that  $\dot{1}_1^T \dot{1}_1 : \mathbf{H} \rightarrow \mathbf{H}$  is one-to-one and onto and has an inverse  $(\dot{1}_1^T \dot{1}_1)^{-1}$  on  $\mathbf{H}$  with  $\|(\dot{1}_1^T \dot{1}_1)^{-1}\| \leq \varepsilon^{-1}$ .

Since  $\dot{1}_1^T \dot{1}_1$  is self-adjoint, by corollary A.1.2 we need only verify that

$$(70) \quad \langle a, \dot{1}_1^T \dot{1}_1 a \rangle \geq \varepsilon \|a\|^2, \quad \text{for all } a \in \mathbf{H}.$$

But from (67), for  $a \in \mathbf{H}$ ,

$$\begin{aligned} \langle a, \dot{1}_1^T \dot{1}_1 a \rangle &= E\{a(Y_1, Y_2) E^Y E^T a(Y_1, Y_2)\} \\ &= E\{a(Y_1, Y_2) E^T a(Y_1, Y_2)\} \\ &= E\{\Delta_1 \Delta_2 a^2(Y_1, Y_2) \\ &\quad + (1 - \Delta_1 \Delta_2) E^T [a(Y_1, Y_2) E^T a(Y_1, Y_2)]\} \\ (71) \quad &= E\{\bar{G}(Y_1, Y_2) a^2(Y_1, Y_2) + (1 - \Delta_1 \Delta_2) [E^T a(Y_1, Y_2)]^2\} \\ &\geq \varepsilon \|a\|^2 \end{aligned}$$

by (68) and (69). This remark (with a different proof) and construction of estimates of  $F$  which are efficient are in a forthcoming Utrecht thesis of Mark van der Laan. It does not appear that the inverse of  $\dot{1}_1^T \dot{1}_1$  can be given in a reasonably explicit fashion in general. □

### 6.7 TRANSFORMATION MODELS

Here we reconsider the models of section 4.7, but with emphasis on estimation of the unknown transformation  $\tau \in \mathbf{T}$ . In each example we first study the model assuming that the parametric part is known, and then later consider the case of unknown parametric part  $\nu$  or  $\theta$ .

**Example 1. Linear regression-transformation model.**

This is just example 4.7.1: thus we suppose that  $Z \sim H$  on  $R^k$  and  $\varepsilon \sim G_0$  are independent, and that  $G_0$  and  $H$  are known. Then, for some  $\nu \in R^k$  and  $\tau \in \mathbf{T}$ ,

$$\tau(Y) = \nu^T Z + \varepsilon,$$

and we observe  $X = (Z, Y)$ . As discussed in section 4.7, some special cases of particular interest are:

$$(1) \quad e^\varepsilon \sim \begin{cases} \text{exponential}(1) & : & \text{Cox model,} \\ \text{Pareto}(\eta) & : & \text{Clayton-Cuzick model,} \\ e^{N(0,1)} & : & \text{generalized Box-Cox model.} \end{cases}$$

□

While the Cox model is a special case of this general class of models, it has a very nice martingale structure which distinguishes it from the other members of the family, and hence merits separate treatment.

**Example 1.A. The Cox proportional hazards model.**

Suppose that  $e^\varepsilon \sim \text{exponential}(1)$  in example 1, and that  $v$  is known. Then, as in example 4.7.1, we can identify  $\tau = \log \Lambda = \log(-\log \bar{G})$  where  $\Lambda$  is a cumulative hazard function with corresponding df  $G$  on  $[0, \infty)$ , and the model is just the proportional hazards, or Cox, model studied in example 3.4.2. As shown there, the score operator for  $g$ , the density associated with  $\Lambda$ , is (cf. (3.4.41), (3.4.49), and (3.4.45))

$$(2) \quad \dot{\mathbf{i}}_2 a(z, t) = (L_r R a)(z, t),$$

or

$$(3) \quad \dot{\mathbf{i}}_2 a(Z, T) = \int_0^\infty R a \, dM_r$$

for  $a \in L_2^0(G)$  satisfying

$$(4) \quad \int_0^\infty (R a)^2 S_0 \, d\Lambda < \infty,$$

where

$$M_r(t) \equiv 1_{[T \leq t]} - r \int_0^t 1_{[T \geq s]} \, d\Lambda(s).$$

Thus, for  $a, b \in L_2^0(G)$  satisfying (4), by a martingale calculation as in (3.4.51)

$$\begin{aligned} \langle \dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2 a, b \rangle_{L_2(G)} &= \langle \dot{\mathbf{i}}_2 a, \dot{\mathbf{i}}_2 b \rangle_0 \\ &= E \left( \int_0^\infty R a \, dM_r \right) \left( \int_0^\infty R b \, dM_r \right) \\ &= E \int_0^\infty R a R b r(vZ) 1_{[T \geq \cdot]} \, d\Lambda \\ &= \int_0^\infty R a R b S_0 \, d\Lambda \\ &= \int_0^\infty \left( \frac{S_0}{\bar{G}} R a \right) R b \, dG \\ &= \langle L \left( \frac{S_0}{\bar{G}} R a \right), b \rangle_{L_2(G)} \quad \text{since } R^T = L, \end{aligned}$$

and hence

$$(5) \quad \dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2 a = L \left( \frac{S_0}{\bar{G}} R a \right).$$

It follows easily that the (formal) inverse is

$$(6) \quad (\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1} b = L \left( \frac{\bar{G}}{S_0} R b \right)$$

(using  $R^{-1} = L, L^{-1} = R$ ). Then by (3), (6), and martingale calculus,

$$\langle \dot{\mathbf{i}}_2 (\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1} a, \dot{\mathbf{i}}_2 (\dot{\mathbf{i}}_2^T \dot{\mathbf{i}}_2)^{-1} b \rangle_0 = E \left( \int_0^\infty \frac{\bar{G}}{S_0} R a \, dM_r \right) \left( \int_0^\infty \frac{\bar{G}}{S_0} R b \, dM_r \right)$$

$$\begin{aligned}
 (7) \quad &= E \int_0^\infty \frac{\bar{G}^2}{S_0^2} Ra Rb 1_{[T \geq \cdot]} r(vZ) d\Lambda \\
 &= \int_0^\infty \frac{\bar{G}^2}{S_0} Ra Rb d\Lambda.
 \end{aligned}$$

By taking  $a = 1_{[0,s]} - G(s)$ ,  $b = 1_{[0,t]} - G(t)$  in (7) (so that  $Ra = (\bar{G}(s)/\bar{G}) 1_{[0,s]}$ ,  $Rb = (\bar{G}(t)/\bar{G}) 1_{[0,t]}$ ), we find that for  $\kappa(P_G) \equiv G \equiv \psi(g)$  (with  $\dot{\psi}^T(\pi_s) = a$  and using (5.4.16) and (5.2.23))

$$(8) \quad I_{\kappa}^{-1}(s, t | P_2) = \bar{G}(s)\bar{G}(t) \int_0^{s \wedge t} \frac{1}{S_0} d\Lambda.$$

This agrees with the calculations of Begun, et. al. (1983, formula (6.8), page 450), for the case when  $v$  is known, and with the known results for the Breslow (1974) estimator of  $\Lambda$  (with known  $v$ ); see, e.g., Tsiatis (1981) and Andersen and Gill (1982).  $\square$

### Example 2. Joint distribution-transformation model.

Now we return to example 4.7.3 with a general parametric core model  $P_0$ . As in propositions 4.7.1 through 4.7.3, we assume that the marginal distribution of  $T$  is *Uniform*(0, 1) under  $P_\theta \in P_0$ , that  $\dot{I}_\theta$  and  $\dot{I}_u$  in (4.7.15) and (4.7.16) exist, and that  $\alpha(u) \leq M(u(1-u))^{-2}$  as in proposition 4.7.1.

Consider the function  $\kappa : P \rightarrow l^\infty(R)$  defined by

$$\kappa(P_{(\theta, G)}) = G \equiv \psi(g).$$

Then, as in example 5.3.1,

$$\begin{aligned}
 \dot{\psi}^T(\pi_t)(u) &= 1_{[0,t]}(u) - G(t) \\
 &= 1_{[0,t]}(u) - t \quad \text{at } G = I \\
 &= \frac{\partial}{\partial u}(u \wedge t - ut).
 \end{aligned}$$

By proposition 4.7.2,  $(\dot{I}_g^T \dot{I}_g)^{-1}$  exists and is bounded. Hence, by (5.4.16), to find the efficient influence function for estimation of  $G$  we want to compute

$$(9) \quad \tilde{I}_\kappa(\pi_t) = \dot{I}_g(\dot{I}_g^T \dot{I}_g)^{-1} \dot{\psi}^T(\pi_t) = \dot{I}_g(\dot{I}_g^T \dot{I}_g)^{-1} (1_{[0,t]} - t);$$

so we want to solve

$$(10) \quad \dot{I}_g^T \dot{I}_g a = 1_{[0,t]} - t$$

for  $a$ . In view of (4.7.31), by defining  $A(u) \equiv \int_0^u a(v) dv$ , this becomes

$$A'(u) + \int_0^1 (1_{[u \leq s]} - s) A(s) \alpha(s) ds = 1_{[0,t]}(u) - t.$$

But this is just the equation satisfied by the Green's function  $\Delta(u, t)$ : by (4.7.39),

$$(11) \quad \frac{\partial}{\partial u} \Delta(u, t) + \int_0^1 (1_{[u \leq s]} - s) \Delta(s, t) \alpha(s) ds = 1_{[u \leq t]} - t,$$

$0 \leq t \leq 1$ . Thus

$$(12) \quad \tilde{I}_{\kappa}(\pi_t)(z, y) = \dot{\mathbf{i}}_g \left( \frac{\partial}{\partial u} \Delta(\cdot, t) \right)(z, y),$$

and, by (5.2.23),

$$(13) \quad \begin{aligned} I_{\kappa}^{-1}(s, t | P_{(\theta, t)}, \mathbf{P}_2) &= E(\dot{\Psi}^T(\pi_s) (\dot{\mathbf{i}}_g^T \dot{\mathbf{i}}_g)^{-1} \dot{\Psi}^T(\pi_t)) \\ &= E(\dot{\Psi}^T(\pi_s) \frac{\partial}{\partial u} \Delta(U, t)) \\ &= \int_0^1 (1_{[0, s]}(u) - s) \frac{\partial}{\partial u} \Delta(u, t) du \\ &= \Delta(s, t), \end{aligned}$$

since  $\Delta(0, t) = \Delta(1, t) = 0$  by lemma 4.7.1. In view of proposition 4.7.1,

$$(14) \quad I_{\kappa}^{-1}(s, t | P_{(\theta, G)}, \mathbf{P}_2) = I_{\kappa}^{-1}(G(s), G(t) | P_{(\theta, t)}, \mathbf{P}_2) = \Delta(G(s), G(t)). \square$$

**Example 3. Copula model with one unknown marginal df.**

Here we calculate bounds for  $\kappa(P_{(\theta, G)})$  in example 4.7.4.A. These calculations are exactly the same as in example 2 above, but now they can be made more explicit in various cases because sometimes we can calculate the Green's function  $\Delta$  explicitly.

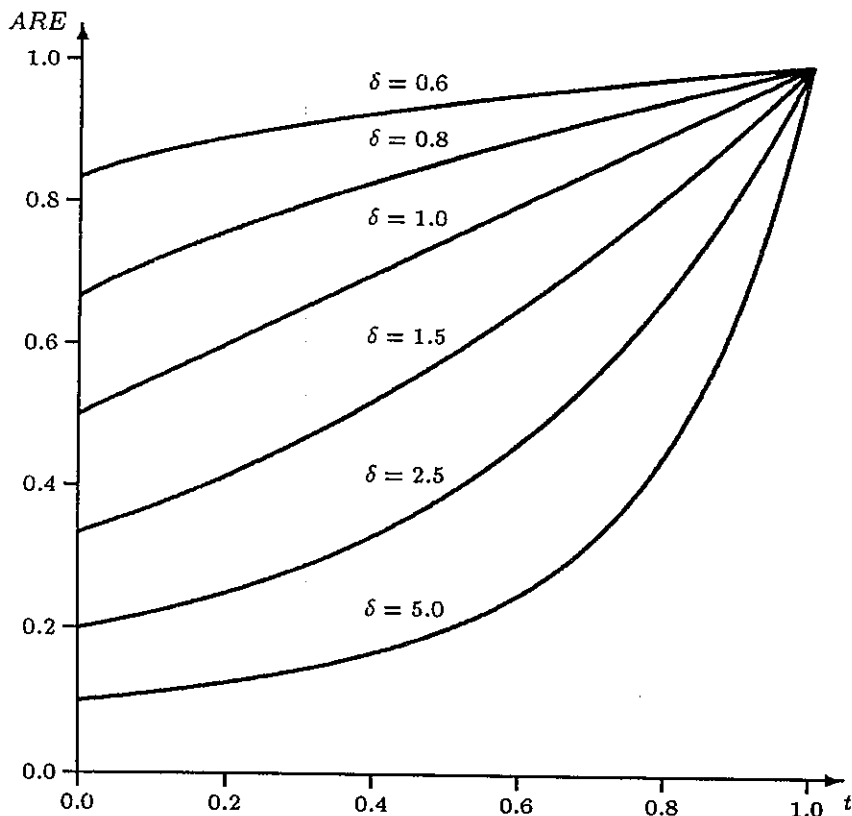


FIGURE 1. Clayton-Oakes copula model: ARE of empirical distribution function to efficient estimator.

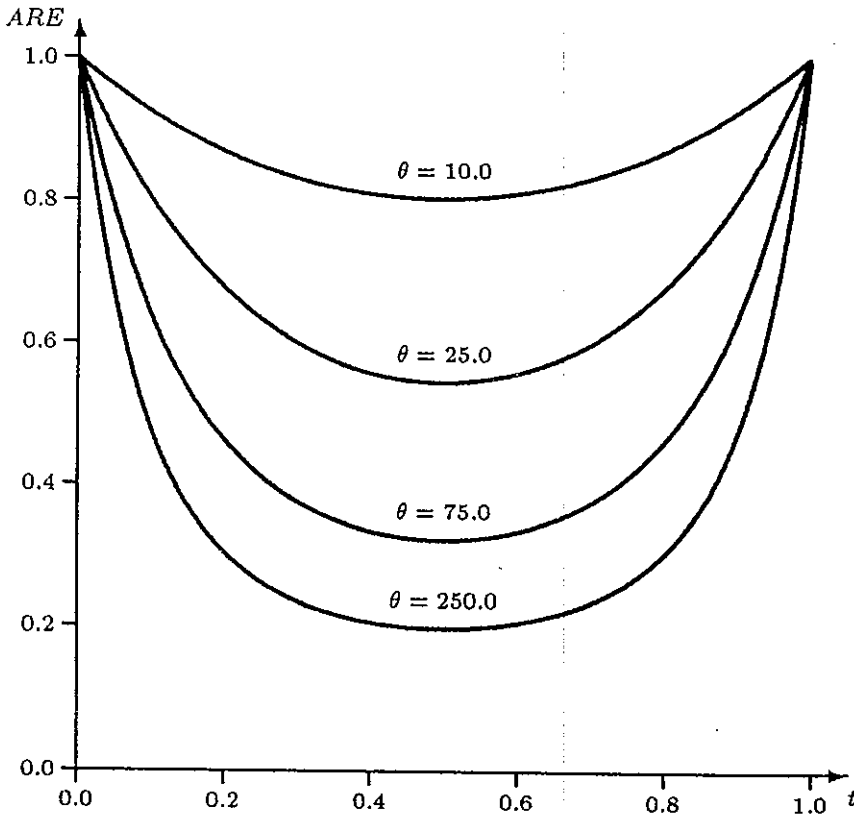


FIGURE 2. Frank copula model: ARE of empirical distribution function to efficient estimator.

Thus we take  $\kappa(P_{(\theta,G)}) = G$  as in example 2, and find that (13) continues to hold with  $\alpha$  being defined by (4.7.52).

Now consider the further special case of the Clayton-Oakes model, example 4.7.4.I.1. Then,  $\Delta(s,t)$  is given by (4.7.84) and hence (13) yields

$$(15) \quad \Gamma_{\kappa}^{-1}(s,t | P_{(\theta,I)}, P_2) = \frac{(s \wedge t)^{1/2+\delta} [(s \vee t)^{1/2-\delta} - (s \wedge t)^{1/2+\delta}]}{(2\delta)}$$

with

$$\delta = \left( \frac{1}{4} + \frac{\theta^2(\theta + 1)}{3\theta + 1} \right)^{1/2}.$$

Consequently, when  $s = t$

$$\Gamma_{\kappa}^{-1}(t,t | P_{(\theta,I)}, P_2) = \frac{1}{2\delta} t(1 - t^{2\delta}).$$

The marginal empirical distribution function estimates  $G$  and has variance  $t(1 - t)$  at  $G = I$ , so the asymptotic relative efficiency of the empirical df relative to an efficient estimator is

$$(16) \quad ARE(\text{empirical}, \text{efficient})(P_2)(t) = \frac{1}{2\delta} \frac{1 - t^{2\delta}}{1 - t},$$

which is plotted as a function of  $t$  for a range of values of  $\delta$  in figure 1.

Similarly, for the special case of the Frank model, example 4.7.4.I.3,

$$I_{\kappa}^{-1}(s, t \mid P_{(\theta, t)}, \mathbf{P}_2) = \Delta(s, t)$$

with  $\Delta$  given by (4.7.88), and hence, when  $s = t$ ,

$$I_{\kappa}^{-1}(t, t \mid P_{(\theta, t)}, \mathbf{P}_2) = \frac{\sinh(dt) \sinh(d(1-t))}{d \sinh(d)},$$

where  $d^2 \equiv 3^{-1} (\log \theta)^2$ . Thus, as before, the asymptotic relative efficiency of the empirical df at  $G = I$  is given by

$$(17) \quad ARE(\text{empirical, efficient})(\mathbf{P}_2)(t) = \frac{\sinh(dt) \sinh(d(1-t))}{dt(1-t) \sinh(d)}$$

which is plotted as a function of  $t$  for several values of  $\theta$  in figure 2. We will construct efficient estimators of  $G$  for example 2 in example 7.6.7.  $\square$

# 7 | Construction of Estimates

## 7.1 INTRODUCTION

In the first seven sections of this chapter we shall discuss a number of important methods for the construction of  $\sqrt{n}$ -consistent estimates of both Euclidean and abstract parameters. Finally, in section 7.8, we shall study when these estimates, or procedures which use these estimates as building blocks, are efficient in the full model considered, or at least on a submodel.

Conceptually, there are essentially only two basic approaches to the construction of estimates of a Euclidean parameter  $v(P)$ ,  $v : \mathbf{P} \rightarrow R^m$ .

- A. Find a "smooth" extension of  $v$ , say  $\bar{v}_n$  where  $\bar{v}_n : M_0 \rightarrow R^m$ ,  $M_0 \supset \mathbf{P}$ , and  $M_0 \supset \{ \text{probability distributions with finite support} \}$ . Then estimate  $v$  by  $\bar{v}_n(IP_n)$ .
- B. Find a "good" estimate  $\hat{P}$  of  $P$  such that  $\hat{P} \in \mathbf{P}$ . Estimate  $v$  by  $v(\hat{P})$ .

If  $v$  is an abstract parameter the same approaches apply. When  $v(P) = P$  for example, we can in most cases extend  $v$  to  $\bar{v}_n$  on  $M_0$  in a natural way and apply method A in estimating  $P$  by  $\bar{v}_n(IP_n) = \hat{P}$ . Then an arbitrary parameter  $\eta(P)$  can be estimated via approach B by  $\eta(\hat{P}) = \eta(\bar{v}_n(IP_n))$ . However, this estimator can be viewed as a type A estimator of  $\eta(\bar{v}_n(P)) = \eta(P)$ . So, really, we only have the extension approach but the distinction between trying to estimate "everything" ( $P$ ) first, and estimating only the piece of interest ( $v(P)$ ), is useful conceptually.

The procedures for estimating  $\theta$  in a regular parametric model  $\{P_\theta : \theta \in \Theta\}$ ,  $\Theta \subset R^k$ , that were discussed in section 2.5 illustrate the two approaches. The general technique to construct  $\sqrt{n}$ -consistent estimators that was described in this section was to estimate  $\theta$  by  $\tilde{\theta}_n = \theta(\Pi(IP_n))$  where  $\Pi : \mathbf{M} \rightarrow \mathbf{P}$  is defined by  $\rho(\Pi(Q), Q) = \min \{ \rho(P, Q) : P \in \mathbf{P} \}$  and  $\rho$  is a metric compatible with  $IP_n$ . Again, we can think of  $\Pi(IP_n)$  as the  $\hat{P}$  required for approach B. On the other hand, the MLE maximizes  $\int l(x, \theta) dIP_n(x)$  and can be considered as a type A estimator of the maximizer of  $\int l(x, \theta) dP(x)$ .

Another simple illustration is provided by the constrained model, example 3.3.2. Here the parameter  $v$  is specified on  $M_0$ , for instance  $\int x dP(x)$ , so that



the natural but inefficient estimate is just  $v(\mathbb{P}_n)$ . With  $\mathbf{P} = \{P \in \mathbf{M}_0 : \gamma(P) = 0\}$ , it is natural to consider the class of extensions,  $\bar{v}_a(P) = v(P) + a(P)\gamma(P)$ . As we noted in (3.3.32) for a special case, if  $v(\mathbb{P}_n)$  has influence function  $\dot{v}(\cdot, P)$  and  $\gamma(\mathbb{P}_n)$  has influence function  $\dot{\gamma}(\cdot, P)$ , then the putatively efficient choice of  $a$  is

$$a_0(P) = - \int \dot{v}(x, P) \dot{\gamma}(x, P) dP(x) / \int \dot{\gamma}^2(x, P) dP(x).$$

It is easy to see that the estimator  $\bar{v}_{a_0}(\mathbb{P}_n)$  is efficient if

$$v(\mathbb{P}_n) = v(P) + n^{-1} \sum_{i=1}^n \dot{v}(X_i, P) + o_p(n^{-1/2}),$$

$$\gamma(\mathbb{P}_n) = \gamma(P) + n^{-1} \sum_{i=1}^n \dot{\gamma}(X_i, P) + o_p(n^{-1/2}),$$

$$\frac{1}{n} \sum_{i=1}^n \dot{\gamma}(X_i, \mathbb{P}_n) \dot{v}(X_i, \mathbb{P}_n) = \int \dot{v}(x, P) \dot{\gamma}(x, P) dP(x) + o_p(1),$$

and

$$\frac{1}{n} \sum_{i=1}^n \dot{\gamma}^2(X_i, \mathbb{P}_n) = \int \dot{\gamma}^2(x, P) dP(x) + o_p(1).$$

Extensions  $\bar{v}$  can be characterized as solutions of optimization problems and/or systems of equations. In the optimization problem formulation we look for a function  $D : \bar{\mathbf{A}} \times \mathbf{M}_0 \rightarrow R$  such that:

- (i)  $D(\bar{v}(\mathbb{P}_n), \mathbb{P}_n) = \min \{D(v, \mathbb{P}_n) : v \in \bar{\mathbf{A}}\}$
- (ii) If  $P \in \mathbf{P}$ ,  $D(v, P)$  has its unique minimum at  $v = v(P)$ .

We call these *generalized minimum contrast estimates* or *GMC-estimates*.

If  $\mathbf{P} = \{P_\theta : \theta \in \Theta\}$ ,  $\Theta \subset R^k$  is parametric and  $\theta$  is identifiable, choices of  $D$  with  $\bar{\mathbf{A}} \equiv R^k$  that we have considered are:

1.  $D(\theta, P) \equiv \rho(P_\theta, P)$  (minimum distance).
2.  $D(\theta, P) \equiv - \int \log p_\theta(x) dP(x)$  (maximum likelihood).

More generally, if  $D(\theta, P) \equiv \int \rho(x, \theta) dP(x)$ , we are led to what Pfanzagl (1979) called *minimum contrast estimates*.

If  $v$  is Euclidean,  $D$  is smooth and  $v(\mathbb{P}_n)$  is a generalized minimum contrast estimate, then  $\bar{v}(\mathbb{P}_n)$  solves

$$(1) \quad W(v, \mathbb{P}_n) = 0,$$

where  $W : R^m \times \mathbf{M}_0 \rightarrow R^m$ ,  $W(\cdot, P) \equiv \nabla D(\cdot, P)$ ,  $\nabla = (\partial/\partial v_1, \dots, \partial/\partial v_m)^T$ , and  $v = (v_1, \dots, v_m)^T$ . Further

$$(2) \quad W(v(P), P) = 0 \quad \text{for all } P \in \mathbf{P}.$$

If  $v(P)$  is the unique solution of (2), we call any solution of the system of equations (1) a *generalized M-estimate* or *GM-estimate*. In particular, if

$D(v, P) = \int \rho(x, v) dP(x)$ ,  $\psi(x, \cdot) \equiv \nabla \rho(x, \cdot)$  is defined for all  $x$ , we are led to estimates which are solutions of

$$(3) \quad \int \psi(x, v) dP_n(x) = 0,$$

that is, to  $M$ -estimates as defined by Huber (1981). These are discussed in section 7.2.

There are natural functions  $D$  for which we cannot pass to  $W$ , for instance,  $D(\theta, P) = d_K(P_\theta, P)$  where  $d_K$  is the Kolmogorov distance. On the other hand, there may be solutions of (1) which may not correspond to minimizers of  $D$  and, of course, most functions  $W$  are not of the form  $\nabla D$ . Finally, any  $\bar{v}$  is both a generalized minimum contrast and a generalized  $M$ -estimate: take  $D(v, P) = (v - \bar{v}(P))^2$ . The utility of these formulations lies in translating the choice of  $\bar{v}$  into a choice of  $D$  or  $\psi$  which it is easier to think about. Further, the generalized minimum contrast formulation enables us to postulate fairly natural conditions for consistency of  $\bar{v}(P_n)$ ; see lemma 2.5.1. The generalized  $M$ -estimate formulation leads naturally to conditions for linearization of  $W(v, \cdot)$ ,  $\sqrt{n}$ -consistency, and asymptotic linearity and normality of  $\bar{v}(P_n)$ . We pursue this theme in section 7.3: Generalized  $M$ -estimates for Euclidean parameters. Here our results are of the type introduced by Cramér (1946): Consistent solutions of (1) are  $\sqrt{n}$ -consistent, asymptotically linear, normal, etc. This conforms to practice where equations such as (1) have to be solved iteratively in any case. If there are multiple roots the "right" one will be obtained only if one starts close enough to the unknown  $v(P)$ . That is, if one is lucky, or has available a crude consistent starting value. Proofs are to be found in section A.10. We apply these results to a number of examples which have arisen in chapters 3 and 4.

In section 7.4, we face the question of existence, uniqueness and consistency of  $GMC$ - and the corresponding  $GM$ -estimates in a very important special case:  $D(\cdot, P)$  convex. This strong condition also simplifies the linearization theorems of 7.3. These theorems, whose proofs may also be found in section A.10, are applied to a pair of important examples, including the Cox estimate.

The definition of generalized minimum contrast estimates obviously extends to abstract  $v$  provided that  $\bar{A}$  in the definition of  $D$  is chosen appropriately. For generalized  $M$ -estimates it turns out that  $W$  now has to be a map from  $\bar{A} \times M_0$  to another abstract space, typically  $\bar{A}^*$ , the dual space of  $\bar{A}$ . If  $A$  is a function space, the passage from  $D$  to  $W$  and, in particular, the meaning of  $\nabla D$ , have to be looked at carefully. Appropriate definitions are given in section 7.5.

The infinite-dimensional nature of  $\bar{A}$  forces us to further generalizations. For instance, take

$$P = \{ \text{all probabilities on } R \text{ dominated by Lebesgue measure} \},$$

take  $\bar{A} = P$ , and define  $D$  as the negative log-likelihood functional

$$D(P, Q) = - \int \log p(x) dQ(x).$$

Then, as we shall see in example 7.5.2,

$$\inf_{P \in \mathbf{P}} D(P, \mathbb{P}_n) = -\infty$$

is approached by  $P^{(h)} \in \mathbf{P}$  such that  $P^{(h)}$  converge to the empirical distribution  $\mathbb{P}_n$  which is not a member of  $\mathbf{P}$ . Various ways out, including nonparametric maximum likelihood, replacement of  $\mathbb{P}_n$  in  $D(\cdot, \mathbb{P}_n)$  by a smoothed version of itself, penalized maximum likelihood, and the method of sieves, are discussed in sections 7.5. They lead to considering generalized minimum contrast estimates as minimizers of maps  $D_n: \bar{\mathbf{A}} \times \mathbf{M}_0 \rightarrow R$  where  $D_n(v, P) \rightarrow D(v, P)$  as  $n \rightarrow \infty$  and  $D$  satisfies (ii). For example, if  $d_H$  is the Hellinger metric and  $\bar{\mathbf{A}} = \mathbf{P}$  as above,  $\inf_P d_H(P, \mathbb{P}_n) = 0$  is not assumed in  $\mathbf{P}$ . However, if we replace  $\mathbb{P}_n$  by any smoothed version, for instance  $\mathbb{P}_n^\# \equiv \mathbb{P}_n * N(0, 1/n)$ , then  $\mathbb{P}_n^\# \in \mathbf{P}$  minimizes  $d_H(P, \mathbb{P}_n^\#)$ . Evidently if we take  $D_n(P, Q) = d_H(P, Q * N(0, 1/n))$ ,  $\mathbb{P}_n^\#$  is the corresponding generalized minimum contrast estimate. The corresponding extension of generalized *M*-estimation based on functions  $W_n: \bar{\mathbf{A}} \times \mathbf{M}_0 \rightarrow \bar{\mathbf{A}}^*$ , with  $W_n \rightarrow W$  satisfying (2), is discussed in section 7.5 and 7.6. Specifically, in section 7.5 we give consistency and rate results using convexity properties of the type discussed in section 7.4. Asymptotic linearity and normality are discussed in section 7.6. The methodology in both sections is applied to a number of important examples from chapters 5 and 6. In these sections we also apply approach B to estimation of Euclidean parameters in examples from chapter 4 after having first constructed "good"  $\hat{P}$  using approach A.

In section 7.7, we combine the themes of sections 7.5 and 7.6 by showing how to construct good estimating equations  $W_n$  for Euclidean parameters using generalized *M*-estimates of abstract "nuisance" parameters. Finally, in section 7.8, we first discuss heuristically when the maximum likelihood related procedures we have introduced in sections 7.5–7.7 should lead to efficient estimates or at least estimates efficient at certain parametric submodels. We note that these heuristic predictions are verified for the chapter 4, 5, and 6 examples we have considered in sections 7.5–7.7. We then show how to construct efficient estimates of Euclidean parameters given a  $\sqrt{n}$ -consistent initial estimate and a suitably consistent and asymptotically unbiased estimate of the efficient influence function. This approach is applied to a number of examples from sections 4.2 and 4.3.

## 7.2 *M*-ESTIMATES FOR EUCLIDEAN PARAMETERS

In this section the familiar *M*-estimates will be discussed. This will serve as a convenient introduction to the generalized *M*-estimates which will be studied in section 7.3.

Suppose that  $\psi: \mathbf{X} \times R^m \rightarrow R^m$  and let

$$(1) \quad W(v, Q) \equiv \int \psi(x, v) dQ(x)$$

be well defined for all  $Q \in \mathbf{M}_0 \supset \mathbf{P}$ . We assume that for each  $P \in \mathbf{P}$ , our

parameter of interest  $v(P) \in R^m$  satisfies

$$(2) \quad W(v(P), P) = 0.$$

We call any root of

$$(3) \quad W(v, P_n) = n^{-1} \sum_{i=1}^n \psi(X_i, v) = 0$$

an *M-estimate*. An estimate  $\hat{v}_n$  is called an *asymptotic M-estimate* or *AM-estimate* if it satisfies

$$(4) \quad W(\hat{v}_n, P_n) = o_p(n^{-1/2}).$$

We assume:

(M1) For any  $\varepsilon_n \downarrow 0$  we have

$$(5) \quad \sup_{|v - v(P)| \leq \varepsilon_n} \frac{|n^{-1/2} \sum_{i=1}^n \{\psi(X_i, v) - \psi(X_i, v(P)) - E_P[\psi(X, v) - \psi(X, v(P))]\}|}{1 + \sqrt{n} |v - v(P)|} = o_p(1).$$

(M3)  $W(v, P) = \int \psi(x, v) dP(x)$  is differentiable and the derivative at  $v(P)$ ,  $\dot{W}(P)$ , is nonsingular.

We obtain:

**Theorem 1.** Suppose that (M1) and (M3) hold,  $\hat{v}_n$  is consistent, and is an AM-estimate. Then  $\hat{v}_n$  is asymptotically linear with influence function  $-\dot{W}^{-1}(P)\psi(\cdot, v(P))$ .

This theorem follows from Huber (1967). In fact, it is an application of Theorem 7.3.1; for a complete discussion see corollary A.10.1.

Condition (M1) is easily seen to be implied by the following familiar condition:

(M1') For some  $\varepsilon > 0$  the random functions  $\sqrt{n}(W(\cdot, P_n) - W(\cdot, P))$  take values in  $l^\infty(B_\varepsilon(v(P)))$  and converge weakly in the sense of appendix A.8 to an a.s. continuous random function; here  $B_\varepsilon(v(P))$  denotes the closed ball of radius  $\varepsilon$  about  $v(P)$ .

Further, if we assume the following differentiability conditions on  $\psi$ , then (M1) and (M3) follow quite easily.

(D1)  $v \rightarrow \psi(X_1, v) = (\psi_1(X_1, v), \dots, \psi_m(X_1, v))^T$  is differentiable a.s. and  $v \rightarrow (\partial/\partial v_j)\psi_i(X_1, v)$  is continuous a.s.,  $i, j = 1, \dots, m$ .

(D2) For all  $v$  with  $|v - v(P)| \leq \varepsilon(P)$  for some  $\varepsilon(P) > 0$  we have

$$(6) \quad \left| \frac{\partial}{\partial v_j} \psi_i(X_1, v) \right| \leq M_{ij}(X_1, P) \quad \text{a.s., where } E_P M_{ij}(X_1, P) < \infty,$$

(D3) The matrix  $\dot{W}(v)$  given by

$$(7) \quad \dot{W}(v) = [E_P(\frac{\partial}{\partial v_j} \psi_i(X_1, v))]_{m \times m}$$

is nonsingular at  $v = v(P)$ . As before, we write  $\dot{W}(P) = \dot{W}(v(P))$ .

Now (M1) is implied by (D1) and (D2), while (M3) comes from (D3). Consequently we have:

**Corollary 1.** Suppose that  $\hat{v}_n$  is a consistent *AM*-estimate, and that (D1) through (D3) hold. Then  $\hat{v}_n$  is asymptotically linear with influence function  $-\dot{W}^{-1}(P)\psi(\cdot, v(P))$ .

Note that if  $P = \{P_\theta : \theta \in \Theta\}$  is a regular parametric model and  $\psi = \dot{\mathbf{i}} = \nabla_\theta \log p(\cdot, \theta)$ , the resulting influence function is efficient in view of the classical identity

$$(8) \quad \int \dot{\mathbf{i}}^T(x, \theta) dP_\theta(x) = [- \int \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p(x, \theta) dP_\theta(x)]_{k \times k} \\ \equiv - \int \ddot{\mathbf{i}}(x, \theta) dP_\theta(x),$$

which holds under the strong conditions of corollary 1. By proposition 3.3.1 the resulting *M*-estimate is regular.

We can go further. If  $T_n$  is an *M*-estimate, the conditions of theorem 1 hold and  $P = \{P_v : v \in N\}$  is regular parametric, then by (2.4.27)  $T_n$  is regular if and only if for all  $v_0$ ,

$$(9) \quad \int \psi(x, v_0) \dot{\mathbf{i}}^T(x, v_0) dP_{v_0}(x) = - \dot{W}(P_{v_0}).$$

But, if interchange of integration and differentiation is permitted in  $\int \psi(x, v) dP_{v_0}(x)$ , (9) is just the identity

$$(10) \quad \left[ \int \psi_i(x, v_0) \dot{\mathbf{i}}_j(x, v_0) dP_{v_0}(x) \right] = - \left[ \int \frac{\partial \psi_i}{\partial v_j}(x, v_0) dP_{v_0}(x) \right],$$

which generalizes (8).

Corollary 1 essentially goes back at least to Cramér (1946). The need for the more general theorem 1 is illustrated by the following example. This example is covered by the conditions of Huber (1967) implying (M1) which are much weaker than the smoothness conditions (D1) and (D2).

**Example 1. The median.**

If  $W(v, P_n) \equiv n^{-1} \sum_{i=1}^n \text{sgn}(X_i - v)$ , then it is easy to see that  $\hat{v} \equiv \text{Median } X_i$  is an *AM*-estimate. Suppose  $F$  has median  $v_0$ , is differentiable at  $v_0$  and has positive density  $f(v_0)$  at  $v_0$ . Then

$$\hat{v}_n = v_0 + (2f(v_0))^{-1} n^{-1} \sum_{i=1}^n \text{sgn}(X_i - v_0) + o_p(n^{-1/2}).$$

To see this apply theorem 1 with  $\psi(x, v) \equiv \text{sgn}(x - v)$ , which is evidently not

differentiable. Then

$$W(v, P) = 1 - 2F(v), \quad \dot{W}(P) = -2f(v_0),$$

and condition (M1') follows from the weak convergence of the empirical process.  $\square$

### *M-estimates in Semiparametric Models*

If  $\mathbf{P} = \{P_{(\theta, G)} : \theta \in \Theta, G \in \mathbf{G}\}$ ,  $\Theta \subset R^k$ , then in order to construct *M*-estimates for  $\theta$  or  $q(\theta) \equiv v(P_{\theta, G})$ , we need to find a function  $\psi(x, \theta)$  such that

$$(11) \quad \int \psi(x, \theta) dP_{(\theta, G)}(x) = 0 \quad \text{for all } \theta \text{ and } G.$$

There is a large class of semiparametric models  $\mathbf{P} = \{P_{(\theta, G)} : \theta \in \Theta, G \in \mathbf{G}\}$ , where, heuristically, we expect that if we apply the method of maximum likelihood (or an approximation) to estimate  $\theta$  in a submodel  $\mathbf{P}_1(G_0) = \{P_{(\theta, G_0)} : \theta \in \Theta\}$ , we obtain estimates  $\tilde{\theta}_n = (\tilde{v}_n, \tilde{\eta}_n)$  such that  $\tilde{v}_n$  is  $\sqrt{n}$ -consistent for the identifiable component  $v$  of  $\theta = (v, \eta)$  in the full model  $\mathbf{P}$ . This is a nice situation since such estimates, in addition to providing starting points for globally efficient estimates, themselves possess a local efficiency property. They are efficient at all points  $P$  of the submodel  $\mathbf{P}_1(G_0)$ .

#### **Example 2. Symmetric location in $k$ dimensions.**

Suppose  $X_{k \times 1} = \theta + e$  where  $e$  has distribution  $G$  symmetric about 0 in the sense that  $G(A) = G(-A)$  for all measurable sets  $A$ . The set of all such  $G$  is  $\mathbf{G}$ . This is one of several possible generalizations of the one-dimensional symmetric location model. If  $G_0$  has smooth density  $g_0$ , the maximum likelihood estimate for  $\theta$  in  $\mathbf{P}_1(G_0)$  is  $\theta(P_n)$  where  $\theta(P)$  satisfies

$$\int \frac{\dot{g}_0}{g_0}(x - \theta) dP(x) = 0$$

and where  $\dot{g}_0 = \nabla g_0$ . But in this case, if  $G \neq G_0$ , but  $G \in \mathbf{G}$ , it is still true (under regularity conditions) that

$$\int \frac{\dot{g}_0}{g_0}(x - \theta) dP_{(\theta, G)}(x) = \int \frac{\dot{g}_0}{g_0}(x) g(x) dx = 0.$$

The reason is that  $(\dot{g}_0/g_0)(x) = -(\dot{g}_0/g_0)(-x)$  and  $g(x) = g(-x)$  for all  $x$ . If  $g_0(x) = \prod_{j=1}^k \phi(x_j)$  and we restrict to symmetric  $G$  with  $\int |x|^2 dG(x) < \infty$ , we are led to  $\theta(P) = \int x dP(x)$  and  $\bar{X}$  as a  $\sqrt{n}$ -consistent *M*-estimate of  $\theta$  in  $\mathbf{P}$ . More generally, if  $G_0$  has i.i.d. components with density  $g_{01}$ , then we formally obtain an *M*-estimate of  $\theta$  solving

$$(12) \quad \sum_{i=1}^n \frac{g_{01}'}{g_{01}}(X_{ij} - \theta_j) = 0, \quad j = 1, \dots, k,$$

where  $X_i = (X_{i1}, \dots, X_{ik})^T$  and  $\theta = (\theta_1, \dots, \theta_k)^T$ . Specializing further to the logistic density  $g_{01}(x) = e^x(1 + e^x)^{-2}$ , (12) has a unique root  $(\hat{\theta}_1, \dots, \hat{\theta}_k)^T$ , and (M1), (M3) are satisfied.  $\square$

If  $\Theta$  is  $k$ -dimensional,  $k > 1$ , then often  $\theta$  is not even identifiable on  $\mathbf{P}$  although it may be well defined on  $\mathbf{P}_1(G_0)$ . Yet if  $q(\theta)$ ,  $q: \Theta \rightarrow R^m$ ,  $m < k$  is identifiable, we can still proceed as in example 2. We illustrate what happens in examples 3 and 4 below.

**Example 3. Regression models.**

We use the notation of example 4.2.2. Fix  $G_0$ . If the error density  $q_0$  is absolutely continuous with derivative  $q'_0$  satisfying the conditions of example 4.2.2, continued, then by (4.2.16)

$$\dot{\mathbf{i}}(x, \theta, g_0) = \sigma^{-1}(z^T \psi(u), \psi(u), u \psi(u) - 1)^T,$$

where  $u \equiv \sigma^{-1}(y - \Delta - \sum_{j=1}^{k-2} v_j z_j)$  and  $\psi = -q'_0/q_0$ . If  $q_0 = \phi$ , we obtain the least squares estimates as the corresponding uniquely defined *M*-estimates. Although  $\sigma$  and  $\Delta$  are not identifiable on  $\mathbf{P}$  it is still true (if  $q$  has second moments and  $Z$  is genuinely  $(k-2)$ -dimensional) that the least squares estimates,  $\tilde{\theta} \equiv (\tilde{v}^T, \tilde{\Delta}, \tilde{\sigma})^T$  are converging to a value  $\theta = (v^T, \Delta, \sigma)^T$  such that  $P = P_{(\theta, G)}$  for some  $G$ . And, of course, the identifiable piece  $v$  is (if  $q$  has second moments) being  $\sqrt{n}$ -consistently estimated.

If  $-\log q_0$  is smooth (with bounded first and second derivatives),  $Z$  is genuinely  $(k-2)$ -dimensional, and  $E|Z|^2 < \infty$  then  $\dot{\mathbf{i}}(x, \theta, g_0)$  satisfies (M1) and (M3). We argue later that if  $-\log q_0$  is strictly convex (for instance,  $q_0$  logistic), then unique *M*-estimates  $\theta^*$  exist. Again  $v^*$  is  $\sqrt{n}$ -consistent, asymptotically linear, etc.  $\square$

The basic feature of these examples is that functions  $\psi(x, \theta) \equiv \dot{\mathbf{i}}(x, \theta, G_0)$  for which necessarily,

$$\int \dot{\mathbf{i}}(x, \theta, G_0) dP_{(\theta, G_0)}(x) = 0$$

end up also satisfying (11). Here is a heuristic explanation of this feature.

**Definition 1.** Let  $G$  be a convex subset of a linear space. The model  $\mathbf{P} = \{P_{(\theta, G)}: \theta \in \Theta, G \in \mathbf{G}\}$  has a *convex parametrization* if for all  $\theta \in \Theta$ ,  $0 \leq \alpha \leq 1$ ,  $G_0, G \in \mathbf{G}$ ,

$$(13) \quad P_{(\theta, \alpha G_0 + (1-\alpha)G)} = \alpha P_{(\theta, G_0)} + (1-\alpha)P_{(\theta, G)}.$$

If  $\mathbf{P}$  is a semiparametric group model as given in section 4.2 and  $\mathbf{G}$  is a convex set of probabilities on  $X$ , then, by definition,  $\mathbf{P}$  is convex. Examples thus include the generalization of the symmetric location model of example 3.2.4, that we considered in example 2, example 3 when the distribution  $H$  of  $Z$  is assumed known, and the elliptic model of example 4 below. Similarly,  $\mathbf{P}$  is convex if it is a mixture model as given by (4.5.1) and  $\mathbf{G}$  is convex. In particular, the errors in variables models are convex. Examples of models which are not convex include regression when both  $H$  and  $G$  are unknown, generalized regression models such as projection pursuit, transformation models including the Cox model of example 3.4.2, and censored models.

The basic property which nice convex models have is:

- (C1) For all  $h$  such that
- (i)  $h \in L_2(P_{(\theta_0, G)})$  for all  $G$ ,
  - (ii)  $h \perp \dot{P}_2(\theta_0, G_0)$ ,
  - (iii)  $\int h dP_{(\theta_0, G_0)} = 0$ ,

we have

$$(14) \quad \int h dP_{(\theta_0, G)} = 0 \quad \text{for all } G \in \mathbf{G}.$$

Here are the heuristics which lead us to expect (14). For given  $G, G_0$ , form a submodel passing through  $(\theta_0, G_0)$ , for  $0 \leq \eta \leq 1$ ,

$$\begin{aligned} Q_\eta &= P_{(\theta_0, (1-\eta)G_0 + \eta G)} \\ &= (1-\eta)P_{(\theta_0, G_0)} + \eta P_{(\theta_0, G)} \end{aligned}$$

by convexity of the parametrization. Formally,

$$(15) \quad \frac{\partial}{\partial \eta} \log q_\eta \Big|_{\eta=0} = \frac{P_{(\theta_0, G)}}{P_{(\theta_0, G_0)}} - 1.$$

Now suppose that (15) actually belongs to  $\dot{P}_2(\theta_0, G_0)$ . Then (ii) implies

$$\int h \left( \frac{P_{(\theta_0, G)}}{P_{(\theta_0, G_0)}} - 1 \right) P_{(\theta_0, G_0)} d\mu = 0,$$

which reduces to

$$\int h P_{(\theta_0, G)} d\mu = 0$$

by (iii).

Now suppose that (C1) and the necessary conditions for adaptive estimation for  $v = \theta$  at  $G = G_0$  hold. Then the efficient score function  $\dot{l}^*(x, P_{(\theta, G_0)} \mid \theta, P_1(G_0))$  for  $\theta$  in  $P_1(G_0)$  equals the score function  $\dot{l}(x, \theta, G_0) \equiv \dot{l}(x, P_{(\theta, G_0)} \mid \theta, P_1(G_0))$  for  $\theta$  and

$$(i) \quad \dot{l}(\cdot, \theta, G_0) \perp \dot{P}_2(\theta, G_0) \text{ for all } \theta$$

and

$$(ii) \quad \int \dot{l}(x, \theta, G_0) dP_{(\theta, G_0)}(x) = 0 \text{ for all } \theta$$

hold. If  $\dot{l}(x, \theta, G_0) \in L_2(P_{(\theta, G)})$  for all  $G$ , then (C1) yields

$$(16) \quad \int \dot{l}(x, \theta, G_0) dP_{(\theta, G)}(x) = 0 \quad \text{for all } G \text{ and } \theta.$$

Therefore if we base  $M$ -estimates  $\hat{\theta}_n$  on  $\dot{l}(\cdot, \theta, G_0)$ , we expect  $\hat{\theta}_n$  to be  $\sqrt{n}$ -consistent on  $\mathbf{P}$  and efficient at all  $P \in P_1(G_0)$ .

Other examples of such partially efficient  $M$ -estimates of  $\theta$  are rare. If we are interested in estimating only the first component  $v$  of  $\theta = (v, \eta)$ , and if (C1), with  $\dot{P}_3$  replacing  $\dot{P}_2$ , and the necessary condition (3.4.24) for adaptive estimation of  $v$  at  $G = G_0$  holds, then the above argument leads to



$$(17) \quad \int I_1^*(x, v, \eta_0, G_0) dP_{(v, \eta_0, G)}(x) = 0,$$

where  $I_1^*(x, v, \eta_0, G_0) = I_1^*(x, P_{(v, \eta_0, G_0)} | v, P_{12}(G_0))$  is the efficient score function for  $v$  in  $P_{12}(G_0)$ , and  $v, \eta_0$  and  $G$  are arbitrary. However, this only yields  $m$  equations for the  $k$  unknown coordinates of  $\theta$ . To extend this system of equations we add two conditions, satisfied in all our examples:

(C2) For all  $\eta_0 \in H$ ,

$$P = \{P_{(v, \eta_0, G)} : v \in N, G \in \mathbf{G}\}.$$

That is  $\eta$  and  $G$  are confounded and varying  $G$  only has the same effect as varying both  $\eta$  and  $G$ .

(C3) For each  $P \in \mathbf{P}$  the equation

$$(18) \quad \int \dot{I}_2(x, v(P), \eta, G_0) dP(x) = 0$$

has a solution  $\eta(P)$ .

Let  $P \in \mathbf{P}$ . Put  $\eta_0 = \eta(P)$ . By (C2) there exists a  $G$  such that  $P = P_{(v, \eta_0, G)}$ . Combine (17) and (18) to get

$$(19) \quad \int \dot{I}(x, v(P), \eta(P), G_0) dP(x) = 0.$$

That is, if  $\theta(IP_n) \equiv (v(IP_n), \eta(IP_n))$  is the  $M$ -estimate corresponding to (19) and the conditions of theorem 1 hold, then  $v(IP_n)$  is  $\sqrt{n}$ -consistent for  $v$  on all of  $\mathbf{P}$  and efficient at all  $P \in P_1(G_0)$ .

#### Example 4. Elliptic models.

We checked in example 4.2.3 that  $v = (\Delta, \Sigma / \text{trace } \Sigma)$  satisfies (C3). We leave it to the reader to check (C2) for  $\eta = \text{trace } \Sigma$ . If we take  $G_0$  to be the spherical  $k$ -dimensional standard Gaussian distribution,  $p(\cdot, \Delta, \Sigma, G_0)$  is the density of  $N_k(\Delta, \Sigma^{-1})$ , the resulting  $M$ -estimate is just

$$\hat{\Delta} = \bar{X}, \quad \hat{\Sigma}^{-1} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T.$$

Evidently  $\bar{X}$  and  $\hat{\Sigma} / \text{trace } \hat{\Sigma}$  are  $\sqrt{n}$ -consistent on  $\mathbf{P}$  as long as  $E |X_1|^4 < \infty$ . It is easy to see that more generally, if for a shape  $G_0$  the system

$$\int \dot{I}(x, \theta, G_0) dP(x) = 0$$

has a unique solution  $\theta(P)$  for all  $P \in \mathbf{P}$ , then in accordance with our heuristics  $v(P_{(\theta, G)}) = v$ . Conditions for unicity of solutions are discussed in Maronna (1976).  $\square$

This approach can be applied to models in which (C3) apparently fails such as errors in variables. We need only add location and scale parameters, adding two redundant parameters to our description of  $\mathbf{P}$ . As we have noted in example 4.5.2, in this parametrization  $\beta$  is adaptively estimable and (C1) and (C2) are easily satisfied. As expected, in the restricted model, taking  $G_0$  to be Gaussian leads to the regular estimate  $\hat{\beta}_p$  which defines a line minimizing the sum of

squared perpendicular distances to the data. We expect that other choices of  $G_0$  will lead to  $\sqrt{n}$ -consistent estimates in both the restricted and Reiersøl models. However, conditions on  $G_0$  for unicity of the solution of (17) are as yet unknown. There are situations other than those in which one can adapt to  $G$  in which functions  $\psi(\cdot, \theta)$  satisfying (10) can be constructed. In the following example a choice of  $\psi(x, \theta)$  is suggested by the form of the efficient score function  $I^*$ .

**Example 5. Mixture models and conditional likelihood.**

Suppose that

$$P = \{P_{(\theta, G)} = \int f(x, \theta, \eta) dG(\eta) : \theta \in \Theta, G \in \mathbf{G}\},$$

where  $f(x, \theta, \eta)$  is given by (4.5.4) with  $T(x, \theta) = T(x)$ ; that is, an exponential family in  $\eta$  with sufficient statistic  $T(X)$ . This is a convex parametrization if  $\mathbf{G}$  is convex. By (4.5.26),

$$I^*(x, \theta) = \dot{I}(x, \theta | T(x)),$$

the derivative of the log-conditional likelihood, which is independent of  $G_0$ . Hence, trivially,

$$\int I^*(x, \theta) dP_{(\theta, G)}(x) = 0 \quad \text{for all } G.$$

In fact,  $I^*$  generates efficient estimates. In general it is easy to see that if model (4.5.2) holds then (D1) through (D3) are valid provided only that if  $P = P_{(\theta_0, G_0)}$ , for some  $\varepsilon > 0$ ,

$$(i) E_{(\theta_0, G_0)} [\sup\{E_{\theta}( |S(X)|^2 | T(X)) : |\theta - \theta_0| \leq \varepsilon\}] < \infty$$

and

$$(ii) I(P | \theta, P) = E_P \text{Var}_{\theta}(S(X) | T(X)) \text{ is nonsingular.}$$

In examples 4.5.1 (model  $S$ ) and 4.5.3 these conditions are satisfied but, of course,  $\hat{\theta}_n$  can be calculated explicitly in both cases.  $\square$

In other situations we can employ symmetries of the problem.

**Example 6. The Has'minskii-Ibragimov model.**

Recall that in this model, first given in section 4.5, we have  $X = (U', Y)$  where  $U', Y$  are independent  $U' \sim G$  and  $Y$  is from the mixture model with density

$$\int q(y, \theta, u) dG(u).$$

The parametrization is convex if  $\mathbf{G}$  is. We can think of  $Y$  as coming from  $(U, Y)$  where  $U \sim G$  and given  $U = u$ ,  $Y$  has density  $q(y, \theta, u)$ . Note that  $P = P_{(\theta, G)} \in \mathbf{P}$  has margins  $G$  and  $\int Q(\cdot, \theta, u) dG(u)$  for  $U'$  and  $Y$  respectively.

Let  $\gamma(y, \theta)$  be given and let

$$\alpha(u, \theta) \equiv E_{\theta}(\gamma(Y, \theta) | U = u) = \int \gamma(y, \theta) q(y, \theta, u) d\mu(y).$$

Then

$$E_G \alpha(U', \theta) = E_{(\theta, G)} \gamma(Y, \theta).$$

Evidently, if

$$\psi(x, \theta) \equiv \gamma(y, \theta) - \alpha(u', \theta),$$

then  $\psi$  satisfies (11).

Satisfaction of (D1) and (D2) for  $\gamma$  and  $v$  separately implies (D1) and (D2) for  $\psi$ . Under conditions permitting interchange of integration and differentiation,

$$E \left( \frac{\partial}{\partial \theta_j} \psi_i(X, \theta) \right) = - E(\gamma_i(Y, \theta) \frac{\partial}{\partial \theta_j} \log q(Y, \theta, U)).$$

(D3) is satisfied at  $G_0$  if, for instance,

$$\gamma_i(y, \theta) \equiv E_{(\theta, G_0)} \left[ \frac{\partial}{\partial \theta_i} \log q(Y, \theta, U) \mid Y \right]$$

and  $\nabla \log q(Y, \theta, U)$  is genuinely  $k$ -dimensional. Has'minskii and Ibragimov (1983) show that if the parametric model with densities  $q(\cdot, \theta, \cdot)$ ,  $\theta \in \Theta$  is regular and  $\theta$  is identifiable in that model, then  $\sqrt{n}$ -consistent estimates of  $\theta$  can be constructed. In fact, by adapting within the class of  $M$ -estimates and using one step estimates starting from  $\sqrt{n}$ -consistent estimates, they show how to construct efficient estimates. It is often not hard in these models to find  $\sqrt{n}$ -consistent starting values  $\bar{\theta}(P_n)$  where  $\bar{\theta}$  is a  $\rho$ -Lipschitz extension of  $\theta$  for  $\rho$  compatible with  $P_n$  in the sense of section 2.5. For example, consider the normal convolution model: then  $Y = U + \varepsilon$  where  $\varepsilon \sim N(\mu, \sigma^2)$  is independent of  $U$ , and  $\theta = (\mu, \sigma^2)$ . Then

$$\mu(P) = \int y dF(y) - \int u dG(u) = E_F Y - E_G(U')$$

and

$$\sigma^2(P) = \text{Var}_F(Y) - \text{Var}_G(U'),$$

where  $G$  and  $F$  are the marginal distributions of  $X = (U', Y)$ , will be satisfactory extensions if, say, we suppose  $G[-M, M] = 1$  for some  $0 < M < \infty$ .

Symmetries are used similarly by Robbins and Zhang (1989) in a related class of models.  $\square$

### 7.3 GENERALIZED $M$ -ESTIMATES FOR EUCLIDEAN PARAMETERS

We begin with a slight extension of the definition of generalized  $M$ -estimates for Euclidean parameters given in sections 7.1.

**Definition 1.** Suppose  $\mathbf{M}_0 \supset \mathbf{P}$  and all distributions with finite support as in sections 7.1 and 7.2. Suppose that  $W_n, W$  map  $R^m \times \mathbf{M}_0 \rightarrow R^m$  and that

- (i)  $W_n(v, P) = W(v, P) + o(1)$  for all  $P \in \mathbf{M}_0$ , all  $v$ ,
- (ii)  $W(v(P), P) = 0$  for all  $P \in \mathbf{P}$ .

For the empirical distribution  $P_n$  we introduce the notation  $W_n(v) = W_n(v, P_n)$ . Then,  $\hat{v}_n$  is a *generalized  $M$ -estimate*, or *GM-estimate* of  $v(P)$  if

$$(1) \quad W_n(\hat{v}_n) = 0.$$

If

$$(2) \quad W_n(\hat{v}_n) = o_p(n^{-1/2}) \quad \text{for all } P \in \mathbf{P},$$

we say  $\hat{v}_n$  is an *asymptotic generalized M-estimate*, or *AGM-estimate*, of  $v$ .

Going beyond  $W_n \equiv W$  enables us to include estimates such as

$$s^2 \equiv (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n}{n-1} \int (x - \int x dP_n)^2 dP_n(x)$$

and, as we have discussed in section 7.1, this extension becomes crucial when we consider estimation of infinite-dimensional parameters.

Estimates of this type were first considered by Filippova (1962) under regularity conditions considerably more stringent than the ones we present. That *some* regularity conditions are necessary to say anything is evident since any estimate  $\hat{v}_n(P_n)$  can be considered as *GM* by taking  $W_n(v, P) = v - \hat{v}_n(P)$ . Introducing *AGM*-estimates is an extension of what Huber (1967) did for *M*-estimates. This extension enables us to include estimates such as the median which only approximately satisfies  $\int \text{sgn}(x - v) dP_n(x) = 0$ . See also example 2 and the discussion of one-step estimates below.

**Example 1. M-estimates.**

With  $\psi$  and  $W$  as in section 7.2, set

$$(3) \quad \begin{aligned} W_n(v) &= W_n(v, P_n) \equiv W(v, P_n) \\ &= \int \psi(x, v) dP_n(x) = n^{-1} \sum_{i=1}^n \psi(X_i, v). \end{aligned}$$

With this  $W_n$  we have called  $\hat{v}_n$  an *M-estimator* if it satisfies (1), or an *asymptotic M-estimator* if it satisfies (2).

As we have seen in section 7.2, under regularity conditions, (asymptotic) *M*-estimates are asymptotically linear estimates of  $v(P)$  with influence function a linear transformation of  $\psi$ .  $\square$

**Example 2. The Hodges-Lehmann estimate.**

Let  $\mathbf{P} = \{\text{all absolutely continuous symmetric distributions on } R\}$ , and let  $v(P)$  be the center of symmetry of  $P \in \mathbf{P}$ . Note that

$$P(X_1 + X_2 > 2v(P)) - \frac{1}{2} = 0$$

for  $X_1, X_2$  i.i.d.  $P$ . Define

$$(4) \quad W(v, Q) \equiv Q(X_1 + X_2 > 2v) - \frac{1}{2},$$

which makes sense for *any* distribution  $Q \in \mathbf{M}$ . This definition (4) of  $W(\cdot, Q)$  for  $Q \in \mathbf{M}$  suggests extensions  $\bar{v}$  of  $v$  to  $\mathbf{M}_0 \supset \mathbf{P}$  as a median of the distribution of the average of a pair  $X_1, X_2$  i.i.d.  $Q$ . For  $Q \in \mathbf{Q} \equiv \mathbf{P} \cup \{\text{all } Q \text{ with strictly positive density}\}$ , the extension  $\bar{v}$  is unique, and satisfies  $W(\bar{v}(Q), Q) = 0$ . Let

$$(5) \quad W_n(v) = W_n(v, P_n) \equiv W(v, P_n) = n^{-2} \sum_{i,j} \left\{ 1_{[X_i + X_j > 2v]} - \frac{1}{2} \right\}.$$

Then the well-known Hodges-Lehmann estimator

$$(6) \quad \hat{v}_n = \text{median}_{i,j} \frac{X_i + X_j}{2}$$

satisfies  $\mathbf{W}_n(\hat{v}_n) = O_p(n^{-2})$ , for  $P \in \mathbf{Q}$ , and hence is an  $AGM$ -estimator of  $\bar{v}(P)$ ,  $P \in \mathbf{Q}$ .  $\square$

**Example 3.  $M$ -estimates for  $v$  with  $\eta$  estimated.**

Suppose  $\mathbf{P} = \{P_\theta : \theta \in \Theta\}$  is a parametric model where  $\theta^T = (v^T, \eta^T)$ , with  $v$  of dimension  $m$ ,  $\eta$  of dimension  $k - m$ , and we want to estimate  $v$ .

It is sometimes convenient to proceed as follows.

(i) Find functionals  $\tilde{\eta}(P)$  or more generally  $\tilde{\eta}(v, P)$  such that  $\tilde{\eta}(v, P_n)$  is a consistent estimate of  $\eta(P)$  under  $\mathbf{P}_2(v)$ .

(ii) Pick  $\psi(\cdot, v, \eta)_{m \times 1}$  such that for all  $\theta$

$$(7) \quad \int \psi(x, v, \eta) dP_\theta(x) = 0,$$

and construct  $GM$ -estimates  $\hat{v}_n$  solving

$$(8) \quad \mathbf{W}_n(v)_{m \times 1} \equiv n^{-1} \sum_{i=1}^n \psi(X_i, v, \tilde{\eta}(v, P_n)) = 0.$$

These estimates correspond to  $\mathbf{W}_n(v, P) \equiv \int \psi(x, v, \tilde{\eta}(v, P)) dP(x)$ . An example of this approach, with  $\tilde{\eta}(P)$  independent of  $v$ , is discussed in Huber (1981, section 6.5, page 140), and occurs in robust  $M$ -estimation of location. Here  $v$ , the center of symmetry, is estimated using

$$\mathbf{W}_n(v) = n^{-1} \sum_{i=1}^n \psi\left(\frac{X_i - v}{\tilde{\eta}(P_n)}\right),$$

where  $\psi$  is an antisymmetric scalar function and  $\tilde{\eta}(P_n)$  is an estimate of the unknown scale  $\eta$ . Other examples are discussed in Gong and Samaniego (1981), who propose this method with  $\psi = \dot{I}_1$ , to obtain estimates in mixture models. See also Parke (1986).  $\square$

The remainder of this section is organized as follows: We first give a very general "master theorem" for  $AGM$ -estimates under a consistency assumption. This generalises the results for  $M$ -estimates and  $AM$ -estimates which were presented in section 7.2. Then we treat one step  $AGM$ -estimates. By considering iteration of the one-step procedures we obtain results for  $GM$ -estimates.

*Asymptotic Generalized  $M$ -estimates: the Master Theorem*

For  $\mathbf{W}_n$ ,  $\mathbf{W}_n$  and  $W$  as in the definition of a  $GM$ -estimate, let  $\mathbf{V}_n$  be the centered and  $\sqrt{n}$ -normalized version of  $\mathbf{W}_n$ ,

$$(9) \quad \mathbf{V}_n(v) \equiv \sqrt{n} \{\mathbf{W}_n(v) - W(v, P)\}$$

for  $v \in N$  open  $\subset R^m$ . We regard  $\mathbf{V}_n$  as a random function with argument  $v$ . Let  $\mathbf{Q}$  be a model with  $\mathbf{Q} \supset \mathbf{P}$ . Here are the key assumptions:

(GM0) There exists  $v : Q \rightarrow R^m$  such that  $v(P)$  satisfies  $W(v(P), P) = 0$  for all  $P \in Q$ .

(GM1) For any  $\varepsilon_n \downarrow 0$  we have

$$(10) \quad \sup \left\{ \frac{|\mathbf{V}_n(v) - \mathbf{V}_n(v(P))|}{1 + \sqrt{n} |v - v(P)|} : |v - v(P)| \leq \varepsilon_n \right\} = o_p(1).$$

(GM2) There is a function  $\psi : X \times Q \rightarrow R^m$  with  $\int \psi(x, P) dP(x) = 0$  and  $|\psi(\cdot, P)| \in L_2(P)$  such that

$$\mathbf{W}_n(v(P)) = n^{-1} \sum_{i=1}^n \psi(X_i, P) + o_p(n^{-1/2}).$$

(GM3)  $W(\cdot, P) = (W_1(\cdot, P), \dots, W_m(\cdot, P))^T$  is differentiable with derivative  $\dot{W}(v, P) \equiv [\partial/\partial v_j W_i(v, P)]_{m \times m}$  and  $\dot{W}(P) \equiv \dot{W}(v(P), P)$  is non-singular.

We introduce the model  $Q$  to emphasize that even though  $W_n$  may be motivated by features of  $P$ , the AGM-estimates corresponding to  $W_n$  can be thought of as estimates of parameters on a larger model. The following weaker version of (GM1) is useful:

(GM1') For all  $M < \infty$ ,

$$\sup \{ |\mathbf{V}_n(v) - \mathbf{V}_n(v(P))| : \sqrt{n} |v - v(P)| \leq M \} = o_p(1).$$

**Theorem 1.** Suppose  $P \in Q$ . Let  $\hat{v}_n$  be an asymptotic generalized  $M$ -estimate of  $v(P)$  on  $Q$ . If  $\hat{v}_n$  is consistent and if (GM0)–(GM3) hold, then  $\hat{v}_n$  is an asymptotically linear estimate of  $v(P)$  with influence function  $-\dot{W}^{-1}(P)\psi(\cdot, P)$ :

$$(11) \quad \sqrt{n}(\hat{v}_n - v(P)) = -\frac{1}{\sqrt{n}} \sum_{i=1}^n \dot{W}^{-1}(P)\psi(X_i, P) + o_p(1). \\ \rightarrow_d N(0, \Sigma(\hat{v}, P)) \quad \text{as } n \rightarrow \infty,$$

where

$$(12) \quad \Sigma(\hat{v}, P) = \dot{W}^{-1}(P)E[\psi(X, P)\psi^T(X, P)][\dot{W}^{-1}(P)]^T.$$

If (GM1) is weakened to (GM1') and  $\hat{v}_n$  is  $\sqrt{n}$ -consistent, then the conclusion remains valid.

Theorem 1 generalizes the results of Huber (1967) and Pollard (1985) to general estimating equations. A detailed proof of this theorem is given in section A.10; see theorem A.10.1. Here are the heuristics behind the (GM) conditions.

Write

$$(13) \quad -\mathbf{W}_n(v(P)) = \mathbf{W}_n(\hat{v}_n) - \mathbf{W}_n(v(P)) + o_p(n^{-1/2}) \quad \text{by (2)}$$

$$= W(\hat{v}_n, P_n) - W(v(P), P_n) + n^{-1/2}(W_n(\hat{v}_n) - W_n(v(P))) + o_P(n^{-1/2}).$$

Condition (GM1) and the consistency of  $\hat{v}_n$  are needed to justify  $n^{-1/2}(W_n(\hat{v}_n) - W_n(v(P))) = o_P(n^{-1/2})$ . Condition (GM2) permits us to replace  $W_n(v(P))$  in (13) by  $n^{-1} \sum_{i=1}^n \psi(X_i, P) + o_P(n^{-1/2})$ . Finally, (GM3) and consistency of  $\hat{v}_n$  are used to replace  $W(\hat{v}_n, P_n) - W(v(P), P_n)$  by  $\dot{W}(P)(\hat{v}_n - v(P))(1 + o_P(1))$ . We end up with

$$\dot{W}(P)(\hat{v}_n - v(P)) = -n^{-1} \sum_{i=1}^n \psi(X_i, P) + o_P(n^{-1/2})$$

as advertised.

It is important to note that special features of the problem can and should be used to verify these conditions. For instance,  $(M1')$  with  $W_n$  no longer restricted to be linear in  $P$  implies (GM1), as does

$$(14) \quad \sup\{|\dot{W}_n(v) - \dot{W}(v(P), P)| : |v - v(P)| \leq \varepsilon_n\} = o_P(1),$$

for  $W_n$  and  $W$  differentiable. This remark underlies corollary 7.2.1 and the argument needed for example 6.

If  $W_n$  is given by (3), condition (GM2) is automatically satisfied. On the other hand, if we have an explicit representation for  $\hat{v}_n$  as  $v_n(P_n)$  where  $v_n: M_0 \rightarrow R^m$  and we take  $W_n(v, Q) \equiv v - v_n(Q)$ , only (GM2) has to be checked. In general (GM2) is often attacked by a method introduced by von Mises (1947) and developed by Filippova (1962). It can be broken down into the following steps.

- (i) Suppose  $M_0$  is convex. Compute formally a Gâteaux derivative,

$$\frac{\partial}{\partial \alpha} W_n(v(P), (1 - \alpha)P + \alpha Q) |_{\alpha=0} \equiv d_P W_n(v(P), P)(Q - P).$$

- (ii) Suppose, as is usually the case, that

$$(15) \quad d_P W_n(v(P), P)(Q - P) = \int \psi_n(x, P) d(Q - P)(x),$$

$$\text{where } \int \psi_n(x, P) dP(x) = 0.$$

- (iii) Show that

$$W_n(v(P), P_n) = W_n(v(P), P) + n^{-1} \sum_{i=1}^n \psi_n(X_i, P) + o_P(n^{-1/2}).$$

- (iv) Check that  $W_n(v(P), P) = o(n^{-1/2})$ ,  $E_P(\psi_n(X, P) - \psi(X, P))^2 \rightarrow 0$ .

The difficult step (iii) is often best done by direct calculation but the points of view of Von Mises (1947), Reeds (1976), Gill (1988), (1989), and Wong and Severini (1991) are sometimes useful.

For  $W_n \equiv W$ , the approach of Reeds is to establish Hadamard (compact) differentiability of  $Q \rightarrow W(v(P), Q)$  when  $M_0$  is viewed as a subset of an

appropriate linear topological space  $S$  and  $W_n$  is defined on  $S$ , and then to use tightness of  $\sqrt{n}(\mathbb{P}_n - P)$ . In practice we need to show that:

For all  $\epsilon > 0$  there exists a compact  $K_\epsilon \subset S$  such that

$$(a) \quad \sup \{ |W_n(v(P), P + \alpha \Delta) - W_n(v(P), P) - \alpha d_P W_n(v(P), P)(\Delta)| : \Delta \in K_\epsilon \} = o(\alpha),$$

$$(b) P(\sqrt{n}(\mathbb{P}_n - P) \in K_\epsilon) \geq 1 - \epsilon.$$

If (15) holds and we put  $\alpha = n^{-1/2}$ ,  $\Delta = \sqrt{n}(\mathbb{P}_n - P)$  we obtain (GM2).

The von Mises approach is to calculate a second Gâteaux derivative (if possible) given by

$$d_P^2 W_n(v(P), P)(Q_1 - P, Q_2 - P) = \frac{\partial^2}{\partial \alpha_1 \partial \alpha_2} W_n(v(P), (1 - \alpha_1 - \alpha_2)P + \alpha_1 Q_1 + \alpha_2 Q_2) |_{\alpha_1 = \alpha_2 = 0},$$

and satisfying

$$d_P^2 W_n(v(P), P)(Q_1 - P, Q_2 - P) = \iint \psi_n(x, y, P) d(Q_1 - P)(x) d(Q_2 - P)(y),$$

where  $\psi_n(x, y, P) = \psi_n(y, x, P)$  and

$$\int \psi_n(x, y, P) dP(x) = 0 \quad \text{a.s.}$$

Then, using Cauchy's form of the Taylor expansion, we can write

$$W_n(v(P), \mathbb{P}_n) = W_n(v(P), P) + \int \psi_n(x, P) d(\mathbb{P}_n - P)(x) + \int_0^1 (1 - \lambda) \iint \psi_n(x, y, (1 - \lambda)P + \lambda \mathbb{P}_n) d(\mathbb{P}_n - P)(x) d(\mathbb{P}_n - P)(y) d\lambda.$$

If

$$(16) \quad \sup_\lambda E_P | \iint \psi_n(x, y, (1 - \lambda)P + \lambda \mathbb{P}_n) d(\mathbb{P}_n - P)(x) d(\mathbb{P}_n - P)(y) | = o_P(n^{-1/2}),$$

(GM2) will follow. If for instance, as in example 2, (when  $v$  is fixed)  $\psi_n(x, y, P) = \psi(x, y)$  doesn't depend on  $n$  and  $P$ , and  $E_P \psi^2(X_1, X_2) < \infty$ , (16) follows trivially from second moment bounds for  $U$  statistics. In our opinion, the most valuable aspect of the von Mises approach is in identifying  $\psi_n(\cdot, P)$  and  $\psi_n(\cdot, \cdot, P)$  as objects to focus on.

The Reeds approach typically has the advantage of simultaneously establishing linearity and regularity via proposition 3.3.1. Suppose  $W_n$  does not depend on  $n$ . Then (a) and (15) suggest that  $W(v(P), \cdot)$  is pathwise differentiable with derivative  $\psi(\cdot, P)$  (in  $L_2(P)$ ). That is, if  $h$  is the tangent corresponding to  $\{P_\eta\}$ , then

$$(17) \quad W(v(P), P_\eta) = \int \psi(x, P) d(P_\eta - P) + o(\eta)$$



$$= \eta \int \psi(x, P) h(x) dP(x) + o(\eta).$$

Now if  $v$  is pathwise differentiable, differentiation of the equation in (GM0) yields

$$(18) \quad \dot{W}(P) \dot{v}(P) + \psi(\cdot, P) = 0.$$

Hence,  $\dot{v}(P) = -\dot{W}^{-1}(P) \psi(\cdot, P)$ , and regularity follows.

**Example 2. The Hodges-Lehmann estimate, continued.**

By (4) and (5)

$$\begin{aligned} \mathbf{V}_n(v) &= \sqrt{n}(\mathbf{W}_n(v) - W(v, P)) \\ &= -\sqrt{n} \left( \int \mathbf{F}_n(2v - x) d\mathbf{F}_n(x) - \int F(2v - x) dF(x) \right) \end{aligned}$$

where  $\mathbf{F}_n$  is the empirical distribution function and  $F$  is the true distribution function. After an integration by parts

$$\begin{aligned} \mathbf{V}_n(v) &= -2\sqrt{n} \int F(2v - x) d(\mathbf{F}_n - F)(x) \\ &\quad - \sqrt{n} \int (\mathbf{F}_n - F)(2v - x) d(\mathbf{F}_n - F)(x). \end{aligned}$$

At  $v = \bar{v}(P)$  the first term equals  $n^{-1/2} \sum_{i=1}^n \psi(X_i, P)$  with

$$\begin{aligned} \psi(x, P) &= -2 \{ F(2\bar{v}(P) - x) - P(X_1 + X_2 \leq 2\bar{v}(P)) \} \\ &= -2 \{ F(2\bar{v}(P) - x) - \frac{1}{2} \}. \end{aligned}$$

We leave it as an exercise to show that, for all  $M > 0$ ,

$$\begin{aligned} &\sqrt{n} \sup \{ | \int (\mathbf{F}_n - F)(2v - x) d(\mathbf{F}_n - F)(x) | : |v| \leq M \} \\ &= \sup \{ | \int_0^1 \{ \sqrt{n}(\mathbf{F}_n - F)(2v - \mathbf{F}_n^{-1}(s)) \\ &\quad - \sqrt{n}(\mathbf{F}_n - F)(2v - F^{-1}(s)) \} ds | : |v| \leq M \} \\ &= o_p(1). \end{aligned}$$

(see, e.g., Shorack and Wellner (1986, chapter 14)) and that

$$\begin{aligned} &\sup \{ \sqrt{n} | \int \{ F(2v - x) - F(2\bar{v}(P) - x) \} d(\mathbf{F}_n - F)(x) | : \\ &\quad |v - \bar{v}(P)| \leq \varepsilon_n \} = o_p(1) \end{aligned}$$

whenever  $\varepsilon_n \downarrow 0$ . Since  $\hat{v}_n$  is easily shown to be a consistent estimate of  $v$  and

$$\dot{W}(P) = -2 \int f(2\bar{v}(P) - x) dF(x),$$

where  $f = F'$ , theorem 1 applies and we conclude that  $\hat{v}_n$  given by (6) has influence function

$$\frac{F(2\bar{v}(P) - x) - \frac{1}{2}}{\int f(2\bar{v}(P) - x) dF(x)}$$

If  $F$  is symmetric ( $P \in \mathbf{P}$ ) about  $v(P)$ , then the influence function reduces to

$$\frac{F(x) - \frac{1}{2}}{\int f^2(x) dx}$$

as is well known; see, e.g., Huber (1981, formula (3.4.20), page 64). Regularity of  $\hat{v}_n$  follows by checking (17) and (18).  $\square$

### One-step Estimates

Unfortunately, equations (1), even if they are likelihood equations, can have, for all  $n$ , multiple (or no) solutions and some of the roots may be inconsistent; see, for instance, Reeds' (1985) treatment of the Cauchy location problem. However if we have a  $\sqrt{n}$ -consistent estimate  $\tilde{v}_n$ , we can construct a consistent AGM-estimate by the *one-step* method.

Suppose (GM0)–(GM3) hold and we have available  $\tilde{v}_n$  which is  $\sqrt{n}$ -consistent and

$$(19) \quad \dot{W}^* = \dot{W}(P) + o_p(1),$$

which is a consistent estimate of  $\dot{W}$ . Define the one-step estimate

$$(20) \quad \hat{v}_n = \tilde{v}_n - (\dot{W}^*)^- \mathbf{W}_n(\tilde{v}_n),$$

where  $(\dot{W}^*)^-$  is a generalized inverse of  $\dot{W}^*$ . If  $\mathbf{W}_n$  is differentiable and  $\dot{W}_n(\tilde{v}_n)$  is a consistent estimate of  $\dot{W}(P)$  (so that (19) and (20) hold with  $\dot{W}^* = \dot{W}_n(\tilde{v}_n)$ ), then  $\hat{v}_n$  is the first Newton-Raphson iteration for solving  $\mathbf{W}_n(v) = 0$  starting from  $\tilde{v}_n$ . We claim that  $\hat{v}_n$  is a consistent AGM-estimate so that the conclusion of the master theorem applies. This estimate is  $\sqrt{n}$ -consistent since  $\tilde{v}_n$  is and

$$(21) \quad \mathbf{W}_n(\tilde{v}_n) = O_p(n^{-1/2})$$

by (GM1) and (GM2).

Further, by (GM1) and the consistency of  $\tilde{v}_n$  and  $\hat{v}_n$ ,

$$\begin{aligned} \mathbf{W}_n(\hat{v}_n) &= \mathbf{W}_n(\tilde{v}_n) + W(\hat{v}_n, P) - W(\tilde{v}_n, P) + o_p(n^{-1/2}) \\ &= \mathbf{W}_n(\tilde{v}_n) - \dot{W}(P)(\dot{W}^*)^- \mathbf{W}_n(\tilde{v}_n) + o_p(n^{-1/2}) \\ &\quad \text{by the } \sqrt{n} \text{-consistency of } \tilde{v}_n, (20), \text{ and (GM3)} \\ &= o_p(n^{-1/2}) \quad \text{by (19) and (21),} \end{aligned}$$

and  $\hat{v}_n$  is an AGM-estimate. With  $\tilde{v}_n$  being  $\sqrt{n}$ -consistent, this argument suffices to apply theorem 1 under (GM1').

Even (GM1') can be weakened if we discretize  $\tilde{v}_n$ . We say that  $\tilde{v}_n$  is *discre-*

tized if it takes as possible values only the vertices of a decomposition of  $R^m$  into cubes of size  $cn^{-1/2}$ ,  $c > 0$ . If  $\tilde{v}_n$  is any  $\sqrt{n}$ -consistent estimate, we can construct a discretized  $\sqrt{n}$ -consistent estimate from it by following the procedure described in section 2.5.

(GM1'')  $\mathbf{V}_n(v_n) - \mathbf{V}_n(v(P)) = o_p(1)$  for any deterministic sequence  $\{v_n\}$  such that  $\sqrt{n}|v_n - v(P)| \leq M$  for some  $M$  and all  $n$ .

**Theorem 2.** Suppose (GM0), (GM2), and (GM3) hold. Suppose that  $\tilde{v}_n$  is  $\sqrt{n}$ -consistent and either: (GM1') holds; or  $\tilde{v}_n$  is discretized and (GM1'') holds. Suppose also that  $\dot{W}^*$  is a consistent estimate of  $\dot{W}(P)$ . Then the conclusion of theorem 1 holds for the one-step estimator  $\hat{v}_n$  given by (20).

**Proof.** The argument under (GM1') is given above. For the discretized preliminary estimate, it can be adapted by the technique of the proof of theorem 2.5.2.  $\square$

We give two examples: a one-step  $M$ -estimate and a one-step  $GM$ -estimate.

**Example 4. Estimation of location for the Cauchy distribution.**

Let

$$p(x, \theta) = \pi^{-1} (1 + (x - \theta)^2)^{-1},$$

$$\psi(x, \theta) \equiv \dot{l}(x, \theta) = 2(x - \theta)(1 + (x - \theta)^2)^{-1}.$$

Take

$$W_n(\theta, P_n) = \int \psi(x, \theta) dP_n(x)$$

so that

$$(22) \quad W(\theta, P) = \int \psi(x, \theta) dP(x) = 2 \int (x - \theta)(1 + (x - \theta)^2)^{-1} dP(x).$$

If  $P$  is symmetric about  $\theta_0$ ,  $v(P) = \theta_0$ . Note that  $W$  is actually defined for all  $P \in \mathbf{M}$ ,

$$\dot{W}(P) = 2 \int ((x - \theta)^2 - 1)(1 + (x - \theta)^2)^{-2} dP(x)$$

and conditions (D1)–(D3) of section 7.2 hold provided  $\dot{W}(P) \neq 0$ . If  $P$  is symmetric about 0 and has a continuous positive density at 0, we can take as preliminary estimate the median  $\tilde{\theta}_n$  and

$$\dot{W}^* \equiv 2n^{-1} \sum_{i=1}^n [(X_i - \tilde{\theta}_n)^2 - 1](1 + (X_i - \tilde{\theta}_n)^2)^{-2}.$$

It is easy to see that (19) is satisfied and

$$\hat{\theta}_n = \tilde{\theta}_n - \frac{\sum_{i=1}^n (X_i - \tilde{\theta}_n)(1 + (X_i - \tilde{\theta}_n)^2)^{-1}}{\sum_{i=1}^n [(X_i - \tilde{\theta}_n)^2 - 1](1 + (X_i - \tilde{\theta}_n)^2)^{-2}}$$

behaves like the MLE for the Cauchy should, and, in particular, is efficient if the Cauchy holds. Note that if  $P$  is not symmetric,  $\theta(P)$  solving  $W(\theta, P) = 0$  is in general not the median of  $P$ . Thus,  $\tilde{\theta}_n$  is not a  $\sqrt{n}$ -consistent estimate of  $\theta(P)$

and the conclusion as well as the conditions of the theorem fail. It is clear heuristically and it can be verified that  $\hat{\theta}_n$  converges to

$$\theta^*(P) \equiv \tilde{\theta} - \dot{W}^{-1}(\tilde{\theta}, P) W(\tilde{\theta}, P),$$

where  $\tilde{\theta} = \text{median}(P)$ , not the root of  $W(\theta, P) = 0$ .  $\hat{\theta}_n$  has influence function which is a linear combination of  $\psi(x, \tilde{\theta}) - W(\tilde{\theta}, P)$  and the score function of  $\tilde{\theta}_n$ , which is  $1_{[x > \tilde{\theta}]} - \frac{1}{2}$ .  $\square$

**Example 5. Minimum Cramér-von Mises distance estimation for the exponential distribution.**

Let

$$D_n(\theta, P_n) \equiv \int_0^\infty (\bar{F}_n(x) - e^{-\theta x})^2 dF_n(x)$$

where  $\bar{F}_n(x) \equiv 1 - F_n(x) \equiv P_n(x, \infty)$  and

$$W_n(\theta, P_n) \equiv 2 \int_0^\infty x (\bar{F}_n(x) - e^{-\theta x}) e^{-\theta x} dF_n(x).$$

$W_n(\theta, P_n) = \dot{D}_n(\theta, P_n)$  is not linear in the observations. But  $W_n$  satisfies (GM0)–(GM3) with

$$(23) \quad W(\theta, P) \equiv 2 \int_0^\infty x (\bar{F}(x) - e^{-\theta x}) e^{-\theta x} dF(x),$$

where  $\theta(P)$  solves  $W(\theta, P) = 0$ . Further, with  $\theta \equiv \theta(P)$ ,

$$\dot{W}(P) = -2 \int_0^\infty x^2 e^{-\theta x} (\bar{F}(x) - 2e^{-\theta x}) dF(x)$$

and

$$\begin{aligned} \psi(x, P) = & 2\{x(\bar{F}(x) - e^{-\theta x})e^{-\theta x} \\ & + \int_0^\infty y(1_{[x > y]} - e^{-\theta y})e^{-\theta y} dF(y)\}. \end{aligned}$$

If  $P$  is exponential,  $\bar{F}(x) = e^{-\theta_0 x}$ , then  $\theta(P) = \theta_0$ , and  $\dot{W}(P) = 4/27\theta_0^2$ . It is easy to see that under these conditions,  $\tilde{\theta}_n \equiv (n^{-1} \sum_{i=1}^n X_i)^{-1}$  is  $\sqrt{n}$ -consistent for  $\theta_0$ , and we can take  $\dot{W}^* \equiv 4/27\tilde{\theta}_n^2$  in the definition of the one step to satisfy (19). Again, the one-step  $\hat{\theta}_n$  is an asymptotically linear estimate of  $\theta(P)$  with influence function given by the master theorem. If  $P$  is not exponential,  $\tilde{\theta}_n$  doesn't estimate  $\theta(P)$ , and  $\dot{W}^*$  doesn't estimate  $\dot{W}(P)$ , although again we can show  $\hat{\theta}_n$  to be an asymptotically linear estimate of a suitable  $\theta^*(P)$ .  $\square$

In these examples regularity follows by calculating (17) and (18) for  $W$  given by (22) and (23) respectively.

### *Newton-Raphson Iteration and Existence of GM-estimates*

Consider the one-step estimate  $\hat{v}_n$  given by (20). Suppose  $\dot{W}_n \equiv [(\partial/\partial v_j) W_{ni}(v)]$  exists where  $W_n \equiv (W_{n1}, \dots, W_{nm})^T$  and that for some  $\varepsilon > 0$

$$\sup\{|\dot{W}_n(v) - \dot{W}(v, P)| : |v - v(P)| \leq \varepsilon\} = o_P(1).$$

It is then possible to take  $\dot{W}^* = \dot{W}_n(\tilde{v}_n)$ . We can iterate this process. Let

$$(24) \quad \begin{aligned} T_n^{(0)} &= \tilde{v}_n, \\ T_n^{(j+1)} &= T_n^{(j)} - \dot{W}_n^-(T_n^{(j)})\mathbf{W}_n(T_n^{(j)}). \end{aligned}$$

Since  $\sqrt{n}$ -consistency of  $T_n^{(j)}$  implies  $\sqrt{n}$ -consistency of  $T_n^{(j+1)}$  under the conditions of theorem 2, it follows that each  $T_n^{(j)}$  is  $\sqrt{n}$ -consistent once  $\tilde{v}_n$  is. Therefore,  $T_n^{(j+1)}$  satisfies the conclusion of the master theorem for all fixed  $j$  if (GM0), (GM1'), (GM2), and (GM3) hold. Of course, we typically want to iterate to the limit  $T_n^{(\infty)}$  which would give us a  $GM$ -estimate. But the limit may not exist, and even if it does, may not satisfy the consistency condition of the master theorem. It is however possible under an additional uniform convergence condition on  $\dot{W}_n(v)$  to establish existence and asymptotic linearity of  $T_n^{(\infty)}$  with probability tending to 1. In fact we shall show a little more. Even if  $\tilde{v}_n$  is not  $\sqrt{n}$ -consistent but only "close" to  $v(P)$ , then  $T_n^{(\infty)}$  exists, is the unique root of  $\mathbf{W}_n(v) = 0$  in a neighborhood of  $v(P)$ , and is asymptotically linear. Unfortunately, we are left with the requirement that, for all  $P$ , the starting value  $\tilde{v}_n$  is "close" to  $v(P)$ , that is essentially with consistency for  $\tilde{v}_n$ . The additional assumption about  $\dot{W}_n(v)$  is:

(U) For some sequence  $\{\epsilon_n\}$  with  $\epsilon_n \downarrow 0, \epsilon_n n^{1/2} \rightarrow \infty$ ,

$$(25) \quad \sup\{|\dot{W}_n(v) - \dot{W}(P)| : |v - v(P)| \leq \epsilon_n\} = o_p(1).$$

Note that assumption (U) is implied by the following two conditions:

(S1)  $\mathbf{W}_n(v(P)) - \dot{W}(P) = o_p(1)$ .

(S2) For some sequence  $\{\epsilon_n\}$  with  $\epsilon_n \downarrow 0, \epsilon_n n^{1/2} \rightarrow \infty$ ,

$$(26) \quad \sup\{\epsilon_n \left| \frac{\partial^2}{\partial v_j \partial v_k} \mathbf{W}_{ni}(v) \right| : |v - v(P)| \leq \epsilon_n, 1 \leq i, j, k \leq m\} = o_p(1).$$

**Theorem 3.** Suppose (GM0), (GM2), (GM3), and (U) hold.

- A. With probability converging to 1,  $\mathbf{W}_n(v)$  has a unique root  $\hat{v}_n$  in  $\{v : |v - v(P)| \leq \epsilon_n\}$ .  $\hat{v}_n$  is  $\sqrt{n}$ -consistent and asymptotically linear with influence function  $-\dot{W}^{-1}(P)\psi(\cdot, v(P))$ .
- B. If there exists an estimator  $\tilde{v}_n$  satisfying  $P(|\tilde{v}_n - v(P)| < \epsilon_n) \rightarrow 1$  for  $\{\epsilon_n\}$  as in (U) and if the Newton-Raphson iteration (24) starts at  $\tilde{v}_n$ , then

$$(27) \quad P(T_n^{(\infty)} \text{ exists and equals } \hat{v}_n) \rightarrow 1, \quad \text{as } n \rightarrow \infty.$$

Here is a major example where the conditions of Theorem 3 are satisfied.

**Example 6. The regression-transformation model.**

Consider the special case of the regression-transformation model of example 4.7.1 in which  $\tau(Y) = T = \theta Z + \epsilon$  and given  $Z = z$ ,  $T$  has density  $g(t - \theta z)$ ,  $\theta$  real. So, we can write the density of  $(Z, Y)$  as

$$p(z, y, \theta, \tau) = g(\tau(y) - \theta z)\tau'(y),$$

with respect to  $H \times$  Lebesgue measure, where  $H$  is a known probability measure on  $Z$  and  $\tau \in \mathbf{T}$ .

We can construct  $GM$ -estimates of  $\theta$  easily once we observe that for fixed  $\theta$ , the model is a nonparametric group model under  $\tau$ . That is, if  $(Z, Y) \sim P \in \mathbf{P}$ ,

$$(Z, \tau(Y)) \sim P\tau^{-1} \in \mathbf{P} \quad \text{for all } \tau \in \mathbf{T},$$

and  $\tau^{-1}$  is identified with the map  $(z, y) \rightarrow (z, \tau^{-1}(y))$ . Let  $R_i$  be the rank of  $Y_i$  among  $Y_1, \dots, Y_n$ . Then  $\{(Z_i, R_i), 1 \leq i \leq n\}$  is a maximal invariant in the sense of Lehmann (1986, section 6.2). Thus, the joint distribution of the  $(Z_i, R_i)$  depends on  $\theta$  only and not on  $\tau$ . Unfortunately, the joint density is very complicated. By Hoeffding's formula

$$(28) \quad \frac{P_\theta(R_i = r_i, 1 \leq i \leq n \mid Z_i = z_i, 1 \leq i \leq n)}{P_{\theta_0}(R_i = r_i, 1 \leq i \leq n \mid Z_i = z_i, 1 \leq i \leq n)} \\ = E_{\theta_0} \left( \prod_{i=1}^n \frac{g(T_i - \theta z_i)}{g(T_i - \theta_0 z_i)} \mid R_i = r_i, 1 \leq i \leq n \right),$$

where  $T_i = \tau(Y_i)$  so that  $R_i$  is also the rank of  $T_i$  among  $T_1, \dots, T_n$ . As far as we know, the conditional likelihood ratio of the ranks has a simple expression only if  $g$  is the density of  $a \log V$  where  $V$  has an exponential distribution. That is just the (uncensored) Cox model. The estimate that we obtain by maximizing the conditional likelihood of the ranks is the Cox estimate discussed in example 7.4.3 below, which is efficient. In general, obtaining efficient estimates in these models is not trivial. However, simple and surprisingly good estimates have been obtained in this context by Cuzick (1988). See also Dabrowska and Doksum (1988), Tsiatis (1990). We follow Cuzick's treatment. In particular, consider linear functions of the ranks,

$$(29) \quad W_n(\theta) \equiv n^{-1} \sum_{i=1}^n Z_i \psi \left( F_{n\theta}^{-1} \left( \frac{R_i}{n+1} \right) - \theta Z_i \right).$$

Since  $W_n(\theta)$  does not depend on  $\tau$  we may without loss of generality in our subsequent analysis as well assume that  $\tau$  is the identity. Then

$$(30) \quad W_n(\theta) = n^{-1} \sum_{i=1}^n Z_i \psi \left( F_{n\theta}^{-1} \left( \frac{n}{n+1} F_n(T_i) \right) - \theta Z_i \right) \\ = \iint z \psi \left( F_{n\theta}^{-1} \left( \frac{n}{n+1} F_n(t) \right) - \theta z \right) dP_n(z, t),$$

where  $F_n$  is the empirical marginal distribution of  $T_1, \dots, T_n$ ,  $P_n$  the empirical distribution of  $(Z_1, T_1), \dots, (Z_n, T_n)$ , and

$$(31) \quad F_{n\theta}(t) \equiv n^{-1} \sum_{i=1}^n G_0(t - \theta Z_i).$$

To motivate this form, consider the choice

$$W_{n0}(\theta) \equiv n^{-1} \sum_{i=1}^n Z_i E_\theta (\psi(T_i - \theta Z_i) \mid (R_j, Z_j), 1 \leq j \leq n),$$

an estimating function which is nonlinear even in the ranks. If  $\psi = -g'/g$ , then  $W_{n0}$  is the derivative of the conditional log-likelihood of the ranks, as can be seen by differentiation of (28) at  $\theta = \theta_0$ . But

$$W_{n0}(\theta) = n^{-1} \sum_{i=1}^n Z_{D_i} E_{\theta}(\psi(T_{(i)} - \theta Z_{D_i}) | (R_j, Z_j), 1 \leq j \leq n),$$

where  $T_{(1)} < \dots < T_{(n)}$  are the order statistics of the  $T_i$  and  $T_{D_i} \equiv T_{(i)}$ . Then  $T_{(i)} \approx F_{\theta}^{-1}(i/(n+1))$  by standard theory, since  $F_{n\theta}$ , given by (31), is close to the marginal distribution  $F_{\theta}(t) \equiv \int G(t - \theta z) dH(z)$  of  $T$ . So if  $\theta$  is true

$$W_{n0}(\theta) \approx W_n(\theta).$$

In fact the approximation is not valid to order  $n^{-1/2}$ ; see Bickel (1986).

However, we proceed with (29) where  $\psi$  satisfies

$$(32) \quad \int \psi(t) g(t) dt = 0.$$

Under suitable conditions, we shall be able to apply Theorem 3. Here is a simplified and strengthened version of Cuzick's (1988) conditions to do this.

- (T1)  $\psi$  is twice continuously differentiable with bounded first and second derivatives.
- (T2)  $Z$  has bounded support.
- (T3)  $[G^{-1}]'(u) \leq C_1 [u(1-u)]^{-\alpha-1}$  for  $0 < \alpha < 1/2$ .
- (T4)  $|(g'/g)(t)| \leq C_2 |t|^{\beta}$  for  $0 < \beta < 1/\alpha$ .
- (T5)  $E\psi'(\epsilon)E(Z^2) - E(E(Z|T)E(\psi'(\epsilon)Z|T)) \neq 0$ .

Conditions (T3) and (T4) are tail conditions on  $G$ , while (T5) is needed for non-singularity of  $\dot{W}$ . It is easy to see that (T3) and (T4) are satisfied if  $g$  is a Gaussian density, and that (T5) holds if  $\psi(t) = t$ .

**Proposition 1.** Suppose (T1) through (T5) hold. Then  $W_n(\theta) = 0$  has a solution which is  $\sqrt{n}$ -consistent for estimating  $\theta(P)$  and asymptotically linear. Moreover, this root can be found by a Newton-Raphson algorithm that starts at a  $\sqrt{n}$ -consistent estimator.

$\sqrt{n}$ -consistent estimators can be constructed in various ways. For example, if  $g$  is positive and  $Z$  is 0 or 1 a.s., the two-sample model, let  $C_j$  be the conditional distribution of  $Y$  given  $Z = j$ , that is

$$C_j(y) = P(\tau^{-1}(T) \leq y | Z = j).$$

Then

$$G^{-1}(C_0(y)) = \tau(y), \quad G^{-1}(C_1(y)) = \tau(y) - \theta,$$

and consequently

$$\theta = G^{-1}C_0C_1^{-1}(\frac{1}{2}) - G^{-1}(\frac{1}{2}).$$

Substituting  $C_{nj}$ , the empirical distribution of the  $T_i$  with  $Z_i = j$ , for  $C_j$ ,  $j = 0, 1$ , we obtain a  $\sqrt{n}$ -consistent estimate of  $\theta$ .

For the proof of proposition 1 we need the following lemma.

**Lemma 1.**

- A. (T3) and (T4) imply  $E_G | \varepsilon |^\beta < \infty$ .  
 B. (T2) and (T3) imply that there exist  $c > 0$  and  $K > 0$  such that for all  $\theta$  with  $|\theta| \leq K$  and all  $u \in (0, 1)$

$$(33) \quad | [F_{n\theta}^{-1}]'(u) | \leq c [u(1-u)]^{-\alpha-1}.$$

**Proof.** To prove A we assume without loss of generality that  $G(0) = \frac{1}{2}$ . It suffices to show then

$$\int_0^{1/2} |G^{-1}(u)|^\beta du < \infty.$$

But the left-hand side equals

$$\begin{aligned} \int_0^{1/2} | \int_u^{1/2} [G^{-1}]'(v) dv |^\beta du &\leq c \int_0^{1/2} | \int_u^{1/2} v^{-\alpha-1} dv |^\beta du \\ &\leq c \int_0^{1/2} (u^{-\alpha} / \alpha)^\beta du < \infty, \end{aligned}$$

where  $c > 0$  and  $\alpha\beta < 1$ . For B, we first note that (T3) is equivalent to

$$c_1 g(t) \geq [G(t)(1-G(t))]^{\alpha+1},$$

which implies, for  $t \leq 0$ ,  $|\theta Z| \leq a$  a.s.

$$\begin{aligned} c_1 f_\theta(t) &= c_1 \int g(t - \theta z) dH(z) \\ &\geq \int [G(t - \theta z)(1 - G(a))]^{\alpha+1} dH(z) \\ &\geq (1 - G(a))^{\alpha+1} [F_\theta(t)]^{\alpha+1}. \end{aligned}$$

Note that  $H$  can be replaced by its empirical here. Together with the same argument for  $t \geq 0$  this yields

$$c f_{n\theta}(t) \geq [F_{n\theta}(t)(1 - F_{n\theta}(t))]^{\alpha+1},$$

which is equivalent to (33). □

**Proof of proposition 1.** We begin by introducing the appropriate  $W$  function. If  $P$  is an arbitrary joint distribution for  $(Z, T)$ , write

$$(a) \quad W(\theta, P) \equiv \iint z \psi(F_\theta^{-1}(F(t)) - \theta z) dP(z, t),$$

where  $F$  is the marginal distribution of  $T$  and  $H$  is the marginal distribution of  $Z$  under  $P$ . If  $F$  is absolutely continuous,  $Z$  is bounded a.s.,  $\psi$  satisfies (T1) and  $g$  satisfies the conditions of lemma 1, then  $W(\theta, P)$  is well defined. To see this note that (T1) implies that  $\psi$  can be bounded by linear functions. Consequently,

$$W(\theta, P) = O(E_P | F_\theta^{-1}(F(T)) |).$$

But



$$E_P |F_\theta^{-1}(F(T))| = \int_0^1 |F_\theta^{-1}(u)| du = E_{F_\theta} |T| \leq E_G |\varepsilon| + |\theta| E_H |Z| < \infty,$$

by lemma 1. Note that if  $P = P_{(\theta, \tau)} \in \mathbf{P}$ , and  $Q_\theta$  is the corresponding distribution of  $(Z_1, T_1)$ , then

$$W(\theta, Q_\theta) = \iint z \psi(t - \theta z) dQ_\theta(z, t) = 0$$

in view of (a). Cuzick (1988) argues conditionally on  $Z_1, \dots, Z_n$ , using results of Van Zuijlen (1978), to show that, if  $P = P_{(\theta, \tau)}$ , the natural expansion

$$(b) \quad W_n(\theta) - W(\theta, P) = \iint z \psi(t - \theta z) d(P_n - P)(z, t) + \iint z \psi'(t - \theta z) f_{n\theta}^{-1}(t) \cdot (F_n(t) - F_{n\theta}(t)) dP(z, t) + o_p(n^{-1/2})$$

is valid. This, of course, may be rewritten in the form of (GM2).

Under (T1) and (T2) it can be verified easily that  $\dot{W}(P) = -(\partial/\partial\theta)W(\theta(P), P)$  equals the left side of (T5). Consequently, (T5) implies (GM3).

We note that

$$(c) \quad \frac{\partial}{\partial\theta} F_{n\theta}^{-1}(v) = \frac{\sum_{i=1}^n Z_i g(F_{n\theta}^{-1}(v) - \theta Z_i)}{\sum_{i=1}^n g(F_{n\theta}^{-1}(v) - \theta Z_i)},$$

and introduce the notation

$$\frac{\partial}{\partial\theta} F_{n\theta}^{-1}\left(\frac{R_i}{n+1}\right) = \bar{Z}_i.$$

We have

$$(d) \quad \frac{\partial}{\partial\theta} W_n(\theta) = -n^{-1} \sum_{i=1}^n Z_i \psi'(F_{n\theta}^{-1}\left(\frac{R_i}{n+1}\right) - \theta Z_i)(Z_i - \bar{Z}_i).$$

In view of  $R_i = n F_n(T_i)$ , (T1), (T2), and lemma 1.B, we can replace  $\frac{n}{n+1} F_n(T_i)$  by  $F_{n\theta}(T_i)$  in (d) with error  $o_p(1)$  for  $\varepsilon n \leq i \leq (1 - \varepsilon)n$ ,  $\varepsilon > 0$ . Since  $Z_i, \bar{Z}_i, \psi'$  are bounded, the remaining part of (d) can be bounded uniformly by  $M\varepsilon$  for  $M$  independent of  $\theta, n$ . We deduce that

$$(e) \quad \frac{\partial}{\partial\theta} W_n(\theta_0) = -n^{-1} \sum_{i=1}^n Z_i \psi'(T_i - \theta_0 Z_i)(Z_i - Z_i^*(\theta_0)) + o_p(1)$$

where

$$Z_i^*(\theta) \equiv \frac{\sum_{j=1}^n Z_j g(T_i - \theta Z_j)}{\sum_{j=1}^n g(T_i - \theta Z_j)}.$$

But

$$\begin{aligned}
 (f) \quad & E \left\{ n^{-1} \sum_{i=1}^n Z_i \psi'(T_i - \theta Z_i) (Z_i - Z_i^*(\theta)) \mid Z_1, \dots, Z_n \right\} \\
 &= n^{-1} \sum_{i=1}^n Z_i^2 \int \psi'(t) g(t) dt - \int \frac{A_n(t) B_n(t)}{C_n(t)} dt
 \end{aligned}$$

where

$$A_n(t) = n^{-1} \sum_{i=1}^n Z_i \psi'(t - \theta Z_i) g(t - \theta Z_i),$$

$$B_n(t) = n^{-1} \sum_{i=1}^n Z_i g(t - \theta Z_i),$$

$$C_n(t) = n^{-1} \sum_{i=1}^n g(t - \theta Z_i).$$

For  $|\theta Z| \leq a$  a.s. we have

$$|f_{n\theta}(t) - g(t)| \leq \frac{1}{n} \sum_{i=1}^n \left| \int_0^{\theta Z_i} g'(t-s) ds \right| \leq \int_{-a}^a |g'(t-s)| ds.$$

Since

$$\int \int_{-a}^a |g'(t-s)| ds dt = 2a \int |g'(t)| dt < \infty$$

by (T4) and lemma 1.A, we have

$$f_{n\theta}(t) = \frac{1}{n} \sum_{i=1}^n g(t - \theta Z_i) \leq M(t),$$

with  $\int M(t) dt < \infty$ . We can use bounded convergence and the law of large numbers to get that (f) converges with probability 1 to

$$E\psi'(\epsilon)E(Z^2) - E(E(Z|T)E(\psi'(\epsilon)Z|T)) = -\frac{\partial}{\partial \theta} W(\theta_0, P).$$

Since, by (T1) and (T2) again,

$$\text{Var} \left\{ n^{-1} \sum_{i=1}^n Z_i \psi'(T_i - \theta Z_i) (Z_i - Z_i^*(\theta)) \mid Z_1, \dots, Z_n \right\} = O_p(n^{-1}),$$

we obtain (S1). Further,

$$\begin{aligned}
 \frac{\partial}{\partial \theta} Z_i^*(\theta) &= \frac{-\sum_{j=1}^n Z_j^2 g'(T_i - \theta Z_j)}{\sum_{j=1}^n g(T_i - \theta Z_j)} \\
 &+ \frac{\sum_{j=1}^n Z_j g(T_i - \theta Z_j) \sum_{k=1}^n Z_k g'(T_i - \theta Z_j)}{[\sum_{j=1}^n g(T_i - \theta Z_j)]^2}.
 \end{aligned}$$

Using (T2) and (T4),

$$\left| \frac{\partial}{\partial \theta} Z_i^*(\theta) \right| \leq M_1 (|\epsilon_i|^\beta + |\theta|^\beta + 1).$$

It is then not hard to bound

$$\left| \frac{\partial^2}{\partial \theta^2} W_n(\theta) \right| \leq M_2 \left( n^{-1} \sum_{i=1}^n |\varepsilon_i|^\beta + |\theta|^\beta + 1 \right),$$

and lemma 1.A yields (S2) and hence (U) for every sequence  $\tilde{\varepsilon}_n \downarrow 0$ ,  $\tilde{\varepsilon}_n n^{1/2} \rightarrow \infty$ . Application of Theorem 3 completes the proof of the proposition. Regularity follows by checking (17) and (18) as usual.  $\square$

### 7.4 GMC-AND GM-ESTIMATES CORRESPONDING TO CONVEX $D$

A key hypothesis of the master theorem and its corollaries is consistency or  $\sqrt{n}$ -consistency of the estimate considered. As we indicated in the introduction, when estimates are defined algorithmically this point is often ignored. The basic technique for establishing consistency for  $MC$ -estimates is due to Wald (1949). This technique and a parallel approach for  $M$ -estimates has been developed by Huber (1967) and Perlman (1972). Extensions of the technique to deal with abstract-valued parameters are due to Kiefer and Wolfowitz (1956), Bahadur (1967), Geman and Hwang (1982), and others. We shall not go into this work here. Instead we focus on the nice situation that the estimates of interest are not only AGM but also, possibly after reparametrization, generalized minimum contrast in the sense of section 1 corresponding to  $D$  convex.

We introduce two classes of functions,  $\mathbf{D}$  and  $\mathbf{W}$ . Let  $N$  be an open convex subset of  $R^m$ . We define  $\mathbf{D}$  to be the class of all functions  $D : N \rightarrow R$  such that  $D$  is convex. Let  $\mathbf{D}_0 \subset \mathbf{D}$  be the subclass of all functions  $D$  which achieve a unique minimum on  $N$ . Consequently, if  $\partial N$  is the boundary of  $N$  in  $\bar{R}^m$ , compactified Euclidean space, then for  $D \in \mathbf{D}_0$

$$(1) \quad \liminf_{v \rightarrow \partial N} D(v, P) > \inf\{D(v, P) : v \in N\}.$$

Suppose we are given  $D : N \times Q \rightarrow R$ , and fix  $Q \in Q$ . Then, with  $D \equiv D(\cdot, Q)$ , the map  $\lambda \rightarrow D(t + \lambda u)$  is convex. If  $D$  is differentiable and  $W \equiv W(\cdot, Q) \equiv \dot{D}(\cdot, Q)$ , then  $\lambda \rightarrow u^T W(t + \lambda u)$  is monotone nondecreasing. We therefore let  $\mathbf{W}$  be the class of all functions  $W : N \rightarrow R^m$  such that for all  $u \in R^m, t \in N$ , the maps  $\lambda \rightarrow u^T W(t + \lambda u)$  from  $\{\lambda \in R : t + \lambda u \in N\}$  to  $R$  are monotone nondecreasing. Let  $\mathbf{W}_0 \subset \mathbf{W}$  be the subclass of functions which have a unique root in  $N$ .

**Theorem 1.** Suppose that  $D : N \times Q \rightarrow R$  belongs to  $\mathbf{D}$  for all  $Q \in Q$  and

$$(2) \quad \sup\{|D(v, P_n) - D(v, P)| : v \in K\} = o_p(1)$$

for all compact  $K \subset N$ . Suppose  $D(\cdot, P) \in \mathbf{D}_0$ , is differentiable and let  $W = \dot{D}$ . Write  $D_n(v)$  for  $D(v, P_n)$ .

Let  $\hat{v}_n$  be the GM-estimate corresponding to  $W = \dot{D}$ . Then

- A.  $W(v, P) = 0$  has a unique solution  $v(P)$ .
- B. With probability tending to 1,  $\hat{v}_n$  exists and minimizes  $D_n(v)$ .
- C.  $\hat{v}_n$  is a consistent estimate of  $v(P)$ .

**Proof.** Part A is standard. For part B, note that by (1), for all  $\delta > 0$ ,

$$\min\{D(v, P) : |v - v(P)| = \delta, v \in N\} > D(v(P), P).$$

Therefore, by (2),

$$P(\min\{D_n(v) : |v - v(P)| = \delta, v \in N\} > D_n(v(P))) \rightarrow 1.$$

Then B follows since

$$(a) \quad P(\min\{D_n(v) : |v - v(P)| \leq \delta\} = \min_N D_n(v)) \rightarrow 1$$

by standard properties of convex functions. Since  $\delta$  in (a) is arbitrary, claim C follows.  $\square$

**Notes.** 1. The *GM*-estimate exists if and only if  $\liminf_{v \rightarrow \partial N} D_n(v) > \min_N D_n(v)$ .

2. If  $D_n$  is strictly convex and the *GM*-estimate exists, it is unique.

If we specialize to *MC*- and *M*-estimation where  $D(v, P) = \int \rho(x, v) dP(x)$  and

$$(3) \quad W(v, P) = \int \psi(x, v) dP(x)$$

for  $\psi = \nabla_v \rho$ , then the hypotheses of theorem 1 are implied by

$$(4) \quad \rho(x, \cdot) \text{ strictly convex}$$

and

$$(5) \quad \liminf_{v \rightarrow \partial N} \int \rho(x, v) dP(x) > \inf_N \int \rho(x, v) dP(x),$$

$$(6) \quad E \sup_K |\rho(X, v) - \rho(X, v(P))| < \infty.$$

We apply these remarks to:

**Example 1. Exponential family.**

Let  $P_v \in \mathbf{P}$  where  $\mathbf{P}$  is an  $m$ -dimensional exponential family with densities (with respect to  $\mu$ ) given by

$$p(x, v) = \exp\{v^T S(x) - b(v)\},$$

$$b(v) = \log \int \exp\{v^T S(x)\} d\mu(x), \quad v \in N.$$

Here  $N \equiv \{v : \int \exp\{v^T S(x)\} d\mu(x) < \infty\}$ , the natural parameter space. Suppose that  $N$  is open and

$$(7) \quad \mu\{v^T S(x) \neq c\} > 0 \quad \text{for all } |v| = 1, \text{ and all } c \in R.$$

We show that, if  $P \in \mathbf{P}$ ,  $\rho(x, v) \equiv -\log p(x, v)$  satisfies (4)–(6) so that the conclusion of theorem 1 holds.

By Lehmann (1983, equation (1.4.12) page 30),

$$(8) \quad \left[ \frac{\partial^2}{\partial v_i \partial v_j} \rho(\cdot, v) \right] = \left[ \frac{\partial^2}{\partial v_i \partial v_j} b(v) \right]$$

$$= E_v [S(X) - E_v S(X)] [S(X) - E_v S(X)]^T$$

which is positive definite in view of (7). Then (4) follows by standard theory. Since  $b$  is analytic, it is easy to see that

$$\dot{D}(v, P) = \int \psi(x, v) dP(x),$$

where

$$(9) \quad \psi(x, v) = -S(x) + \dot{b}(v).$$

If  $P \equiv P_{v_0} \in \mathbf{P}$ ,

$$W(v_0, P) = -E_{v_0} S(X) + \dot{b}(v_0) = 0.$$

So  $D \in \mathbf{D}_0$  and, in fact, (5) follows from the strict convexity of  $D(v, P)$ . Finally (6) follows since

$$E \sup_K |\psi(X, v)| \leq E |S(X)| + \sup_K |\dot{b}(v)| < \infty.$$

Note that without additional conditions (see Barndorff-Nielsen (1978), for example) existence of the MLE is only guaranteed with probability tending to 1. For example, suppose  $\mathbf{P}$  is the binomial(1,  $p$ ) family,  $0 < p < 1$ , reparametrized by  $p = e^v(1 + e^v)^{-1}$ ,  $v \in R$ . Then, the MLE does not exist if  $\bar{X} = 0$  or 1, but this happens with probability tending to 0 if  $P \in \mathbf{P}$ .

Now condition (D2) of section 7.2 follows from (9) and the analyticity of  $b$ , and (D3) follows from (8). Finally, since by proposition 1, on  $\mathbf{P}$ ,  $\hat{v}_n$  is consistent, we conclude that  $\hat{v}_n$  is asymptotically linear and, of course, efficient.

Suppose  $P \notin \mathbf{P}$  but  $S(X) \in L_2(P)$  and

$$\dot{b}(v) = E_P S(X)$$

has a solution  $v(P) \in N$ . Then it is easy to check as above that  $\hat{v}_n$  is an asymptotically linear estimate of  $v(P)$  with influence function,  $\overset{\dots-1}{b}(v(P))(S(x) - \dot{b}(v(P)))$ , where  $\ddot{b}(v) = [(\partial^2/\partial v_i \partial v_j) b(v)]$ . Compare Berk (1972).  $\square$

Having  $D \in \mathbf{D}$  and  $D(\cdot, P) \in \mathbf{D}_0$  not only guarantees existence, unicity, and consistency of GMC-estimates, but also yields asymptotic linearity more easily. Although the following key hypotheses involve only  $W: N \times \mathbf{Q} \rightarrow R^m$ , there is typically a function  $D$  in the background.

Here are the hypotheses we will use: Assume that (GM0) and (GM3) hold, and let  $v_0 \equiv v(P)$ .

(C1)  $P(W(\cdot, P_n) \in W_0) \rightarrow 1$  as  $n \rightarrow \infty$ .

(C2) For each fixed  $\tau \in R^m$

$$\sqrt{n} \{W(v_0 + n^{-1/2}\tau, P_n) - W(v_0, P_n)\} = \dot{W}(P)\tau + o_p(1).$$

(C3)  $W(v_0, P_n) = \int \psi(x, P) dP_n(x) + o_p(n^{-1/2})$  where  $|\psi| \in L_2(P)$ ,  $\int \psi(x, P) dP(x) = 0$ .

We will also use a strengthening of (C1):

(C1')  $W(\cdot, Q) \in W_0$  for all  $Q \in M_0$  where  $M_0 \supset Q \cup \{\text{all realizations of the empirical measures } P_n, n \geq 1\}$ .

The following useful result is due to Brown (1985) and Ritov (1987).

**Theorem 2.** Suppose that  $X_1, \dots, X_n$  are i.i.d.  $P \in Q$ . Suppose that (GM0), (GM3), (C1'), (C2), and (C3) hold. Then  $\hat{v}_n$  corresponding to  $W_n(v) \equiv W(v, P_n)$  is uniquely defined and asymptotically linear with influence function  $-\dot{W}^{-1}(P)\psi(\cdot, P)$ . If the hypothesis (C1') is replaced by (C1), then the AGM-estimate  $\hat{v}_n$  exists and asymptotic linearity continues to hold.

**Proof.** See theorem A.10.3. □

**Example 2. M-estimates in the linear model.**

Consider estimation of  $\theta = (v^T, \sigma)^T$  in the semiparametric linear model  $P$  described by

$$(10) \quad Y = v^T Z + \sigma \varepsilon,$$

where  $\varepsilon \sim G$  is independent of  $Z \sim H$  on  $R^m$  based on observation of ( $n$  i.i.d. copies of)  $X \equiv (Z^T, Y)$ . This is just as in example 4.2.2, but with the roles of  $(G, g)$  and  $(Q, q)$  interchanged.

We view  $P$  as a subset of  $Q = M$ , the collection of all distributions on  $R^{m+1}$ .

Suppose that  $\rho: R \rightarrow R^+$  is twice continuously differentiable, and strictly convex and positive with minimum at 0. For  $\theta = (v^T, \sigma)^T \in R^m \times R^+$ ,  $Q \in Q$  on  $R^{m+1}$ , define

$$(11) \quad D(\theta, Q) = \int \rho\left(\frac{y - v^T z}{\sigma}\right) dQ(z, y) + \log \sigma \\ \equiv E_Q \rho(\varepsilon(\theta)) + \log \sigma.$$

If the density  $g$  is twice differentiable, symmetric about 0 and strongly unimodal, then  $\rho = -\log g$  satisfies our hypotheses. The minimum contrast estimates are then just the maximum likelihood estimates for  $v$  and  $\sigma$  in the linear regression model of example 4.2.2 with  $\varepsilon \sim g$  known.

In order to investigate the behavior of the minimum contrast estimate based on  $D$ , we introduce a new parametrization  $\tau = (\delta^T, \gamma)^T = (\sigma^{-1}v^T, \sigma^{-1})^T$ , and define accordingly

$$D_1(\tau, Q) \equiv D(\theta(\tau), Q) \\ = \int \rho(\gamma y - \delta^T z) dQ(y, z) - \log \gamma$$

where  $\theta(\tau) = (\gamma^{-1}\delta^T, \gamma^{-1})^T$ . Now, it can be shown that if

$$(12) \quad J(\theta, Q) \equiv E_Q [\rho''(\varepsilon(\theta))(|Z|^2 + \varepsilon^2(\theta)) \\ + |\rho'(\varepsilon(\theta))| (|\varepsilon(\theta)| + |Z|)]$$

is finite and continuous as a function of  $\theta$ , then  $D_1(\tau, Q)$  can be differentiated twice under the integral sign. Therefore, if  $\tilde{\varepsilon}(\tau) \equiv \varepsilon(\theta(\tau))$ ,

$$W^T(\tau, Q) = \nabla_{\tau} D_1(\tau, Q) = E_Q [(-Z^T, Y)\rho'(\tilde{\epsilon}(\tau))] - \left(0, \frac{1}{\gamma}\right)$$

and

$$\dot{W}(\tau, Q) = E_Q [(-Z^T, Y)^T(-Z^T, Y)\rho''(\tilde{\epsilon}(\tau))] + \begin{pmatrix} 0 & 0 \\ 0 & \gamma^{-2} \end{pmatrix}.$$

Hence  $D_1(\cdot, Q)$  is strictly convex if for all  $\tau \in R^m \times R^+$ ,  $(a^T, b)^T \in R^{m+1}$ ,  $(a^T, b) \neq 0$

$$E_Q (|bY - a^T Z|^2 \rho''(\tilde{\epsilon}(\tau))) + b^2 \gamma^{-2} > 0,$$

or, equivalently, if and only if, for all  $(a, b^T) \neq 0$ ,  $Q(a^T Z \neq bY) + b^2 > 0$ , or, for all  $a \neq 0$ ,

$$(13) \quad Q(a^T Z = 0) < 1.$$

Moreover, if, in addition, for all  $a \neq 0$ ,

$$(14) \quad Q(Y = a^T Z) < 1,$$

then  $D_1(\tau, Q) \rightarrow \infty$  as  $|\tau| \rightarrow \infty$ . Evidently,  $D_1(\tau, Q) \rightarrow \infty$  as  $\gamma \rightarrow 0$ . Together, conditions (13) and (14) imply that  $D_1(\cdot, Q) \in \mathbf{D}$ . A standard Taylor expansion shows that if (12) holds, then condition (C2) of theorem 2 holds. Hence, if (12)–(14) are satisfied for  $Q = P$ , then  $(\hat{v}_n, \hat{\sigma})$  is asymptotically linear. If  $P \in \mathbf{P}$ , the semiparametric linear model of example 4.2.2, then the parameter  $v(P)$  defined by minimizing (11) equals  $v$  for all  $G, H, \Delta, \sigma$ . Therefore  $\hat{v}_n$  is  $\sqrt{n}$ -consistent for  $v$  in this model. In example 4.2.2,  $\sigma$  is not identifiable, so all that we can say of  $\hat{\sigma}_n$  is that it is a  $\sqrt{n}$ -consistent estimate of  $\sigma(P)$  defined by (11).

Thus, for example, if  $g$  is the logistic density, the corresponding

$$\rho(t) = 2 \log(1 + e^{-t}) + t$$

yields asymptotically linear equivariant estimates of  $v$  in the linear model. □

**Example 3. Linear regression with right censoring.**

The model is as given in example 4.6.4. We observe  $X = (Z, Y \wedge C, 1_{[Y \leq C]}) \equiv (Z, V, \Delta)$  where  $\epsilon \equiv Y - v^T(P)Z$  is independent of  $(Z, C)$  and its distribution is otherwise arbitrary. Let,  $\hat{\epsilon}(v) \equiv V - v^T Z$ . Tsiatis (1990) and Ritov and Fygenon (1990) suggested estimating  $v$  using the GMC-estimator based on

$$D(v, P) \equiv E[\Delta_1(\hat{\epsilon}_2(v) - \hat{\epsilon}_1(v)) 1_{[\hat{\epsilon}_2(v) \geq \hat{\epsilon}_1(v)]}],$$

where, for  $i = 1, \dots, n$ ,  $\epsilon_i, \hat{\epsilon}_i(v), \Delta_i$  correspond to  $X_i$ , which are i.i.d.  $P$ . Then  $D \in \mathbf{D}$  with gradient given a.e. by

$$W(v, P) = E\{(Z_1 - Z_2)\Delta_1 1_{[\hat{\epsilon}_2(v) \geq \hat{\epsilon}_1(v)]}\}.$$

Minimization of  $D(v, P_n)$  leads to the AGM-estimate based on

$$W(v, P_n) = \frac{1}{n^2} \sum_i \sum_j (Z_i - Z_j) \Delta_i 1_{[\hat{\epsilon}_i(v) \geq \hat{\epsilon}_j(v)]}.$$

Note that  $W(v, P_n)$  is neither continuous nor strictly monotone, so that  $W(v, P_n) = 0$  may have no solution, but  $\hat{v}_n$  such that  $W(\hat{v}_n, P_n) = O_p(n^{-2})$  exists.

**Proposition 1.** Suppose that  $|Z| \leq M$ ,  $f$  the density of  $\varepsilon$  has finite Fisher information,  $P(Y \leq C) > 0$ , and the distribution of  $Z$  given  $Y \leq C$  is not concentrated on a hyperplane of  $R^m$ . Then:

- A.  $W(v(P), P) = 0$   
 B.  $\hat{v}_n$  is asymptotically linear.

**Proof.** We verify the conditions of theorem 2. It is not too difficult to verify (C1)–(C3) by standard empirical process arguments as in example 7.3.2, or alternatively  $U$  statistic theory as in Serfling (1980), and we leave this to the reader. We do, however, verify (GM0) which is just part A and (GM3). Note that without loss of generality we may take  $\gamma(P) = 0$ . Since if  $X^* = (Z, Y - v^T(P)Z, \Delta)$  with corresponding  $P^*$ , then

$$(a) \quad W(v, P) = W(v - v(P), P^*).$$

In the case  $v(P) = 0$  we can write

$$(b) \quad W(v, P) = E\{(Z_1 - Z_2) 1_{[\varepsilon_1 \leq C_1]} 1_{[\varepsilon_1 - v_2 \leq v^T(Z_1 - Z_2)]}\}.$$

Then

$$(c) \quad W(0, P) = E(Z_1 - Z_2) 1_{[\varepsilon_1 \leq C_1 \wedge C_2 \wedge \varepsilon_2]}.$$

But, given  $Z_1, Z_2, C_1, C_2$ , the indicators  $1_{[\varepsilon_1 \leq C_1 \wedge C_2 \wedge \varepsilon_2]}$  and  $1_{[\varepsilon_2 \leq C_1 \wedge C_2 \wedge \varepsilon_1]}$  have the same distribution. Part A follows by interchanging indices in  $W(0, P)$ .

Further, by conditioning on  $Z_1, Z_2, C_1, C_2, \varepsilon_2$ , we get from (b),

$$(d) \quad W(v, P) = E\{(Z_1 - Z_2)F(C_1 \wedge V_2 + v^T(Z_1 - Z_2))\},$$

where  $F$  is the distribution function of  $\varepsilon_1$ . Then

$$\begin{aligned} & u^T \frac{\partial}{\partial \lambda} W(\lambda u, P) |_{\lambda=0} \\ &= u^T E\{(Z_1 - Z_2)(Z_1 - Z_2)^T F(V_2) 1_{[V_2 < C_1]}\} u \\ &\geq u^T E\{(Z_1 - Z_2)(Z_1 - Z_2)^T \int_{-\infty}^{C_1 \wedge C_2} F^2(t) dt\} u > 0 \end{aligned}$$

under our assumptions. Thus  $\dot{W}(P)$  is nonsingular, (GM3) holds, and part B follows.  $\square$

#### Example 4. The Cox estimate.

Suppose that  $X = (Z_{m \times 1}, Y, \Delta)$  where  $Y = T \wedge C$ ,  $\Delta = 1_{[T \geq C]}$ , and, given  $Z$ ,  $T$  and  $C$  are independent. Assume also that the conditional distribution of  $T$  given  $Z$  is  $1 - \bar{G}^r$  where  $r = r(v^T z) = \exp(v^T z)$ ,  $\bar{G} = 1 - G$ . Then  $X \sim P$  belongs to the generalization  $\mathbf{P}$  of the Cox model, example 3.4.2, in which right censoring of  $T$  by  $C$  is allowed and the distribution of the censoring variable  $C$



can depend on the covariate. For simplicity in the sequel we take  $m = 1$ . The argument for  $m > 1$  is only notationally more complicated.

Our goal is to use theorem 2 to give asymptotic properties of the Cox partial likelihood estimator even when the model fails. Suppose that  $P$  is any distribution of  $(Z, Y, \Delta)$  on  $R \times R \times \{0, 1\}$  satisfying:

- (i)  $Z$  is not degenerate a.s.  $P$ ; i.e.,  $P(Z = z_0) \neq 1$  for all  $z_0$ .
- (ii)  $Z$  is bounded a.s.  $P$ ; i.e.,  $P(|Z| \leq C) = 1$ .

Let  $\mathbf{Q}$  denote the collection of all such distributions  $P$ . We suppose that  $X_1, \dots, X_n$  are i.i.d.  $P \in \mathbf{Q}$ ; note that  $P$  is not necessarily in the Cox model  $\mathbf{P} \subset \mathbf{Q}$ .

We define the Cox estimate of  $v$  using the notation of Tsiatis (1981). For  $Q \in \mathbf{Q}$ , let

$$S_j(t, v, Q) = \int z^j r(vz) 1_{[s \geq t]} dQ^{(12)}(z, s), \quad \text{for } j = 0, 1, 2,$$

where  $Q^{(12)}$  is the marginal distribution of  $(Z, Y)$ . Let

$$(15) \quad D_0(v, Q) \equiv -v \int z dQ^{(1)}(z) + \int \log S_0(t, v, Q) dQ^{(2)}(t)$$

for  $v \in N \equiv R$  where

$$Q^{(1)}(z) \equiv Q(Z \leq z, \Delta = 1)$$

and

$$Q^{(2)}(t) \equiv Q(Y \leq t, \Delta = 1).$$

It is easy to check that  $D_0(\cdot, Q) \in \mathbf{D}$ . Note that

$$\begin{aligned} D_0(v, P_n) &= -\frac{1}{n} v \sum_{i=1}^n \Delta_i Z_i + \frac{1}{n} \sum_{i=1}^n \Delta_i \log \left\{ \frac{1}{n} \sum_{j \in R_i} e^{vZ_j} \right\} \\ &= -\frac{1}{n} \sum_{i=1}^n \Delta_i \log \left\{ \frac{e^{vZ_i}}{n^{-1} \sum_{j \in R_i} e^{vZ_j}} \right\}, \end{aligned}$$

where the "risk set"  $R_i \equiv \{j : T_j \geq T_i\}$ . Hence the resulting minimum contrast estimate  $\bar{v}(P_n)$  is the Cox partial likelihood estimator of  $v$ . It solves

$$W_0(v, P_n) = 0$$

for

$$W_0(v, Q) = - \int z dQ^{(1)}(z) + \int \frac{S_1}{S_0}(t, v, Q) dQ^{(2)}(t).$$

Unfortunately, the resulting estimator using all the data is difficult to analyze because of the contributions of observations in the right tail. For discussions of the problems associated with the estimator based on all the data, see, e.g., Tsiatis (1981, page 95) (his  $T_0 =$  our  $\tau$  defined below), and Andersen and Gill (1982, pages 1109–1110) (their 1 = our  $\tau$ ). Although Andersen and Gill (1982) give some results for the estimator based on all the data, here we take the approach of

the majority of the literature by defining a modified estimator based on most, but not all, of the data as follows.

Let  $\tau_Y \equiv \inf \{s : P(Y \geq s) = 0\}$ , and suppose that  $\tau < \tau_Y$  is fixed. Thus  $P(Y \geq \tau) > 0$ . Define

$$(16) \quad D(v, Q) \equiv -v \int_0^\tau \int z dQ(z, t, 1) + \int_0^\tau \log S_0(t, v, Q) dQ^{(2)}(t).$$

Then we base our estimate on minimization of

$$D(v, P_n) \equiv -\frac{1}{n} \sum_{i=1}^n \Delta_i 1_{[Y_i \leq \tau]} \log \left\{ \frac{e^{vZ_i}}{n^{-1} \sum_{j \in R_i} e^{vZ_j}} \right\}.$$

The resulting minimum contrast estimate  $\hat{v}(P_n, \tau)$  is also commonly called the Cox estimate. It solves  $W(v, P_n, \tau) = 0$  for

$$(17) \quad W(v, Q) = -\int_0^\tau \int z dQ(z, t, 1) + \int_0^\tau \frac{S_1}{S_0}(t, v, Q) dQ^{(2)}(t).$$

We prove, for the GMC-estimate  $\hat{v}_n$  corresponding to  $D, W$  given by (16) and (17):

**Proposition 2.** Suppose conditions (i) and (ii) of this example hold for  $Q = P$ . If  $v(Q)$  minimizes (16), then

$$(18) \quad \dot{W}(v, Q) = \int_0^\tau \left\{ \frac{S_2}{S_0}(t, v, Q) - \left( \frac{S_1}{S_0} \right)^2(t, v, Q) \right\} dQ^{(2)}(t) > 0,$$

$$(19) \quad \begin{aligned} \psi(x, P) = & -\{\Delta z 1_{[y \leq \tau]} - E(\Delta Z 1_{[Y \leq \tau]})\} \\ & + \left\{ \Delta \frac{S_1}{S_0}(y) 1_{[y \leq \tau]} - E\left(\Delta \frac{S_1}{S_0}(Y) 1_{[Y \leq \tau]}\right) \right\} \\ & + e^{vz} \left\{ z \int_0^{Y \wedge \tau} \frac{1}{S_0}(s) dP^{(2)}(s) - \int_0^{Y \wedge \tau} \frac{S_1}{S_0^2}(s) dP^{(2)}(s) \right\} \end{aligned}$$

and  $\hat{v}_n$  is asymptotically linear with influence function  $-\dot{W}^{-1}(P)\psi(\cdot, P)$ . Further,  $\hat{v}_n$  is regular on  $\mathbf{P}$ .

**Proof.** By the Cauchy-Schwarz inequality we have  $S_1^2(t) \leq S_0(t)S_2(t)$  for all  $t$ , with strict inequality unless  $Z$  is constant, and hence (18). Also note that  $(S_1/S_0)(t, v, Q)$  varies from  $Q$ -essinf  $Z$  to  $Q$ -esssup  $Z$  as  $v$  varies from  $-\infty$  to  $+\infty$ . Hence, for  $Q \in \mathbf{Q}$ , we have  $D \equiv D(\cdot, Q) \in \mathbf{D}$ .

We verify the conditions of theorem 2 under the assumption  $P \in \mathbf{Q}$ . Note first that by (ii), with probability 1,

$$\left| \frac{\partial^2}{\partial v^2} W(v, P_n) \right| \leq 6C^3.$$

So to verify (C2) and (C3) of theorem 2 it suffices to show that

$$(20) \quad \dot{W}(v_0, P_n) = \dot{W}(v_0, P) + o_p(1) \equiv \dot{W}(P) + o_p(1),$$

and

$$(21) \quad W(v_0, P_n) = W(v_0, P) + \int \psi(x, P) dP_n + o_p(n^{-1/2}).$$

Let  $R(t, Q)$  be the integrand in  $\dot{W}(v, Q)$  given in (17) with  $v$  replaced by  $v_0 = v(P)$ , and define

$$\begin{aligned} \mathcal{E}_{nj}(t) &\equiv S_j(t, v_0, P_n), & S_j(t) &\equiv S_j(t, v_0, P), \\ A_j(t) &\equiv \sqrt{n}[\mathcal{E}_{nj}(t) - S_j(t)], & j &= 0, 1, 2. \end{aligned}$$

Using a standard tightness argument (for details see Tsiatis (1981), lemma A.3), we get

$$(22) \quad \sup_t |A_j(t)| = O_p(1) \quad \text{for } j = 0, 1.$$

Since  $P(Y \geq \tau) > 0$ , we have

$$(23) \quad \inf\{S_0(t, v, P) : t \leq \tau\} = S_0(\tau, v, P) \geq e^{-v|C|} P(Y \geq \tau) > 0.$$

Hence

$$(24) \quad \sup\{|R(t, P_n) - R(t, P)| : t \leq \tau\} = O_p(n^{-1/2}).$$

By using  $|R(t, P)| \leq 2C^2$ , the weak law of large numbers, and (23),

$$\begin{aligned} (25) \quad &|\dot{W}(v_0, P_n) - \dot{W}(v_0, P)| \\ &= \left| \int_0^\tau R(t, P_n) dP_n^{(2)}(t) - \int_0^\tau R(t, P) dP^{(2)}(t) \right| \\ &\leq \left| \int_0^\tau R(t, P) (dP_n^{(2)} - dP^{(2)})(t) \right| \\ &\quad + \sup\{|R(t, P_n) - R(t, P)| : t \leq \tau\} \\ &= o_p(1) \end{aligned}$$

so (20) holds. To prove (21), let

$$\begin{aligned} (26) \quad B_n(t) &\equiv \sqrt{n} \left\{ \frac{\mathcal{E}_{n1}}{\mathcal{E}_{n0}}(t) - \frac{S_1}{S_0}(t) \right\} - \frac{A_1(t)}{S_0(t)} + A_0(t) \frac{S_1}{S_0^2}(t) \\ &= \frac{A_0}{S_0}(t) \left\{ \frac{S_1}{S_0}(t) - \frac{\mathcal{E}_{n1}}{\mathcal{E}_{n0}}(t) \right\} \\ (27) \quad &= \frac{A_0}{S_0}(t) \left\{ \frac{n^{-1/2} A_0 S_1}{\mathcal{E}_{n0} S_0}(t) - \frac{n^{-1/2} A_1}{\mathcal{E}_{n0}}(t) \right\}. \end{aligned}$$

We claim that

$$(28) \quad \sup\{|B_n(t)| : t \leq \tau\} = o_p(1).$$

To see this, note that if  $Y_{(n)}$  is the largest uncensored observation, then

$$(29) \quad \sup_{t \leq Y_{(n)}} \frac{S_0(t)}{\mathcal{E}_{n0}(t)} = \sup_{t \leq Y_{(n)}} \frac{S_0(t, v_0, P)}{S_0(t, v_0, P_n)}$$

$$\leq e^{2\nu_0 c} \sup_{t \leq Y_0} \frac{P[t, \infty]}{IP_n[t, \infty]} = O_p(1)$$

by inequality (10.3.1), page 412, Shorack and Wellner (1986). Then (28) follows from (22), (23), (27), and (29).

Now we can write, using (26) and  $W(\nu_0, P) = 0$  at the last step,

$$\begin{aligned} W(\nu_0, IP_n) &= - \int_0^\tau \int z dIP_n(z, t, 1) + \int_0^\tau \frac{S_1}{S_0}(t, \nu_0, IP_n) dIP_n^{(2)}(t) \\ &= W(\nu_0, P) - \int_0^\tau \int z d(IP_n - P)(z, t, 1) \\ &\quad + \int_0^\tau \frac{S_1}{S_0}(t, \nu_0, P) d(IP_n^{(2)} - P^{(2)})(t) \\ &\quad + \int_0^\tau \left\{ \frac{S_{n1}}{S_{n0}}(t) - \frac{S_1}{S_0}(t) \right\} dIP_n^{(2)}(t) \\ (30) \quad &= - \int_0^\tau \int z d(IP_n - P)(z, t, 1) \\ &\quad + \int_0^\tau \frac{S_1}{S_0}(t, \nu_0, P) d(IP_n^{(2)} - P^{(2)})(t) \\ &\quad + n^{-1/2} \int_0^\tau \left\{ \frac{A_1}{S_0}(t) - \frac{S_1 A_0}{S_0^2}(t) \right\} dP^{(2)}(t) \\ &\quad + n^{-1/2} \int_0^\tau \left\{ \frac{A_1}{S_0}(t) - \frac{S_1 A_0}{S_0^2}(t) \right\} d(IP_n^{(2)} - P^{(2)})(t) \\ &\quad + n^{-1/2} \int_0^\tau B_n(t) dIP_n^{(2)}(t). \end{aligned}$$

The first three terms on the right side in (30) are linear, the fifth is  $o_p(n^{-1/2})$  by (28). The fourth term may be written as the  $V$ -statistic

$$(31) \quad V_n = n^{-2} \sum_{i,j} K(X_i, X_j),$$

where the kernel

$$\begin{aligned} K(x_1, x_2) &= (z_1 - \frac{S_1}{S_0}(y_2)) \frac{e^{\nu_0 z_1}}{S_0(y_2)} 1_{[y_1 \geq y_2]} \Delta_2 1_{[y_2 \leq \tau]} \\ &\quad - \int_0^\tau (z_1 - \frac{S_1}{S_0}(t)) \frac{e^{\nu_0 z_1}}{S_0(t)} 1_{[y_1 \geq t]} dP^{(2)}(t) \end{aligned}$$

is bounded in absolute value by  $4Ce^{\nu_0 c} / S_0(\tau)$  and satisfies  $E K(X_1, x_2) = E K(x_1, X_2) = 0$ . Consequently

$$E V_n^2 = O(n^{-2})$$

and  $V_n = O_p(n^{-1})$ . Hence, (20), (21) and the first part of the proposition follow. Regularity on the model  $P$  can be verified as usual by checking that  $v(P)$  is pathwise differentiable on  $P$  and  $\dot{v} = \psi$ .  $\square$

If  $P$  belongs to the uncensored Cox model, then

$$S_0(s) = \int_{-\infty}^{\infty} r(v_0z) \bar{G}^{r(v_0z)}(s) h(z) dz = \bar{G}(s) \frac{dP^{(2)}}{dG}(s),$$

so that (19) becomes, with  $\tau = \infty$ ,

$$(32) \quad \psi(x, P) = - \left( z - \frac{S_1}{S_0}(y) \right) + e^{yz} \left\{ -z \log \bar{G}(y) - \int_0^y \frac{S_1}{S_0}(s) \frac{dG}{G}(s) \right\}$$

in agreement with (3.4.34) and (3.4.44) (after noting the  $-$  sign in theorem 2).

To obtain efficiency both in the uncensored and censored models we need to use  $\tau = \tau_Y$  in (18) and (19). This can be achieved by employing the Andersen-Gill (1982) approach under their conditions. More generally we believe, though this remains to be shown, that this can be achieved by taking estimates  $\hat{v}_n(\tau_n)$  corresponding to (16) with  $\tau_n \rightarrow \infty$ . Proposition 2 has been extended to a large class of inefficient estimates of  $v$  by Sasieni (1989).  $\square$

### 7.5 ESTIMATION OF $P$ AND OTHER INFINITE-DIMENSIONAL PARAMETERS: METHODS, CONSISTENCY, AND RATES OF CONVERGENCE

We consider abstract parameters  $v$  taking their values in a subset  $A$  of an appropriate Banach space  $B$  as in section 5.2. The most important such parameter is  $v(P) = P$  where  $P$  is suitably identified with  $A$ . For instance, if  $P$  is identified with its distribution function on the sample space  $R^d$ , we can consider  $P$  as a subset of  $l^\infty(R^d)$ . Or if  $P$  is identified with the map  $f \rightarrow \int f dP$  for  $f \in F \subset L_2(P)$ , then we can consider  $P$  as a subset of  $l^\infty(F)$ . This is the point of view of Wellner (1989), and Sheehy and Wellner (1988) as a natural generalization of the work of Gill (1989). If  $P$  is a semiparametric model  $\{P_{(\theta, G)} : \theta \in \Theta, G \in G\}$ , we can think of  $P$  as identified with  $\Theta \times G$  where  $G$  is a subset of a Banach space  $B_1$ . If  $\Theta \subset R^k$  we are naturally led to  $B = R^k \times B_1$ . In this situation we may, of course, just focus on a part of  $P$ , for example,  $v(P_{(\theta, G)}) = G$ .

If  $P$  is not regular parametric there are, of course, important natural representations of  $P$  which we do not expect to be able to estimate at rate  $n^{-1/2}$ . For instance, identify  $P$  with  $s(P) \equiv \sqrt{dP/d\mu} = \sqrt{p}$  viewed as an element of  $L_2(\mu)$ . Let  $M_\mu = \{P : P \ll \mu\}$  and if  $\mu$  is Lebesgue measure on  $(0, 1)$ ,  $\varepsilon \in (0, 1)$ ,

$$P = \{P \in M_\mu : \varepsilon \leq p \leq \varepsilon^{-1}, |p'| \leq \varepsilon^{-1}, |p''| \leq \varepsilon^{-1}\}.$$

Then it is easy to show from results of Farrell (1972) or Bretagnolle and Huber (1979) (using the second inequality in (A.6.3)) that for any estimate  $\hat{s}$  and any  $\delta > 0$

$$n^{2/5+\delta} \sup_{P \in \mathcal{P}} E_P \| \hat{s} - s(P) \|_{\mu} \rightarrow \infty;$$

see for instance, Prakasa Rao (1983, page 31), or Devroye and Györfi (1985, page 38).

In this section we consider the estimation of infinite-dimensional parameters which are not necessarily estimable at rate  $n^{-1/2}$ . Using generalizations of the GM and GMC methodology developed in sections 3 and 4, we give a number of examples illustrating the scope and limitations of various methods and then discuss consistency and rates of convergence for such procedures in norms which are typically weaker than the norm on  $\mathbf{B}$ . In fact the usual  $\mathbf{B}$  topology is typically used only for the definition of tangent sets. In the next section we discuss asymptotic linearity for the AGM-estimates proposed here and, following in part Wong and Severini (1991), show how this may be applied to obtain Gaussianity and regularity of estimates of (possibly Banach-valued) parameters  $v(P)$  which are estimable at rate  $n^{-1/2}$ .

The GMC method of construction generalizes formally without any difficulty. Given  $D_n : \mathbf{A} \times \mathbf{M}_0 \rightarrow R$  with  $\mathbf{A} \subset \mathbf{B}$ , the estimate  $\hat{v}_n \in \mathbf{A}$  is GMC if it minimizes  $D_n(\cdot, \mathcal{P}_n)$ .

The generalization of the GM method that follows owes much to Wong and Severini (1991) and Gill (1989). As in section 3.2, let  $\dot{\mathbf{A}}^0$  be the tangent set of  $\mathbf{A}$  at a point  $a \in \mathbf{A}$ , and let  $\dot{\mathbf{A}}$  be the tangent space. (Note that definitions 3.2.1 and 3.2.2 can be used for Banach spaces too; see also appendix 5). Suppose that  $q : \mathbf{A} \rightarrow \Gamma$  where  $\Gamma$  is a Banach space. The following notion of differentiability is discussed in section A.5 and will be used here and later. Let  $\| \cdot \|_{\mathbf{B}}$  and  $\| \cdot \|_{\Gamma}$  be norms on  $\mathbf{B}$ , and  $\Gamma$  respectively (but not necessarily the Banach space norms). If  $T$  is a linear operator from  $\mathbf{B}$  to  $\Gamma$ , let  $\| T \|_{\mathbf{B}\Gamma} \equiv \sup\{ \| T(b) \|_{\Gamma} : \| b \|_{\mathbf{B}} \leq 1 \}$ .

**Definition 1.**  $A. q$  is  $\mathbf{B}, \Gamma$  differentiable at  $v_0 \in \mathbf{A}$  if and only if there exists a bounded linear operator  $\dot{q}(v_0)$  from  $\dot{\mathbf{A}}$  to  $\Gamma$  such that for all  $b, \{b_t\}$  with  $b_t \rightarrow b$  as  $t \rightarrow 0$  and  $v_0 + t b_t \in \mathbf{A}$  for  $|t|$  sufficiently small,

$$\| q(v_0 + t b_t) - q(v_0) - \dot{q}(v_0)b \|_{\Gamma} = o(t) \quad \text{as } t \rightarrow 0.$$

Here  $\dot{q}(v_0)$  bounded means that  $\| \dot{q}(v_0) \|_{\mathbf{B}\Gamma} < \infty$ .

We now motivate the definition of GM that we give below. Suppose  $D_n(\cdot, \mathcal{P}_n)$  is Hadamard differentiable on  $\mathbf{A}$  and let  $W_n(v, \mathcal{P}_n) : \dot{\mathbf{A}} \rightarrow R$  denote a derivative. If  $\hat{v}_n$  is GMC,

$$(1) \quad W_n(\hat{v}_n, \mathcal{P}_n)(h) = 0 \quad \text{for all } h \in \dot{\mathbf{A}}^0.$$

If  $\| \cdot \|_{\mathbf{B}}$  is the Banach norm, since  $W_n(\hat{v}_n, \mathcal{P}_n)$  is linear (as a map from  $\dot{\mathbf{A}}$  to  $R$ ) and bounded, (1) holds for all  $h \in \dot{\mathbf{A}}$  and  $W_n(\hat{v}_n, \mathcal{P}_n)$  is an element of the dual

space  $\dot{A}^*$ . So we can think of (1) with  $h \in \dot{A}$  as saying that  $\hat{v}_n$  is a zero of the map  $W_n(\cdot, P_n) : A \rightarrow \dot{A}^*$ . This leads to:

**Definition 2.** If  $W_n : A \times M_0 \rightarrow \Gamma$  where  $\Gamma$  is Banach, then  $\hat{v}_n$  is a *GM-estimate (respectively an AGM-estimate) if it solves*

$$(2) \quad W_n(\hat{v}_n, P_n) = 0 \quad (\text{respectively } o_p(n^{-1/2}) \text{ in the } \|\cdot\|_\Gamma \text{ norm}).$$

Of course, as in the finite-dimensional case, not all solutions of (1) are GMC. Moreover, (1) is not well defined when  $\dot{A}$  depends on  $\hat{v}_n$ . A further problem is that, by definition  $b \in \dot{A}_0$  implies that  $-b \in \dot{A}$ , and hence, directions  $b$ , such that for any  $b_t \rightarrow b$ ,  $v_0 + tb_t \in A$  only for  $t \geq 0$ , do not belong to  $\dot{A}_0$ . These problems appear in example 1 below, which is analyzed by the GMC method without the introduction of a  $W_n$  function.

The main difficulties hidden in these formal extensions are:

- (i) Existence and unicity of GMC- or GM-estimates  $\hat{v}_n$ .
- (ii) Computability of  $\hat{v}_n$ .

Note that by taking  $A = A_n$  to be a finite-dimensional manifold as in the "method of sieves" discussed below, these two problems become no more difficult than they were for  $A$  Euclidean.

As we have seen, even in the Euclidean case  $\hat{v}_n$  may not exist (or be unique). In the infinite-dimensional case without careful choice of  $D_n$  and/or its domain  $A \times M_0$ , nonexistence becomes the rule rather than the exception. Example 2 below is typical.

If  $v(P) = P$ , parametric experience suggests we consider  $D(P, Q) = \rho(P, Q)$  where  $\rho$  is a metric compatible with  $P_n$ . Unfortunately, such procedures seem not to lead to elegant or easily interpretable solutions and are not efficient even though, as we have seen in the parametric case, they can be comparatively easy to analyze. On the other hand, the extensions of maximum likelihood which we consider in most of this chapter often do lead to computable efficient procedures.

Suppose  $P$  is, as usual, dominated by  $\mu$ , and  $v(P) = P$ . Naive extension of maximum likelihood leads us to consider GMC-estimates corresponding to  $D$  given by

$$(3) \quad D(P, Q) \equiv - \int_{-\infty}^{\infty} \log p(x) dQ(x), \quad \text{where } p \equiv \frac{dP}{d\mu}.$$

Here is an example where this method yields computable estimates.

**Example 1. Exponential mixture model.**

We follow Jewell (1982) in considering the model

$$p(x, G) = \int_0^{\infty} \eta e^{-\eta x} dG(\eta) \quad \text{for } x > 0,$$

where  $G$  is an arbitrary distribution on  $R^+$ . Given  $X_1, \dots, X_n$ , we want to estimate  $G$  via (3). Thus we have to find  $G$  such that

$$\frac{1}{n} \sum_{i=1}^n \log \int_0^{\infty} \eta e^{-\eta X_i} dG(\eta)$$

is maximal. Consider the set  $E \subset (R^+)^n$  which is the convex hull of the curve

$$E_0 \equiv \{(\eta e^{-\eta X_1}, \dots, \eta e^{-\eta X_n}) : \eta > 0\}.$$

$E$  is closed since  $E_0$  is bounded; see Rockafellar (1970). Note that  $e \in E$  iff

$$e = \left( \int_0^{\infty} \eta e^{-\eta X_1} dG(\eta), \dots, \int_0^{\infty} \eta e^{-\eta X_n} dG(\eta) \right)$$

for some  $G$ . So we have to find that point  $e_0$  in  $E$ , for which  $\sum_{i=1}^n \log e_{0i}$  is maximal. Since the logarithm is increasing,  $e_0$  has to be a boundary point of  $E$ . Because  $E$  is convex, the supporting hyperplane theorem yields the existence of a hyperplane  $H$  with  $e_0 \in H$  and  $E$  at one side of  $H$ . Consequently  $e_0 \in E \cap H$  and  $E \cap H$  is the convex hull within  $H$  of  $E_0 \cap H$ . By a theorem of Carathéodory (see Theorem 17.1 of Rockafellar (1970))  $e_0$  is a convex combination of at most  $n$  points in  $E_0 \cap H \subset E_0$ , since  $H$  is  $(n-1)$ -dimensional. But this means that there is an MLE of  $G$  which is discrete and puts its mass at  $n$  points at most. The uniqueness of this MLE follows from the strict concavity in  $\lambda$  of  $\sum_{i=1}^n \log \int_0^{\infty} \eta e^{-\eta X_i} d(\lambda G_1 + (1-\lambda) dG_2)(\eta)$ . This reduces the problem of characterizing the GMC-estimate to the  $2n-1$  parameter maximum likelihood problem in which we restrict to  $G$ 's concentrating on  $n$  (unknown) points. The solution of this problem is discussed further in Laird (1978) and Jewell (1982), and extended to general mixtures of exponential families by Lindsay (1983a, 1983b, 1983c). Here, though not explicit, the GMC-estimate is in principle computable. However, little beyond consistency is known of its behavior; see Pfanzagl (1990) for example.  $\square$

The next example shows the limited scope of the naive extension of maximum likelihood.

**Example 2. Absolutely continuous distributions on  $R$ .**

Suppose  $\mu$  is Lebesgue measure and let  $P = P_\mu$ . Given  $X_1, \dots, X_n$ , take, for instance,  $P_{n\sigma}$  with density

$$p_{n\sigma}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\sigma} \phi\left(\frac{x - X_i}{\sigma}\right)$$

where  $\phi$  denotes the standard normal density function. Then

$$\begin{aligned} & \int \log p_{n\sigma}(x) dP_n(x) \\ &= -\log [n\sigma] + \log [\phi(0)] + \frac{1}{n} \sum_{j=1}^n \log \left\{ 1 + \sum_{i \neq j} \frac{\phi((X_j - X_i)/\sigma)}{\phi(0)} \right\} \\ & \rightarrow \infty \quad \text{as } \sigma \rightarrow 0. \end{aligned}$$

Of course  $P_{n\sigma}$  tends to  $P_n$  as  $\sigma \rightarrow 0$ , but  $P_n \notin P$ .  $\square$

There are several generic ways of modifying maximum likelihood which can



give usable procedures in situations such as these where the naive extension fails.

The first way leads to nonparametric maximum likelihood estimation. Other modifications to which we will return later in this section, involve sieves, regularized maximum likelihood, and penalized maximum likelihood.

*Nonparametric Maximum Likelihood*

Take  $v(P) = P$ . If we let  $A = M$  rather than just  $P$ , we would expect  $\hat{v}_n = IP_n$ . Unfortunately,  $D(P, Q)$  is now not well defined since  $M$  is not dominated by  $\mu$  and densities  $p$  as in (3) can no longer be specified. Kiefer and Wolfowitz (1956) proposed a way out of this difficulty. Their idea is to enlarge  $P$  to an undominated family  $\bar{P}$  and interpret (3) suitably so as to obtain reasonable estimates  $\hat{P}$  in interesting cases. Their estimate satisfies  $\hat{P} \in \bar{P}$  and

$$(4) \quad \int \log \left[ \frac{d\hat{P}}{d\mu} \right] dIP_n \geq \int \log \left[ \frac{dP}{d\mu} \right] dIP_n$$

for all  $P \in P$ , some measure  $\mu$  dominating both  $P$  and  $\hat{P}$ , and some determination of  $d\hat{P} / d\mu$  and  $dP / d\mu$ . This method is called *nonparametric maximum likelihood*. In most interesting cases, initial consideration of this method yields a model  $P_{IP_n}$ , which depends on the data  $IP_n$ , and is dominated by counting measure on  $\{X_1, \dots, X_n\}$ . This model includes most candidates for the title nonparametric maximum likelihood estimator (NPMLE). Finding the member of  $P_{IP_n}$  that is the NPMLE needs only straightforward calculation. We can think of this as a GMC method, but only formally, by defining

$$D(P, Q) = - \int \log \frac{dP}{d\mu_Q} dQ, \quad P \in P_Q.$$

The only arbitrary feature of this approach is the choice of  $\bar{P}$  and the versions of the Radon-Nikodym derivatives. This second difficulty can sometimes be remedied by using the definition of Scholz (1980), and the choice of  $\bar{P}$  is often dictated by the initial choice of  $P$ .

Here are some examples of the method of nonparametric maximum likelihood:

**Example 2. Absolutely continuous distributions on  $R$ , continued.**

It is natural to take  $\bar{P} = M$  in this case. Note that if  $\hat{P}$  satisfies (4), it must be carried by  $X_1, \dots, X_n$ . To see this, write

$$P = (1 - \epsilon)\hat{P} + \epsilon\tilde{P},$$

where  $\hat{P}\{X_1, \dots, X_n\} = 1$ ,  $\tilde{P}\{X_1, \dots, X_n\} = 0$  and  $0 \leq \epsilon \leq 1$ . Choose  $\mu = \hat{P} + \tilde{P}$ . Then

$$\prod_{i=1}^n \frac{d\hat{P}}{d\mu}(X_i) = 1 > (1 - \epsilon)^n = \prod_{i=1}^n \frac{dP}{d\mu}(X_i),$$

unless  $\varepsilon = 0$ . We restrict  $P$  to  $\mathbf{P}_{P_n}$ , the set of all distributions dominated by  $\mu_{P_n}$ , the counting measure concentrated on  $x_1, \dots, x_m$ , the distinct values of the  $X_i$ , let  $p_j = P\{x_j\}$ , and maximize the likelihood as a function of  $p_1, \dots, p_m$ . Thus we maximize

$$\sum_{j=1}^m n_j \log p_j \quad \text{with } n_j = \#\{X_i = x_j : i = 1, \dots, n\}$$

and we of course obtain  $\hat{p}_j = n_j/n$ . We conclude that the nonparametric maximum likelihood estimate  $\hat{P}$  is just the empirical measure  $IP_n$ .  $\square$

**Example 3. Absolutely continuous distributions symmetric about zero.**

Here  $\mathbf{P} = \{\text{all absolutely continuous distributions symmetric about zero}\}$ . For  $\bar{\mathbf{P}} = \{\text{all distributions symmetric about zero}\}$ , we can argue as in example 2 that we can restrict ourselves to discrete distributions concentrating on  $\pm x_1, \dots, \pm x_m$ . If  $p_1, \dots, p_m$  are the masses assigned to  $x_1, \dots, x_m$ , the likelihood is proportional to

$$\prod_{j=1}^m p_j^{(n_{j+} + n_{j-})},$$

where  $n_{j+} \equiv \#\{X_i = x_j\}$ ,  $n_{j-} \equiv \#\{X_i = -x_j\}$ . We conclude that the NPMLE is given by

$$\hat{P}(A) = \frac{1}{2} \{IP_n(A) + IP_n(-A)\}.$$

$\square$

Here is an important example in which the NPMLE is not explicit but still "finite dimensional" and hence, in principle, computable.

**Example 4. Biased sampling.**

The model as given in example 4.4.4 has densities for  $(I, X)$  with respect to counting  $\times$  Lebesgue measure given by

$$(5) \quad p(i, x) = \frac{\lambda_i w_i(x) q(x)}{W_i(Q)} \quad \text{for } i = 1, \dots, s \text{ and } x \in \mathbf{X}$$

where  $w_i \geq 0$ ,  $\sum_{i=1}^s \lambda_i = 1$ , and

$$W_i(Q) = \int w_i(x) dQ(x).$$

We begin by assuming that the  $\lambda_i$ 's are unknown.

As in example 2, if we drop the absolute continuity requirement and apply nonparametric maximum likelihood, we need only consider distributions on  $\{i_1, \dots, i_J\} \times \{x_1, \dots, x_K\}$  where the  $i_j$  are the distinct values of  $I_1, \dots, I_n$  and the  $x_k$  are the distinct values of  $X_1, \dots, X_n$ . For simplicity, take  $i_j = j$ ,  $1 \leq j \leq J$ . The resulting discrete analog of (5) can be rewritten as

$$p(j, x_k) = \frac{\theta_j q_k w_j(x_k)}{\sum_{a,b} \theta_a q_b w_a(x_b)}$$

$$= \exp\{v_j + \mu_k - b(v, \mu)\} w_j(x_k),$$

an exponential family where

$$v_j = \log \theta_j, \quad \theta_j = \frac{\lambda_j}{W_j(Q)}, \quad \mu_k = \log q_k,$$

vary freely over  $R$  and

$$b(v, \mu) = \log\left\{\sum_{j,k} \theta_j q_k w_j(x_k)\right\}, \quad q_k = q(x_k).$$

To make these new parameters identifiable, we can require that  $\sum_j \theta_j = \sum_k q_k = 1$ . The old parameters can be regained via

$$\lambda_j = \sum_k p(j, x_k) = \theta_j \sum_k q_k w_j(x_k) \exp(-b(v, \mu))$$

and

$$Q(x) = \sum_{x_i \leq x} q_k.$$

By standard exponential family theory (for example Lehmann (1983, page 438)) if MLE's  $\hat{\theta}$ ,  $\hat{q}$  exist, they satisfy the likelihood equations, and any solution of these equations is an MLE. The log-likelihood expressed in terms of  $\theta_1, \dots, \theta_J$ ,  $q_1, \dots, q_K$ , is up to a factor  $n$  and a term independent of the parameters  $\theta_j$  and  $q_k$ ,

$$\sum_j \hat{\lambda}_j^* \log \theta_j + \sum_k \hat{q}_k^* \log q_k - \log \sum_{a,b} \theta_a q_b w_a(x_b),$$

where  $\hat{\lambda}_j^*$  and  $\hat{q}_k^*$  are the empirical frequencies of stratum  $j$  and  $x_k$  respectively. Maximization of this expression over  $\theta_j$  yields

$$(6) \quad \hat{\lambda}_j^* = \hat{\theta}_j \sum_k \hat{q}_k w_j(x_k) / \sum_{a,b} \hat{\theta}_a \hat{q}_b w_a(x_b),$$

while maximization over  $q_k$  yields

$$(7) \quad \hat{q}_k^* = \hat{q}_k \sum_j \hat{\theta}_j w_j(x_k) / \sum_{a,b} \hat{\theta}_a \hat{q}_b w_a(x_b).$$

Vardi (1985) gives necessary and sufficient conditions for MLE's to exist and be unique. Substituting (7) into (6) we can exhibit  $\hat{\theta}$  as a  $GM$ -estimate and  $\hat{G}$  as an explicit function of  $\hat{\theta}$  and  $P_n$ , and hence also in  $GM$ -estimate form. For details see example 7.6.3.

A standard algorithm for solving the system (6) and (7) subject to the normalization  $\sum_j \hat{\theta}_j = \sum_k \hat{q}_k = 1$ , is to iterate between the two sets of equations. In each step the log-likelihood is maximized with respect to one set of parameters with the other set of parameters kept fixed. Standard convex function theory yields convergence of this procedure to the MLE. Formally, we start at an arbitrary point with all parameters different from 0 and iterate

$$(8) \quad \begin{aligned} \text{Step 1: } \hat{\theta}_j^{\text{NEW}} &= \frac{\hat{\lambda}_j^* / \sum_k \hat{q}_k^{\text{OLD}} w_j(x_k)}{\sum_a \hat{\lambda}_a^* / \sum_b \hat{q}_b^{\text{OLD}} w_a(x_b)} \\ \text{Step 2: } \hat{q}_k^{\text{NEW}} &= \frac{\hat{q}_k^* / \sum_j \hat{\theta}_j^{\text{NEW}} w_j(x_k)}{\sum_b \hat{q}_b^* / \sum_a \hat{\theta}_a^{\text{NEW}} w_a(x_b)}. \end{aligned}$$

Note that the  $\hat{\theta}_j^{\text{NEW}}$  from step 1 satisfy (6) with  $\hat{q}_k = \hat{q}_k^{\text{OLD}}$  and the side condition  $\sum_j \theta_j = 1$ , and a similar remark holds for step 2. The convergence can be slow if  $K$  is large. An alternative implemented by Cosslett (1981) and Vardi (1985) is to substitute (7) into (6), iteratively solve the resulting system of linear equations, and then use (7). This alternative can be much faster if  $K$  is much larger than  $J$  as is typically the case. We will discuss the asymptotic theory of these estimates in example 7.6.3. □

Here is an important example where the NPMLE method leads to an explicit estimate.

**Example 5. Random censoring.**

Extend the model as given by (6.6.1) to  $\bar{P}$  where we permit  $F, G$  to be arbitrary and even to put mass at  $\infty$ .

We can then write

$$(9) \quad dP(t, \Delta) = \Delta(1 - G(t-)) dF(t) + (1 - \Delta)(1 - F(t)) dG(t).$$

We claim, extending the discussion in example 6.6.1, that  $\bar{P} = \{\text{all distributions on } [0, \infty] \times \{0, 1\}\}$ . To prove this, given any  $P$  on  $[0, \infty] \times \{0, 1\}$ , let

$$\bar{H}(t-) = P(T \geq t) = 1 - (H_1(t-) + H_2(t-)),$$

where

$$H_1(t) = P(T \leq t, \Delta = 1), \quad H_2(t) = P(T \leq t, \Delta = 0).$$

As usual, given  $F$  on  $[0, \infty]$ , let

$$(10) \quad \Lambda_1(t) = \int_{[0,t]} \frac{1}{1 - F(t-)} dF(t)$$

be its hazard rate. Motivated by (9), let

$$(11) \quad \Lambda_1(t) \equiv \int_{[0,t]} \frac{1_{[\bar{H}(s-) > 0]}}{1 - H(t-)} dH_1(s)$$

and hence define  $1 - F(s)$  on  $[0, \tau]$  with  $\tau \equiv \sup \{s : \bar{H}(s) > 0\}$  by (6.6.6) and  $dF(\infty) = 1 - F(\tau)$ . Further, let

$$(12) \quad 1 - G(t-) = \frac{\bar{H}(t-)}{1 - F(t-)}, \quad t \leq \tau,$$

and  $dG(\tau) = 1 - G(\tau-)$ . By (10) and (11)

$$(13) \quad \begin{aligned} dH_1(t) &= (1 - F(t-))^{-1} \bar{H}(t-) dF(t) \\ &= (1 - G(t-)) dF(t). \end{aligned}$$

From (12), by a straightforward passage to the limit,

$$(14) \quad (1 - G(t-)) dF(t) + (1 - F(t)) dG(t) = d(H_1(t) + H_2(t)).$$

Now (13) and (14) yield the representation (9) for  $P$ . Indeed,  $\bar{\mathbf{P}}$  is all distributions on  $[0, \infty] \times \{0, 1\}$ . Hence, by the continuation of example 2 the NPMLE of  $P$  is just  $\hat{P}_n$ . By (6.6.6), any NPMLE  $\hat{F}_n$  of  $F$  must satisfy

$$(15) \quad 1 - \hat{F}_n(t) = \prod_{s \leq t} \left(1 - \frac{dH_{n1}(s)}{\bar{H}_n(s-)}\right) \quad \text{for } t \leq T_{(n)}$$

where  $T_{(1)} < \dots < T_{(n)}$  are the ordered  $T_i$ ,

$$\begin{aligned} \bar{H}_n(t-) &\equiv 1 - (H_{n1}(t-) + H_{n2}(t-)), \\ d\bar{H}_{n1}(t) &\equiv dP_n(t, 1), \quad dH_{n2}(t) = dP_n(t, 0). \end{aligned}$$

Of course, (15) is just the Kaplan-Meier estimate. Since  $P_n\{\infty\} = 0$ , (15) defines  $\hat{F}_n$  uniquely as a probability concentrating on the uncensored  $T_i$  if  $T_{(n)}$  is uncensored and placing some mass beyond  $T_{(n)}$  if  $T_{(n)}$  is censored.  $\square$

After these encouraging examples, we present a familiar example in which the method fails.

**Example 6. Symmetric location.**

Let  $p(x, \theta, G) = g(x - \theta)$ ,  $g$  symmetric about 0, and  $g$  the density of  $G$ . We naturally take  $\bar{\mathbf{P}}$  to be all distributions symmetric about some point. If  $\theta$  were fixed, it follows from example 3 that our estimate of  $P$ , say  $\hat{P}_\theta$ , would be the distribution assigning mass  $1/2n$  to  $X_i$  and  $1/2n$  to  $2\theta - X_i$ ,  $i = 1, \dots, n$ , with obvious modifications if there are ties under  $X_1, \dots, X_n$  or if  $\theta$  is some  $X_j$  or the midpoint of some pair of  $X_j$ 's. Note that, if there are no ties, if  $\theta_1 = (X_i + X_j)/2$  for some  $i$  and  $j$  (possibly  $i = j$ ) and if  $\theta_2$  is not of this type, then

$$\frac{d\hat{P}_{\theta_1}}{d\hat{P}_{\theta_2}}(X_i) = 2, \quad \prod_{h=1}^n \frac{d\hat{P}_{\theta_1}}{d\hat{P}_{\theta_2}}(X_h) = \begin{cases} 4 & \text{if } i \neq j, \\ 2 & \text{if } i = j. \end{cases}$$

This shows that with probability 1 under  $P \in \mathbf{P}$ , the midpoint  $\hat{\theta}$  of any pair of distinct  $X_j$ 's yields  $\hat{P}_{\hat{\theta}}$  as an NPMLE, so that maximum likelihood gives us no guidance in estimating  $P$ . Note however that we can construct efficient estimates of "nice" parameters in this model: first estimate  $\theta$  efficiently by  $\hat{\theta}$ —see section 7.8—and then estimate  $G$  by the maximum likelihood estimate assigning mass  $1/2n$  to  $\pm(X_i - \hat{\theta})$ ,  $1 \leq i \leq n$ .  $\square$

*Modifications: Sieves, Regularization, and Regularized MLEs*

Three modifications of maximum likelihood which deal with the difficulties that examples 2 and 6 pose have been considered.

One approach is the *method of sieves*; see, e.g., Grenander (1981). We select a sieve of submodels  $\{P_m\}_{m \geq 1}$ , such that  $P = \overline{\bigcup_m P_m}$ , the weak closure of  $\bigcup_m P_m$ . Suppose further that each of the  $P_m$ 's admits a maximum likelihood estimate  $\hat{P}_m \in P_m$  such that

$$\int \log \hat{p}_m(x) dP_n = \max\left\{ \int \log p dP_n : P \in P_m \right\}.$$

Finally, we designate  $\hat{P}_{m(n)}$  (for a suitable, possibly data determined sequence  $m(n)$ ) as the *method of sieves estimate* of  $P$ . The  $\hat{P}_m$  should be readily computable so that the  $P_m$  are either regular parametric submodels or have explicitly computable nonparametric maximum likelihood estimates. This method has appeared in other contexts in applied mathematics, for instance see the "selection method" in Tikhonov and Arsenin (1977). Heuristically the method should produce efficient estimates provided that we can choose  $m(n)$  so that the bias incurred by replacing  $P$  by  $P_m$  is  $o(n^{-1/2})$ . Efficient methods of estimation on  $P_m$  other than maximum likelihood can also be considered; see example 10. We shall find it useful to consider a slight extension, which we also call the method of sieves. Drop the requirement that  $P_m \subset P$ , but suppose  $P$  is the limit of  $P_m$  in the sense that if  $P_m \in P_m$ ,  $P_m \rightarrow P$ , then  $P \in P$  and any  $P \in P$  is a limit point of  $P_m \in P_m$ .

The method of *regularization*, also discussed in this context by Grenander, consists of replacing the functional  $p \rightarrow \int \log p dP_n$ , whose maximum is not assumed, by a more "regular" functional whose maximum is uniquely assumed at a point which "behaves well." The most important class of methods of this type have been studied under the name of *penalized maximum likelihood*. Usually the functional is of the form

$$-\int \log p dP_n + \lambda_n J(p).$$

For instance, in example 2,  $J(p)$  is usually taken as a penalty for lack of smoothness such as

$$J(p) = \int \frac{[p']^2}{p}(x) dx.$$

This method can be viewed formally as an extension of choosing the posterior mode as an estimate in a Bayesian context. In this formulation,  $\exp(-\lambda_n J(p))$  plays the role of prior density.

A third approach is to consider formally the maximum likelihood functional,  $Q \rightarrow P_Q \equiv P(Q)$  where  $P_Q \in P$  with density  $p_Q$  such that

$$(16) \quad \int \log p_Q(x) dQ(x) = \max_{P \in P} \int \log p(x) dQ(x).$$

If  $Q = I_{P_n}$ , the empirical distribution, this is just the maximum likelihood estimate. By choosing  $Q = I_{P_n^\#}$ , a regularized version of  $I_{P_n}$ , we may have  $P(I_{P_n^\#})$  defined even though  $P(I_{P_n})$  is not. For instance, in example 2, if  $I_{P_n^\#}$  is  $I_{P_n}$  kernel smoothed, its density is

$$p_n^\#(x) = \frac{1}{nb_n} \sum_{i=1}^n \phi\left(\frac{x - X_i}{b_n}\right),$$

and the resulting estimate of  $P$  is just  $\mathbb{P}_n^\#$ . We shall call this the *regularized maximum likelihood method*. It is not hard to see that at least in this special case, regularized maximum likelihood is a special case of the sieve method. An example of this approach is given by Staniswalis (1989).

We limit our further discussion to these extensions of maximum likelihood and some modifications. The main class of methods we do not discuss are non-parametric Bayesian methods and Bayesian sieves. Only the first of these approaches has been studied extensively, with a primary focus on Dirichlet priors. A good introduction to the issues involved is in Diaconis and Freedman (1986). See also Hjort (1990).

The difficulty with all three extensions is the introduction of additional choices that have to be made. For sieves, we must choose the dimension of the member of the sieve for which we calculate, as well as the sieve itself; for penalized maximum likelihood both  $\lambda_n$  and the penalty functional  $J$  must be specified; and for our third approach, the method of smoothing and the bandwidth must be selected. Moreover, questions of existence and computability remain. In the examples we discuss in this section and section 7.6, we shall often study modifications of these methods which lead to procedures which have formally the same asymptotic properties but are more tractable. We now briefly present examples of each of the three approaches. These examples will be pursued in detail later in this section and in section 7.6.

**Example 7. Density estimation by the method of sieves.**

Suppose

$\mathbf{P} = \{\text{all distributions on } [0, 1] \text{ with densities } p \text{ such that } \log p \in \mathbf{L}\}$ ,  
where

$$\mathbf{L} \equiv \{s : s \text{ is twice differentiable, } \|s\|_\infty + \|s''\|_\infty < \infty\},$$

a Sobolev space of functions on  $[0, 1]$  with  $\|s\| \equiv \|s\|_\infty + \|s''\|_\infty$ . Let  $v(P) = \log p - \int_0^1 \log p(y) dy$ . This peculiar centering is introduced for technical reasons. Clearly we can recover  $P$  from  $v(P) = v$  by  $p(x) = e^{v(x)} / \int_0^1 e^{v(y)} dy$ . The following sieve has been proposed by Stone (1990) among others. Let

$$\mathbf{C}_m = \{\text{cubic splines with knots at } \frac{j}{m}, 0 \leq j \leq m,\}$$

which are continuous and have

continuous first and second derivatives}

and let

$$\mathbf{P}_m = \{P \in \mathbf{P} : v(P) \in \mathbf{C}_m\}.$$

Note that  $\mathbf{P}_m$  is an exponential family which we can parametrize as follows:

$$(17) \quad p(x, \alpha_1, \dots, \alpha_{m+4}) = \exp \left\{ \sum_{j=2}^{m+4} \alpha_j b_j(x) - c(\alpha) \right\},$$

where  $b_1, \dots, b_{m+4}$  is a  $B$ -spline basis for  $C_m$ ,  $b_1 = 1$ , and

$$(18) \quad \exp c(\alpha) = \int_0^1 \exp \left\{ \sum_{j=2}^{m+4} \alpha_j b_j(x) \right\} dx.$$

See De Boor (1978) for a discussion of splines and their properties.

Since the natural parameter space is  $R^{m+4}$ , standard theory, see for example Brown (1986), leads to the existence of unique MLE's  $\hat{\alpha}$  of  $\alpha$  satisfying

$$n^{-1} \sum_{i=1}^n b_j(X_i) = \int_0^1 b_j(x) p(x, \alpha) dx, \quad j = 2, \dots, m+4,$$

and the corresponding estimate of  $v$

$$(19) \quad \hat{v}_n = \sum_{j=2}^{m+4} \hat{\alpha}_j \{ b_j(\cdot) - \int_0^1 b_j(y) dy \}.$$

We put this estimate in GMC and GM forms below. □

Here is an example where it is useful to marry sieves and an efficient method of estimation.

**Example 8. Estimating a joint distribution with one or both marginals known.**

We follow the theme and notation of examples 6.2.2 and 6.2.3. We observe  $X = (U, V) \in [0, 1]^2$ , and suppose  $P$  is restricted only by:

Case (i)

$$(20) \quad P(U \leq u) = F_0(u) \quad \text{for all } u;$$

or

Case (ii) both

$$(21) \quad \begin{aligned} P(U \leq u) &= F_0(u) && \text{for all } u, \text{ and} \\ P(V \leq v) &= G_0(v) && \text{for all } v. \end{aligned}$$

We take  $v(P) = P$ .

Let

$$I_{jkm} \equiv \left[ \frac{j-1}{m}, \frac{j}{m} \right) \times \left[ \frac{k-1}{m}, \frac{k}{m} \right) \quad 1 \leq j, k \leq m.$$

A natural sieve to consider is

$$\mathbf{P}_m \equiv \{ P : \text{the conditional distribution of } X \text{ given } X \in I_{jkm} \text{ is arbitrary, and if } p_{jkm} \equiv P(X \in I_{jkm}), \text{ the } p_{jkm} \text{ satisfy the marginal restrictions for case (i), or case (ii)} \}.$$

This is a nonparametric sieve with unknown parameters,  $p_{jkm}$  and the conditional distribution of  $X$  given  $X \in I_{jkm}$ , where



$$(22) \quad \sum_k p_{jkm} = f_{jmo}$$

for case (i), and, in addition,

$$(23) \quad \sum_j p_{jkm} = g_{kmo}$$

for case (ii), are fixed and known.

If  $N_{jkm} \equiv \sum_{i=1}^n 1[X_i \in I_{jkm}]$ , the NPMLE (relative to the sieve) of  $P(\cdot | X \in I_{jkm})$  is the distribution assigning mass  $N_{jkm}^{-1}$  to all  $X_i \in I_{jkm}$ . Now consider  $\{p_{jkm}\}$ . It seems plausible that we can get efficiency with a method known to be asymptotically equivalent to maximum likelihood for fixed  $m$ . In particular, consider modified minimum  $\chi^2$  estimation where the  $\hat{p}_{jkm}$  minimize, subject to (22) for case (i), and (22) and (23) for case (ii),

$$\sum_{j,k} \frac{(N_{jkm} - np_{jkm})^2}{N_{jkm}} 1[N_{jkm} > 0].$$

Let a + subscript denote summation over the appropriate index. For case (i)

$$(24) \quad \hat{p}_{jkm} = \frac{N_{jkm}}{N_{j+m}} f_{jmo}$$

is also the maximum likelihood estimate. For case (ii),

$$(25) \quad \hat{p}_{jkm} = \frac{N_{jkm}}{n} (1 - \hat{\alpha}_j - \hat{\beta}_k),$$

where  $\hat{\alpha}$  and  $\hat{\beta}$  solve

$$(26) \quad \begin{aligned} \hat{\beta}_k \frac{N_{+km}}{n} + \sum_j \hat{\alpha}_j \frac{N_{jkm}}{n} &= \frac{N_{+km}}{n} - g_{kmo}, \quad 1 \leq k \leq m, \\ \hat{\alpha}_j \frac{N_{j+m}}{n} + \sum_k \hat{\beta}_k \frac{N_{jkm}}{n} &= \frac{N_{j+m}}{n} - f_{jmo}, \quad 1 \leq j \leq m. \end{aligned}$$

This example can be put in GMC form but is more easily dealt with directly. We pursue case (i) in section 7.6. The form (25), (26) of the estimate in case (ii) is exploited in Bickel, Ritov and Wellner (1991). □

**Example 9. Regression estimation by penalized maximum likelihood.**

Consider the nonparametric regression model in which we observe  $X = (Z, Y) \in [0, 1] \times R$  where

$$(27) \quad Y = t(Z) + \varepsilon$$

and  $Z$  and  $\varepsilon$  are independent,  $\varepsilon \sim N(0, \sigma^2)$ ,  $t \in \mathbf{T}$ , a Sobolev space with norm  $\|\cdot\|_s$ . We could let  $\mathbf{T} = \mathbf{L}$  given in example 6, but prefer in view of (29) below to take

$$(28) \quad \mathbf{T} = \{t: t \text{ twice differentiable, } \|t\|_s^2 \equiv \|t\|_2^2 + \|t''\|_2^2 < \infty\},$$

where  $\|h\|_2^2 = \int_0^1 h^2(x) dx$ . Any function  $\hat{t}$  with  $\hat{t}(Z_i) = Y_i$  is an MLE of  $t$ . Such estimates are inconsistent. The penalized MLE by definition minimizes

$$(29) \quad D_n(t, \mathbb{P}_n) \equiv \int (y - t(z))^2 d\mathbb{P}_n(z, y) + \lambda_n \|t''\|_2^2,$$

where the penalty  $\lambda_n$  is  $> 0$ . We show later that a unique minimizer  $\hat{t}_n$  exists and is in fact a cubic spline with knots at the data points; see Silverman (1985) for a discussion of such procedures.  $\square$

**Example 10. Regularized MLE in the Has'minskii-Ibragimov model.**

We follow the notation of our discussion in section 4.5. Suppose that  $q(y, \theta, \eta)$  does not depend on  $\theta$  and write  $q(y, \theta, \eta) \equiv q_0(y, \eta)$ . We are interested in estimating  $v = \log g$ . The joint density of  $X = (U', Y)$  is

$$p(u', y, g) = e^{v(u')} q(y, v),$$

where  $q(y, v) \equiv \int q_0(y, \eta) \exp\{v(\eta)\} d\eta$ . We suppose  $G$  with density  $g$  concentrates on  $[0, 1]$ . View  $v$  as an element of  $\mathbf{B} \equiv L_2([0, 1])$ . Let

$$\mathbf{A} = \{v \in \mathbf{B} : \int_0^1 \exp\{v(u)\} du = 1, \|v\|_\infty + \|v''\|_\infty < \infty\},$$

$$\mathbf{P} \equiv \{P_g : v \in \mathbf{A}\}.$$

Let  $p_{n1}^\#$  be an estimate of the density  $g$  of  $U'$  based on  $\mathbb{P}_{n1}$ , the empirical distribution of  $U'_1, \dots, U'_n$ . The estimate has to be chosen with care; see section 7.6. Let  $\mathbb{P}_n^\# = \mathbb{P}_{n1}^\# \times \mathbb{P}_{n2}$  where  $\mathbb{P}_{n1}^\#$  has density  $p_{n1}^\#$  and  $\mathbb{P}_{n2}$  is the empirical distribution of  $Y_1, \dots, Y_n$ . That is we regularize  $\mathbb{P}_{n1}$  only.

Formally, maximize

$$(30) \quad \int \log p(u', y, v) d\mathbb{P}_n^\#(u', y) \\ = \int v(u') p_{n1}^\#(u') du' + \int \log \int q_0(y, \eta) \exp\{v(\eta)\} d\eta d\mathbb{P}_{n2}(y)$$

subject to  $\int \exp\{v(u)\} du = 1$ . Differentiation of (30) with respect to  $v$  yields, after standard manipulations,

$$(31) \quad W_n(v, \mathbb{P}_n) \equiv p_{n1}^\#(\cdot) + \exp\{v(\cdot)\} \left\{ \int \frac{q_0(y, \cdot)}{q(y, v)} d\mathbb{P}_{n2}(y) - 2 \right\} \\ = 0$$

for the maximizing  $\hat{v}_n$ . If  $\hat{v}_n$  exists it is evidently GMC. The form of (31) indicates that under regularity conditions we can represent  $\hat{v}_n$  as a GM-estimate. We do this and establish asymptotic existence, unicity, as well as consistency and asymptotic linearity in section 7.6.  $\square$

*Asymptotic Existence, Consistency, and Rates in the Convex Case*

Suppose that, for  $v \in \mathbf{A}$  and  $P \in \mathbf{P}$ , the following two conditions hold:

(C1)  $D_n(v, \mathbb{P}_n) \rightarrow_p D(v, P),$

(C2)  $D(v, P) > D(v(P), P)$  unless  $v = v(P).$

We then hope that the GMC-estimate  $\hat{v}_n$  asymptotically exists and is a consistent estimate of  $v(P)$  in some sense. This is, of course, not true without further conditions such as those of Wald (1949) and Bahadur (1960). An application of this approach to estimation by sieves is Geman and Hwang (1982). We do not go into the Wald-Bahadur approach but discuss, with examples, two elementary approaches for the convex case. Our goal will be to obtain consistency and rate results. It is important to note that these will typically *not* be results for the “natural” norm on  $\mathbf{B}$  but rather for weaker norms. Such norms are discussed more generally in section 7.6.

We begin with a theorem from convex analysis.

**Theorem 1.** Let  $P$  be given and suppose that  $N \subset A \subset \mathbf{B}$  is open and convex. Suppose that

- (i)  $D(v, P)$  is convex in  $v \in N$ ,
- (ii)  $v \rightarrow D(v, P)$  is continuous in the norm topology of  $\mathbf{B}$ ,
- (iii)  $\liminf_{v \rightarrow \partial N} D(v, P) > \inf_{v \in N} D(v, P)$ ,  $\liminf_{\|v\|_{\mathbf{B}} \rightarrow \infty} D(v, P) > \inf_{v \in N} D(v, P)$ , where  $\partial N$  is the boundary of  $N$ .

Then the minimum of  $D(\cdot, P)$  is assumed on  $N$ . If  $D(\cdot, P)$  is strictly convex, the minimum is uniquely attained.

**Proof.** The crucial point here is that (i) and (ii) imply weak lower semicontinuity of  $D(\cdot, P)$  by corollary 2.2 of chapter I of Ekeland and Temam (1976). The result follows from proposition 1.2 of chapter II of Ekeland and Temam (1976). □

We have used this result successfully in section 4 where  $\mathbf{B}$  was Euclidean. For  $\mathbf{B}$  infinite-dimensional it is common that (iii) fails; see examples 6 and 7 below. However, the following refinement of theorem 1 is useful.

**Theorem 2.** Suppose that (i)–(ii) of theorem 1 hold for each member of a sequence  $D_n$ . Let  $\{\mathbf{B}_n\}$  be a sequence of Banach subspaces of  $\mathbf{B}$ . Let  $v_n \equiv v_n(P_0)$  minimize  $D_n(v, P_0)$  and  $\hat{v}_n$  minimize  $D_n(v, P_n)$  on  $\mathbf{B}_n$ . Let  $N_n$  be a closed convex neighborhood of  $v_n$  in  $\mathbf{B}_n$ . Suppose that  $D_n(v, P_0)$  is strictly convex for all  $n$  and, if  $\partial N_n \equiv$  boundary of  $N_n$  in  $\mathbf{B}_n$ ,

$$(32) \quad \sup \left\{ \left| \frac{D_n(v, P_n) - D_n(v_n, P_n)}{D_n(v, P_0) - D_n(v_n, P_0)} - 1 \right| : v \in \partial N_n \right\} = o_p(1).$$

Then

$$P_0(\hat{v}_n \in N_n) \rightarrow 1.$$

**Proof.** By strict convexity

$$\begin{aligned} 0 &< \inf \{ D_n(v, P_0) - D_n(v_n, P_0) : v \in \partial N_n \} \\ &= (1 + o_p(1)) \inf \{ D_n(v, P_n) - D_n(v_n, P_n) : v \in \partial N_n \} \end{aligned}$$

by (32), and the result follows from theorem 1. □

We note again that  $N_n$  is typically *not* a sphere in the norm on  $\mathbf{B}$  but one in a weaker norm. Here is an application of theorem 2.

**Example 7. Density estimation by the method of sieves, continued.**

We put this and similar sieve examples into the GMC framework. Let

$$\mathbf{B} \equiv \mathbf{L} \cap \{b : \int b(x) dx = 0\}.$$

a Sobolev space endowed with the norm  $\|b\|_\infty + \|b''\|_\infty$ . Let  $\mathbf{C}_m$ , be the linear space of cubic splines defined previously,  $m = m(n)$ ,  $\mathbf{B}_n = \mathbf{C}_m \cap \mathbf{B}$ , and

$$D(b, P) \equiv - \int b(x) dP(x) + \log \int e^{b(x)} dx.$$

Then  $D(\cdot, P)$  is convex, and  $\hat{v}_n$  defined in (19) minimizes  $D(\cdot, P_n)$  over  $\mathbf{B}_n$ .

**Proposition 1.** If in the above situation  $v_0 = v(P_0)$ ,  $m = m(n) \rightarrow \infty$  and  $m = o(n^{1/2})$ , then  $\hat{v}_n$  is consistent in both the  $l^\infty$  ( $[0,1]$ ) and  $L_2$ ( $[0,1]$ ) norms. In fact:

- A.  $\|\hat{v}_n - v_0\|_\infty = O_p(mn^{-1/2} + m^{-3/2})$ ,
- B.  $\|\hat{v}_n - v_0\|_2 = O_p((m/n)^{1/2} + m^{-2})$ .

More general results of this type may be found in Stone (1990).

**Proof.** We begin by recording some well known basic properties of  $\mathbf{B}_n$ .

For any  $v \in \mathbf{B}_n$ , write

$$(a) \quad v(x) = \sum_{i=0}^{m-1} \sum_{j=0}^3 a_{ji}(v) b_{mij}(x),$$

where  $b_{mij}(x) = m^j(x - i/m)^j 1_{[i/m < x \leq (i+1)/m]}$ .

Let  $M_i = \sup\{|v(x)| : i/m < x \leq (i+1)/m\}$ . We see that

$$(b) \quad \int_{i/m}^{(i+1)/m} v^2(x) dx \geq \int_{i/m}^{(i+1)/m} \{M_i m^3(x - \frac{i}{m})^3\}^2 dx = \frac{M_i^2}{7m}.$$

Consequently, for  $v \in \mathbf{B}_n$ , we have

$$(c) \quad \|v\|_\infty^2 \leq 7m \|v\|_2^2.$$

Let  $a_i(v) = \max_j |a_{ji}(v)|$ . Minimizing  $\int_0^1 (a_0 + a_1 y + a_2 y^2 + a_3 y^3)^2 dy$  over the  $a_i$ 's keeping  $a_k$  fixed and doing this for  $k = 0, \dots, 3$ , we notice that there exists a positive constant  $c$  such that

$$(d) \quad \int_{i/m}^{(i+1)/m} v^2(x) dx \geq \frac{c a_i^2(v)}{m}, \quad i = 1, \dots, m-1.$$

This yields

$$(e) \quad \|v\|_2^2 \geq \frac{c}{m} \sum_i a_i^2(v).$$

Let  $v_n$  minimize  $D(v, P_0)$  on  $\mathbf{B}_n$ . We can argue for the existence of  $v_n$  in the same way as for the existence of  $\hat{v}_n$ . Next we derive the rates at which  $v_n$

approaches  $v_0$ . Given  $b \in \mathbf{B}$ , let  $\pi_m b$  be the cubic spline such that  $\pi_m b^{(k)}(j/m) = b^{(k)}(j/m)$  for  $k = 0, 1, 2, j = 0, \dots, m$ . It is well known and easy to see that for some constant  $C$

$$\|\pi_m b\|_\infty + \|(\pi_m b)'\|_\infty \leq C(\|b\|_\infty + \|b'\|_\infty),$$

$$(f) \quad \|\pi_m b - b\|_\infty \leq \frac{C \|b\|_\infty}{m^2}.$$

As a consequence

$$(g) \quad D(\pi_m v_0, P_0) - D(v_0, P_0) = O(m^{-2}).$$

For any  $v \in \mathbf{B}$ , let  $Tv = v - \log \int e^{v(x)} dx$ . Note that

$$(h) \quad \int e^{Tv(x)} dx = 1.$$

Consider  $f(\lambda) = D(v_0 + \lambda(v_n - v_0), P_0)$ . Its minimum is achieved at  $\lambda = 0$  and hence

$$(i) \quad 0 = f'(0)$$

$$\begin{aligned} &= - \int (v_n(x) - v_0(x)) dP_0(x) + \frac{\int (v_n(x) - v_0(x)) e^{v_0(x)} dx}{\int e^{v_0(x)} dx} \\ &= - \int (Tv_n(x) - Tv_0(x)) dP_0(x) + \int (Tv_n(x) - Tv_0(x)) e^{Tv_0(x)} dx. \end{aligned}$$

Now

$$\begin{aligned} (j) \quad D(v_n, P_0) - D(v_0, P_0) &= - \int (Tv_n(x) - Tv_0(x)) dP_0(x) \\ &= - \int (Tv_n(x) - Tv_0(x)) e^{Tv_0(x)} dx, \quad \text{by (i)} \\ &= \int (e^{Tv_n(x) - Tv_0(x)} - 1 - (Tv_n(x) - Tv_0(x))) e^{Tv_0(x)} dx, \end{aligned}$$

by (h). We claim that  $\|Tv_n - Tv_0\|_\infty \rightarrow 0$ . Suppose that for some  $\varepsilon > 0$ ,  $\|Tv_{n'} - Tv_0\|_\infty > \varepsilon$ ,  $n' \rightarrow \infty$ . Then  $|Tv_{n'}(x) - Tv_0(x)| > \varepsilon/2$  on an interval of length at least  $c/m(n')$  for some universal  $c > 0$ , since  $Tv_0$  is fixed and  $Tv_n$  are polynomials on intervals of length  $1/m$ . But the right side of (j) is then of order at least  $m^{-1}$ , a contradiction to (g). Hence  $\|Tv_n - Tv_0\|_\infty \rightarrow 0$  and (j) yields

$$(k) \quad D(v_n, P_0) - D(v_0, P_0) = \left(\frac{1}{2} + o(1)\right) \int (Tv_n(x) - v_0(x))^2 e^{Tv_0(x)} dx.$$

The same argument proves that

$$\begin{aligned} (l) \quad D(\pi_m v_0, P_0) - D(v_0, P_0) &= \left(\frac{1}{2} + o(1)\right) \int (T\pi_m v_0(x) - Tv_0(x))^2 e^{Tv_0(x)} dx \\ &= O(m^{-4}) \end{aligned}$$

by (f). Since  $D(v_n, P_0) \leq D(\pi_m v_0, P_0)$  we conclude that the expressions in (k) and (l) are of order  $O(m^{-4})$ ; i.e.,

$$(m) \quad \|Tv_n - Tv_0\|_2^2 = O(m^{-4}),$$

which entails

$$(n) \quad \|v_n - v_0\|_2^2 = O(m^{-4}).$$

We are now in a position to find suitable  $N_n$  for which (32) holds. If  $v \in B_n$ , by (a),

$$\begin{aligned} D(v, IP_n) - D(v_n, IP_n) - D(v, P_0) + D(v_n, P_0) \\ = \int (v_n(x) - v(x)) d(IP_n - P_0)(x) \\ = \sum_{i=0}^{m-1} \sum_{j=0}^3 a_{ji}(v_n - v) W_{jin}, \end{aligned}$$

where

$$W_{jin} = \int b_{mij}(x) d(IP_n - P_0)(x).$$

Suppose that  $\sqrt{m} \varepsilon_n \rightarrow 0$  and  $\varepsilon_n \sqrt{n/m} \rightarrow \infty$  and  $N_n = \{v \in B_n : \|v - v_n\|_2 \leq \varepsilon_n\}$ . Then

$$\begin{aligned} (o) \quad \sup \left\{ \frac{\left| \int (v_n(x) - v(x)) d(IP_n - P_0)(x) \right|}{\|v_n - v\|_2^2} : v \in \partial N_n \right\} \\ \leq \varepsilon_n^{-2} \max_{c \sum a_i^2/m \leq \varepsilon_n^2} \sum_{i=0}^{m-1} |a_i| \sum_{j=0}^3 |W_{jin}|, \quad \text{by (e)} \\ \leq \varepsilon_n^{-1} \sqrt{\frac{m}{c}} \left( \sum_{i=0}^{m-1} \left( \sum_{j=0}^3 |W_{jin}| \right)^2 \right)^{1/2}, \quad \text{by Cauchy Schwarz.} \end{aligned}$$

Simple expectation considerations yield that

$$(p) \quad \varepsilon_n^{-1} \sqrt{\frac{m}{c}} \left( \sum_{i=0}^{m-1} \left( \sum_{j=0}^3 |W_{jin}| \right)^2 \right)^{1/2} = O_p \left( \sqrt{\frac{m}{n}} \varepsilon_n^{-1} \right) = o_p(1).$$

By the same argument as for (k) and (l), for  $v$  with  $\|v - v_n\|_2 = \varepsilon_n$ ,

$$\begin{aligned} (q) \quad D(v, P_0) - D(v_n, P_0) &= \left( \frac{1}{2} + O(\varepsilon_n) \right) \int (Tv(x) - Tv_n(x))^2 e^{Tv_n(x)} dx \\ &\geq C_0 \|v - v_n\|_2^2, \quad \text{for some } C_0. \end{aligned}$$

Since (o), (p), and (q) imply (32), theorem 2 yields

$$(r) \quad \|\hat{v}_n - v_n\|_2 = O_p(\varepsilon_n).$$

Together with (n) this implies part B because  $\varepsilon_n \sqrt{n/m}$  diverges arbitrarily slowly. Part A follows from part B and (c).

Note: Part B with  $m = n^{1/5}$  yields the best possible rate for  $L_2$  estimation of  $v$ , but part A with the optimal choice  $m = n^{1/5}$  is not sharp, yielding only a

best rate of  $n^{-3/10}$  rather than the optimal  $n^{-1/3}(\log n)^{1/2}$ . This is a defect of our argument rather than the method of estimation, see Stone (1990).  $\square$

To get consistency and rate results it may be convenient to work directly with the likelihood equations. In such cases, the following theorem can be useful.

**Theorem 3.** Suppose  $D_n(\cdot, P_n)$  is convex, and suppose that  $W_n(\cdot, P_n) = \dot{D}_n(\cdot, P_n)$  is as in (1). Define  $v_n, \hat{v}_n, B_n, N_n$  as in theorem 2, but make  $N_n$  open. Suppose that

$$(33) \quad P(\inf \{W_n(v, P_n)(v - v_n) : v \in \partial N_n\} > 0) \rightarrow 1.$$

Then  $P(\hat{v}_n \in N_n) \rightarrow 1$ .

**Proof.** Since  $D_n(v_n + \lambda(\hat{v}_n - v_n), P_n)$  is convex,

$$g_n(\lambda) \equiv W_n(v_n + \lambda(\hat{v}_n - v_n), P_n)(\hat{v}_n - v_n)$$

is nondecreasing in  $\lambda$ . Let  $\hat{\lambda}_n$  be such that  $v_n + \hat{\lambda}_n(\hat{v}_n - v_n) \in \partial N_n$ . Then

$$\begin{aligned} P(\hat{v}_n \notin N_n) &= P(\hat{\lambda}_n \leq 1) \\ &\leq P(g_n(\hat{\lambda}_n) \leq 0), \quad \text{by (1),} \\ &\leq P(\inf \{W_n(v, P_n)(v - v_n) : v \in \partial N_n\} \leq 0) \rightarrow 0 \end{aligned}$$

by (33).  $\square$

**Example 9. Regression estimation by penalized maximum likelihood, continued.**

In this case

$$(34) \quad \begin{aligned} W_n(t, P_n)(h) &\equiv \frac{\partial}{\partial \eta} D_n(t + \eta h, P_n) |_{\eta=0} \\ &= 2 \{n^{-1} \sum_{i=1}^n (t(Z_i) - Y_i) h(Z_i) + \lambda_n \int t''(z) h''(z) dz\}, \end{aligned}$$

which is indeed a bounded linear functional on  $T$  since the maps  $h \rightarrow h(Z_i)$  and  $h \rightarrow \int t'' h''$  are continuous linear functionals. Note that (34) reveals that  $\hat{t}_n$  is a cubic spline. For  $z \neq Z_i$ , we need only see that for smooth  $h$  with support an arbitrarily small neighborhood of  $z$ ,  $\int \hat{t}_n''(z) h''(z) dz = 0$ , and hence by standard arguments  $\hat{t}_n^{(iv)}(z) = 0$ . Smooth fit at the knots  $Z_{(1)}, \dots, Z_{(n)}$  is now forced by  $\hat{t}_n \in T$ . Computation of  $\hat{t}_n$  reduces to a finite dimensional optimization problem, in fact, the solution of a system of linear equations.

Let  $A = T$  and suppose  $Y_i = t_0(Z_i) + \varepsilon_i, i = 1, \dots, n$ , where  $t_0 \in T$ , under  $P_0$ . Let  $t_n \in T$  be the minimizer of

$$E_0(t(Z) - t_0(Z))^2 + \lambda_n \int [t''(z)]^2 dz.$$

**Proposition 2.** Suppose  $\lambda_n \rightarrow 0$  and  $n^{2/3} \lambda_n / \log n \rightarrow \infty$  in the above situation.

A. Then

$$\begin{aligned} &\int (\hat{t}_n(z) - t_n(z))^2 dP_0(z) + \lambda_n \int (\hat{t}_n''(z) - t_n''(z))^2 dz \\ &= O_p(\lambda_n^{-1/2} n^{-1} \sqrt{\log n}). \end{aligned}$$

B. Furthermore, if  $P_0$  has a density bounded away from 0 then  $\hat{t}_n$  is consistent in the  $L_2([0,1])$  norm. In fact,

$$\|\hat{t}_n - t_0\|_2^2 = O_p(\lambda_n).$$

**Proof.** Note that  $t_n$  is the solution of

$$(a) \quad 0 = E_0\{(t_n(Z) - t_0(Z))b(Z)\} + \lambda_n \int t_n''(z)b''(z) dz, \quad b \in \mathbf{T}.$$

With  $b = h_n \equiv t_n - t_0$ , this yields

$$(b) \quad E_0 h_n^2(Z) = -\lambda_n \int (t_n''(z) - \frac{1}{2} t_0''(z))^2 dz + \frac{1}{4} \lambda_n \int (t_0''(z))^2 dz \\ \leq \frac{1}{4} \lambda_n \|t_0''\|_2^2$$

and

$$(c) \quad 0 \geq \|h_n''\|_2^2 + \int t_0''(z) h_n''(z) dz,$$

which implies by Cauchy-Schwarz,

$$(d) \quad \|h_n''\|_2 \leq \|t_0''\|_2.$$

Hence, by Sobolev embedding inequalities (see, e.g., Corollary 5.16 of Adams (1975)), for a universal constant  $C$

$$(e) \quad \|h_n'\|_\infty + \|h_n\|_\infty \leq C \|h_n''\|_2 \leq C \|t_0''\|_2.$$

Let

$$(f) \quad N_n = \{t: E_0\{(t(Z) - t_n(Z))^2\} + \lambda_n \|t'' - t_n''\|_2^2 < \gamma_n^2\}$$

for  $\gamma_n$  to be specified later. Now

$$(g) \quad \frac{1}{2} W_n(t, P_n)(t - t_n) \\ = \frac{1}{n} \sum_{i=1}^n (t(Z_i) - t_0(Z_i))(t(Z_i) - t_n(Z_i)) \\ - \frac{1}{n} \sum_{i=1}^n \varepsilon_i (t(Z_i) - t_n(Z_i)) \\ + \lambda_n \int t''(z)(t''(z) - t_n''(z)) dz.$$

We obtain from this and (a) with  $b = h \equiv t - t_n$  that

$$(h) \quad \frac{1}{2} W_n(t, P_n)(t - t_n) \\ = \frac{1}{n} \sum_{i=1}^n h^2(Z_i) - \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(Z_i) \\ + \frac{1}{n} \sum_{i=1}^n h_n(Z_i) h(Z_i) + \lambda_n \int t_n''(z) h''(z) dz$$



$$\begin{aligned}
 &+ \lambda_n \int (h''(z))^2 dz \\
 &= E_0(h^2(Z)) + \lambda_n \|h''\|_2^2 + \int h^2(z) d(P_n - P_0)(z) \\
 &\quad - \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(Z_i) + \int h_n(z) h(z) d(P_n - P_0)(z).
 \end{aligned}$$

But on  $\partial N_n$  the first two terms on the right side of (h) sum to  $\gamma_n^2$ . Moreover, for  $h \in \bar{N}_n$  as in (e)

(i)  $\|h'\|_\infty + \|h\|_\infty \leq C \gamma_n \lambda_n^{-1/2}$ .

Hence, for  $h \in \partial N_n$ , by (f)

(j)  $Var_0 h^2(Z) \leq E_0 h^4(Z) \leq \|h\|_\infty^2 E_0 h^2(Z) \leq C^2 \gamma_n^4 \lambda_n^{-1}$

and by (b)

(k)  $Var_0 h_n h(Z) \leq E_0 h_n^2 h^2(Z) \leq \|h\|_\infty^2 E_0 h_n^2(Z) \leq \frac{1}{4} C^2 \|t_0''\|_2^2 \gamma_n^2$ .

If  $h \in \partial N_n$ , it follows from (e), (i), (j), and (k), that we can apply remark A.6.1 to see that up to multiplicative constants  $\lambda_n \gamma_n^{-2} h^2$  and  $\lambda_n^{1/2} \gamma_n^{-1} h h_n$  belong to  $F_n(\alpha_n, r)$  satisfying the conditions of corollary A.6.3 with  $d = m = r = 1$ ,  $\alpha_n = \lambda_n$ ,  $\mu_n = \alpha_n^{1/4}$  provided that

(l)  $n^{1/2} \lambda_n^{3/4} \rightarrow \infty$ .

Note that remark A.6.1 with  $d = m = r = 1$  is applicable also to the functions  $\lambda_n^{1/2} \gamma_n^{-1} h$ ,  $h \in \partial N_n$ . Let  $w : R \rightarrow [0, 1]$  be a smooth delimiter, say  $w(t) = 1 \wedge [2 - |t|]^+$ . It easily follows that up to multiplicative constants  $\gamma_n^{-1} \lambda_n^{1/2} \sigma^{-1} (2 \log n)^{-1/2} \varepsilon h(x) w(\varepsilon \sigma^{-1} (2 \log n)^{-1/2})$  belongs to  $F_n(\alpha_n, r)$  with  $r = 1$ ,  $\alpha_n = \lambda_n / \log n$  and satisfies the conditions of corollary A.6.3 with  $\mu_n = \alpha_n^{1/4}$ , provided

(m)  $n^{1/2} \lambda_n^{3/4} (\log n)^{-3/4} \rightarrow \infty$ .

Since  $P(\max_i |\varepsilon_i| \geq \sigma \sqrt{2 \log n}) = o(1)$ , we can apply corollary A.6.3 (with a truncation argument) to conclude that

$$\begin{aligned}
 &\sup\{|\int h^2 d(P_n - P_0)| + |\int h_n h d(P_n - P_0)| + |n^{-1} \sum_{i=1}^n \varepsilon_i h(X_i)| : \\
 &\quad \|h\|_2^2 + \lambda_n \|h''\|_2^2 = \gamma_n^2\} \\
 &= O_p(\gamma_n^2 \lambda_n^{-3/4} n^{-1/2} + \gamma_n \lambda_n^{-1/4} n^{-1/2} (\log n)^{1/4}).
 \end{aligned}$$

This term is  $o_p(\gamma_n^2)$  if

(n)  $n^{1/2} \lambda_n^{3/4} \rightarrow \infty$ ,  $\gamma_n \lambda_n^{1/4} n^{1/2} (\log n)^{-1/4} \rightarrow \infty$ .

Therefore

$$\inf \left\{ \frac{1}{2} W_n(t, P_n)(t - t_n) : t \in \partial N_n \right\} = \gamma_n^2(1 + o_p(1))$$

if (l), (m), and (n) hold. We apply theorem 3 and part A follows from (f) with  $\gamma_n = \lambda_n^{-1/4} n^{-1/2} (\log n)^{1/4} d_n$ ,  $d_n \rightarrow \infty$  arbitrarily slowly. Part B follows from part A and (b) since  $n^{-1} \sqrt{\log n} \lambda_n^{-1/2} = o(\lambda_n)$ .  $\square$

### 7.6 ESTIMATION OF INFINITE-DIMENSIONAL PARAMETERS: ASYMPTOTICS AND APPLICATIONS

In this section we deal with asymptotic linearity and Gaussianity of asymptotic generalized *M*- (*AGM*-), and asymptotic generalized minimum contrast (*AGMC*-) estimates. We use the notation of the previous section.

We saw in examples 6 and 7 of section 7.5 that we can often only establish convergence results under weaker norms than that belonging to **B**. Analogously, in this section, we may only be able to justify linear approximations to estimates under weaker norms and obtain asymptotic normality for restricted classes of linear functionals of estimates. Here are the notions we need.

Let  $B_0^*$  be a bounded subset of  $B^*$ , the dual space of **B** such that  $b^*b \equiv b^*(b) = 0$  for all  $b^* \in B_0^*$  implies  $b = 0$ . Define

$$\|b\|_* \equiv \|b\|_{*B} \equiv \sup \{ |b^*(b)| : b^* \in B_0^* \}.$$

Then  $\|\cdot\|_*$  is a norm. As examples we mention:

- (i) If  $B_0^*$  is the unit ball in  $B^*$ ,  $\|b\|_* = \|b\|$ .
- (ii) Take  $B = l^\infty([0, 1])$  and identify  $B_0^*$  with the unit sphere in  $L_p([0, 1])$  via  $f \rightarrow b_f^*$  where  $b_f^*(b) \equiv \int fb \, d\mu$ . Then,  $\|b\|_* = \left( \int |b|^q \, d\mu \right)^{1/q}$  where  $1/p + 1/q = 1$ .

Now we give further appropriate definitions of asymptotic linearity:

**Definition 1.**  $\hat{v}_n$  is weakly  $B_0^*$  linear at  $P$  if and only if for each  $b^* \in B_0^* \subset B^*$  there exists  $\psi(\cdot, b^*, P) \in L_2(P)$  with  $E\psi(X, b^*, P) = 0$  such that

$$(1) \quad b^* \hat{v}_n = b^* v(P) + n^{-1} \sum_{i=1}^n \psi(X_i, b^*, P) + o_p(n^{-1/2}).$$

It is  $B_0^*$  linear at  $P$  if and only if there exists  $\psi(\cdot, P) : X \rightarrow B$  such that

$$(2) \quad \|\hat{v}_n - v(P) - n^{-1} \sum_{i=1}^n \psi(X_i, P)\|_* = o_p(n^{-1/2}).$$

It is linear at  $P$  if and only if (2) holds with  $\|\cdot\|_* = \|\cdot\|$ . Note that for economy of description we have dropped the qualifier ‘‘asymptotically’’ in our definition.

**Definition 2.** The estimator sequence  $\{\hat{v}_n\}$  is weakly  $B_0^*$  Gaussian regular if and only if, for all  $b^* \in B_0^*$ ,  $b^* \hat{v}_n$  is a Gaussian regular estimate of  $b^*v$ . The sequence  $\{\hat{v}_n\}$  is  $B_0^*$  (locally) Gaussian regular at  $P_0 \in P$ , if and only if there

exists a  $\mathbb{Z}_0$  such that for every regular curve  $\{P_\eta\}$  in  $\mathbf{P}$  through  $P_0$  and  $\eta_n = O(n^{-1/2})$ ,

$$\mathbb{Z}_n \equiv \sqrt{n}(\hat{v}_n - v(P_{\eta_n})) \Rightarrow \mathbb{Z}_0.$$

Here  $\mathbb{Z}_0$  is a Gaussian,  $\mathbf{B}$  valued random element which is tight and measurable with respect to the Borel  $\sigma$ -field generated by the  $\|\cdot\|_*$  topology.  $\mathbb{Z}_0$ , of course, depends on  $P_0$ . Again if  $\|\cdot\|_* = \|\cdot\|$ , we say  $\hat{v}_n$  is *Gaussian regular*.

Here are some useful relations between these properties.

**Proposition 1.**

A. If  $\hat{v}_n$  is weakly  $\mathbf{B}_0^*$  linear at  $P_0 \in \mathbf{P}$ ,  $v$  is weakly pathwise differentiable and

$$(3) \quad \psi(\cdot, b^*, P_0) - b^* \dot{v}(P_0) \perp \dot{\mathbf{P}}(P_0)$$

for all  $b^* \in \mathbf{B}_0^*$ , then  $\hat{v}_n$  is weakly  $\mathbf{B}_0^*$  Gaussian regular at  $P_0$ .

B. If  $\hat{v}_n$  is  $\mathbf{B}_0^*$  linear at  $P_0 \in \mathbf{P}$ , the conditions from A hold, and (for  $\{P_\eta\}$ ,  $\{\eta_n\}$ , and  $\mathbb{Z}_0$  as in definition 2)

$$(4) \quad n^{-1/2} \sum_{i=1}^n \psi(X_i, P_{\eta_n}) \Rightarrow \mathbb{Z}_0,$$

then  $\hat{v}_n$  is  $\mathbf{B}_0^*$  Gaussian regular at  $P_0$ .

**Proof.** A follows from proposition 3.3.1; B follows from A and the fact that tightness under  $P_0$  implies tightness under  $P_{\eta_n}$ . □

Linearity and Gaussianity for nonlinear estimates  $q(\hat{v}_n)$  where  $q$  maps  $\mathbf{B}$  to another Banach space  $\Gamma$  require linearity and Gaussianity of  $\hat{v}_n$  as well as the following alternate additional conditions proposed respectively by Gill (1989) and Wong and Severini (1991).

Let  $q, \hat{v}_n$  be as above. Let  $\Gamma_0^* = \{g^* \in \Gamma^* : g^* \dot{q}(v_0) \in \mathbf{B}_0^*\}$ . Let  $\|\cdot\|_{\mathbf{B}}$  be a norm on  $\mathbf{B}$  and  $\|\cdot\|_{*\Gamma}$  be the  $\Gamma_0^*$  norm on  $\Gamma$ .

(GD1)  $q$  is  $\|\cdot\|_{\mathbf{B}} - \|\cdot\|_{*\Gamma}$ -differentiable at  $v_0 \equiv v(P_0)$  with derivative  $\dot{q}(v_0)$ .

(GD2')  $\hat{v}_n$  is  $\mathbf{B}_0^*$  linear and Gaussian regular.

(GD2)  $\hat{v}_n$  is weakly  $\mathbf{B}_0^*$  linear and Gaussian regular.

(GD3) There exists  $\dot{q}(v_0)$  as above such that

$$(5) \quad q(v) = q(v_0) + \dot{q}(v_0)(v - v_0) + R(v, v_0)$$

where  $\|R(v, v_0)\|_{*\Gamma} = O(\|v - v_0\|_{\mathbf{B}}^2)$ . Note that  $\dot{q}(v_0)$  may not be unique.

(GD4)  $\|\hat{v}_n - v_0\|_{\mathbf{B}} = o_p(n^{-1/4})$ .

**Proposition 2. The delta method**

A. If (GD1) and (GD2') hold, then  $q(\hat{v}_n)$  is  $\Gamma_0^*$  linear and Gaussian regular.

B. If (GD2), (GD3), and (GD4) hold, then  $q(\hat{v}_n)$  is weakly  $\Gamma_0^*$  linear and Gaussian regular.

**Proof.** If (GD2') holds, then  $n^{1/2}(\hat{v}_n - v_0)$  is tight in the  $\|\cdot\|_B$  topology and hence  $\|\hat{v}_n - v_0\|_B = o_p(1)$ . But (GD1) implies that

$$\|q(\hat{v}_n) - q(v_0) - \dot{q}(v_0)(\hat{v}_n - v_0)\|_{*\Gamma} = o_p(n^{-1/2}),$$

see, for instance, Fernholz (1983, page 19). Our first claim follows. The second claim is immediate since  $\|R(\hat{v}_n, v_0)\|_{*\Gamma} = o_p(n^{-1/2})$  by (GD3) and (GD4).  $\square$

### Explicitly Defined Estimates

If  $\hat{v}_n$  is given explicitly via a functional  $v_n$  on a convex set  $M_0 \supset P \cup \{\text{distributions with finite support}\}$ ,  $\hat{v}_n = v_n(\mathcal{I}P_n)$ , we can try to linearize using the delta method as in the Euclidean case. That is:

(i) We calculate formally,  
 $\psi_n(x, b^*, P) \equiv (\partial/\partial \varepsilon) b^* v_n((1 - \varepsilon)P + \varepsilon \delta_x) |_{\varepsilon=0}$ .

(ii) We show that

$$(6) \quad b^* \hat{v}_n = b^* v(P) + n^{-1} \sum_{i=1}^n \psi_n(X_i, b^*, P) + o_p(n^{-1/2}),$$

using the techniques discussed in sections 7.2 and 7.3.

(iii) We determine  $\psi(\cdot, b^*, P)$ , the  $L_2(P)$  limit of  $\psi_n(\cdot, b^*, P)$ .

(iv) We check the conditions of proposition 1.

With (i)–(iv) we have established weak  $B_0^*$  linearity, regularity, and Gaussianity. As we shall see in the following examples it is usually just as easy to apply steps (i)–(iii) directly to  $\hat{v}_n$  rather than  $b^* \hat{v}_n$  to obtain  $\psi_n(\cdot, P)$ ,  $\psi(\cdot, P)$  and establish linearity as well as weak Gaussian regularity. To obtain Gaussian regularity, we need to involve a functional central limit theorem as we indicated in (4) of proposition 1.

#### Example 1. The Nelson-Aalen and Kaplan-Meier estimators.

In the situation of example 7.5.5 we assume that  $\tau_0 > 0$  satisfies  $P_0(T > \tau_0) = (1 - G_{01}(\tau_0))(1 - G_{02}(\tau_0)) > 0$ . Let  $v(P) = \Lambda_1$ , be the cumulative hazard function on  $[0, \tau_0]$ . The subscript 0 is used consistently to refer to quantities calculated under  $P_0$ .

View  $\Lambda_1$  (and  $H_{n1}$ , etc.) as an object in  $B \equiv l^\infty(F)$  where  $F = \{1_{[0,s]} : s \leq \tau_0\}$  via the correspondence  $\Lambda_1(1_{[0,s]}) = \Lambda_1(s)$ . Then, the Nelson-Aalen estimate of  $\Lambda_1$  is  $\hat{\Lambda}_1 = q(\mathcal{I}P_n)$  where, in view of (7.5.11),

$$(7) \quad q(P) \equiv \int_0^\cdot (1 - H(t-))^{-1} dH_1(t).$$

**Proposition 3.** The estimate  $\hat{\Lambda}_1$  is linear, Gaussian regular with influence function given by (6.6.7). The Kaplan-Meier estimate  $\hat{G}_{KM}$  of  $G_1$ , is linear Gaussian regular with influence function given by (6.6.14).

**Proof.** Let  $\bar{H}_n(s) \equiv 1 - H_n(s)$ . On the event  $[\bar{H}_n(y) > 0]$

$$\sqrt{n}(\hat{\Lambda}_1(y) - \Lambda_{01}(y)) = \sqrt{n} \int_0^y \{1_{[\bar{H}_n(t-)>0]} \frac{1}{\bar{H}_n(t-)} dH_{n1}(t) - d\Lambda_{01}(t)\},$$

where the process on the right side is a martingale. Hence it converges weakly to a Gaussian process by a martingale central limit theorem, see, e.g., Shorack and Wellner (1986, theorem 7.1.1, page 298). A bivariate central limit theorem for martingales and a contiguity argument yield Gaussian regularity. We can also establish linearity and give an alternative proof of Gaussian regularity.

Write

$$\begin{aligned}
 (a) \quad & \sqrt{n} \left\{ \int_0^y \frac{dH_{n1}(t)}{1 - H_n(t-)} - \int_0^y \frac{dH_{01}(t)}{1 - H_0(t-)} \right\} \\
 &= \sqrt{n} \int_0^y \frac{d(H_{n1} - H_{01})(t)}{1 - H_0(t-)} \\
 &\quad + \int_0^y \frac{\sqrt{n}(H_n - H_0)(t-)}{(1 - H_0(t-))^2} dH_{01}(t) \\
 &\quad + \sqrt{n} \int_0^y \frac{(H_n(t-) - H_0(t-))^2}{(1 - H_n(t-))(1 - H_0(t-))^2} dH_0(t) \\
 &\quad + \sqrt{n} \int_0^y \frac{H_n(t-) - H_0(t-)}{(1 - H_n(t-))(1 - H_0(t-))} d(H_{n1} - H_{01})(t).
 \end{aligned}$$

The sum of the first two terms is a continuous function of the empirical process  $(\sqrt{n}(H_{n1} - H_{01}), \sqrt{n}(H_n - H_0))$  into  $l^\infty(\mathbb{F})$  and hence it converges weakly by Donsker's theorem. The first of the two remainder terms is  $O_p(\sqrt{n} \|H_n - H_0\|_\infty^2) = O_p(n^{-1/2})$ .

Write

$$\begin{aligned}
 (b) \quad & \sqrt{n} \int_0^y \frac{H_n(t-) - H_0(t-)}{(1 - H_n(t-))(1 - H_0(t-))} d(H_{n1}(t) - H_{01}(t)) \\
 &= \int_0^{H_{n1}(y)} W_n(H_{n1}^{-1}(s)) ds - \int_0^{H_{01}(y)} W_n(H_{01}^{-1}(s)) ds \\
 &= \int_0^{H_{n1}(y)} (W_n(H_{n1}^{-1}(s)) - W_n(H_{01}^{-1}(s))) ds \\
 &\quad - \int_{H_{n1}(y)}^{H_{01}(y)} W_n(H_{01}^{-1}(s)) ds,
 \end{aligned}$$

where

$$W_n(t) \equiv \sqrt{n} \frac{H_n(t-) - H_0(t-)}{(1 - H_n(t-))(1 - H_0(t-))}.$$

Therefore, the second remainder term is

$$\begin{aligned}
 & O_p(\sup\{|W_n(H_{n1}^{-1}(s)) - W_n(H_{01}^{-1}(s))| : |s| \leq H_{n1}(\tau_0)\}) \\
 & + o_p(\sup\{|W_n(H_{01}^{-1}(s))| : |s| \leq \max\{H_{n1}(\tau_0), H_{01}(\tau_0)\}\}) = o_p(1)
 \end{aligned}$$

since, by Donsker's theorem,  $W_n(\cdot)$  converges weakly to a continuous process and  $\sup\{|H_n^{-1}(s) - H_{01}^{-1}(s)| : 0 \leq s \leq \tau_0\} = O_p(n^{-1/2})$ . That the influence function is as given by (6.6.7) follows from the calculations of Example 6.6.1.A. In this example we have checked that the influence function of this estimate is efficient. Regularity follows by proposition 3.3.1.

Apply proposition 2 to see that  $\hat{G}_1(\cdot) \equiv \exp(-\hat{\Lambda}_1(\cdot))$  is linear and Gaussian since  $\Lambda \rightarrow e^{-\Lambda}$  is Fréchet differentiable. Breslow and Crowley (1974) show by an elementary argument that  $\|\hat{G}_1 - \hat{G}_{KM}\|_\infty = O_p(n^{-1})$ , and linearity and Gaussian regularity of  $\hat{G}_{KM}$  follow.  $\square$

**Example 2. Estimating a joint distribution with one marginal known.**

This is a continuation of example 7.5.8. We take  $\mathbf{B} \equiv l^\infty(\mathbf{F})$  for  $\mathbf{F}$  a bounded subset of  $l^\infty([0, 1]^2)$  whose members are measurable and identify  $P(f) \equiv \int f dP$  for  $f \in \mathbf{F}$ . Identify  $\mathbf{B}_0^*$  with  $\mathbf{F}$  via the usual  $f(b) \equiv b(f)$ , so  $\|b\|_* = \sup\{|b(f)| : \|f\| \leq 1\}$ .

Without loss of generality we take  $F_0$  as the uniform distribution on  $[0, 1]$ , since we can always "probability transform" ourselves to that case if  $F_0$  is continuous. Consider  $\hat{P}_n$  defined as follows. The conditional distribution of  $X = (U, V)$  given  $(U, V) \in I_{jkm}$  under  $\hat{P}_n$  is the conditional empirical distribution and  $\hat{P}_n(I_{jkm})$  is given by (7.5.24) with  $m = m_n$ . Let  $\hat{P}_n(f) = \int f d\hat{P}_n$  for  $f \in \mathbf{F}$ .

**Proposition 4.** If  $\hat{P}_n$  is as above and  $m_n^2/n = o(1)$ , then  $\hat{P}_n$  is weakly  $\mathbf{B}_0^*$  linear and Gaussian regular with influence function

$$(8) \quad \psi(X, f, P_0) = f(U, V) - E_0(f(U, V) | U).$$

**Proof.** Write  $m \equiv m_n$ ,  $J_{ij} \equiv 1_{[(j-1)/m \leq U_i < j/m]}$  and note that  $f_{jm0}$  from (7.5.22) equals  $m^{-1}$ . Then

$$\int f(u, v) d\hat{P}_n(u, v) = \sum_{j=1}^m \frac{1}{mN_{j+m}} \sum_{i=1}^n f(U_i, V_i) J_{ij}$$

where  $0/0 = 0$ . Let

$$f_m(u) = \sum_{j=1}^m 1_{[(j-1)/m \leq u < j/m]} E_0(f(U, V) | U \in [\frac{j-1}{m}, \frac{j}{m})).$$

Then

$$\begin{aligned} & \int f(u, v) d\hat{P}_n(u, v) - \int f(u, v) dP_0(u, v) \\ &= \sum_{j=1}^m \frac{1}{mN_{j+m}} \sum_{i=1}^n (f(U_i, V_i) - f_m(U_i)) J_{ij} \\ &= \frac{1}{n} \sum_{i=1}^n (f(U_i, V_i) - f_m(U_i)) \\ &\quad - \sum_{j=1}^m \left(\frac{N_{j+m}}{n} - \frac{1}{m}\right) \frac{1}{N_{j+m}} \sum_{i=1}^n (f(U_i, V_i) - f_m(U_i)) J_{ij} \end{aligned}$$

$$= \frac{1}{n} \sum_{i=1}^n [f(U_i, V_i) - E_0(f(U_i, V_i) | U_i)] - \frac{1}{n} \sum_{i=1}^n [f_m(U_i) - E_0(f(U_i, V_i) | U_i)] + R_n.$$

We claim that

(a) 
$$\frac{1}{\sqrt{n}} \sum_{i=1}^n [f_m(U_i) - E_0(f(U_i, V_i) | U_i)] = o_p(1),$$

(b)  $|R_n| = o_p(n^{-1/2}).$

The left side of (a) has mean 0 and variance

$$\int [f_m(u) - E_0(f(U, V) | U=u)]^2 du,$$

which tends to zero as  $n \rightarrow \infty$  since the integrand does for Lebesgue almost all  $u$  (cf. Theorems 18.4 and 18.5 of Hewitt and Stromberg (1965)). To verify (b), check that

$$\begin{aligned} |R_n| &= \left| \sum_{j=1}^m \left( \frac{N_{j+m}}{n} - \frac{1}{m} \right) \frac{1}{N_{j+m}} \sum_{i=1}^n (f(U_i, V_i) - f_m(U_i)) J_{ij} \right| \\ &\leq \left\{ \sum_{j=1}^m \left( \frac{N_{j+m}}{n} - \frac{1}{m} \right)^2 \sum_{j=1}^m \left( \frac{1}{N_{j+m}} \sum_{i=1}^n (f(U_i, V_i) - f_m(U_i)) J_{ij} \right)^2 \right\}^{1/2}. \end{aligned}$$

Now

$$E \sum_{j=1}^m \left( \frac{N_{j+m}}{n} - \frac{1}{m} \right)^2 = \frac{1}{n} \left( 1 - \frac{1}{n} \right).$$

Furthermore,

(c) 
$$\begin{aligned} E_0 \left\{ \sum_{j=1}^m \left( \frac{1}{N_{j+m}} \sum_{i=1}^n (f(U_i, V_i) - f_m(U_i)) J_{ij} \right)^2 \mid N_{1+m}, \dots, N_{m+m} \right\} \\ = \sum_{j=1}^m \frac{m}{N_{j+m}} \int 1_{[(j-1)/m \leq u < j/m]} (f(u, v) - f_m(u))^2 dP_0(u, v) 1_{[N_{j+m} > 0]}. \end{aligned}$$

But,

$$P \left( \inf_{1 \leq j \leq m} \frac{N_{j+m}}{n} \leq \frac{1}{m} - \frac{\log m}{\sqrt{mn}} \right) \leq m \exp \left[ -\frac{1}{2} \log^2 m (1 + o(1)) \right]$$

by Bernstein's inequality. Therefore

$$\inf_j \frac{N_{j+m}}{n} = \frac{1}{m} (1 + O_p(\sqrt{\frac{m}{n}} \log m)).$$

Hence, the conditional expectation in (c) is of the order

$$\frac{m^2}{n} (1 + o_p(1)) E_0(f(U, V) - f_m(U))^2 = o_p(1),$$

and (b) follows.

We conclude that the estimator is weakly  $B_0^*$  linear with influence function (8), and by proposition 1.A is weakly  $B_0^*$  Gaussian regular.  $\square$

Note that by (6.2.10)  $\hat{I}_n^P$  is efficient in the sense of theorem 5.2.2. Strong regularity and Gaussianity hold if  $F$  is reduced to a class such that the theory of Pollard (1984) can be applied. See also theorem B of Wellner (1989).

It is shown in Bickel, Ritov, and Wellner (1991) that weak linearity holds for case (ii) of example 7.5.8 as well.

Note that the results depend rather crucially on the nature of the sieve. For instance, if we used the "natural" sieve which specifies that the conditional distribution of  $V$  given  $U = u$  is the same for all  $u \in [(j-1)/m, j/m)$ ,  $1 \leq j \leq m$ , and require the appropriate marginal condition, then linearity would hold in general only for  $f$  which are functions of  $V$  only.  $\square$

**Example 3. Biased sampling.**

We use the notation of example 7.5.4 and consider estimation of  $Q$  viewed as an element of  $l^\infty(F)$ ,  $F$  a bounded subset of  $l^\infty(R)$  whose members are measurable,  $B_0^*$  as in the previous example so that  $\|b\|_* = \sup\{|b(f)| : f \in F\}$ .

Equations (7.5.7) and (7.5.8) define  $\hat{Q}_n$  implicitly.

**Proposition 5.** Suppose that

$$(9) \quad P_0(a(I) + b(X) = c) < 1$$

for every constant  $c$  and all functions  $a$  and  $b$  which are not both constant, and

$$(10) \quad P_0(\sum_a w_a(X) \theta_{0a} \geq \varepsilon) = 1$$

for some  $\varepsilon > 0$ . Then,  $\hat{Q}_n$  is a weakly  $B_0^*$  linear, Gaussian regular estimate of  $Q$ .

**Proof.** Substitute  $\hat{q}_k$  from (7.5.7) into (7.5.6) to obtain the  $s - 1$  equations

$$(a) \quad \hat{\lambda}_j^* - \hat{\theta}_j \sum_k \hat{q}_k^* w_j(x_k) \{ \sum_a \hat{\theta}_a w_a(x_k) \}^{-1} = 0.$$

Thus  $\hat{\theta} = \theta(\hat{I}_n)$  is a GM-estimate where

$$W_j(\theta, P) \equiv \int \{ 1_{[I=j]} - \theta_j w_j(x) [\sum_a \theta_a w_a(x)]^{-1} \} dP(j, x).$$

It is shown by Gill, Vardi, and Wellner (1988) that if (9) and (10) hold then,  $W \equiv (W_1, \dots, W_{s-1})$  satisfies the conditions of theorem 7.3.3. Hence, with probability tending to 1,  $\hat{\theta}$  exists and is linear and Gaussian regular.

Write  $\theta_0$  for  $\theta(P_0)$  and  $c(\theta, Q) \equiv \int [\sum_a \theta_a w_a(x)]^{-1} dQ(x)$ . Now (7.5.8) gives us an explicit representation of  $\hat{Q}_n$  in terms of  $\hat{\theta}$ ,

$$(b) \quad \int f(x) d\hat{Q}_n(x) = \hat{c}^{-1} \int f(x) [\sum_a \hat{\theta}_a w_a(x)]^{-1} d\hat{Q}_n^*(x),$$

where  $\hat{Q}_n^*$  is the empirical d.f. of  $X$ ,  $\hat{c} = c(\hat{\theta}, \hat{Q}_n^*)$ . Expansion around  $(\theta_0, Q_0^*)$  yields

$$(c) \quad \int f(x) [\sum_a w_a(x) \hat{\theta}_a]^{-1} d\hat{Q}_n^*(x)$$



$$\begin{aligned}
 &= \int f(x) dQ_0(x) + \int f(x) \left[ \sum_a \theta_{0a} w_a(x) \right]^{-1} d(\hat{Q}_n^* - Q_0^*)(x) \\
 &\quad - \sum_b (\hat{\theta}_b - \theta_{0b}) \left( \int f(x) w_b(x) \left[ \sum_a \theta_{0a} w_a(x) \right]^{-1} dQ_0(x) \right) \\
 &\quad + \|f\|_\infty O_p(n^{-1}),
 \end{aligned}$$

where  $Q_0^*$  is the marginal distribution of  $X$  under  $P_0$ . In particular for  $f \equiv 1$ ,

$$\begin{aligned}
 (d) \quad \hat{c} &= 1 + \int \left[ \sum_a \theta_{0a} w_a(x) \right]^{-1} d(\hat{Q}_n^* - Q_0^*)(x) \\
 &\quad - \sum_b (\hat{\theta}_b - \theta_{0b}) \left\{ \int w_b(x) \left[ \sum_a \theta_{0a} w_a(x) \right]^{-1} dQ_0(x) \right\} \\
 &\quad + O_p(n^{-1}).
 \end{aligned}$$

Linearity follows as do weak  $B_0^*$  regularity, Gaussianity, and efficiency. For details see Gill, Vardi, and Wellner (1988).  $\square$

*AGM-estimates: Linearity and Gaussianity*

Our formulation here owes much to Filippova (1962) and Wong and Severini (1991). In examples 7.5.7 and 7.5.10 we have seen that the representation of  $P$  or  $G$  which appears naturally in the definition of  $W_n$  is  $v = \log p$  or  $\log g$ , quantities which can typically *not* be estimated at rate  $n^{-1/2}$ . Yet, equivalent parameters  $h(v)$  may be so estimated. For instance, consider the regularized maximum likelihood estimate of  $p$  given by  $p_n^\#(x) \equiv b_n^{-1} \int \phi((x-y)/b_n) dP_n(y)$ . Then,  $n^{1/2} \|p_n^\# - p\| \rightarrow \infty$  in probability if  $b_n \rightarrow 0$  at an appropriate rate for any reasonable norm. Yet if  $h \in L_2(P)$ ,

$$\int h p_n^\# = \int h p_n + n^{-1} \sum_{i=1}^n \{h_n(X_i) - \int h p_n\},$$

where

$$p_n(x) \equiv b_n^{-1} \int \phi\left(\frac{x-y}{b_n}\right) p(y) dy, \quad h_n(x) = \frac{1}{b_n} \int h(y) \phi\left(\frac{x-y}{b_n}\right) dy.$$

If

$$(11) \quad \|h'\|_\infty + \|h''\|_\infty < \infty,$$

then

$$\begin{aligned}
 (12) \quad \int h(x) p_n(x) dx &= \iint h(y + b_n z) p(y) \phi(z) dy dz \\
 &= \int h(x) p(x) dx + O(b_n^2).
 \end{aligned}$$

Also,

$$E(h_n(X_1) - h(X_1))^2 = o(1).$$

We deduce that if  $b_n = o(n^{-1/4})$  and  $h$  satisfies (11), then  $\int h p_n^\#$  is an efficient estimate of  $\int h p$  in example 7.5.2. This is an example of an application of corollary 2 below. To deal with such cases we focus on linearization of  $W_n$  in  $v$  and  $P_n$  and then deduce linearity and Gaussianity of suitable  $h(\hat{v}_n)$ .

We shall follow our approach of section 7.5 by first giving general assumptions leading to a trivial master theorem and then giving, in examples 4–6, various ways of verifying these conditions.

We write  $v_0 \equiv v(P_0)$  and require that statements hold for all  $P_0 \in \mathbf{P}$ . Let  $\|\cdot\|_{\mathbf{B}}$ , be a norm on  $\mathbf{B}$ . Let  $\Gamma_0^*$  be a bounded subset of  $\Gamma^*$  such that  $g^*(g) = 0$  for all  $g^* \in \Gamma_0^*$  implies  $g = 0$  and denote the corresponding norm on  $\Gamma$  by  $\|\cdot\|_{*\Gamma}$ .

The assumptions are given in two parts, the first being implied by the second.

(GGM0) For all  $g^* \in \Gamma_0^*$ ,  $g^* W_n(v, P) = g^* W(v, P) + o(1)$  for all  $v \in \mathbf{A}$ ,  $P \in \mathbf{M}$ . Furthermore,

$$(13) \quad g^* W_n(v_0, P_0) = o(n^{-1/2}).$$

(GGM0') (GGM0) with (13) replaced by  $\|W_n(v_0, P_0)\|_{*\Gamma} = o(n^{-1/2})$ .

$$(GGM1) \quad \|W_n(\hat{v}_n, P_n) - W_n(\hat{v}_n, P_0) - W_n(v_0, P_n) + W_n(v_0, P_0)\|_{*\Gamma} = o_p(n^{-1/2}).$$

(GGM2) For all  $g^* \in \Gamma_0^*$ ,

$$(14) \quad g^* W_n(v_0, P_n) = g^* W_n(v_0, P_0) + n^{-1} \sum_{i=1}^n g^* w(X_i, P_0) + o_p(n^{-1/2}).$$

$$(GGM2') \quad \|W_n(v_0, P_n) - W_n(v_0, P_0) - n^{-1} \sum_{i=1}^n w(X_i, P_0)\|_{*\Gamma} = o_p(n^{-1/2}).$$

Note that if (GGM0) (respectively (GGM0')) holds, then  $W_n(v_0, P_0)$  can be absorbed into  $o_p(n^{-1/2})$  here.

(GGM3)  $W_n(\cdot, P_0)$  (respectively  $W(\cdot, P_0)$ ) is Hadamard differentiable on  $\mathbf{A}$  at  $v_0$  with derivative  $\dot{W}_n \equiv \dot{W}_n(v_0, P_0)$  (respectively  $\dot{W}$ ).

(GGM4) For all  $g^* \in \Gamma_0^*$ ,

$$(15) \quad g^* W_n(\hat{v}_n, P_0) - g^* W_n(v_0, P_0) = g^* \dot{W}_n(\hat{v}_n - v_0) + o_p(n^{-1/2})$$

and

$$g^* (\dot{W}_n - \dot{W})(\hat{v}_n - v_0) = o_p(n^{-1/2}).$$

$$(GGM4') \quad \|W_n(\hat{v}_n, P_0) - W_n(v_0, P_0) - \dot{W}_n(\hat{v}_n - v_0)\|_{*\Gamma} = o_p(n^{-1/2})$$

and

$$(16) \quad \|(\dot{W}_n - \dot{W})(\hat{v}_n - v_0)\|_{*\Gamma} = o_p(n^{-1/2}).$$

These conditions correspond to (GM0)–(GM3) in section 7.3. The only assumption that is missing is the invertibility of  $\dot{W}_n$  (or its limit  $\dot{W}$ ). This assumption can fail in the infinite-dimensional case; see examples 4 and 5. Wong and Severini (1991) noted that useful conclusions could still be drawn even if  $\dot{W}_n$  is singular. We can think of (17) below as linearization of the estimate  $\dot{W}_n(\hat{v}_n - v_0)$  of the locally identifiable part of the parameter  $\dot{W}_n(v - v_0)$ .

**Theorem 1.** Suppose that  $\hat{v}_n$  satisfies  $W_n(\hat{v}_n, P_n) = o_p(n^{-1/2})$  and (GGM1) and (GGM3) hold.

A. Under (GGM0), (GGM2), and (GGM4),

$$(17) \quad g^* \dot{W}(\hat{v}_n - v_0) = -n^{-1} \sum_{i=1}^n g^* w(X_i, P_0) + o_p(n^{-1/2}).$$

B. Under (GGM0'), (GGM2'), and (GGM4'),

$$(18) \quad \|\dot{W}(\hat{v}_n - v_0) + n^{-1} \sum_{i=1}^n w(X_i, P_0)\|_{*\Gamma} = o_p(n^{-1/2}).$$

**Proof.** The proof is trivial. We interpret  $o_p$  in terms of the appropriate norms. For instance, for part B, write

$$\begin{aligned} -W_n(v_0, P_n) &= W_n(\hat{v}_n, P_n) - W_n(v_0, P_n) + o_p(n^{-1/2}), && \text{by definition of } \hat{v}_n \\ (a) \quad &= W_n(\hat{v}_n, P_0) - W_n(v_0, P_0) + o_p(n^{-1/2}), && \text{by (GGM1)} \\ &= \dot{W}(\hat{v}_n - v_0) + o_p(n^{-1/2}) && \text{by (GGM3) and (GGM4')}. \end{aligned}$$

Apply (GGM0') and (GGM2') to obtain the conclusion.  $\square$

The following two corollaries are immediate. They give conditions under which we can obtain asymptotic linearity and Gaussianity of linear and nonlinear functions of  $\hat{v}_n$ .

**Corollary 1** (Linearity and Gaussianity). Specify  $\Gamma_0^* \subset \Gamma^*$  and let  $B_0^* = \{b^* \in B^*: b^* = g^* \dot{W} \text{ for some } g^* \in \Gamma_0^*\}$ . Then:

- If the conditions of theorem 1.A hold, where here and subsequently  $b^* = g^* \dot{W}$ , then  $\hat{v}_n$  is weakly  $B_0^*$  linear with  $\psi(\cdot, b^*, P) = -g^* w(\cdot, P)$ .
- If the conditions of theorem 1.B hold, then  $\hat{v}_n$  is  $B_0^*$  linear.
- If, in addition, for each  $b^* \in B_0^*$ ,  $g^* \hat{v}_n$  satisfies the conditions of proposition 1.A then  $\hat{v}_n$  is weakly  $B_0^*$  Gaussian regular.

**Corollary 2.** Suppose the conditions of corollary 1.B hold. Let  $\|\cdot\|_B = \|\cdot\|_{*B}$  with  $B_0^*$  as above. Suppose that:

- $q: B \rightarrow R$  satisfies (GD3).
- $\dot{q}(v_0) \in B_0^*$ .
- $\hat{v}_n$  satisfies (GD4).

Then  $q(\hat{v}_n)$  is weakly  $\Gamma_0^*$  linear and Gaussian regular.

We can evidently extend corollary 2 to  $q$  Banach-valued. We leave this to the reader.

**Example 4. Density estimation by sieves.**

In the notation we have established in example 7.5.7,

$$W_n(\hat{v}_n, P_n) = 0,$$

where  $W_n : \mathbf{B} \times \mathbf{M} \rightarrow \mathbf{B}^* \equiv \Gamma$  is given by

$$(19) \quad \langle W_n(b_1, P), b_2 \rangle \equiv - \int \pi_m b_2 dP \\ + \int \pi_m b_2(x) \exp[\pi_m b_1(x)] dx \left( \int \exp[\pi_m b_1(x)] dx \right)^{-1},$$

where  $\langle b^*, b \rangle \equiv b^*(b)$  and  $\mathbf{B} \subset L^\infty([0,1])$  is a Sobolev space with norm  $\|b\|_{\mathbf{B}} = \|b\|_\infty + \|b''\|_\infty$ . Evidently,  $W$  is defined by replacing  $\pi_m b_1$ ,  $\pi_m b_2$  by  $b_1$ ,  $b_2$  throughout. If  $P_0$  has density  $p_0$ ,  $\langle W(v_0, P_0), b_2 \rangle = 0$  for all  $b_2 \in \mathbf{B}$  if and only if  $p_0(x) = \exp[v_0(x)] / \int \exp[v_0(y)] dy$ . In particular

$$v_0 = \log p_0 - \int \log p_0(x) dx$$

is a solution. Further,  $\dot{W}_n(b, P)$  doesn't depend on  $P$  and

$$\langle \dot{W}_n(b_1), b_2 \rangle = \text{Cov}_m(\pi_m b_1, \pi_m b_2),$$

where the subscript  $m$  of the covariance denotes computation under the density proportional to  $\exp[\pi_m v_0(x)]$ . Finally,

$$(20) \quad \langle \dot{W}(b_1), b_2 \rangle = \int (b_2(x) - E_0 b_2) b_1(x) p_0(x) dx = \text{Cov}_0(b_1, b_2).$$

**Proposition 6.** Let  $m = o(n^{1/2})$ ,  $mn^{-1/4} \rightarrow \infty$ . Then

$$(21) \quad \langle \dot{W}(\hat{v}_n - v_0), b \rangle = n^{-1} \sum_{i=1}^n [b(X_i) - E_0 b(X_i)] + o_p(n^{-1/2})$$

for all  $b \in \mathbf{B}$ . That is, if we identify  $\mathbf{B}_0^*$  with  $\{b \in \mathbf{B} : \|b\| \leq 1, \int b p_0 = 0\}$  with the convention  $b(h) = \int hb p_0$ , then  $\hat{v}_n$  is weakly  $\mathbf{B}_0^*$  linear Gaussian regular.

**Proof.** Let  $\Gamma_0^*$  = unit ball in  $\mathbf{B}$  viewed as usual as a subset of  $\Gamma^* = \mathbf{B}^{**}$ , by the identification  $g^* \leftrightarrow b_{g^*}$  where  $g^*(b^*) \equiv b^*(b_{g^*})$ . Then, indeed,  $g^* \dot{W}(b) = \langle \dot{W}(b), b_{g^*} \rangle$ , and we see that the proposition follows from corollary 1.A once we have verified (GGM1), (GGM3) and (GGM0), (GGM2') and (GGM4'). Since  $\|v_0 - \pi_m v_0\|_\infty = O(m^{-2})$  by (f) of the proof of proposition 7.5.1, (GGM0) holds. (GGM3) holds because of

$$\|W(v, P_0) - W(v_0, P_0) - \dot{W}(v - v_0)\|_{\Gamma^*} \\ = \sup\{ | \text{Cov}_0(b, \frac{e^{v-v_0}}{E_0 e^{v-v_0}} - (v - v_0)) | : b \in \Gamma_0^* \} \\ = (E_0 e^{v-v_0})^{-1} \sup\{ | \text{Cov}_0(b, e^{v-v_0} - 1 - (v - v_0)) |$$

$$+ |Cov_0[b, (v - v_0)(1 - E_0 e^{v-v_0})]| : b \in \Gamma_0^*$$

$$= O(\|v - v_0\|_\infty^2) = o(\|v - v_0\|_B).$$

(GGM1) is also immediate since

$$W_n(\hat{v}_n, P_n) - W_n(v_0, P_n) = W_n(\hat{v}_n, P_0) - W_n(v_0, P_0).$$

We now establish (GGM4').

Let the subscript  $n\lambda$  indicate that computation is carried out under the density  $p_{n\lambda} \equiv \exp(v_{n\lambda}) / \int \exp(v_{n\lambda}(y)) dy$  where  $v_{n\lambda} \equiv v_0 + \lambda(\hat{v}_n - v_0)$ . Then calculate

$$\begin{aligned} \frac{\partial^2}{\partial \lambda^2} \langle W_n(v_{n\lambda}, P_0), b \rangle &= Cov_{n\lambda}(\pi_m b, [\pi_m(\hat{v}_n - v_0)]^2) \\ &\quad - 2Cov_{n\lambda}(\pi_m b, \pi_m(\hat{v}_n - v_0)) E_{n\lambda} \pi_m(\hat{v}_n - v_0) \\ &= Cov_{n\lambda}(\pi_m b, \Delta_n^2) \end{aligned}$$

where

$$\Delta_n \equiv \pi_m(\hat{v}_n - v_0) - E_{n\lambda} \pi_m(\hat{v}_n - v_0).$$

Since, by proposition 7.5.1.A,  $\sup_\lambda \|p_{n\lambda} - p_{n0}\|_\infty = o(1)$ , we deduce that

$$| \frac{\partial^2}{\partial \lambda^2} \langle W_n(v_{n\lambda}, P_0), b \rangle | \leq (1 + o(1)) \|b\|_\infty \int \Delta_n^2 p_0.$$

Hence,

$$\begin{aligned} &| \langle W_n(\hat{v}_n, P_0) - W_n(v_0, P_0) - \dot{W}_n(\hat{v}_n - v_0), b \rangle | \\ &= | \int_0^1 (1 - \lambda) \frac{\partial^2}{\partial \lambda^2} \langle W_n(v_{n\lambda}, P_0), b \rangle d\lambda | \\ &\leq \frac{1}{2} (1 + o(1)) \|b\|_\infty \int \Delta_n^2 p_0 \\ &= O_p(m/n) + O_p(m^{-4}) = o_p(n^{-1/2}) \end{aligned}$$

by part B of proposition 7.5.1. Further, since  $\pi_m \hat{v}_n = \hat{v}_n$

$$\begin{aligned} &\|(\dot{W}_n - \dot{W})(\hat{v}_n - v_0)\|_{*\Gamma} \\ &= \sup_{b \in \Gamma_0} |Cov_m(\hat{v}_n - \pi_m v_0, \pi_m b) - Cov_0(\hat{v}_n - v_0, b)| \\ &= \sup_{b \in \Gamma_0} |Cov_m(\hat{v}_n - \pi_m v_0, \pi_m b) - Cov_0(\hat{v}_n - \pi_m v_0, \pi_m b) \\ &\quad + Cov_0(v_0 - \pi_m v_0, \pi_m b) - Cov_0(\hat{v}_n - v_0, b - \pi_m b)| \end{aligned}$$

$$\begin{aligned}
&= O_p(\|\hat{v}_n - \pi_m v_0\|_\infty \|v_0 - \pi_m v_0\|_\infty + \|v_0 - \pi_m v_0\|_\infty \\
&\quad + \|\hat{v}_n - v_0\|_\infty m^{-2}) \\
&= O_p(m^{-2}) = o_p(n^{-1/2}),
\end{aligned}$$

by proposition 7.5.1.A and (f) of its proof. (GGM4') follows.

Finally, (GGM2) holds since

$$\begin{aligned}
\langle W_n(v_0, P_n), b_2 \rangle - \langle W_n(v_0, P_0), b_2 \rangle &= - \int \pi_m b_2 d(P_n - P_0) \\
&= - \int b_2 d(P_n - dP_0) \\
&\quad + O_p(\|b_2\|_\infty m^{-2}).
\end{aligned}$$

The proposition is proved.  $\square$

**Remark 1.** Since the map  $v \rightarrow \log p = v - \log \int e^{v(x)} dx$  from  $\mathbf{B}$  to  $\mathbf{B}$  is Fréchet differentiable at  $v_0$ , it is easy to see that (21) continues to hold if  $v = \log p - \int \log p$  is replaced by  $\log p$ .

**Remark 2.** By proposition 7.5.1 the sequences  $\{m(n)\}$  satisfying the conditions of proposition 6 are such that (GD4) holds:

$$(22) \quad \|\hat{v}_n - v_0\|_2 = o_p(n^{-1/4}).$$

**Remark 3.**  $\dot{W}$  does not have a bounded inverse on  $\{b \in \mathbf{B} : E_0 b = 0\}$ . To see this note that

$$\begin{aligned}
&\inf\{|\langle \dot{W}(b_1), b_2 \rangle| : \|b_1\| = \|b_2\| = 1\} \\
&= \inf\{\text{Var}(b) : \|b\| = 1\} = 0.
\end{aligned}$$

Indeed, we can find  $b_m$  with  $\|b_m\| = 1$  such that  $E b_m^2 \rightarrow 0$ , for instance  $b_m(x) \equiv (2\pi m)^{-2} \cos(2\pi m x) / \{1 + (2\pi m)^{-2}\}$ .

**Remark 4.** Despite remark 3, we can apply corollary 2 and remark 2 to obtain linearity and Gaussian regularity of other interesting functionals. For instance, let  $q(v) = \int_0^1 x \exp v(x) dx / \int_0^1 \exp v(x) dx$  be the mean. It is easy to see that if  $\|\cdot\|_{\mathbf{B}}$  is the  $L_2(P_0)$  norm then  $q$  satisfies (GD1), (GD3), and (GD4) with

$$\dot{q}(v_0)(b) = \int (x - q(v_0)) b(x) dP_0(x) = \text{Cov}_0(X, b(X)).$$

So we can apply corollary 2. Not surprisingly,  $q(\hat{v}_n)$  is efficient.

In this example  $\hat{v}_n$  is *not* (strongly) Gaussian. But, a "smoothed" version  $q(v)$ ,  $q: \mathbf{B} \rightarrow \mathbf{B}_1$  where  $\mathbf{B}_1$  is a "smaller" Banach space, might be. This is true since, for instance,  $\sup\{n^{-1/2} \sum_{i=1}^n [b(X_i) - E_0(b(X_i))]: b \in \mathbf{B}\} = \infty$ , so that the limit of the finite-dimensional cylinder measures of  $n^{1/2}(\hat{v}_n - v_0)$  is not tight. Yet if we define  $q(b) = \int \exp b(s) ds / \int_0^1 \exp b(s) ds$ ,  $q: \mathbf{B} \rightarrow l^\infty(\{1_{[0,s]} : s \leq 1\})$ , then  $\sqrt{n}(q(\hat{v}_n) - q(v_0))$  is Gaussian regular. For a

systematic discussion of this phenomenon, see Millar (1983). However, as corollary 2 reveals, if interest centers on Euclidean functionals this is immaterial.  $\square$

**Example 5. Regression by penalized maximum likelihood.**

Recall that in the situation of example 7.5.9,  $\mathbf{B}$  is the Sobolev space of twice differentiable functions on  $[0, 1]$  with norm  $\|b\|_{\mathbf{B}} = \{\int h^2 + \int (h'')^2\}^{1/2}$ . By (7.5.33) and (a) of the proof of proposition 7.5.2,  $\hat{t}_n$  is GM for  $W_n: \mathbf{B} \times \mathbf{M} \rightarrow \mathbf{B}^*$  given by:

$$\langle W_n(t, P), b \rangle = \int (t(z) - y) b(z) dP(z, y) + \lambda_n \int t''(z) b''(z) dz.$$

Here

$$\langle \dot{W}_n(t, P)(b_1), b \rangle = \int b_1 b(z) dP(z, y) + \lambda_n \int b_1'' b''(z) dz$$

and, of course,

$$(23) \quad \langle W(t, P), b \rangle = \int (t(z) - y) b(z) dP(z, y),$$

$$(24) \quad \langle \dot{W}(t, P)(b_1), b \rangle = \int b_1 b(z) dP(z, y).$$

**Proposition 7.** Let  $\lambda_n = o(n^{-1/2})$  and  $n^{2/3} \lambda_n / \log n \rightarrow \infty$ , and suppose  $Z$  has a positive marginal density  $p_0(\cdot)$ . Then:

A.  $\langle \dot{W}(\hat{t}_n - t_0), b \rangle = n^{-1} \sum_{i=1}^n b(Z_i) \varepsilon_i + o_p(n^{-1/2})$ ,

B.  $\|\hat{t}_n - t_0\|_2 = o_p(n^{-1/4})$ .

**Proof.** If  $\Gamma_0^*$  is as in the proof of proposition 6, for part A we again just need to check the conditions of theorem 1. Again, (GGM0) and (GGM3) are immediate. We check (GGM1). Now

$$(a) \quad \langle W_n(\hat{t}_n, \mathbb{P}_n) - W_n(t_0, \mathbb{P}_n) - W_n(\hat{t}_n, P_0) + W_n(t_0, P_0), b \rangle \\ = \int (\hat{t}_n - t_0)(z) b(z) d(\mathbb{P}_n - P_0)(z, y).$$

Note that

$$(b) \quad \|b(\hat{t}_n - t_n)\|_{\infty} \leq \|b\|_{\infty} \|\hat{t}_n - t_n\|_{\infty},$$

$$(c) \quad \|[b(\hat{t}_n - t_n)]'\|_{\infty} \leq \|b\|_{\infty} \|\hat{t}_n' - t_n'\|_{\infty} + \|b'\|_{\infty} \|\hat{t}_n - t_n\|_{\infty}.$$

But, by proposition 7.5.2.A and Sobolev's inequalities in (e) of its proof

$$(d) \quad \|\hat{t}_n - t_n\|_{\infty} + \|\hat{t}_n' - t_n'\|_{\infty} = O_p(\lambda_n^{-3/4} n^{-1/2} (\log n)^{1/4}).$$

We deduce from (b) through (d), remark A.6.1, and proposition 7.5.2.A, that  $\lambda_n^{3/4} n^{1/2} (\log n)^{-1/4} b(\hat{t}_n - t_n)$  belongs to  $\mathbf{F}_n$  which satisfies the conditions of corollary A.6.3 with  $r = 1, \alpha = \alpha_n = \lambda_n, \lambda = \mu_n = K\alpha_n^{1/4}$ . Hence

$$\sup \{ |\int (\hat{t}_n - t_0) b(z) d(\mathbb{P}_n - P_0)(z, y)| : \|b\|_{\mathbf{B}} \leq 1 \} \\ = O_p(\lambda_n^{-3/4} (\log n)^{1/4} n^{-1} \sup \{ n^{1/2} |\int f d(\mathbb{P}_n - P_0)| : f \in \mathbf{F}_n \}) \\ = O_p(\lambda_n^{-1/2} n^{-1} (\log n)^{1/4})$$

$$= o_p(n^{-1/2}),$$

and (GGM1) holds.

The first part of (GGM4') is obvious since  $W_n(t, P)$  is linear in  $t$ . In view of the definition of  $t_n$  just before proposition 7.5.2 and in view of part A of this proposition we have

$$\begin{aligned} \|(\dot{W}_n - \dot{W})(\hat{t}_n - t_0)\|_{*\Gamma} &= \lambda_n \sup\{|\int (\hat{t}_n - t_0)'' b''(z) dz| : \|b\|_{\mathbf{B}} \leq 1\} \\ &\leq \lambda_n \|\hat{t}_n'' - t_0''\|_2 \\ &\leq \lambda_n \{\|\hat{t}_n'' - t_n''\|_2 + \|t_n''\|_2 + \|t_0''\|_2\} \\ &\leq \lambda_n \{\|\hat{t}_n'' - t_n''\|_2 + 2\|t_0''\|_2\} \\ &= O_p(\lambda_n) = o_p(n^{-1/2}). \end{aligned}$$

Finally (GGM2) is immediate and part A of the proposition follows. Part B is a restatement of part B of proposition 7.5.2. □

Again it is easy to check that  $\dot{W}$  does not have a bounded inverse. However, as in the previous examples if

$$(25) \quad q(t) = q(t_0) + \int (t - t_0)(z) a(z, t_0) dz + O(\|t - t_0\|_2^2)$$

and  $\|a\|_2 < \infty$ , then  $q(\hat{t}_n)$  is asymptotically linear, Gaussian regular with influence function  $a(z, t_0) \in$ . In particular, if  $q(t) \equiv \int t(z) a(z) p_0(z) dz$ , then

$$(26) \quad q(t_0) = E(E(Y | Z) a(Z)) = E(Y a(Z)).$$

The natural efficient estimate of  $q(t_0)$  is then  $n^{-1} \sum_{i=1}^n a(Z_i) Y_i$  which we have just seen agrees with  $q(\hat{t}_n)$  to order  $o_p(n^{-1/2})$ . So, in this case, and generally  $q(\hat{t}_n)$  is efficient. □

### Asymptotic Existence, Consistency, and Linearity of GM-Estimates

The problems of lack of existence of unique GM-estimates (which are exact roots of their defining equations) that we encountered in section 7.3 hold a fortiori in the infinite-dimensional case. Again, if we have available a suitably consistent estimate  $\tilde{v}_n$  and an estimate of  $\dot{W}_n$  we can construct a one-step GM-estimate. However, this construction requires inversion of an operator. If we iterate this procedure and the iterates converge we expect to obtain a GM-estimate which is the unique solution of  $W_n(v, P_n) = 0$  in a neighborhood of  $v(P)$ . We have already seen in the finite-dimensional case that this estimate behaves well and we shall see that under suitable conditions the same is true in the infinite-dimensional case. It may be possible, as we illustrate in example 6, to obtain this estimate through iterations which do not involve inversion of  $\dot{W}_n$ . This is evidently useful.

We begin with a rate of convergence result for the one step.

For each  $P_0 \in \mathbf{P}$  and corresponding  $v(P_0) \equiv v_0$  we suppose that there exists a nested (decreasing) sequence of convex subsets  $N_n$  of  $\mathbf{B}$  containing  $v_0$



and norms  $\|\cdot\|_{\mathbf{B}}$  on  $\mathbf{B}$  and  $\|\cdot\|_{\Gamma}$  on  $\Gamma$  such that

(O1) There exists an estimate  $\tilde{v}_n$  of  $v$ , which satisfies

$$P_0(\tilde{v}_n \in N_n) \rightarrow 1,$$

$$\|\tilde{v}_n - v_0\|_{\mathbf{B}} = o_p(n^{-1/4}).$$

(O2) Write  $W_n(v)$  for  $W_n(v, P_n)$ ,  $W(v)$  for  $W(v, P_0)$ . Suppose  $W_n, W$  are continuously differentiable on  $N_1$  when  $\mathbf{B}, \Gamma$  are endowed with  $\|\cdot\|_{\mathbf{B}}, \|\cdot\|_{\Gamma}$  respectively, and write  $\|\cdot\|_{\mathbf{B}\Gamma}$  for the operator norm. Then

$$\sup\{\|\dot{W}_n(v) - \dot{W}(v_0)\|_{\mathbf{B}\Gamma} : \|v - v_0\|_{\mathbf{B}} \leq n^{-1/4}, v \in N_n\} = O_p(n^{-1/4}).$$

(O3)  $\dot{W}(v_0)$  is onto and has a bounded inverse  $\dot{W}^{-1}(v_0)$ .

(O4)  $\|W_n(v_0)\|_{\Gamma} = o_p(n^{-1/4})$ .

Define the *one-step estimate*,

$$(27) \quad \hat{v}_n \equiv \tilde{v}_n - \hat{W}_n^- W_n(\tilde{v}_n),$$

where  $\hat{W}_n^- \equiv \dot{W}_n^{-1}(\tilde{v}_n)$  if  $\dot{W}_n(\tilde{v}_n)$  has a bounded inverse and  $\hat{v}_n = \tilde{v}_n$  otherwise.

**Theorem 2.** If (O1)–(O4) hold then

A.  $\|\hat{v}_n - v_0\|_{\mathbf{B}} = o_p(n^{-1/4})$  and  $W_n(\hat{v}_n) = o_p(n^{-1/2})$ .

B. If, in addition (14) holds with  $g^* W_n(v_0, P_0)$  deleted, then  $\hat{v}_n$  is weakly  $\mathbf{B}^*$  linear, Gaussian regular with influence function  $-\dot{W}^{-1} w(\cdot, P_0)$ , where  $w$  is given in (14).

**Proof.** By (O2) and (O3) and a result from functional analysis (Kantorovich and Akilov (1982, theorem 5.4, page 155)),

$$P_0(\dot{W}_n(v) \text{ has a bounded inverse for all } \|v - v_0\|_{\mathbf{B}} \leq n^{-1/4}, v \in N_n) \rightarrow 1,$$

and

$$(a) \sup\{\|\dot{W}_n^{-1}(v) - \dot{W}^{-1}(v_0)\|_{\mathbf{B}\Gamma} : \|v - v_0\|_{\mathbf{B}} \leq n^{-1/4}, v \in N_n\} = O_p(n^{-1/4}).$$

By (O1),  $\|\tilde{v}_n - v_0\|_{\mathbf{B}} \leq n^{-1/4}$  for large  $n$ , and hence (a) and (O2) yield

$$(b) \sup\{\|I - \hat{W}_n^- (\dot{W}_n(\lambda v_0 + (1 - \lambda)\tilde{v}_n))\|_{\mathbf{B}} : \lambda \in (0, 1)\} = O_p(n^{-1/4}).$$

We now write

$$\begin{aligned} (c) \quad \hat{v}_n - v_0 &= \tilde{v}_n - v_0 - \hat{W}_n^- (W_n(\tilde{v}_n) - W_n(v_0)) - \hat{W}_n^- W_n(v_0) \\ &= (I - \hat{W}_n^- \dot{W}_n(v^*)) (\tilde{v}_n - v_0) \\ &\quad - (\hat{W}_n^- - \dot{W}^{-1}(v_0)) W_n(v_0) - \dot{W}^{-1}(v_0) W_n(v_0), \end{aligned}$$

where  $v^*$  is an intermediate point between  $v_0$  and  $\tilde{v}_n$ . The  $\|\cdot\|_{\mathbf{B}}$  norm of the first term on the right side of (c) is  $o_p(n^{-1/2})$  by (b) and (O1), and that of the second term is  $o_p(n^{-1/2})$  by (a) and (O4). The first result of part A follows from (O3) and (O4), and the second one from (O2) by expanding  $W_n(\hat{v}_n)$  around  $\tilde{v}_n$ .

For part B, note that part A and (O2) enable us to replace the third member of (a) in the proof of theorem 1 by

$$\dot{W}_n(\lambda \hat{v}_n + (1 - \lambda)v_0)(\hat{v}_n - v_0) + o_p(n^{-1/2}), \quad \text{for some } \lambda \in [0, 1],$$

thus saving the validity of (a) of that proof and of the conclusion of theorem 1.A for the present situation. □

As in section 7.3 we can define

$$(28) \quad T_n^{(0)} \equiv \tilde{v}_n$$

$$T_n^{(j+1)} \equiv T_n^{(j)} - \dot{W}_n^-(T_n^{(j)})W_n(T_n^{(j)}), \quad \text{for } j \geq 0,$$

and conclude that  $T_n^{(j+1)}$  behaves like  $T_n^{(j)}$  for all  $j \geq 1$  if the conditions of theorem 2 hold for  $T_n^{(j)}$ . We go further in the following theorem.

**Theorem 3.** Suppose (O1)–(O4) hold with  $N_n = B$ . Then:

A. If  $T_n^{(\infty)}$  is the limit of the iterations,

$$P(T_n^{(\infty)} \text{ is defined and } W_n(T_n^{(\infty)}, P_n) = 0) \rightarrow 1, \quad \text{as } n \rightarrow \infty.$$

Further,

$$(29) \quad \|T_n^{(\infty)} - v_0\|_B = o_p(n^{-1/4}).$$

B. If, in addition, (GGMO) and (GGM2) hold, then  $T_n^{(\infty)}$  is weakly  $B_0^*$  linear, Gaussian regular with influence function  $-\dot{W}_n^{-1} w(\cdot, P_0)$  where  $w$  is given in (14).

**Proof.** Part A is proved in essentially the same way as theorem 7.3.3 (see A.10). The differences are:

(i)  $\sqrt{n}$ -consistency of  $v_n^{(\infty)}$  is replaced by  $\|v_n^{(\infty)} - v_0\|_B = o_p(n^{-1/4})$  since (O4) has replaced  $W_n(v_0) = O_p(n^{-1/2})$  which holds automatically in the Euclidean case (under (GM2) from section 7.3).

(ii) The general contraction mapping fixed point theorem (Kantorovich and Aki-  
lov (1982, page 474)) needs to be used to conclude that  $v_n^{(\infty)}$  exists.

Finally part B follows by considering  $T_n^{(\infty)}$  as a one-step estimate starting from itself and applying theorem 2.B. □

In fact, as in the Euclidean case, (O2), (O3) and  $\|W_n(v_0)\|_\Gamma = o_p(1)$  suffice for the asymptotic existence, for some  $\delta > 0$ , of a unique root  $v_n^{(\infty)}$  of  $W_n(v)$  in  $\{v: \|v - v_0\|_B < \delta\}$  such that  $\|v_n^{(\infty)} - v_0\|_B = o_p(n^{-1/4})$ . To compute the root we, of course, need some convergent iteration scheme and a good starting point. Theorem 3 gives conditions under which we can do this with a Newton-Raphson scheme. However, we may use other algorithms.

We illustrate this remark and theorem 2 with an important example.

**Example 6. The Has'minskii-Ibragimov model.**

We continue the study of the regularized maximum likelihood estimate of example 7.5.10. We shall impose brutal and far from necessary conditions on the model to facilitate our computations. We shall show elsewhere how this method can be refined in this class of examples, to yield procedures which work under minimal conditions.

Recall that in example 7.5.10 we defined  $\mathbf{B} = L_2([0,1])$  and  $\mathbf{A} = \{v \in \mathbf{B} : \int e^v = 1, \|v\|_\infty + \|v''\|_\infty < \infty\}$ . For simplicity of notation we take  $q_0(\cdot, \eta)$  and hence  $q(\cdot, v)$  as densities with respect to Lebesgue measure. This is immaterial to our conclusions. We impose the following restriction on  $q_0$ ,

$$(30) \quad 0 < \varepsilon \leq \frac{q_0(y, u)}{q_0(y, u')} \leq \frac{1}{\varepsilon} < \infty$$

for all  $y, 0 \leq u, u' \leq 1$ . This condition rules out many interesting models, such as the normal convolution model. However, e.g., convolution models with  $Y$  Cauchy or double exponential and  $U$  a location parameter satisfy (30). Further, we require, under  $P_0 = P_{g_0} \in \mathbf{P}$  with  $v_0 = \log g_0$ ,

$$(31) \quad \sup_u \left| n^{-1} \sum_{i=1}^n \left[ \frac{q_0(Y_i, u)}{q(Y_i, v_0)} - 1 \right] \right| = o_p(n^{-1/4}).$$

This condition holds if (30) does and  $\{q_0(\cdot, u) : u \in [0, 1]\}$  is a regular parametric model. To see this note that, (30) yields for all  $u_1, u_2$

$$\begin{aligned} \int \frac{(q_0(y, u_1) - q_0(y, u_2))^2}{q(y, v_0)} dy &\leq \frac{4}{\varepsilon} \int \{q^{1/2}(y, u_1) - q^{1/2}(y, u_2)\}^2 dy \\ &\leq c(u_1 - u_2)^2, \end{aligned}$$

and our claim follows from a standard fluctuation inequality; see, for example, Billingsley (1968, theorem 12.3, page 95).

Let  $p_{n1}^\#$  be the sieve estimate of examples 4 and 7.5.7 based on  $U_1', \dots, U_n'$  with  $mn^{-1/4} \rightarrow 0, mn^{-1/6} \rightarrow \infty$ . Define  $W_n(v, P_n)$  by (7.5.31) for  $v \in \mathbf{A}$ , and let

$$(32) \quad W(v, P_0) = g_0 + e^v \left\{ \int \frac{q_0(y, \cdot)}{q(y, v)} q(y, v_0) dy - 2 \right\}.$$

$W$  is the limit of  $W_n$  and evidently,  $W(v_0, P_0) = 0$ . Note that  $W_n(\cdot, P_n)$  and  $W(\cdot, P_0)$  are well defined as maps from  $\mathbf{A}$  to  $\Gamma = L_2([0,1])$ , since

$$\begin{aligned} \|W(v, P_0)\|_\infty &\leq e^{1v_0} + e^{1v_1} \left\{ \frac{1}{\varepsilon} + 2 \right\}, \\ \|W_n(v, P_n)\|_\Gamma &= O_p(1 + E_0 \int e^{2v(u)} \left\{ \int \frac{q_0(y, u)}{q(y, v)} dP_{n2}(y) \right\}^2 du) \\ &= O_p(1), \end{aligned}$$

by (30) and the boundedness of  $p_{n1}^*$  as noted in proposition 7.5.1.

We shall check the conditions of theorem 2 to prove:

**Proposition 8.** Let  $v_0 \in \mathbf{A}$ , (30) and (31) hold, and  $p_{n1}^\#, W_n$  and  $W$  be as above. Denote  $\hat{v}_n^{(0)} \equiv \log p_{n1}^\#$ . Then:

- A.  $F(\dot{W}_n(\hat{v}_n^{(0)}, P_n)$  exists and has a bounded inverse)  $\rightarrow 1$ .
- B. Let  $\hat{v}_n^{(1)} \equiv \hat{v}_n^{(0)} - \dot{W}_n^{-1}(\hat{v}_n^{(0)}, P_n) W_n(\hat{v}_n^{(0)}, P_n)$ .

Then, for all  $b \in \mathbf{B}$

$$(33) \quad \begin{aligned} \|\hat{v}_n^{(1)} - v_0\|_2 &= o_p(n^{-1/4}), \\ \int b \hat{v}_n^{(1)}(u) du - \int b v_0(u) du &= n^{-1} \sum_{i=1}^n \psi(X_i, b, P_0) + o_p(n^{-1/2}), \end{aligned}$$

where

$$\psi(u', y, b, P_0) = \Pi_0(b | \dot{P}(g_0)) = \psi_1(u', b) + \psi_2(y, b)$$

given in proposition 4.5.1. That is, if we take  $\Gamma_0^*$   $\equiv$  unit ball in  $\mathbf{B}$  viewed as a subset of  $\mathbf{B}^* = L_2([0, 1])$ ; then  $\hat{v}_n^{(1)}$  is weakly  $\Gamma_0^*$  linear and Gaussian.

**Proof.** We verify (14) without the first (undefined) term of its right side, and (O1) through (O4) with

$$N_n = \{v \in \mathbf{A} : \|v - v_0\|_\infty \leq n^{-1/4}\} \quad \text{and} \quad \tilde{v}_n = \hat{v}_n^{(0)}.$$

By proposition 7.5.1 with  $\hat{v}_n = \hat{v}_n^{(0)}$  we have  $\|\hat{v}_n^{(0)} - v_0\|_\infty = o_p(n^{-1/4})$  and  $\|\hat{v}_n^{(0)} - v_0\|_2 = o_p(n^{-1/3})$  and hence (O1).

The Gâteaux derivative of  $W_n$  is

$$(a) \quad \begin{aligned} \dot{W}_n(v)(b)(w) &\equiv \frac{\partial}{\partial \lambda} W_n(v + \lambda b) |_{\lambda=0}(w) \\ &= e^{v(w)} \left\{ b(w) \left[ \int \frac{q_0(y, w)}{q(y, v)} dP_{n2}(y) - 2 \right] \right. \\ &\quad \left. - \int \frac{q_0(y, w)}{q^2(y, v)} \left[ \int q_0(y, u) e^{v(u)} b(u) du \right] dP_{n2}(y) \right\}. \end{aligned}$$

We claim

$$(b) \quad \sup\{\|\dot{W}_n(v)(b)\|_\Gamma : v \in N_1, \|b\|_2 \leq 1\} = O_p(1).$$

To see this check that by (30) and for  $v \in N_1, \|b\|_2 \leq 1$ ,

$$(c) \quad \int e^{2v(w)} \left\{ b^2(w) \left[ \int \frac{q_0(y, w)}{q(y, v)} dP_{n2}(y) \right]^2 \right\} dw \leq e^{2(1v_0)_+ + 1} \epsilon^{-2}$$

and

$$(d) \quad \begin{aligned} \int e^{2v(w)} \left\{ \int \frac{q_0(y, w)}{q^2(y, v)} \left[ \int q_0(y, u) e^{v(u)} b(u) du \right] dP_{n2}(y) \right\}^2 dw \\ \leq \exp(2\|v_0\|_\infty + 1) \epsilon^{-4}. \end{aligned}$$

From (a), (c), and (d), claim (b) follows. Next we calculate

$$(e) \quad \begin{aligned} (\dot{W}_n(v) - \dot{W}_n(v_0))(b)(w) \\ = -2e^{v_0(w)} b(w) \{e^{(v-v_0)(w)} - 1\} \end{aligned}$$

$$\begin{aligned}
 & + e^{v_0(w)} b(w) \int \frac{q_0(y, w)}{q(y, v_0)} \left\{ e^{(v-v_0)(w)} \frac{q(y, v_0)}{q(y, v)} - 1 \right\} d\mathbb{P}_{n2}(y) \\
 & - e^{v_0(w)} \int \int \frac{q_0(y, w) q_0(y, u)}{q^2(y, v_0)} e^{v_0(u)} b(u) \Delta(u, w, y) du d\mathbb{P}_{n2}(y) \\
 & = \text{I}' + \text{II}' + \text{III}',
 \end{aligned}$$

where

$$\Delta(u, w, y) = \exp[(v - v_0)(u) + (v - v_0)(w)] \frac{q^2(y, v_0)}{q^2(y, v)} - 1.$$

Since

$$\text{(f)} \quad \left\| \frac{q(y, v)}{q(y, v_0)} - 1 \right\|_\infty = O(\|v - v_0\|_\infty),$$

we have for  $\|b\|_2 \leq 1, v \in N_1$ , in view of (30),

$$\text{(g)} \quad \|\text{I}'\|_2^2 = O(\|v - v_0\|_\infty^2),$$

$$\begin{aligned}
 \text{(h)} \quad \|\text{II}'\|_2^2 & \leq \exp(2\|v_0\|_\infty) \int b^2(w) \int \frac{q_0^2(y, w)}{q^2(y, v_0)} \\
 & \quad \cdot \left\{ e^{(v-v_0)(w)} \frac{q(y, v_0)}{q(y, v)} - 1 \right\}^2 d\mathbb{P}_{n2}(y) dw \\
 & = O_p(\|v - v_0\|_\infty^2).
 \end{aligned}$$

Furthermore, we see,

$$\sup_{u, w, y} |\Delta(u, w, y)| = O(\|v - v_0\|_\infty)$$

and hence

$$\text{(i)} \quad \|\text{III}'\|_2^2 = O_p(\|v - v_0\|_\infty^2).$$

From (e), (g), (h), (i) and the definition of  $N_n$ , we obtain

$$\text{(j)} \quad \sup\{\|\dot{W}_n(v) - \dot{W}_n(v_0)\|_{\text{BF}} : v \in N_n\} = O_p(n^{-1/4}).$$

By proposition A.5.1.E it follows from (b) and (a refinement of) (j) that  $W_n$  is continuously Fréchet differentiable. The same holds for  $W$ . Next, we write

$$\begin{aligned}
 \text{(k)} \quad & (\dot{W}_n(v_0) - \dot{W}(v_0))(b)(w) \\
 & = e^{v_0(w)} b(w) \frac{1}{n} \sum_{i=1}^n \left\{ \frac{q_0(Y_i, w)}{q(Y_i, v_0)} - 1 \right\} \\
 & \quad - \frac{1}{n} \sum_{i=1}^n \left\{ \frac{q_0(Y_i, w)}{q^2(Y_i, v_0)} \int q_0(Y_i, u) e^{v_0(u)} b(u) du \right.
 \end{aligned}$$

$$- E_0 \left\{ \frac{q_0(Y_1, w)}{q^2(Y_1, v_0)} \int q_0(Y_1, u) e^{v_0(u)} b(u) du \right\} \Bigg\}.$$

But

$$\begin{aligned} (l) \quad & \int e^{2v_0(w)} b^2(w) \left( \frac{1}{n} \sum_{i=1}^n \left[ \frac{q_0(Y_i, w)}{q(Y_i, v_0)} - 1 \right] \right)^2 dw \\ & \leq \exp(2 \|v_0\|_\infty) \|b\|_2^2 \sup_w \left| \frac{1}{n} \sum_{i=1}^n \left[ \frac{q_0(Y_i, w)}{q(Y_i, v_0)} - 1 \right] \right|^2 \\ & = \|b\|_2^2 o_p(n^{-1/2}) \end{aligned}$$

by (31). Also, the  $L_2$  norm of the second term in (k) is bounded by

$$\begin{aligned} (m) \quad & \exp(\|v_0\|_\infty) \|b\|_2 \left[ \iint \left( \frac{1}{n} \sum_{i=1}^n \left( \frac{q_0(Y_i, w)}{q^2(Y_i, v_0)} q_0(Y_i, u) \right. \right. \right. \\ & \left. \left. \left. - E \left[ \frac{q_0(Y_1, w)}{q^2(Y_1, v_0)} q_0(Y_1, u) \right] \right) \right)^2 du dw \right]^{1/2} \\ & = \|b\|_2 O_p(n^{-1/2}). \end{aligned}$$

(O2) follows from (j)–(m).

Since

$$\begin{aligned} W_n(v_0)(w) &= p_{n1}^\#(w) - e^{v_0(w)} \\ &+ e^{v_0(w)} \left\{ \int \left[ \frac{q_0(y, w)}{q(y, v_0)} - 1 \right] dP_{n2}(y) \right\}, \\ \|p_{n1}^\# - e^{v_0}\|_F^2 &= O_p(\|\hat{v}_n^{(0)} - v_0\|_2^2) = o_p(n^{-2/3}), \end{aligned}$$

and by (l)

$$\|e^{v_0(\cdot)} \int \left[ \frac{q_0(y, \cdot)}{q(y, v_0)} - 1 \right] dP_{n2}(y)\|_F^2 = o_p(n^{-1/2}),$$

we have (O4).

With

$$\begin{aligned} (n) \quad \langle w(X, P_0), b \rangle &= b(U') - \int b(u) g_0(u) du \\ &+ \int b(u) g_0(u) \left\{ \frac{q_0(Y, u)}{q(Y, v_0)} - 1 \right\} du, \end{aligned}$$

we obtain by proposition 6,

$$\langle W_n(v_0, P_n), b \rangle = \langle \frac{1}{n} \sum_{i=1}^n w(X_i, P_0), b \rangle + o_p(n^{-1/2}).$$

To complete the proof of (14) with the first term of the right side deleted we note that

$$\int e^{v_0(u)} b(u) \left[ \frac{q_0(Y, u)}{q(Y, v_0)} - 1 \right] du = O_p(1), \quad \|b\|_2 \leq 1.$$

In order to prove (O3), note that by (a) with  $n = \infty$ ,

$$\begin{aligned} (o) \quad & \dot{W}(v_0)(b)(U) \\ &= -g_0(U) \left\{ b(U) + \int \frac{q_0(y, U)}{q(y, v_0)} \int b(w) q_0(y, w) g_0(w) dw dy \right\} \\ &= -D_{g_0}(I + E_0^U E_0^Y)(b)(U), \end{aligned}$$

where  $I$  is the identity,  $E_0^U$  and  $E_0^Y$  are the conditional expectation operators of proposition 4.5.1 defined by

$$\begin{aligned} (E_0^Y b)(y) &= E_0(b(U) | Y = y), \\ (E_0^U a)(u) &= E_0(a(Y) | U = u), \end{aligned}$$

and  $D_g : L_2(g_0) \rightarrow L_2(g_0)$  via  $D_g b = g b$ . Here  $L_2(g_0)$  is the Hilbert space of square integrable functions with respect to the density  $g_0$  on  $[0,1]$  and inner product  $\langle \cdot, \cdot \rangle_{g_0}$  and norm  $\|\cdot\|_{g_0}$ . Note that  $g_0$  is bounded and bounded away from 0. Hence it suffices to prove that

$$I + E_0^U E_0^Y : L_2(g_0) \rightarrow L_2(g_0)$$

has a bounded inverse. Now, for  $b \in L_2(g_0)$ ,

$$\begin{aligned} & \|b\|_{g_0} \|(I + E_0^U E_0^Y)b\|_{g_0} \\ & \geq \langle b, (I + E_0^U E_0^Y)b \rangle_{g_0} \\ & = \|b\|_{g_0}^2 + \iint b(u) \left\{ \int \frac{q_0(y, u) q_0(u, w)}{q(y, v_0)} dy \right\} b(w) g_0(u) g_0(w) du dw \\ & = \|b\|_{g_0}^2 + \int \left\{ \int \frac{q_0(y, u)}{q(y, v_0)} b(u) g_0(u) du \right\}^2 q(y, v_0) dy \\ & \geq \|b\|_{g_0}^2. \end{aligned}$$

Consequently, by corollary A.1.2,  $I + E_0^U E_0^Y$  and hence  $\dot{W}(v_0)$  have bounded inverses as maps from  $L_2(g_0)$  to  $L_2(g_0)$ , and

$$\dot{W}^{-1}(v_0) = -(I + E_0^U E_0^Y)^{-1} D_{g_0}^{-1}.$$

The proposition follows from theorem 2. □

In the next proposition we will present an iteration scheme to get  $\hat{v}_n$  which

solves  $W_n(v, P_n) = 0$  and which does not require computation of  $\hat{W}_n^{-1}$ . For this scheme to work we need the following condition:

$$(34) \sup \left\{ \left\| \int \int \frac{q_0(y, \cdot)}{q^2(y, v_0)} q_0(y, w) b(w) g_0(w) dw d(P_{n2} - P_{02})(y) \right\|_\infty : \|b\|_\infty \leq 1 \right\},$$

$$= o_p(n^{-1/4}),$$

where  $P_{02}$  is the marginal of  $Y$  under  $P_0$ . This condition is not as onerous as it seems. If  $q_0(\cdot, u)q_0(\cdot, w)q^{-2}(\cdot, v_0)$  obey uniform Lipschitz conditions so do the inner integrals and we can apply corollary A.6.3. Note that (34) implies (31) by taking  $b \equiv 1$ .

**Proposition 9.** Let  $\hat{v}_n^{(0)} = \log p_{n1}^\#$ , and for  $j = 1, 2, \dots$  set

$$(35) \quad \exp(\hat{v}_n^{(j)}) = p_{n1}^\# + \exp(\hat{v}_n^{(j-1)}) \left( \int \frac{q_0(y, \cdot)}{q(y, \hat{v}_n^{(j-1)})} dP_{n2}(y) - 1 \right).$$

If (30) and (34) hold, then with probability tending to 1 under  $P_0$ , as  $n \rightarrow \infty$ ,  $W_n(v, P_n) = 0$  has a solution  $\hat{v}_n$  in  $\{v: \|v - v_0\|_\infty \leq n^{-1/4}\}$  and  $\hat{v}_n^{(j)}$  converges in  $j$  to  $\hat{v}_n$ .

**Proof.** Let  $B_n = \{b \in B: \int b g_0 = 1, \|b - 1\|_\infty \leq n^{-1/4}/2\}$ . Recall that  $B = L_2([0, 1])$  and note that  $B_n \subset l^\infty([0, 1])$ . For  $b \in B_n$ ,  $v \equiv \log(bg_0)$  we define

$$(a) \quad T_n(b) = \frac{p_{n1}^\#}{g_0} + b \left( \int \frac{q_0(y, \cdot)}{q(y, v)} dP_{n2}(y) - 1 \right).$$

Since by proposition 7.5.1

$$(b) \quad \left\| \frac{p_{n1}^\#}{g_0} - 1 \right\|_\infty = O_p \left( \left\| \log \left( \frac{p_{n1}^\#}{g_0} \right) \right\|_\infty \right) = o_p(n^{-1/4}),$$

$\hat{v}_n^{(0)} \in B_n$  for  $n$  large and it is enough to show that, with probability tending to 1,  $T_n$  maps  $B_n$  into  $B_n$  and is an  $l^\infty$  contraction. To see this note that if  $\hat{b}$  is the unique fixed point of  $T_n$  then  $\hat{v}_n^{(j)}$  converges to  $\hat{v}_n^{(\infty)} = \log(\hat{b} g_0)$  and  $\hat{v}_n$  exists, equals  $\hat{v}_n^{(\infty)}$  and satisfies  $\|v - v_0\|_\infty \leq n^{-1/4}$ . Write, for  $b, b + \delta \in B_n$ ,  $v_\delta = \log((b + \delta)g_0)$ ,

$$\begin{aligned} & \| \delta \|_\infty^{-1} \{T_n(b + \delta) - T_n(b)\} \\ &= \frac{\delta}{\| \delta \|_\infty} \left\{ \int \frac{q_0(y, \cdot)}{q(y, v_\delta)} dP_{n2}(y) - 1 \right\} \\ & \quad + b \int \frac{q_0(y, \cdot)}{q(y, v_\delta)q(y, v)} \int q_0(y, w) \frac{\delta(w)}{\| \delta \|_\infty} g_0(w) dw dP_{n2}(y) \end{aligned}$$

$$(c) \quad = I(b, \delta) + II(b, \delta).$$



By using (31) and  $\|v_\delta - v_0\|_\infty \leq n^{-1/4}$ , we see that

$$\begin{aligned}
 & \sup\{\|I(b, \delta)\|_\infty : b, b + \delta \in \mathbf{B}_n\} \\
 (d) \quad & \leq \sup\left\{\left\|\int \frac{q_0(y, \cdot)}{q(y, v_0)} dP_{n2}(y) - 1\right\|_\infty \right. \\
 & \quad \left. + \left\|\frac{q(\cdot, v_0)}{q(\cdot, v_\delta)} - 1\right\|_\infty \left\|\int \frac{q_0(y, \cdot)}{q(y, v_0)} dP_{n2}(y)\right\|_\infty \right. \\
 & \quad \left. : b + \delta \in \mathbf{B}_n\right\} \\
 & = o_p(n^{-1/4}) + O_p(n^{-1/4}) = o_p(1).
 \end{aligned}$$

Similarly, by (34),

$$\begin{aligned}
 (e) \quad & \sup\{\|II(b, \delta)\|_\infty : b, b + \delta \in \mathbf{B}_n\} \\
 & = \sup\left\{\left\|\int \frac{q_0(y, \cdot)}{q(y, v_0)} \int q_0(y, w) \frac{\delta}{\|\delta\|_\infty}(w) g_0(w) dw dy\right\|_\infty : \right. \\
 & \quad \left. b + \delta \in \mathbf{B}_n\right\} + o_p(1).
 \end{aligned}$$

Let  $m(y)$  be a median of  $q_0(y, U')$  under  $g_0$ . Since  $\int \delta g_0 = 0$ , the first term on the right side of (e) equals

$$\begin{aligned}
 (f) \quad & \sup\left\{\left\|\int \frac{q_0(y, \cdot)}{q(y, v_0)} \int [q_0(y, w) - m(y)] \frac{\delta}{\|\delta\|_\infty}(w) g_0(w) dw dy\right\|_\infty : \right. \\
 & \quad \left. b + \delta \in \mathbf{B}_n\right\}
 \end{aligned}$$

$$\leq 1 - \inf_{0 \leq u \leq 1} \int \frac{q_0(y, u)}{q(y, v_0)} \int [q_0(y, w) - |q_0(y, w) - m(y)|] g_0(w) dw dy.$$

Note that by (30) for all  $y$ , 0 is not a median of  $q_0(y, U')$  and hence  $m(y) > 0$ . Consequently

$$\int [q_0(y, w) - |q_0(y, w) - m(y)|] g_0(w) dw > 0$$

for all  $y$  and the expressions in (f) are bounded from above by  $1 - c_0$  for some  $c_0 > 0$ . From this and (c) through (e) it follows that  $T_n$  is a contraction with  $\|T_n(b_1) - T_n(b_2)\|_\infty \leq (1 - C_n) \|b_1 - b_2\|_\infty$  and  $P(C_n > c) \rightarrow 1$  for some  $c > 0$ . Furthermore,

$$\int T_n(b) g_0 = \iint \frac{b(u) q_0(y, u) g_0(u) du}{q(y, v)} dP_{n2}(y) = 1.$$

Finally, if  $b \in \mathbf{B}_n$ ,

$$\|T_n(b) - 1\|_\infty \leq \|T_n(b) - T_n(1)\|_\infty + \left\|\frac{p_{n1}^\#}{g_0} - 1\right\|_\infty$$

$$\begin{aligned}
 & + \left\| \int \frac{q_0(y, \cdot)}{q(y, v_0)} dP_{n2}(y) - 1 \right\|_\infty (1 + o(1)) \\
 & \leq \frac{(1 - C_n)n^{-1/4}}{2} + o_p(n^{-1/4}) \\
 & = \frac{(1 - c + o_p(1))n^{-1/4}}{2}
 \end{aligned}$$

by (b) and (31). Therefore,  $T_n$  maps  $\mathbf{B}_n$  into  $\mathbf{B}_n$  and the proposition follows.  $\square$

**Remark 5.** From  $\hat{v}_n$  we can get  $\hat{g}_n \equiv e^{\hat{v}_n} / \int e^{\hat{v}_n}$  and, as in examples 4 and 5, asymptotically linear Gaussian estimates of functionals of the type discussed in remark 4 of example 4. By proposition 4.5.1, the estimates are efficient.

**Remark 6.** An interesting special case in which the NPMLLE exists and is asymptotically linear and Gaussian under mild conditions appears in Vardi and Zhang (1992).

Our next and final example points again to the idea we have been advocating. Find a tractable procedure using whatever heuristic principles are appropriate rather than sticking to an “optimal” method of estimation whose optimality can only be guaranteed under conditions which are both difficult to check and often do not apply.

**Example 7. The joint distribution-transformation model.**

We consider estimation of  $\tau$  in the joint distribution transformation model of examples 4.7.3 and 6.7.2. We suppose  $\theta = \theta_0$  and the distribution  $m_0$  of  $Z$  known and write the model in standardized form with density

$$p(z, y, \tau) = p_0(z, \tau(y))\tau'(y)$$

with respect to Lebesgue measure  $\times m_0$ . For simplicity, we suppose  $m_0$  has finite support and without loss of generality we take the marginal distribution of  $Y$ , when  $\tau = \text{identity}$ , as  $Uniform(0, 1)$ . Let  $\mathbf{B} = l^\infty([0, 1])$ . We take  $b = \tau'$  as parameter and identify  $\mathbf{A}$  with  $\{b : b \geq 0, \int_0^1 b(s) ds = 1\}$ . Thus, if  $v_0 > 0$ ,  $\dot{\mathbf{A}} = \{b : \int_0^1 b(s) ds = 0\}$ . Let  $p_n^\#(z, y)$  be an estimate of  $p(z, y, \tau)$  and

$$p_n^\#(y) \equiv \int p_n^\#(z, y) dm_0(z),$$

$$r_0(z, y) \equiv - \frac{p_0'}{p_0}(z, y),$$

where the prime denotes differentiation with respect to  $y$ . Formally apply the principle of regularized maximum likelihood and maximize  $\iint \log p(z, y, \tau) p_n^\#(z, y) dm_0(z) dy$  over  $\mathbf{A}$  using a Lagrange multiplier  $c(P_n^\#)$  for  $\int \tau'(y) dy = 1$ . We obtain

$$- \iint_s^1 r_0(z, \hat{\tau}(y)) p_n^\#(z, y) dm_0(z) dy + \frac{p_n^\#(s)}{\hat{\tau}'(s)} = c(P_n^\#),$$

and hence, since  $\hat{\tau}(0) = 0$ ,

$$(36) \quad \hat{\tau}(u) = \int_0^u \{ c(\mathbb{P}_n^\#) + \int_s^1 \int r_0(z, \hat{\tau}(y)) p_n^\#(z, y) dm_0(z) dy \}^{-1} p_n^\#(s) ds,$$

where  $c(\mathbb{P}_n^\#)$  is "determined" by  $\hat{\tau}(1) = 1$ . Formally we can think of

$$\hat{\tau}(u) \equiv \tau(u, \mathbb{P}_n^\#),$$

where  $\tau(u, Q)$  "solves"

$$(37) \quad \tau(u, Q) = \int_0^u \{ c(Q) + \int_s^1 \int r_0(z, \tau(y, Q)) dQ(z, y) \}^{-1} dH(s, Q),$$

where  $H(\cdot, Q)$  is the marginal distribution of  $Y$  under  $Q$ , and  $c(Q)$  is defined by  $\tau(1, Q) = 1$ . Unfortunately existence of the right side of (37) and existence and unicity of  $\tau(\cdot, Q)$  solving (37) is at this point an open question. Rewrite (37) in terms of the (function-valued) parameter  $v : [0, 1] \rightarrow [0, 1]$  by

$$v(\cdot, Q) \equiv \tau(H^{-1}(\cdot, Q), Q)$$

as

$$(38) \quad v(u, Q) = \int_0^u \{ c(Q) + \int_s^1 \int r_0(z, v(y, Q)) d\tilde{Q}(z, y) \}^{-1} ds$$

where  $\tilde{Q}(z, y) \equiv Q(z, H^{-1}(y, Q))$  is the distribution of  $(Z, H(Y, Q))$  under  $Q$ . If  $Q = P_\tau \in \mathbf{P}$ , since  $\tau(Y) \sim \text{Uniform}(0, 1)$ ,

$$H(\cdot, P_\tau) = \tau(\cdot).$$

Note that if we take  $c(P_\tau) = 1$ ,  $\tau(u, Q) = u$  is a solution of (37) for  $Q = P_\tau$ . This follows by (38) since then  $\int p_0(z, y) dm_0(z) \equiv 1$  and

$$\begin{aligned} \int_s^1 \int r_0(z, y) dQ(z, H^{-1}(y, Q)) &= \int_s^1 \int r_0(z, y) p_0(z, y) dm_0(z) dy \\ &= - \int (p_0(z, 1) - p_0(z, s)) dm_0(z) \\ &= 0. \end{aligned}$$

This suggests that we apply the one-step method to obtain the functional  $W(v, c, \mathbb{P}_n^\#)$  where  $\mathbb{P}_n^\#$  has density  $p_n^\#(\cdot, \cdot)$  and  $W : \mathbf{A} \times \mathbf{R} \rightarrow \mathbf{A} \times \mathbf{R}$  is defined by

$$(39) \quad W(v, c, Q) \equiv \begin{pmatrix} v(\cdot) - \int_0^1 \{ c + \int_s^1 \int r_0(z, v(y)) d\tilde{Q}(z, y) \}^{-1} ds \\ 1 - \int_0^1 \{ c + \int_s^1 \int r_0(z, v(y)) d\tilde{Q}(z, y) \}^{-1} ds \end{pmatrix}^T$$

The natural  $(\tilde{v}(\cdot), \tilde{c})$  is (identity, 1). It is possible under suitable conditions to show that theorem 2 applies. In fact, see Bickel (1986), it is easier simply to formally expand the first component of  $W$  around (identity, 1) and use the boundary conditions  $v(1, \mathbb{P}_n^\#) = 1$  to obtain as an approximate equation for  $\Delta(u) \equiv v(u, \mathbb{P}_n^\#) - u$ ,

$$(40) \quad \hat{\Delta}(u) = - \int_0^1 K(s, u) \alpha(s) \hat{\Delta}(s) ds - \int \int_0^1 K(s, u) r_0(z, s) d\tilde{\mathbb{P}}_n^\#(z, s),$$

where

$$K(s, u) \equiv s \wedge u - su,$$

$$\alpha(y) = \int r_0'(z, y) p_0(z, y) dm_0(z) = \int r_0^2(z, y) p_0(z, y) dm_0(z).$$

Note that (40) is meaningful even if  $\tilde{P}_n^\#$  is replaced by  $\tilde{P}_n$ , the empirical distribution of  $(Z, H_n(Y))$  where  $H_n$  is the marginal of  $P_n$ .

**Proposition 10.** Suppose that for all  $z$ ,  $p_0(z, \cdot)$  is twice continuously differentiable and  $Z$  has finite support. Further suppose  $\alpha(\cdot)$  is continuous on  $[0, 1]$ . Then

A. The equation

$$\Delta(u) + \int_0^1 K(s, u) \alpha(s) \Delta(s) ds = - \int \int_0^1 K(s, u) r_0(z, s) d\tilde{P}_n(z, s)$$

has a unique solution  $\tilde{\Delta}(\cdot)$  which is continuous, piecewise differentiable, and given by

$$\tilde{\Delta}(u) = - \int \int_0^1 \Delta(s, u) r_0(z, s) d\tilde{P}_n(z, s)$$

B. Here  $\Delta(s, u)$  is the Green's function of the Sturm-Liouville equation

$$y''(t) - \alpha(t)y(t) = 0$$

with boundary conditions  $y(0) = 0, y(1) = 0$ .

C. Furthermore if

$$\hat{\tau}(u) = H_n(u) + \tilde{\Delta}(H_n(u)) \quad \text{for } u = Y_j, \quad j = 1, \dots, n$$

and linearly interpolated in between, then  $P(\tau \text{ is a homeomorphism}) \rightarrow 1$ . In any case,  $\hat{\tau}$  is linear, regular, Gaussian with influence function given by (6.7.9) and (6.7.12).

The proof of this proposition involves establishing that the linear operator

$$L : b \rightarrow b + \int_0^1 K(s, \cdot) \alpha(s) b(s) ds$$

from  $\mathbf{B}$  to  $\mathbf{B}$  has a bounded inverse, characterizing its inverse using Sturm-Liouville theory and then applying standard weak convergence techniques to  $\sqrt{n}(P_n - P_0)$ . See Bickel (1986) for the type of argument needed. This proposition unfortunately does not cover the Cox model or the normal transformation model. The conditions on  $\alpha$  can be considerably weakened to include these transformation models; see Klaassen (1993) and Wellner (1993) for the necessary techniques.  $\square$

## 7.7 JOINT ESTIMATION OF EUCLIDEAN AND INFINITE-DIMENSIONAL PARAMETERS

Strictly speaking, since a Euclidean and an infinite-dimensional parameter can always be combined into one grand infinite-dimensional parameter this section and the last deal with the same question. In fact, we will only deal with the

following kind of situation:  $\mathbf{P} = \{P_{(\theta, G)}: \theta \in \Theta, G \in \mathbf{G}\}$ ,  $\Theta \subset R^k$ ,  $\mathbf{G} \subset l^\infty(T)$  for  $T$  a compact subset of  $R^q$ . We want to estimate  $(\theta, G)$  and have available for each  $\theta \in \Theta$ , an estimate  $\hat{G}_\theta \in l^\infty(T)$  of a parameter  $G_\theta(P) \in l^\infty(T)$  such that the following conditions hold:

- (R1)  $G_\theta(P_{(\theta, G)}) = G$  for all  $(\theta, G)$ .
- (R2) For every  $P_0 = P_{(\theta_0, G_0)} \in \mathbf{P}$  the map  $\theta \rightarrow G_\theta(P_0)$  from  $\Theta$  to  $l^\infty(T)$  is Hadamard differentiable at  $\theta_0$  with derivative  $\dot{G}_0: R^k \rightarrow l^\infty(T)$ .
- (R3) If  $\hat{G}, G(P)$  are viewed as elements of  $l^\infty(\Theta \times T)$  then  $\hat{G}$  is asymptotically linear and  $\sqrt{n}(\hat{G} - G(P))$  converges weakly to a continuous Gaussian process on  $\Theta \times T$ .

In our applications  $\hat{G}_\theta$  will be an NPML or other explicit estimate of  $G$  in the model  $\mathbf{P}_2(\theta)$ , that is, under the assumption that  $\theta$  is known. With these estimates we create AGM-estimates  $\hat{\theta}$  of  $\theta$  using the natural extension of the method outlined in example 7.3.3. That is, we find an appropriate function  $\psi: \mathbf{X} \times \Theta \times \mathbf{G} \rightarrow R^k$  such that, for all  $\theta, G$ ,

$$(1) \quad E_{(\theta, G)} \psi(X, \theta, G) = 0$$

and base our estimates on the function

$$(2) \quad W_n(\theta) \equiv W(\theta, P_n) \equiv \int \psi(x, \theta, \hat{G}_\theta) dP_n(x).$$

By (R1) and (1) we expect  $\hat{\theta}$  to be consistent and, under further conditions, asymptotically linear, regular, and Gaussian. As in the parametric case it is trivial to construct such functions  $\psi(\cdot, \theta, G)$ . For instance, if we want efficiency of  $\hat{\theta}$  at all  $P \in \mathbf{P}_1(G_0)$ , it is natural to try  $\psi(x, \theta, G) = \dot{I}_1(x, \theta, G_0) - E_{(\theta, G)} \dot{I}_1(X, \theta, G_0)$ . Given  $\hat{\theta}$  we estimate  $G$  by  $\hat{G}_{\hat{\theta}}$ , and use (R2) and (R3) to apply the delta method. If  $P_0 \equiv P_{(\theta_0, G_0)}$ , and, as in (R2),  $\dot{G}_0$  is the Hadamard derivative of  $\theta \rightarrow G_\theta(P_0)$  at  $\theta_0$ , then

$$(3) \quad \begin{aligned} \sqrt{n}(\hat{G}_{\hat{\theta}} - G_0) &= \sqrt{n}((\hat{G}_{\hat{\theta}} - G_{\hat{\theta}}) + \dot{G}_0(\hat{\theta} - \theta_0)) + o_p(1) \\ &= \sqrt{n}(\hat{G}_{\theta_0} - G_0) + \dot{G}_0 \sqrt{n}(\hat{\theta} - \theta_0) + o_p(1), \end{aligned}$$

and we obtain linearity of  $\hat{G}_{\hat{\theta}}$  from the linearity of  $\hat{\theta}$  and  $\hat{G}_{\theta_0}$ .

Here is a simple example.

**Example 1. Biased sampling regression.**

This model extends the discussion in example 4.4.3. It is treated in detail in Bickel and Ritov (1991). We observe  $X = (I, Z, Y)$  where  $I$  is a stratum indicator taking values  $1, \dots, s$ ,  $Z$  is discrete taking values  $z_1, \dots, z_K$  and

$$(4) \quad \begin{aligned} p(i, z_j, y) &= \lambda_i^* h_j f(y - v^T z_j) w(i, z_j, y) / W_i(v, f, h), \\ W_i(v, f, h) &= \sum_j h_j \int f(y) w(i, z_j, y + v^T z_j) dy \end{aligned}$$

and all  $\lambda_i^* > 0$ ,  $\sum_{i=1}^s \lambda_i^* = 1$ , all  $h_j > 0$ ,  $\sum_{j=1}^K h_j = 1$ . These parameters and  $v \in R^m$  and  $f$  are all unknown but  $w(i, z_j, y)$  is a known weight function. For

example, if  $s = 1$ ,  $w(1, z_j, y) = 1_{[y \leq y_0]}$  corresponds to truncated regression as in example 4.4.3. It can be shown that we can reparametrize (4) as

$$(5) \quad p(i, z_j, y) = \lambda_i h_j f(y - v^T z_j) w(i, z_j, y) / W(v, G)$$

where

$$W(v, G) = \sum_{i,j} \lambda_i h_j \int f(y) w(i, z_j, y + v^T z_j) dy$$

and  $\lambda_i > 0$ ,  $\sum_{i=1}^s \lambda_i = 1$ . Here  $G \leftrightarrow (\lambda, h, F)$  where  $f = F'$  and we think of  $G$  as a function on  $\{1, \dots, s\} \times \{1, \dots, K\} \times R$ . It has been shown in example 7.6.3 that, for fixed  $v$ , the NPMLE  $\hat{G}_v$  can be obtained by first getting  $\hat{\lambda}_v, \hat{h}_v$  as ordinary  $M$ -estimates based on  $s + K - 2$  equations and then getting  $\hat{F}_v$  as an explicit function of  $\hat{\lambda}_v, \hat{h}_v$ , and  $P_n$ . It is easy to see that

$$(6) \quad \dot{l}_1(x, v, G_0) = -z \frac{f_0'}{f_0}(y - v^T z) + E_{G_0}(Z \frac{f_0'}{f_0}(Y - v^T Z))$$

which leads us to consider, for a given function  $q$ ,

$$(7) \quad \psi(X, v, G) = Zq(Y - v^T Z) - \frac{1}{W(v, G)} \sum_{i,j} \lambda_i h_j z_j \int q(y) w(i, z_j, y + v^T z_j) dF(y).$$

Further manipulations reveal that setting  $W_n(v) = 0$  with  $W_n$  given by (2) and  $\psi$  as in (7) yields an additional set of  $m$  equations of the  $M$ -type (with  $\hat{\lambda}_v, \hat{h}_v$  appearing in their  $\psi$  functions). Hence  $\hat{v}$  is part of an  $M$ -estimate for the  $m + s + K - 2$  parameters  $(v, \lambda, h)$  and theorem 7.2.1 applies directly. We obtain as our estimate of  $G$ ,  $\hat{G} = (\hat{v}, \hat{\lambda}_{\hat{v}}, \hat{F}_{\hat{v}})$ .

Bickel and Ritov (1991) provide mild conditions on this model which guarantee the conditions required for the validity of theorems 7.2.1 and 7.3.2 as well as (R1)–(R3) and carry through the program of (1)–(3) above for  $(\hat{v}, F_{\hat{v}})$ . In fact, they establish regularity by verifying the conditions of corollary 7.2.1 and extending the delta method of example 7.6.3 along the lines of (3). The form of the influence functions and details of the argument are given in their paper.  $\square$

In general we cannot hope that we get an ordinary  $M$  equation. However, we can and shall give reasonably simple conditions for the validity of the “discretized” version of theorem 7.3.2 for

$$(8) \quad W_n(\theta) = \int \psi(x, \theta, \hat{G}_\theta) dP_n(x).$$

We then apply our methods to censored regression.

Here are the conditions. Fix  $P \in \mathcal{P}$ . Suppose  $\mathbf{X} \subset R^d$  is compact.

(D1) Let

$$(9) \quad U_n(x, \theta) = \psi(x, \theta, \hat{G}_\theta) - \psi(x, \theta, G_\theta)$$

where  $G_\theta \equiv G_\theta(P)$ . For every deterministic sequence  $\theta_n = \theta_0 + o(1)$

the processes  $n^{1/2} U_n(\cdot, \theta_n)$  converge weakly as elements of  $l^\infty(\mathbf{X})$  to a process concentrating on  $C(\mathbf{X})$ , the continuous functions on  $\mathbf{X}$ .

(D2) For each  $x$ ,

$$(10) \quad n^{1/2}(U_n(x, \theta_n) - U_n(x, \theta_0)) = o_p(1).$$

(D3) Define

$$(11) \quad V_n(\theta) \equiv \int \psi(x, \theta, G_\theta) d(\mathbb{P}_n - P)(x).$$

Then

$$(12) \quad n^{1/2}(V_n(\theta_n) - V_n(\theta_0)) = o_p(1).$$

(D4)  $\int \psi(\cdot, \theta_0, \hat{G}_{\theta_0}) dP(x)$  is an asymptotically linear estimate of  $\int \psi(x, \theta_0, G_0) dP(x)$ ,

$$(13) \quad \int U_n(x, \theta_0) dP(x) = \int w(y, P) d\mathbb{P}_n(y) + o_p(n^{-1/2}),$$

where  $w(y, P) \in L_2(P)$ ,  $\int w(y, P) dP(y) = 0$ .

Let

$$(14) \quad W(\theta) \equiv \int \psi(x, \theta, G_\theta) dP(x)$$

so that  $W(\theta_0) = 0$ .

(D5)  $\theta \rightarrow W(\theta)$  is differentiable at  $\theta_0$  with nonsingular derivative  $\dot{W}(\theta_0)$ .

(D6) There exists  $\tilde{\theta}_n$  a  $\sqrt{n}$ -consistent regular estimate of  $\theta$  which is discretized, i.e., with values in an  $n^{-1/2}$  grid of  $R^k$ .

**Lemma 1.** Suppose (D1)–(D4) hold. Then

$$(15) \quad W_n(\theta_n) - W_n(\theta_0) = W(\theta_n) - W(\theta_0) + o_p(n^{-1/2}),$$

$$(16) \quad W_n(\theta_0) = W(\theta_0) + \int \lambda(x, P) d\mathbb{P}_n(x) + o_p(n^{-1/2}),$$

where

$$(17) \quad \lambda(x, P) = \psi(x, \theta_0, G_0) + w(x, P).$$

**Proof.** Write

$$\begin{aligned} (a) \quad & W_n(\theta_n) - W(\theta_n) - W_n(\theta_0) + W(\theta_0) \\ &= \int (\psi(x, \theta_n, G_{\theta_n}) - \psi(x, \theta_0, G_0)) d(\mathbb{P}_n - P)(x) \\ &\quad + \int [\psi(x, \theta_n, \hat{G}_{\theta_n}) - \psi(x, \theta_n, G_{\theta_n}) \\ &\quad \quad - \psi(x, \theta_0, \hat{G}_{\theta_0}) + \psi(x, \theta_0, G_0)] dP(x) \\ &\quad + \int [\psi(x, \theta, \hat{G}_\theta) - \psi(x, \theta, G_\theta) \\ &\quad \quad - \psi(x, \theta_0, \hat{G}_{\theta_0}) + \psi(x, \theta_0, G_0)] d(\mathbb{P}_n - P)(x) \end{aligned}$$

$$= V_n(\theta_n) - V_n(\theta_0) + \int (U_n(x, \theta_n) - U_n(x, \theta_0)) dP(x) \\ + \int (U_n(x, \theta) - U_n(x, \theta_0)) d(\mathbb{P}_n - P)(x).$$

Now,

$$(b) \quad V_n(\theta_n) = V_n(\theta_0) + o_p(n^{-1/2})$$

by (12). Furthermore,

$$(c) \quad \left| \int (U_n(x, \theta_n) - U_n(x, \theta_0)) dP(x) \right| \\ \leq \left| \int (U_n(x, \theta_n) - U_n(x, \theta_0)) 1_{\{\|U_n(\cdot, \theta_n) - U_n(\cdot, \theta_0)\|_\infty \leq Mn^{-1/2}\}} dP(x) \right| \\ + \|U_n(\cdot, \theta_n) - U_n(\cdot, \theta_0)\|_\infty 1_{\{\|U_n(\cdot, \theta_n) - U_n(\cdot, \theta_0)\|_\infty > Mn^{-1/2}\}}.$$

The first term on the right-hand side of (c) is  $o_p(n^{-1/2})$  under  $P$  by (10) and the dominated convergence theorem. On the other hand,

$$P(\text{second term} \geq \varepsilon n^{-1/2}) \leq P(\|U_n(\cdot, \theta_n) - U_n(\cdot, \theta_0)\|_\infty > Mn^{-1/2})$$

and by (D1)

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} P(\|U_n(\cdot, \theta_n) - U_n(\cdot, \theta_0)\|_\infty > Mn^{-1/2}) = 0.$$

So

$$(d) \quad \int (U_n(x, \theta_n) - U_n(x, \theta_0)) dP(x) = o_p(n^{-1/2}).$$

By (D1),  $Z_n = \sqrt{n}U_n(\cdot, \theta_n)$  converges weakly in  $l^\infty(\mathbf{X})$  to  $Z$  which takes its values in  $\mathbf{C}(\mathbf{X})$ . Hence, by the Skorokhod construction theorem A.8.2, there exists a probability space and versions of  $Z_n$  and  $Z$  such that

$$\|Z_n - Z\|_\infty \rightarrow 0 \quad \text{a.s.}$$

Since

$$(e) \quad \left| \int Z_n d(\mathbb{P}_n - P) \right| \leq 2 \|Z_n - Z\|_\infty + \left| \int Z d(\mathbb{P}_n - P) \right|,$$

it suffices to prove

$$(f) \quad \left| \int Z d(\mathbb{P}_n - P) \right| \rightarrow_p 0.$$

Now, since the class  $BL(\mathbf{X}, |\cdot|)$  of bounded Lipschitz functions is dense in  $\mathbf{C}(\mathbf{X})$  for  $\|\cdot\|_\infty$  by Theorem 11.2.4 of Dudley (1989) and the compactness of  $\mathbf{X}$ , there exist for every  $\varepsilon > 0$  a finite constant  $M$  and a process  $Z_\varepsilon$  taking values in the class  $BL_M$  of bounded Lipschitz functions which are bounded in norm by  $M$  such that

$$(g) \quad P(\|Z - Z_\varepsilon\|_\infty \geq \varepsilon) \leq \varepsilon.$$

Now

$$(h) \quad \left| \int Z_\varepsilon d(\mathbb{P}_n - P) \right| \leq \sup\{\left| \int h d(\mathbb{P}_n - P) \right| : h \in BL_M\} \rightarrow_{a.s.} 0$$



by corollary 11.3.4 of Dudley (1989), and together with (g) this proves (f) and hence (15).

Write

$$(i) \quad W_n(\theta_0) = \int \psi(x, \theta_0, G_0) dP_n(x) + \int U_n(x, \theta_0) dP(x) \\ + \int U_n(x, \theta_0) d(P_n - P)(x).$$

By (e) through (h) the last term in (i) is  $o_p(n^{-1/2})$ , while the rest of (16) follows from (13).  $\square$

Let  $\lambda_n \downarrow 0$ ,  $n^{1/2} \lambda_n \rightarrow \infty$ . For a fixed basis vector  $e_i \in R^k$  let  $\tilde{\theta}_{\lambda_i} = \tilde{\theta}_n + \lambda_n e_i$  and

$$\tilde{W}_{ij} = \lambda_n^{-1} \int (\psi_j(x, \tilde{\theta}_{\lambda_i}, \hat{G}_{\tilde{\theta}_{\lambda_i}}) - \psi_j(x, \tilde{\theta}_n, \hat{G}_{\tilde{\theta}_n})) dP_n(x),$$

where  $\psi = (\psi_1, \dots, \psi_k)^T$ .

**Lemma 2.** If (D1)–(D3) and (D5), (D6) hold, then

$$(18) \quad \hat{W} \equiv [\hat{W}_{ij}] = \dot{W}(\theta_0) + o_p(1).$$

**Proof.** Since  $\tilde{\theta}_n$  is discretized, it suffices (see the proof of theorem 2.5.2) to prove (18) for  $\tilde{W}_{ij}$  which is  $\hat{W}_{ij}$  with  $\tilde{\theta}_n$  and  $\tilde{\theta}_{\lambda_i}$  replaced by arbitrary sequences of constants  $\theta_n$  and  $\theta_{\lambda_i} = \theta_n + \lambda_n e_i$  say. Let  $U_{nj}, V_{nj}$  denote the  $j$ th components of  $U_n, V_n$  respectively. Then

$$\tilde{W}_{ij} = \lambda_n^{-1} \int (\psi_j(x, \theta_{\lambda_i}, G_{\theta_{\lambda_i}}) - \psi_j(x, \theta_n, G_{\theta_n})) dP(x) \\ + \lambda_n^{-1} \{ \int (U_{nj}(x, \theta_{\lambda_i}) - U_{nj}(x, \theta_n)) dP_n(x) \\ + V_{nj}(\theta_{\lambda_i}) - V_{nj}(\theta_n) \} \\ = \lambda_n^{-1} \int (\psi_j(x, \theta_{\lambda_i}, G_{\theta_{\lambda_i}}) - \psi_j(x, \theta_n, G_{\theta_n})) dP(x) \\ + o_p(\lambda_n^{-1} n^{-1/2})$$

by (e) and (b) of the proof of lemma 1. The lemma follows from (D5).  $\square$

Lemmas 1 and 2 and theorem 7.3.2 yield the following theorem:

**Theorem 1.** Suppose (D1)–(D6) hold.

A. Then  $\hat{\theta}_n \equiv \tilde{\theta}_n - \hat{W}^{-1} \dot{W}_n(\tilde{\theta}_n)$  is asymptotically linear and regular with influence function  $-\dot{W}^{-1}(\theta_0) \lambda(x, P_0)$ .

B. If further (R1)–(R3) hold and, by (R3), we write

$$(19) \quad \hat{G}_{\hat{\theta}_n}(\cdot) = G_0(\cdot) + \int \gamma(x, \cdot, P_0) dP_n(x) + o_p(n^{-1/2}),$$

then  $\hat{G}_{\hat{\theta}_n}$  is linear and regular with influence function

$$\gamma(x, \cdot, P_0) + \dot{G}_0 \dot{W}^{-1}(\theta_0) \lambda(x, P_0) \in I^\infty(T).$$

**Remark 1.** Suppose that (D1)–(D3) are strengthened as follows. For some  $\Delta_n \rightarrow 0$  with  $n^{1/2} \Delta_n \rightarrow \infty$ :

(D1')  $\lim_{\delta \rightarrow 0} \limsup_{n \rightarrow \infty} P(\omega(\delta, \Delta_n) > \varepsilon) = 0$ , where

$$\omega(\delta, \Delta_n) \equiv \sup\{\sqrt{n} |U_n(x, \theta) - U_n(x', \theta_0)| : |x - x'| \leq \delta, |\theta - \theta_0| \leq \Delta_n\}.$$

(D2')  $\sup\{n^{1/2} |U_n(x, \theta) - U_n(x, \theta_0)| : |\theta - \theta_0| \leq \Delta_n\} = o_p(1)$ .

(D3')  $\sup\{n^{-1/2} |V_n(\theta) - V_n(\theta_0)| : |\theta - \theta_0| \leq \Delta_n\} = o_p(1)$ .

Then we do not need to restrict  $\tilde{\theta}_n$  to be discretized.

**Remark 2.** If  $\theta = (\nu, \eta)$  where  $\eta$  becomes unidentifiable if  $G$  varies, we can treat  $H = (\eta, G)$  as a new identifiable nuisance parameter and proceed as before, replacing  $H$  by, for example,  $\hat{H}_\nu$ , the NPMLE for  $H$  with  $\nu$  fixed.

**Remark 3.** Under subsidiary conditions (see theorem 7.3.3) we can establish the existence of a solution  $\hat{\theta}_n$  of  $W_n(\theta) = 0$  which is asymptotically linear with influence function  $-\dot{W}^{-1}\lambda(\cdot, P)$ . We do not go into this further here.

**Example 2. Censored regression.**

In this example we follow Ritov (1990) in exhibiting (locally) efficient estimates for the coefficients of a linear regression and the distribution function of the error when right censoring is present. We employ the notation of example 4.6.4. Buckley and James (1979), who were the first to apply the method of example 7.3.3 in a semiparametric setting, did so in this model.

We begin with a function  $w : R \rightarrow R$  which is bounded, and note that, if  $P \in \mathbf{P}$ ,

$$(20) \quad E(Z - EZ)w(\varepsilon) = 0.$$

Suppose for simplicity that the distribution of  $Z$  is known. Then let

$$\begin{aligned} (21) \quad \psi(x, \nu, F, H) &= E((Z - EZ)w(\varepsilon) | X = x) \\ &= (z - EZ)[\Delta w(\nu - \nu^T z) \\ &\quad + (1 - \Delta)E(w(\varepsilon) | V = \nu, Z = z, \Delta = 0)] \\ &= (z - EZ)(w(e_\nu) - (1 - \Delta)R_F w(e_\nu)), \end{aligned}$$

where, as defined in (A.1.19),  $R_F$  is the  $R$  operator corresponding to  $F$  given by

$$R_F w(u) = w(u) - \frac{\int_u^\infty w(t) dF(t)}{1 - F(u)},$$

and  $e_\nu \equiv \nu - \nu^T Z$ . From (20), it follows that  $\psi(\cdot, \nu, F, H)$  satisfies (1). It has the attractive feature of not depending on the conditional distribution of  $C$  given  $Z$ . If  $\nu$  is fixed, note that by example 7.6.1 the NPMLE of  $F$  is just the Kaplan-Meier estimate based on observing a sample distributed as  $((Y - \nu^T Z) \wedge (C - \nu^T Z), 1_{[Y - \nu^T Z \leq C - \nu^T Z]})$ , while  $\hat{H}_\nu$  does not depend on  $\nu$  and is just the empirical distribution of the  $Z_i$ .

Let

$$w_1(u) = -\{1 - F_0(u)\}R_{F_0} w(u)\mathcal{K}_0(u)P(C - \nu^T Z \geq u),$$

and

$$w_2(u) = f_0(u) \{R_{F_0} w(u)\}^2 P(C - v^T Z \geq u),$$

where

$$\lambda'_0(u) = \left\{ \frac{f'_0(u)}{1 - F_0(u)} + \frac{f_0^2(u)}{(1 - F_0(u))^2} \right\}.$$

Ritov (1990) proves a generalization of the following theorem.

**Theorem 2.** Suppose that  $F_0$  has finite Fisher information,  $Z$  has compact support, and that  $s(\cdot)$  is a function with  $\int s^2(t) dF_0(t) < \infty$ ,  $s(\cdot)$  has a derivative finitely integrable over compact sets, and let

$$(22) \quad w(t) \equiv s(t \wedge c)$$

where  $c$  is fixed such that  $\bar{G}(c) \equiv P((Y \wedge C) - v_0^T Z > c) > 0$ . Suppose further that  $A \equiv \int \text{Var}(Z | C - v_0^T Z \geq u) w_1(u) du$  is nonsingular where  $\text{Var}(Z | E)$  is the conditional covariance matrix of  $Z$  given the event  $E$ . Then, with probability converging to one:

- A. If  $W_n$  is defined by (8) and (21),  $W_n(v, P_n) = 0$  has a solution  $\hat{v}_n$ .
- B.  $\hat{v}_n$  is  $\sqrt{n}$ -consistent and  $\sqrt{n}(\hat{v}_n - v_0) \rightarrow_d N(0, A^{-1}VA^{-1})$  where  $V = \int \text{Var}(Z | C - v_0^T Z \geq u) w_2(u) du$ .

We sketch an argument using lemmas 1 and 2 which establishes the same result for the one step estimator  $\hat{v}_n$  given in theorem 1. Using further arguments of Tsiatis (1990) it is possible to establish (D1')–(D3') for  $n^{1/2} \Delta_n$  bounded. If  $n^{1/2} \Delta_n$  is not bounded the errors become  $o_p(\Delta_n)$ . This is still quite adequate to establish lemma 2 and hence theorem 1 without the discreteness of  $\tilde{\theta}_n$ .

We identify

$$\bar{F}_v(t) = \exp\left\{-\int_0^t \frac{dEQ_v(u)}{E\bar{Q}_v(u)}\right\},$$

where

$$(23) \quad \begin{aligned} EQ_v(u) &= P(\epsilon \leq [u + (v - v_0)^T Z] \wedge [C - v_0^T Z]) \\ &= E\{F_0([u + (v - v_0)^T Z] \wedge [C - v_0^T Z])\}, \end{aligned}$$

$$(24) \quad E\bar{Q}_v(u) = E\{\bar{F}_0(u + (v - v_0)^T Z)P(C \geq u + (v - v_0)^T Z)\}.$$

We showed in example 7.6.1 that if  $\hat{F}_v$  is the Kaplan-Meier estimate of  $F_v$ , then  $\sqrt{n}(\hat{F}_{v_0} - F_0)$  converges weakly in  $l^\infty[0, c]$  to a process with a.s. continuous sample functions and

$$(25) \quad \|\hat{F}_{v_0}(\cdot) - F_0(\cdot) - n^{-1} \sum_{i=1}^n \mathcal{Y}(\cdot, X_i, P)\|_\infty = o_p(n^{-1/2})$$

where, with  $x = (z, v, \Delta)$ ,

$$(26) \quad \gamma(t, x, P) = \bar{F}_0(t) [\Delta 1_{[v \leq t]} / \bar{G}(v) - E_0(\Delta 1_{[V \leq t]} / \bar{G}(V))].$$

The same argument readily yields that weak convergence to the same limit continues to hold if  $v_0$  is replaced by  $v_n = v_0 + o(1)$ , and then

$$(27) \quad \|\hat{F}_{v_n}(\cdot) - F_{v_n}(\cdot) - n^{-1} \sum_{i=1}^n \gamma_n(t, X_i, P)\|_\infty = o_p(n^{-1/2})$$

where

$$(28) \quad \gamma_n(t, x, P) = \bar{F}_{v_n}(t) \left\{ \Delta \frac{[1_{v \leq (v_n - v_0)^T z + t}]}{\bar{G}(v - (v_n - v_0)^T z)} - E \left\{ \Delta \frac{[1_{V \leq (v_n - v_0)^T Z + t}]}{\bar{G}(V - (v_n - v_0)^T Z)} \right\} \right\}.$$

So, if  $U_n$  is given by (9),

$$U_n(x, v) = (z - \bar{Z})(1 - \Delta)(R_{F_v}(e_v) - R_{\hat{F}_v}(e_v)) \\ + (EZ - \bar{Z})(\Delta s(e_v) + (1 - \Delta)[w(e_v) - R_{F_v}(e_v)]).$$

By the assumption on  $w$  and  $s$ , and the boundedness of  $|Z|$  (by  $M$  say), it is evidently enough to establish (D1) and (D2) with  $\psi(x, v, F, H)$  replaced by  $w(e_v) - R_{F_v}(e_v)$  in the definition of  $U_n$ . Further, by our discussion and the continuity of  $F_0$  we have  $\hat{F}_{v_n}(t + |v_n| M) \rightarrow_p F_0(t) > 0$ . Since  $w(t)$  is constant for  $t \geq c$ ,  $R_{\hat{F}_v}(t) = R_{F_v}(t)$  for  $t \geq c$ . By further elementary manipulations we see that it is enough to establish (D1) and (D2) with  $U_n$  replaced by

$$(29) \quad \tilde{U}_n(z, v, v) \equiv \int_{e_v}^{\infty} \frac{w(t) d(\hat{F}_v(t) - F_v(t))}{\bar{F}_v(e_v)} \\ + (\hat{F}_v(e_v) - F_v(e_v)) \int_{e_v}^{\infty} \frac{w(t) dF_v(t)}{\bar{F}_v^2(e_v)} \\ = - \int_{e_v}^{\infty} \frac{(\hat{F}_v(t) - F_v(t)) w'(t) dt}{\bar{F}_v(e_v)} \\ - (\hat{F}_v(e_v) - F_v(e_v)) \frac{R_{F_v}(e_v)}{\bar{F}_v(e_v)},$$

viewed as a random element with values in  $I^\infty([-\infty, c] \times [z : |z| \leq M])$ . Claim (D1) now follows readily from the weak convergence of  $n^{1/2}(\hat{F}_v - F_v)$ . For (D2) note that for  $e_0 \leq c$  where we write  $e_0$  for  $e_{v_0}$ ,

$$\tilde{U}_n(z, v, v_n) - \tilde{U}_n(z, v + (v_n - v_0)^T z, v_n) = o_p(n^{-1/2})$$

by weak convergence. By (22)–(29)

$$(30) \quad \tilde{U}_n(z, v + (v_n - v_0)^T z, v_n) - \tilde{U}_n(z, v, v_0)$$

$$\begin{aligned}
 &= - \int_{e_0}^{\infty} w'(t) \int (\gamma_n(t, y, P) - \gamma(t, y, P)) dP_n(y) dt / \bar{F}_0(e_0) \\
 &\quad - \int (\gamma_n(e_0, y, P) - \gamma(e_0, y, P)) dP_n(y) \frac{R_{F_0}(e_0)}{\bar{F}_0(e_0)} \\
 &\quad + R_n(z, v, v_n),
 \end{aligned}$$

where  $\|R_n\|_{\infty} = o_p(n^{-1/2})$ . It is easy to see that (10) follows by a second moment calculation. We get (12) by a similar  $L_2$  calculation. To complete the proof of our simplified version of Ritov's (1990) theorem, we need to exhibit a  $\sqrt{n}$ -consistent estimate  $\tilde{v}_n$  and establish regularity of  $\hat{v}_n$ . Example 7.4.3 and proposition 7.4.1 provide such an estimate.  $\square$

### 7.8 EFFICIENT ESTIMATION

In sections 7.3–7.7 we have shown how GM estimating equations can be used for the estimation of both Euclidean and abstract parameters. In some examples, 7.6.1, 7.6.2, 7.6.3, 7.5.5, we noted that the influence functions so obtained agreed with efficient influence functions computed in chapters 3 to 6. In other examples, 7.6.4, 7.7.1 and 7.7.2 we did not arrive at efficient, but merely at regular estimates. This section has two parts:

- A brief heuristic discussion (suggested by Gill (1988)) of when GM estimating equations can be expected to lead to efficient estimates of abstract parameters and of Euclidean parameters which are functions of the latter.
- An extended and rigorous account of how efficient estimates of a Euclidean parameter  $v$  can be constructed given a  $\sqrt{n}$ -consistent estimate of  $v$  and a suitable estimate of the efficient influence function. We apply these constructions to group and mixture models.

#### *When Do GM Estimating Equations Lead to Efficient Estimates?*

Suppose we are given functions  $W_n : A \times M_0 \rightarrow \Gamma$  which satisfy the hypotheses of theorem 7.6.1. When will a GM-estimate  $\hat{v}_n$  based on  $W_n$  be weakly  $B_0^*$  efficient? That is, when will  $b^* \hat{v}_n$  be an efficient estimate of  $b^* v$  for all  $b^* \in B_0^*$ ? The trivial answer from proposition 3.3.1 is that the influence function of  $b^* \hat{v}_n$  has to belong to  $\dot{P}$ . From corollary 7.6.1 this means (under regularity conditions)

$$(1) \quad \gamma^* w(\cdot, P) \in \dot{P}$$

for all  $\gamma^* \in \Gamma_0^*$ , where  $\Gamma_0^*$  is related to  $B_0^*$  as in the corollary.

We give heuristic arguments that the maximum likelihood related methods should lead to efficient procedures.

**Ordinary MLE.**

Suppose  $\mathbf{P} = \{P_\nu : \nu \in \mathbf{A}\}$ ,  $\mathbf{A} \subset \mathbf{H}$ , a Hilbert space. Let  $s : \mathbf{A} \rightarrow L_2(\mu)$  and  $\dot{s}(\nu) : \dot{\mathbf{A}} \rightarrow L_2(\mu)$  be as in chapter 5, and let  $\dot{\mathbf{I}}(\nu) : \dot{\mathbf{A}} \rightarrow L_2(P_\nu)$  be the bounded linear operator as in (5.4.3). The likelihood equations for  $\hat{\nu}_n$  are then

$$(2) \quad \int \dot{\mathbf{I}}(\hat{\nu}_n) h(x) dP_n(x) = 0$$

for all  $h \in \dot{\mathbf{A}}$ . That is,  $\Gamma_0^* = \Gamma^* = \Gamma = \dot{\mathbf{A}}$  and  $W(\nu, P)$  with  $W(\nu, P)(h) = \int \dot{\mathbf{I}}(\nu) h(x) dP(x)$  maps  $\dot{\mathbf{A}}$  into  $R$ , if we assume  $(1/s(\nu)) dP/d\mu \in L_2(\mu)$ . Note that  $\dot{s}(\nu)(h) \in L_2(\mu)$ . The function  $w$  corresponding to  $W$  defined in (7.6.14) satisfies

$$(3) \quad w(x, P_\nu)(h) = \dot{\mathbf{I}}(\nu)h(x).$$

By definition  $w(\cdot, P)(h) \in \dot{\mathbf{P}}$  for all  $h \in \dot{\mathbf{A}}^0$ . Since the linear map  $h \rightarrow w(\cdot, P)(h)$  is continuous in view of the boundedness of  $\dot{\mathbf{I}}$ , we can conclude  $w(\cdot, P)(h) \in \dot{\mathbf{P}}$  for all  $h \in \dot{\mathbf{A}}$ , and we have efficiency. If, further,  $\hat{\nu}_n$  and  $W$  satisfy the conditions of theorem 7.6.1 we have

$$(4) \quad \dot{W}(\nu_0)(\hat{\nu}_n - \nu(P)) = n^{-1} \sum_{i=1}^n w(X_i, P) + o_p(n^{-1/2}).$$

Hence, if the estimate  $\hat{q}(\hat{\nu}_n)$  of the real parameter  $q(\nu)$  satisfies the conditions of corollary 7.6.2, its influence function will be of the form  $w(X_i, P)(h)$ , so that (1) holds and we have efficiency. Unfortunately, this heuristic argument can typically be made rigorous only for parametric models. In non- and semiparametric cases, as we have seen,  $\hat{\nu}_n$  typically doesn't exist. The heuristic suggests, however, that NPMLE's should be efficient. Indeed they are in examples 7.6.1–7.6.3.

The argument for the regularized MLE is the same if we assume  $\int w(x, P)(h) dP_n^\#(x) = \int w(x, P)(h) dP_n(x) + o_p(n^{-1/2})$ . Indeed this is the case for the Has'minskii-Ibragimov-model of example 7.6.6.

**Penalized MLE.**

We simply replace (2) by

$$(5) \quad \int \dot{\mathbf{I}}(\hat{\nu}_n)(h)(x) dP_n(x) = \lambda_n \dot{J}(\hat{\nu}_n)(h)$$

with  $\dot{J}(\nu) : \dot{\mathbf{A}} \rightarrow R$ , so that

$$(6) \quad W_n(\nu, P)(h) = \int \dot{\mathbf{I}}(\nu)(h)(x) dP(x) - \lambda_n \dot{J}(\nu)(h).$$

To satisfy (7.6.13) we must have  $\lambda_n = o(n^{-1/2})$  and

$$(7) \quad W(\nu, P)(h) = \int \dot{\mathbf{I}}(\nu)(h)(x) dP(x)$$

as before. Example 7.6.5 is an instance of such efficiency.

**Sieves.**

Suppose  $\mathbf{P}_m$  is parametrized by  $\mathbf{A}_m \subset \mathbf{A}$  with corresponding  $\dot{\mathbf{A}}_m \subset \dot{\mathbf{A}}$ . We are led to

$$\int \dot{\mathbf{i}}(\hat{v}_n)(\pi_m h)(x) dP_n(x) = 0,$$

for  $\hat{v}_n \in \mathbf{A}_m, h \in \dot{\mathbf{A}}$  where  $\pi_m h \in \dot{\mathbf{A}}_m$  is, for example, the projection of  $h$  into  $\dot{\mathbf{A}}_m$ . Thus, we can formally define

$$(8) \quad W_n(v, P)(h) = \int \dot{\mathbf{i}}(v)(\pi_m h)(x) dP(x).$$

Again (7.6.13) suggests we need  $\sup\{\|\pi_m h - h\|: \|h\| = 1\} = o(n^{-1/2})$  and  $W$  given by (7). Example 7.6.4 is an instance.

The examples suggest that the essential difficulty in efficient estimation by likelihood-based methods is not in establishing efficiency but rather in showing that the linearization assumptions (GGM2), (GGM4) hold for  $W_n$  satisfying the bias control assumption (7.6.13).

It is also interesting to see when the method outlined in section 7.6 leads to efficient estimates. Suppose  $\mathbf{P} = \{P_{(\theta, G)}: \theta \in \Theta, G \in \mathbf{G}\}$ . Let  $\psi(x, \theta) = \dot{\mathbf{i}}_1(x, \theta, G_0)$  for  $G_0$  fixed and  $\psi(x, \theta, G) \equiv \psi(x, \theta) - \int \psi(x, \theta) dP_{(\theta, G)}(x)$ . Suppose that  $\hat{G}_\theta, \psi(\cdot, \theta, G)$  satisfy all the assumptions of theorem 7.7.1. Then, by (7.7.19),

$$(9) \quad \hat{G}_{\theta_0}(y) - G_0(y) = \int \gamma(x, y, P_0) dP_n(x) + o_p(n^{-1/2})$$

and by (D4) of section 7.7

$$(10) \quad \int \psi(x, \theta_0, \hat{G}_{\theta_0}) dP_0(x) = \int w(y, P_0) dP_n(y) + o_p(n^{-1/2}).$$

Suppose we have a linear expansion in  $G$  of the form

$$(11) \quad \begin{aligned} \psi(x, \theta_0, \hat{G}_{\theta_0}) - \psi(x, \theta_0, G_0) \\ = \int c(x, y, P_0)(\hat{G}_{\theta_0}(y) - G_0(y)) dm(y) + o_p(n^{-1/2}). \end{aligned}$$

Then we would expect after substituting (9) in (11) that, for  $w$  given in (10),

$$\begin{aligned} \int \int \int c(x, y, P_0) \gamma(u, y, P_0) dm(y) dP_n(u) dP_0(x) \\ = \int w(u, P_0) dP_n(u), \end{aligned}$$

and hence,

$$(12) \quad w(u, P_0) = \int \int c(x, y, P_0) \gamma(u, y, P_0) dm(y) dP_0(x).$$

If  $\hat{G}_{\theta_0}$  is efficient in estimating  $G$  in  $\mathbf{P}_2(\theta_0)$ , then  $\gamma(\cdot, y, P_0) \in \dot{\mathbf{P}}_2(P_0)$  for all  $y$ . Suppose that the right side of (12) can be approximated in  $L_2(P_0)$  by Riemann-Stieltjes sums or equivalently that it is the  $L_2(P_0)$  limit of elements of the form

$$\sum_{j=1}^J m_j \left( \int c(x, y_j, P_0) dP_0(x) \right) \gamma(\cdot, y_j, P_0)$$

for scalar  $m_j, y_j \in \mathbf{X}, 1 \leq j \leq J$ . This is known as *Bochner integrability*. If  $\gamma(\cdot, y, P_0) \in \dot{\mathbf{P}}_2(P_0)$  for all  $y$  this yields  $w(\cdot, P_0) \in \dot{\mathbf{P}}_2(P_0)$ . Since  $\psi(\cdot, \theta_0) \in \dot{\mathbf{P}}_1(P_0)$  we conclude that, under these assumptions,  $\lambda(\cdot, \cdot)$  from (7.7.17) satisfies

$$(13) \quad \lambda(\cdot, P_0) \in \dot{P}_1(P_0) + \dot{P}_2(P_0) \subset \dot{P}(P_0).$$

Hence,  $\hat{\theta}_n$  given by theorem 7.7.1 is efficient within  $\mathbf{P}$  at  $P_{(\theta_0, G_0)}$  for all  $\theta_0$  or equivalently at all points of the parametric submodel  $\mathbf{P}_1(G_0)$ . This is shown by Ritov (1990) for censored regression; see example 7.7.2.

**Remark 1.** If either  $\hat{G}_{\theta_0}$  is not efficient in  $\mathbf{P}_2(\theta_0)$  or  $\psi(\cdot, \theta_0) \notin \dot{P}_1(G_0)$  even efficiency at the points of  $\mathbf{P}_1(G_0)$  fails.

**Remark 2.**  $\hat{G}_{\hat{\theta}}$  is under the same conditions efficient at all points of  $\mathbf{P}_1(G_0)$ . Thus we expect and it can be shown that the estimates of examples 7.7.1 and, if  $w(\epsilon) = \epsilon$ , those of 7.7.2 are efficient at all points of the Gaussian submodel. In order to have efficiency on  $\mathbf{P}$  it is necessary to use  $\psi(x, \theta)$  itself depending on  $\hat{G}_{\hat{\theta}}$ , that is, we need to estimate  $\dot{I}_1(\cdot, \theta, G)$ . This has been done by Ritov (1984) for censored regression. We now discuss this issue in a more general context.

### *Necessary and Sufficient Conditions for the Construction of Efficient Estimators*

For regular parametric models a general construction of efficient estimators of the parameter indexing the model has been given in Theorem 2.5.2. According to (2.5.6) this construction is based on a  $\sqrt{n}$ -consistent estimator and knowledge of the efficient influence function. This suggests that for estimation of Euclidean parameters in nonparametric and semiparametric models existence of a  $\sqrt{n}$ -consistent estimator, computation of the efficient influence function, and a suitable estimator thereof should suffice for the construction of an efficient estimator. This is essentially true as we shall show in corollary 1 below.

Suppose that

$$(14) \quad \mathbf{P} = \{P_{(\theta, G)} : \theta \in \Theta \subset R^k, G \in \mathbf{G}\}$$

is a semiparametric model as in section 3.4. For simplicity, suppose that  $v = \theta$  (so  $m = k$ ). Under the conditions of corollary 3.4.1 the efficient influence function for estimation of  $\theta$  is well defined and equals

$$(15) \quad \tilde{I}(x, P_0 | \theta, \mathbf{P}) = I^{-1}(P_0 | \theta, \mathbf{P}) I_0^*(x, P_0 | \theta, \mathbf{P})$$

at  $P_0 = P_{(\theta_0, G_0)}$ . Now we want to emphasize the dependence of  $\tilde{I}$  on  $\theta_0$  and  $G_0$ , so (dropping the 0 subscript) we write

$$(16) \quad \tilde{I}(x; \theta; G) \equiv \tilde{I}(x, P_{(\theta, G)} | \theta, \mathbf{P}).$$

Since  $\tilde{I} \in \dot{P}$ , we have

$$(17) \quad \int \tilde{I}(x; \theta; G) dP_{(\theta, G)}(x) = E\tilde{I} = 0$$

and

$$(18) \quad \int |\tilde{I}(x; \theta; G)|^2 dP_{(\theta, G)}(x) = E|\tilde{I}|^2 < \infty$$

for all  $\theta \in \Theta$ .



More generally, we consider an influence function  $\psi(x; \theta; G)$  which is not necessarily efficient, but still satisfies (17) and (18). That is:

- (i)  $\int \psi(x; \theta; G) dP_{(\theta, G)}(x) = 0$ , and  
 (ii)  $\int |\psi(x; \theta; G)|^2 dP_{(\theta, G)}(x) < \infty$  for all  $(\theta, G) \in \Theta \times G$ .

Our first theorem gives conditions for existence of an estimator which is locally asymptotically linear on  $\mathbf{P}_1 = \mathbf{P}_1(G) = \{P_{(\theta, G)} : \theta \in \Theta\}$  with influence function  $\psi$ . It requires (i) and (ii) and the following additional conditions:

- (iii) (Contiguity)  $\{P_{(\theta_n, G)}^n\} \triangleleft \triangleright \{P_{(\theta, G)}^n\}$  for all  $(\theta, G) \in \Theta \times G$  and all sequences  $\{\theta_n\}$  with  $\sqrt{n} |\theta_n - \theta| = O(1)$ .

- (iv) (Smoothness).

$$\sqrt{n} \left\{ \theta_n - \theta + n^{-1} \sum_{i=1}^n [\psi(X_i; \theta_n; G) - \psi(X_i; \theta; G)] \right\} = o_{P_{(\theta, G)}}(1) \quad \text{for all } (\theta, G) \in \Theta \times G \text{ and all sequences } \{\theta_n\} \text{ with } \sqrt{n} |\theta_n - \theta| = O(1).$$

- (v) (Preliminary estimator) There exists a  $\sqrt{n}$ -consistent preliminary estimator  $\tilde{\theta}_n$  of  $\theta$ , i.e.,

$$\sqrt{n}(\tilde{\theta}_n - \theta) = o_{P_{(\theta, G)}}(1)$$

for all  $(\theta, G) \in \Theta \times G$ .

- (vi) ( $\sqrt{n}$ -unbiased consistent estimator of  $\psi$ ). There exists an estimator  $\tilde{\Psi}_n(\cdot; \cdot; \underline{X}) \equiv \tilde{\Psi}_n(\cdot; \cdot; X_1, \dots, X_n)$  of the function  $\psi(\cdot; \cdot; G)$  satisfying

$$(19) \quad \sqrt{n} \int \tilde{\Psi}_n(x; \theta_n; \underline{X}) dP_{(\theta_n, G)}(x) = o_{P_{(\theta_n, G)}}(1)$$

and

$$(20) \quad \int |\tilde{\Psi}_n(x; \theta_n; \underline{X}) - \psi(x; \theta_n; G)|^2 dP_{(\theta_n, G)}(x) = o_{P_{(\theta_n, G)}}(1)$$

for all  $(\theta, G) \in \Theta \times G$  and all sequences  $\{\theta_n\}$  with  $\sqrt{n} |\theta_n - \theta| = O(1)$ .

**Theorem 1.** Suppose that (i)–(vi) hold. Then, there exists a sequence of estimators  $\hat{\theta}_n$  of  $\theta$  satisfying

$$(21) \quad \sqrt{n} \left\{ \hat{\theta}_n - \theta_n - \frac{1}{n} \sum_{i=1}^n \psi(X_i; \theta_n; G) \right\} = o_{P_{(\theta_n, G)}}(1)$$

for every  $(\theta, G) \in \Theta \times G$  and every sequence  $\{\theta_n\}$  with  $\sqrt{n} |\theta_n - \theta| = O(1)$ .

If the model  $\mathbf{P}_1(G)$  is regular for each  $G \in G$ , then the contiguity assumption (iii) holds in view of proposition 2.1.3. If  $\psi = \tilde{I}$ , the smoothness condition (iv) also follows from regularity of  $\mathbf{P}_1(G)$  if  $I_0^* = \tilde{I}_1$ , or from existence of a regular least favorable family; see (2.1.15) or Bickel (1982, (6.43) page 670).

We see that conditions (i) through (iv) are automatically fulfilled, when  $\mathbf{P}$  is a regular parametric model and  $\psi = \tilde{I}$ . Furthermore, (vi) is satisfied then since there is nothing to estimate;  $\tilde{I}(\cdot; \cdot; G)$  is known. Consequently, for a regular parametric model theorem 1 implies theorem 2.5.2.

Theorem 1 with  $\psi = \tilde{I}$  has been proved by Bickel (1982) for the case where there exists  $\tilde{\psi}_n$  for which the left side of (19) vanishes. Schick (1986) has a proof of essentially theorem 1 with  $\psi = \tilde{I}$ . The particular construction of  $\hat{\theta}_n$  which we use is from Klaassen (1987). All these authors base their proofs on "data splitting."

Here is the estimator we will use in theorem 1 which we will prove after a discussion and the statement of theorem 2.

Suppose that  $\{\lambda_n\}$ ,  $\{\mu_n\}$ , and  $\{v_n\}$  are sequences of positive integers with  $0 < \lambda_n < \mu_n < v_n < n$  and

$$\frac{\lambda_n}{n} \rightarrow \lambda, \quad \frac{\mu_n}{n} \rightarrow \mu, \quad \frac{v_n}{n} \rightarrow \nu,$$

where  $0 < \lambda < \mu < \nu < 1$ . Thus  $\lambda_n$ ,  $\mu_n$ , and  $v_n$  split the data into four independent blocks. Define two independent preliminary estimators of  $\theta$ , using the first and third blocks, by

$$\tilde{\theta}_{n1} \equiv \tilde{t}_{\lambda_n}(X_1, \dots, X_{\lambda_n}), \quad \tilde{\theta}_{n2} \equiv \tilde{t}_{v_n - \mu_n}(X_{\mu_n+1}, \dots, X_{v_n}),$$

where  $\tilde{\theta}_n = \tilde{t}_n(X_1, \dots, X_n)$  for all  $n \geq 1$ . Define two independent influence function estimators, using the second and fourth blocks, by

$$\tilde{\psi}_{n1}(x; \theta) \equiv \tilde{\psi}_{\mu_n - \lambda_n}(x; \theta; X_{\lambda_n+1}, \dots, X_{\mu_n})$$

and

$$\tilde{\psi}_{n2}(x; \theta) \equiv \tilde{\psi}_{n - v_n}(x; \theta; X_{v_n+1}, \dots, X_n).$$

Then, as we show below, an estimator satisfying (21) is defined by

$$(22) \quad \hat{\theta}_n \equiv \frac{\mu_n}{n} \left\{ \tilde{\theta}_{n2} + \frac{1}{\mu_n} \sum_{i=1}^{\mu_n} \tilde{\psi}_{n2}(X_i; \tilde{\theta}_{n2}) \right\} \\ + \frac{n - \mu_n}{n} \left\{ \tilde{\theta}_{n1} + \frac{1}{n - \mu_n} \sum_{i=\mu_n+1}^n \tilde{\psi}_{n1}(X_i; \tilde{\theta}_{n1}) \right\}.$$

The sufficient conditions (iv) through (vi) of theorem 1 are also necessary in a strong sense. If there exists an estimator, which is locally asymptotically linear in  $\psi$  on  $\mathbf{P}_1$ , and hence satisfies (21), and if the contiguity condition (iii) holds, then subtracting the left side of (21) from the left side of (21) with  $\theta_n = \theta$  yields (iv). Furthermore, (i), (ii), and (21) imply (v) by the central limit theorem. Finally, the following theorem of Klaassen (1987) says that if (21) holds and  $\psi$  satisfies a uniform integrability condition, then it is possible to construct estimators  $\tilde{\psi}_n(\cdot; \cdot; \underline{X})$  of the influence function  $\psi(\cdot; \cdot; G)$  satisfying (19) and (20). While this theorem is not practically useful (because the hypothesis that (21) holds implies that we already have an estimator  $\hat{\theta}_n$  in hand), it gives theoretical reassurance that we are not asking for too much in the hypotheses of theorem 1.

**Theorem 2.** Suppose that (i) and (ii) hold, that there exists an estimator  $\hat{\theta}_n$  satisfying (21), and that  $\psi_n(\cdot) \equiv \psi(\cdot; \theta_n; G)$  are  $P_n \equiv P_{(\theta_n, G)}$ -uniformly

square integrable for all  $(\theta, G) \in \Theta \times G$  and all sequences  $\{\theta_n\}$  with  $\sqrt{n} |\theta_n - \theta| = O(1)$ , that is,

$$(23) \quad \limsup_{n \rightarrow \infty} E_{P_n} |\psi_n|^2 1_{[|\psi_n| \geq \lambda]} \rightarrow 0 \quad \text{as } \lambda \rightarrow \infty.$$

Then there is a sequence of estimators  $\tilde{\psi}_n(\cdot; \cdot; \underline{X}) = \tilde{\psi}_n(\cdot; \cdot; X_1, \dots, X_n)$  of the influence function  $\psi(\cdot; \cdot; G)$  which satisfies both (19), and (20).

**Proof of theorem 1.** In view of conditions (iii) and (iv), it suffices to show that  $\{\hat{\theta}_n\}$  defined by (22) satisfies (21) with  $\theta_n = \theta$ . Therefore we will prove that

$$(a) \quad \frac{n - \mu_n}{\sqrt{n}} \left\{ \tilde{\theta}_{n1} - \theta + \frac{1}{n - \mu_n} \sum_{i=\mu_n+1}^n [\tilde{\psi}_{n1}(X_i; \tilde{\theta}_{n1}) - \psi(X_i; \theta; G)] \right\} = o_{P_{(\theta, G)}}(1).$$

Now by (v) and the independence of  $\tilde{\theta}_{n1}$  (which is based on the first block of  $X_i$ 's,  $X_1, \dots, X_{\lambda_n}$ ) and  $(X_{\lambda_n+1}, \dots, X_n)$ , (a) holds if for all  $\{\theta_n\}$  with  $\sqrt{n} |\theta_n - \theta| = O(1)$ , we have

$$(b) \quad \sqrt{n} \left\{ \theta_n - \theta + \frac{1}{n - \mu_n} \sum_{i=\mu_n+1}^n [\tilde{\psi}_{n1}(X_i; \theta_n) - \psi(X_i; \theta; G)] \right\} = o_{P_{(\theta, G)}}(1).$$

Now note that for i.i.d. random  $k$ -vectors  $Y_1, \dots, Y_m$  we have

$$(c) \quad E \left| \frac{1}{\sqrt{m}} \sum_{i=1}^m Y_i \right|^2 = E |Y_1|^2 + (m - 1) |E(Y_1)|^2.$$

Further, (i), (19), and (20) imply, by using (c) conditionally on  $X_{\lambda_n+1}, \dots, X_{\mu_n}$ , that

$$(d) \quad E_n \left\{ \left| \frac{1}{\sqrt{n - \mu_n}} \sum_{i=\mu_n+1}^n [\tilde{\psi}_{n1}(X_i; \theta_n) - \psi_n(X_i)] \right|^2 \mid X_{\lambda_n+1}, \dots, X_{\mu_n} \right\} = o_{P_{(\theta_n, G)}}(1),$$

where  $E_n$  is computed under  $P_n$ . Hence

$$(e) \quad \frac{1}{\sqrt{n - \mu_n}} \sum_{i=\mu_n+1}^n [\tilde{\psi}_{n1}(X_i; \theta_n) - \psi_n(X_i)] = o_{P_{(\theta_n, G)}}(1).$$

By the contiguity assumption (iii), this implies that (e) holds with  $o_{P_{(\theta_n, G)}}(1)$  replaced by  $o_{P_{(\theta, G)}}(1)$ . Combining (b) and (e) shows that it suffices to have that

$$(f) \quad \sqrt{n} \left\{ \theta_n - \theta + \frac{1}{n - \mu_n} \sum_{i=\mu_n+1}^n [\psi_n(X_i) - \psi(X_i; \theta; G)] \right\} = o_{P_{(\theta, G)}}(1).$$

But (f) follows from assumption (iv). □

Here is a more symmetric variant of the data splitting scheme which led to the estimator (22). Suppose that  $\{\lambda_n\}$  and  $\{\mu_n\}$  are sequences of positive integers with  $0 < \lambda_n < \mu_n < n$  and

$$\frac{\lambda_n}{n} \rightarrow \lambda, \quad \frac{\eta_n}{n} \rightarrow \eta,$$

where  $0 < \lambda < \mu < 1$ . Thus  $\lambda_n$  and  $\mu_n$  split the data into three independent blocks. Define preliminary estimators

$$\begin{aligned} \tilde{\theta}_{n3} &\equiv \tilde{t}_{\lambda_n}(X_1, \dots, X_{\lambda_n}), \\ \tilde{\theta}_{n2} &\equiv \tilde{t}_{\mu_n - \lambda_n}(X_{\lambda_n+1}, \dots, X_{\mu_n}), \\ \tilde{\theta}_{n1} &\equiv \tilde{t}_{n - \mu_n}(X_{\mu_n+1}, \dots, X_n), \end{aligned}$$

where, as before,  $\tilde{\theta}_n \equiv \tilde{t}_n(X_1, \dots, X_n)$  for all  $n \geq 1$ . Then define independent influence function estimators by

$$\begin{aligned} \tilde{\psi}_{n1}(x, \theta) &\equiv \tilde{\psi}_{\lambda_n}(x; \theta; X_1, \dots, X_{\lambda_n}), \\ \tilde{\psi}_{n2}(x, \theta) &\equiv \tilde{\psi}_{\mu_n - \lambda_n}(x; \theta; X_{\lambda_n+1}, \dots, X_{\mu_n}), \\ \tilde{\psi}_{n3}(x, \theta) &\equiv \tilde{\psi}_{n - \mu_n}(x; \theta; X_{\mu_n+1}, \dots, X_n). \end{aligned}$$

Then under the hypothesis of theorem 1, another estimator satisfying (21) is given by

$$\begin{aligned} (24) \quad \hat{\theta}_n &\equiv \frac{\lambda_n}{n} \left\{ \tilde{\theta}_{n2} + \frac{1}{\lambda_n} \sum_{i=1}^{\lambda_n} \tilde{\psi}_{n3}(X_i; \tilde{\theta}_{n2}) \right\} \\ &\quad + \frac{\mu_n - \lambda_n}{n} \left\{ \tilde{\theta}_{n3} + \frac{1}{\mu_n - \lambda_n} \sum_{i=\lambda_n+1}^{\mu_n} \tilde{\psi}_{n1}(X_i; \tilde{\theta}_{n3}) \right\} \\ &\quad + \frac{n - \mu_n}{n} \left\{ \tilde{\theta}_{n1} + \frac{1}{n - \mu_n} \sum_{i=\mu_n+1}^n \tilde{\psi}_{n2}(X_i; \tilde{\theta}_{n1}) \right\}. \end{aligned}$$

As a corollary to theorem 1 we obtain a generalization of theorem 2.5.2 to semiparametric models.

**Corollary 1.** Let  $P$  be a semiparametric model as in (14) with  $P_1$  regular,  $\dot{P} = \dot{P}_1 + \dot{P}_2$ , and  $v = \theta$  pathwise differentiable on  $P$ . Then the efficient influence function  $\tilde{l}(\cdot; \theta; G)$  is well defined by (15) and (16). If (iv), (v), and (vi) hold, with  $\psi = \dot{l}$ , then there exists an estimator (defined by (22) for example) which is efficient in the sense of definition 3.3.3.

**Proof.** In view of corollary 3.4.1,  $\tilde{l}(\cdot; \theta; G)$  is well defined, and theorem 1 and proposition 3.3.1 yield the result.  $\square$

Corollary 1 can be used to construct efficient estimates in a wide range of examples. In this section we begin by considering estimation of location in detail. We then sketch the application of this approach to other models.

**Example 1. Symmetric location.**

This is the model of examples 3.2.4 and 3.4.1. The efficient influence function for estimation of  $\theta$  is

$$(25) \quad \tilde{l}(x, \theta) = - \frac{1}{I(G)} \frac{g'}{g}(x - \theta) \quad \text{for } x \in R$$

where

$$(26) \quad I(G) \equiv \int \frac{[g'(x)]^2}{g(x)} dx$$

is the Fisher information for location in  $g$ . Since  $g$  is symmetric (even), the score function for location,  $-g'/g$ , is antisymmetric (odd). □

**Example 2. Regression.**

This is the model of example 4.2.2. The efficient influence function for estimation of the regression parameter  $v \in R^{k-2}$  is

$$(27) \quad \tilde{I}(x, \theta) = \frac{\sigma^2}{I_q} \{E(Z - EZ)(Z - EZ)^T\}^{-1} \frac{z - E(Z)}{\sigma} \left(-\frac{q'}{q}(u)\right)$$

where

$$(28) \quad u \equiv \frac{1}{\sigma} (y - \Delta - \sum_{j=1}^{k-2} v_j z_j), \quad \theta = (v, \Delta, \sigma),$$

and  $I_q$  is the Fisher information for location in  $q$  as in (26). □

**Example 3. Elliptic model.**

This is the model of example 4.2.3. The efficient score function for estimation of  $\theta \equiv (\Delta, \Sigma / \text{trace}(\Sigma))$  is

$$(29) \quad \Gamma^*(x, \Delta, \Sigma) = \dot{\omega}[\dot{I}(u)]$$

where

$$(30) \quad u \equiv S^{-1}(x - \Delta), \quad S = \Sigma^{1/2},$$

$$(31) \quad \dot{\omega}[d\theta] = (S^{-1}[d\Delta], S^{-1}[d\Sigma](S^{-1})^T),$$

and  $\dot{I}(u)$  is given by

$$(32) \quad \dot{I}_i(u) = -2u_i \frac{\tilde{g}'}{\tilde{g}}(|u|^2),$$

$$(33) \quad \dot{I}_{ij}(u) = -(2 - \delta_{ij}) \{u_i u_j \frac{\tilde{g}'}{\tilde{g}}(|u|^2) + \frac{1}{2} \delta_{ij}\}.$$

□

**Example 4. Paired exponential mixture model.**

This is example 4.5.4. Here the efficient influence function is

$$(34) \quad \tilde{I}(x, \theta) = \left\{ t \frac{p_T'}{p_T}(t) - 1 \right\} \frac{1}{t} \frac{x_1 - \theta x_2}{2\theta} \frac{3\theta^2}{1 + \frac{1}{4} I_{\text{scale}}(p_T)},$$

where

$$t = -(x_1 + \theta x_2), \quad T = -(X^{(1)} + \theta X^{(2)}),$$

and  $p_T$  is the density of  $T$  given by

$$p_T(t) = -t \int_0^\infty \eta^2 \exp(\eta t) dG(\eta) \quad \text{for } t < 0. \quad \square$$

A common feature of all of these four examples is that estimation of the efficient influence function  $\tilde{I}$  of the parametric part of the model essentially reduces to estimation of the score for location  $-g'/g$  for some density  $g$ :  $g$  itself in example 1,  $q$  in example 2,  $\tilde{g}$  in example 3, and  $p_T$  in example 4. The following proposition shows that this is possible in considerable generality.

**Proposition 1.** Suppose that  $X_1, X_2, \dots$  are i.i.d. real-valued random variables with absolutely continuous density  $g(x) \equiv \int_{-\infty}^x g'(t) dt$  and df  $G(x) = \int_{-\infty}^x g(t) dt$  satisfying

$$(35) \quad I_k \equiv \int x^{2k} \frac{[g'(x)]^2}{g(x)} dx < \infty$$

for some integer  $k$ . Let

$$(36) \quad h_k(x) \equiv x^k \frac{g'(x)}{g(x)}.$$

Then there are estimators  $\hat{h}_k \equiv \hat{h}_k(\cdot; \underline{X}) \equiv \hat{h}_k(\cdot; X_1, \dots, X_n)$  and  $\hat{g} \equiv \hat{g}(\cdot; \underline{X})$  of  $h_k$  and  $g$  satisfying

$$(37) \quad \int |\hat{h}_k(x) - h_k(x)|^2 g(x) dx = o_p(1)$$

and

$$(38) \quad \hat{I}_k \equiv \int \hat{h}_k^2(x) \hat{g}(x) dx \rightarrow_p I_k.$$

Note that, for  $\hat{h}_k(x) = x^k h_0(x)$ ,

$$\int |\hat{h}_k(x) - h_k(x)|^2 g(x) dx = \int |\hat{h}_0(x) - h_0(x)|^2 x^{2k} g(x) dx$$

so that proving (37) for  $k \geq 1$  is equivalent to proving (37) with  $k = 0$  and a weighted metric.

The proof of proposition 1 will be deferred to the end of this section. We now show how proposition 1 can be used in example 1 to verify the hypothesis (vi) of theorem 1. We also exhibit in example 1 an estimate  $\hat{\theta}_n$  satisfying (v). Hence we can obtain asymptotically efficient estimates via (22) or (24) with  $\tilde{\psi}_n = \tilde{I}_n$ , as shown in corollary 1. We also sketch the extension and construction of appropriate  $\tilde{\theta}_n$  in examples 2–4. Proposition 1 is due to Bickel (1982).

**Example 1. Symmetric location, continued.**

We verify the hypotheses of corollary 1.  $\mathbf{P}_1$  is regular by corollary A.5.1,  $\dot{\mathbf{P}} = \dot{\mathbf{P}}_1 + \dot{\mathbf{P}}_2$  by (3.2.20) of example 3.2.4 and  $v = \theta$  is pathwise differentiable on  $\mathbf{P}$  with  $\dot{v} = \tilde{I}$ , given by (25). The smoothness hypothesis (iv) with  $\psi = \tilde{I}$  follows from (2.1.15) of proposition 2.1.2 in this example. There are lots of preliminary estimators satisfying (v). One of the simplest examples is the Hodges-Lehmann estimator of example 7.3.2, which is asymptotically linear with influence function

$$\left(\int g^2\right)^{-1} (G(\cdot) - \frac{1}{2}),$$

as shown in example 7.3.2 (continued). Note that  $0 < \int g^2 < I^{1/2}(G) < \infty$ , since by Cauchy-Schwarz

$$g(x) = \int_{-\infty}^x g'(y) dy \leq \int |g'| g^{-1/2} g^{1/2} \leq I^{1/2}(G).$$

To obtain an estimator of  $\tilde{I}$  satisfying (vi) we apply proposition 1.

Define  $Y_i \equiv Y_i(\theta) \equiv X_i - \theta$ , so that  $Y_1, \dots, Y_n$  are i.i.d.  $g$ . By proposition 1 there is an estimator  $\hat{h}_0$  of  $g'/g$  based on the  $Y_i$ 's satisfying (37). Now we use  $\hat{h}_0$  to estimate  $\tilde{I}$ , but in a way that forces (anti) symmetry. Set

$$(39) \quad \tilde{I}_n(x; \theta; \underline{X}) = \frac{1}{2} \{-\hat{h}_0(x - \theta; \underline{Y}(\theta)) + \hat{h}_0(-x + \theta; \underline{Y}(\theta))\} / \hat{I}_n$$

where

$$\hat{I}_n \equiv \int \hat{h}_0^2(y) \hat{g}(y) dy.$$

Note that  $\tilde{I}_n$  is antisymmetric (odd) about  $\theta$ . Since  $g(\cdot - \theta)$  is symmetric (even) about  $\theta$ , it follows immediately that (19) holds exactly:

$$\sqrt{n} \int \tilde{I}_n(x; \theta_n; \underline{X}) g(x - \theta_n) dx = 0.$$

Furthermore, since  $\tilde{I}$  is antisymmetric, and  $g(\cdot - \theta)$  is symmetric,

$$\begin{aligned} & \int |\tilde{I}_n(x; \theta_n; \underline{X}) - \tilde{I}(x; \theta_n; G)|^2 g(x - \theta_n) dx \\ & \leq \frac{3}{4\hat{I}_n^2} \left\{ \int |\hat{h}_0(x - \theta_n) - \frac{g'}{g}(x - \theta_n)|^2 g(x - \theta_n) dx \right. \\ & \quad \left. + \int |\hat{h}_0(-x + \theta) - \frac{g'}{g}(-x + \theta_n)|^2 g(-x + \theta_n) dx \right\} \\ & \quad + 3 \left| \frac{1}{\hat{I}_n} - \frac{1}{I_g} \right|^2 \int \frac{[g']^2}{g}(y) dy \\ & = o_p(1) \quad \text{by proposition 1.} \end{aligned}$$

Thus (20) holds with  $\psi = \tilde{I}$ , and corollary 1 applies. □

*Examples 2-4, continued.*

Preliminary estimators can be constructed in example 2 fairly easily. For instance, the maximum likelihood estimate  $(\hat{\Delta}, \hat{\beta}, \hat{\sigma})$  when  $g$  is the logistic density is known to be unique by a theorem of Scholz (1974). The asymptotic normality of  $n^{1/2}(\hat{\Delta} - \Delta(P), \hat{\beta} - \beta(P), \hat{\sigma} - \sigma(P))$  follows readily from theorem 7.2.1. Suppose the preliminary estimator is discretized to an  $n^{-1/2}$  grid yielding  $(\tilde{\Delta}, \tilde{v}, \tilde{\sigma})$ . It is fairly easy to show by a contiguity argument that if  $\hat{h}_0$  is com-

puted by replacing  $X_1, \dots, X_n$  by the residuals  $Y_i - \tilde{\Delta} - \sum_{j=1}^{k-2} \tilde{v}_j Z_{ij}$  and  $\hat{I}$  is defined by substituting this  $\hat{h}_0, \tilde{\sigma}, \bar{Z}$  for  $h_0, \sigma, E(Z)$  into  $\tilde{I}$ , then the conditions of corollary 1 are satisfied. In example 3, we can construct preliminary  $M$ -estimates  $\hat{S}$  and  $\hat{\Delta}$  such that  $n^{1/2}(\hat{S} - S(P), \hat{\Delta} - \Delta(P))$  is asymptotically normal, and if  $P_{(\Delta, \Sigma)} \in \mathcal{P}$ , then  $S^2(P) = \lambda(P)\Sigma$  for  $\lambda$  a scalar,  $\Delta(P) = \Delta$ . This can be done using results of Maronna (1976). Discretize to  $\tilde{S}, \tilde{\Delta}$  and form  $\hat{R}_i^2 = (X_i - \tilde{\Delta})' \tilde{S}^{-2} (X_i - \tilde{\Delta}), i = 1, \dots, n$ . Then compute  $\hat{h}_0$  using the  $\hat{R}_i^2$ , and substitute in  $\tilde{I}, \tilde{S}$  for  $S, \tilde{\Delta}$  for  $\Delta$ , and  $\hat{h}_0(u) - (d/2 - 1)u^{-1}$  for  $(\tilde{g}'/\tilde{g})(u)$ . Again, the conditions of corollary 1 are satisfied. Finally, in example 4 we can use the simple preliminary estimate  $\hat{\theta} = \{\sum_{i=1}^n X_i^{(1)}\} / \{\sum_{i=1}^n X_i^{(2)}\}$ , discretize to  $\tilde{\theta}$ , form  $\hat{T}_i = -(X_{1i} + \tilde{\theta} X_{2i})$ , use  $\hat{h}_0$  based on the  $\hat{T}_i$  to estimate  $p_T'/p_T$ , etc. Details of these applications to examples 1-3 of proposition 1 may be found in Bickel (1982). □

*Symmetric Constructions*

The estimators  $\hat{\theta}_n$  in (22) or (24), which were shown in theorem 1 to satisfy (21) and in corollary 1 to be efficient if  $\psi = \tilde{I}$ , are not completely satisfactory because of their lack of symmetry in the observations. We would like to show that more natural estimators  $\hat{\theta}_n$  of  $\theta$  using all the data symmetrically to estimate  $\psi$  and  $\tilde{I}$  respectively continue to satisfy (21) under the hypotheses (19) and (20) of theorem 2.

We proceed by assuming existence of an estimator

$$(40) \quad \hat{\psi}_n(\cdot; \cdot; \underline{X}_j) \equiv \tilde{\psi}_{n-1}(\cdot; \cdot; X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$$

of the influence function  $\psi(\cdot; \cdot; G)$  as before but with  $\tilde{\psi}_n(\cdot; \cdot; \underline{X}_j)$  invariant under permutations of  $\underline{X}_j$ . Our estimator of  $\theta$  will be

$$(41) \quad \hat{\theta}_n = \tilde{\theta}_n + \frac{1}{n} \sum_{i=1}^n \hat{\psi}_n(X_i; \tilde{\theta}_n; \underline{X}_i),$$

where  $\tilde{\theta}_n$  is a  $\sqrt{n}$ -consistent preliminary estimator. Note that  $\hat{\theta}_n$  is invariant under permutations of  $X_1, \dots, X_n$  if  $\tilde{\theta}_n$  is, and that the influence function estimators involve all the observations but the one for which the influence is scored. We want to show that (21) continues to hold, so we rewrite

$$\begin{aligned} & \sqrt{n} \{ \hat{\theta}_n - \theta_n - \frac{1}{n} \sum_{i=1}^n \psi(X_i; \theta_n; G) \} \\ &= \sqrt{n} \{ \tilde{\theta}_n - \theta_n - \frac{1}{n} \sum_{i=1}^n [\psi(X_i; \tilde{\theta}_n; G) - \psi(X_i; \theta_n; G)] \} \\ & \quad + \frac{1}{\sqrt{n}} \sum_{i=1}^n [\hat{\psi}_n(X_i; \tilde{\theta}_n; \underline{X}_i) - \psi(X_i; \tilde{\theta}_n; G)]. \end{aligned}$$

By (iv) and the technique of the proof of theorem 2.5.2 the first term will be  $o_p(1)$  if  $\tilde{\theta}_n$  is a discretized  $\sqrt{n}$ -consistent estimator. Then we want to show that



the second term is  $o_p(1)$ : again, if  $\tilde{\theta}_n$  is discretized this will hold if, for all  $\theta_n$  with  $\sqrt{n} |\theta_n - \theta| = O(1)$ , we have

$$(42) \quad \frac{1}{\sqrt{n}} \sum_{i=1}^n [\hat{\psi}_n(X_i; \theta_n; \underline{X}_i) - \psi(X_i; \theta_n; G)] = o_p(1).$$

This can be proved using the approach of Schick (1987), which we outline below. This leads to extensions of theorem 1 and corollary 1. We formulate the generalization of theorem 1, using the notation  $P_n$  and  $E_n$  introduced in theorem 2 and the proof of theorem 1.

**Theorem 3.** Suppose that (i) through (v) hold and that there exists an estimator  $\hat{\psi}_n(\cdot; \cdot; \underline{X}_n)$  of  $\psi(\cdot; \cdot; G)$  which is symmetric in the observations, is defined as in (40), and satisfies

$$(43) \quad \sqrt{n} \int \hat{\psi}_n(x; \theta_n; \underline{X}_n) dP_n(x) = o_{P_{(\theta_n, G)}}(1),$$

$$(44) \quad E_n \int |\hat{\psi}_n(x; \theta_n; \underline{X}_n) - \psi(x; \theta_n; G)|^2 dP_n(x) = o(1),$$

and

$$(45) \quad n E_n \int |\hat{\psi}_n(x; \theta_n; \underline{X}_n) - \hat{\psi}_n(x; \theta_n; \underline{X}_1)|^2 dP_n(x) = o(1)$$

for all  $(\theta, G) \in \Theta \times G$  and all sequences  $\{\theta_n\}$  with  $\sqrt{n} |\theta_n - \theta| = O(1)$ . Then  $\hat{\theta}_n$  defined by (41) with  $\tilde{\theta}_n$  discretized satisfies (21).

The proof of theorem 3 is deferred to the end of this section. Note that, although  $\hat{\theta}_n$  can be chosen to be symmetric in the observations, it still has the drawback of being based on a discretized preliminary estimator. Furthermore, note that (43) and (44) imply the conditions (19) and (20) of theorem 2. The hypothesis (45) together with the expectation in (44) represent the price which is paid by using simultaneously all the observations but one, to construct the influence function estimators. Equation (45) ensures that no single observation contributes excessively to  $\hat{\psi}_n$ .

Despite the theoretical backing for these estimates, practical implementation is often difficult. For a study of their implementation in the context of example 4.2.2, see Hsieh and Manski (1987). Successful implementations for examples 1 and 2 are given by Faraway (1992) and Jin (1992). For other types of examples, such as the transformation models of section 4.7, several studies suggest that other estimates, such as those based on direct maximization of a likelihood or partial likelihood, actually perform much better for finite sample sizes even though they are more difficult to study theoretically.

**Example 1. Symmetric location, continued.**

In this example, efficient estimators which are permutation invariant in the observations, have been given by Stone (1975) and Beran (1978). To construct such an estimator via theorem 3 we need a slight addition to proposition 1.

**Proposition 2.** Under the assumptions of proposition 1 with  $k = 0$ , there exists an estimator  $\hat{h}(\cdot; \underline{X}_n)$  of  $h \equiv h_0$  satisfying

$$(46) \quad E \int |\hat{h}(x; \underline{X}_n) - h(x)|^2 g(x) dx = o(1)$$

and

$$(47) \quad n E \int |\hat{h}(x; \underline{X}_n) - \hat{h}(x; \underline{X}_1)|^2 g(x) dx = o(1).$$

Similarly as in the first continuation of example 1 we define, applying the same conventions to  $\underline{Y}(\theta) \equiv (Y_1(\theta), \dots, Y_n(\theta))$  as we did to  $\underline{X}(\theta)$ ,

$$(48) \quad \hat{\psi}_n(x; \theta; \underline{X}_n) = \frac{I_0}{2} \{-\hat{h}(x - \theta; \underline{Y}_n(\theta)) + \hat{h}(-x + \theta; \underline{Y}_n(\theta))\}$$

as an "estimator" of  $\psi(x; \theta; G) = -h(x - \theta)/I_0 = \tilde{l}(x; \theta; G)$ . By (anti) symmetry we have

$$\sqrt{n} \int \hat{\psi}_n(x; \theta_n; \underline{X}_n) g(x - \theta_n) dx = 0$$

and, by (46),

$$\begin{aligned} E_n \int |\hat{\psi}_n(x; \theta_n; \underline{X}_n) - \psi(x; \theta_n; G)|^2 g(x - \theta_n) dx \\ \leq I_0^{-2} E_n \int |\hat{h}(x - \theta_n; \underline{Y}_n(\theta_n)) - h(x - \theta_n)|^2 g(x - \theta_n) dx \\ = o(1). \end{aligned}$$

Furthermore, (45) follows from (47). Consequently, applying theorem 3 we obtain, for discretized  $\hat{\theta}_n$ ,

$$(49) \quad \sqrt{n} \{\tilde{\theta}_n + \frac{1}{n} \sum_{i=1}^n \hat{\psi}_n(X_i; \tilde{\theta}_n; \underline{X}_i) - \theta_n - \frac{1}{n} \sum_{i=1}^n \tilde{l}(X_i; \tilde{\theta}_n; G)\} = o_{P_{(\theta_n, G)}}(1).$$

For the efficiency of

$$(50) \quad \hat{\theta}_n = \tilde{\theta}_n + \frac{1}{n} \sum_{i=1}^n \hat{l}_n(X_i; \tilde{\theta}_n; \underline{X}_i)$$

with

$$\hat{l}_n(x; \theta; \underline{X}_i) = \frac{I_0 \hat{\psi}_n(x; \theta; \underline{X}_i)}{\hat{I}_n}, \quad \hat{I}_n \text{ as in (39),}$$

it remains to be shown that

$$(51) \quad \frac{1}{\sqrt{n}} \sum_{i=1}^n [\hat{\psi}_n(X_i; \tilde{\theta}_n; \underline{X}_i) - \hat{l}_n(X_i; \tilde{\theta}_n; \underline{X}_i)] = o_{P_{(\theta_n, G)}}(1).$$

But the central limit theorem and (49) show

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\psi}_n(X_i; \tilde{\theta}_n; \underline{X}_i) = o_{P_{(\theta_n, G)}}(1),$$

which together with (38) implies (51).  $\square$

It seems plausible that this type of construction is also applicable in examples 2-4 modulo some rather nasty technicalities. This technique has been applied to two-sample location-scale problems by Park (1990).

*General Parameters*

There seem to be no simple generalizations of theorems 1 and 3 available for general parameters. To illustrate the difficulty suppose  $\mathbf{P} = \{P_{(G,H)} \mid G \in \mathbf{G}, H \in \mathbf{H}\}$  where  $\mathbf{G} \subset l^2$ , the Hilbert space of square summable sequences. For instance, if  $G$  is a distribution on the circle, identify  $G$  with its Fourier coefficients. Suppose we have available a preliminary weakly regular estimate  $\tilde{G} = (\tilde{g}_1, \tilde{g}_2, \dots) \in l^2$  of  $G$  and  $\tilde{H}$  of  $H$  and for each  $j$  an estimate  $\hat{I}_j(\cdot; g_j; \tilde{g}_i, i \neq j, \tilde{H})$  of  $I_j(\cdot; g_j; g_i, i \neq j, H)$ , the influence function of the parameter  $g_j$  which satisfies (19) and (20). Then, we can use theorem 1 to construct an efficient estimate  $\hat{g}_j$  of  $g_j$  by our methods using  $\tilde{g}_j$  as a preliminary estimate and  $\hat{I}_j$  as above. But to complete the argument we now need to put the estimates together into  $(\hat{g}_1, \hat{g}_2, \dots)$  as an estimate of  $G$ . However, there is no guarantee that  $(\hat{g}_1, \hat{g}_2, \dots) \in l^2$ ! If we suppose not only  $\tilde{g} \in l^2$ , but also that the  $\hat{I}_j$  are compatible, that is,  $\hat{I}(\cdot; \tilde{g}_i, i \geq 1, \tilde{H}) \in l^2$ , and let  $\hat{I}_j(\cdot; \tilde{g}_i; \tilde{g}_i, i \neq j, \tilde{H})$  be the  $j$ th component of  $\hat{I}(\cdot; \tilde{g}_i, i \geq 1, \tilde{H})$ , this objection is overcome. We can conclude that  $(\hat{g}_1, \hat{g}_2, \dots)$  are weakly efficient with respect to  $\mathbf{B}_0^* = [\text{coordinate maps}]$ . However, efficiency in the strong sense as defined in definition 5.2.7 requires tightness of  $n^{-1/2} \sum_{i=1}^n \hat{I}(X_i; \tilde{g}_j, j \geq 1, \tilde{H})$  as Banach valued random elements. Evidently, this can be checked in special cases, but hardly qualifies as a simple condition.

*Proofs of Propositions 1 and 2, and Theorem 3*

**Proof of proposition 1.** We give the proof only for the case  $k = 0$ . The proof for general  $k$  is similar. Let  $w(x)$  be the logistic density function

$$w(x) \equiv \frac{e^{-x}}{[1 + e^{-x}]^2},$$

and let  $b_n$  be a sequence of positive numbers with  $b_n \rightarrow 0$ . Set

$$(a) \quad \hat{g}(x; \underline{X}) \equiv \hat{g}(x) \equiv \frac{1}{nb_n} \sum_{i=1}^n w\left(\frac{x - X_i}{b_n}\right)$$

so that

$$(b) \quad \hat{g}'(x; \underline{X}) = \hat{g}'(x) = \frac{1}{nb_n^2} \sum_{i=1}^n w'\left(\frac{x - X_i}{b_n}\right).$$

Define

$$(c) \quad \hat{h}_0(x; \underline{X}) \equiv \hat{h}_0(x) \equiv \begin{cases} \hat{g}'(x) & \text{if } |x| \leq d_n, |\hat{g}'(x)| \leq c_n \hat{g}(x), \\ \hat{g} & \\ 0 & \text{otherwise,} \end{cases}$$

where  $c_n \rightarrow \infty, d_n \rightarrow \infty, b_n \rightarrow 0$ ,

$$(d) \quad b_n c_n^2 \rightarrow 0 \quad (\text{so } b_n c_n \rightarrow 0 \text{ also})$$

and

$$(e) \quad d_n b_n^{-3} n^{-1} \rightarrow 0,$$

e.g.,  $c_n = b_n^{-1/3}$ ,  $d_n = b_n^{-1}$  and  $n b_n^4 \rightarrow \infty$ . Also, set

$$\begin{aligned} g_n(x) &\equiv E \hat{g}(x; \underline{X}) \\ &= \frac{1}{b_n} \int w\left(\frac{x-y}{b_n}\right) g(y) dy \end{aligned}$$

$$(f) \quad \equiv \int w(z) g(x - b_n z) dz$$

and

$$\begin{aligned} (g) \quad g'_n(x) &\equiv E \hat{g}'(x; \underline{X}) = \frac{1}{b_n^2} \int w'\left(\frac{x-y}{b_n}\right) g(y) dy \\ &= \frac{1}{b_n} \int w'(z) g(x - b_n z) dz \\ &= \int w(z) g'(x - b_n z) dz. \end{aligned}$$

Now the proof proceeds by the decomposition

$$\begin{aligned} &\int |\hat{h}_0(y) - \frac{g'}{g}(y)|^2 g(y) dy \\ &\leq 3 \left\{ \int |\hat{h}_0(y) - \hat{h}_0\left(\frac{g_n}{g}\right)^{1/2}(y)|^2 g(y) dy \right. \\ &\quad + \int \left| \hat{h}_0\left(\frac{g_n}{g}\right)^{1/2}(y) - \left(\frac{g_n'}{g_n}\right)\left(\frac{g_n}{g}\right)^{1/2}(y) \right|^2 g(y) dy \\ &\quad \left. + \int \left| \frac{g_n'}{g_n}\left(\frac{g_n}{g}\right)^{1/2}(y) - \frac{g'}{g}(y) \right|^2 g(y) dy \right\} \\ &= 3 \left\{ \int \hat{h}_0^2(y) |\sqrt{g_n(y)} - \sqrt{g(y)}|^2 dy \right. \\ &\quad + \int \left| \hat{h}_0(y) - \frac{g_n'}{g_n}(y) \right|^2 g_n(y) dy \\ &\quad \left. + \int \left| \frac{g_n'}{\sqrt{g_n}}(y) - \frac{g'}{\sqrt{g}}(y) \right|^2 dy \right\} \end{aligned}$$

$$(h) \quad \equiv 3(I + II + III).$$

Note that this decomposition avoids integrals like  $\int (g_n'/g_n)^2 g$  which may not be bounded in  $n$ . The three terms in (h) are handled by the following three lemmas.

**Lemma 1.** If  $b_n c_n \rightarrow 0$ , then

$$I \equiv \int \hat{h}_0^2(y) [\sqrt{g_n(y)} - \sqrt{g(y)}]^2 dy \rightarrow_p 0.$$

**Lemma 2.** If the conditions of the proposition hold, and  $\hat{h}_0(y)$  is defined by (c)–(e), then

$$II \equiv \int \left| \hat{h}_0(y) - \frac{g'_n(y)}{g_n} \right|^2 g_n(y) dy \rightarrow_p 0.$$

**Lemma 3.** If  $b_n \rightarrow 0$ , then

$$III \equiv \int \left| \frac{g'_n(y)}{\sqrt{g_n}} - \frac{g'}{\sqrt{g}} \right|^2 dy \rightarrow 0.$$

**Proof of lemma 3.** First, by a well-known inequality for Fisher information (see, e.g., Hájek and Šidák (1967, page 17)), it follows from (f) that

$$(i) \quad I(G_n) \equiv \int \frac{[g'_n]^2}{g_n}(y) dy \leq \int \frac{[g']^2}{g}(y) dy \equiv I(G) < \infty$$

so that

$$(j) \quad \limsup_{n \rightarrow \infty} \int \frac{[g'_n]^2}{g_n}(y) dy \leq \int \frac{[g']^2}{g}(y) dy = I(G) < \infty.$$

To apply Vitali's theorem, lemma A.7.5, to III, it remains only to show that

$$(k) \quad \frac{g'_n(y)}{\sqrt{g_n(y)}} \rightarrow_\lambda \frac{g'(y)}{\sqrt{g(y)}},$$

where  $\lambda$  is Lebesgue measure and  $\rightarrow_\lambda$  denotes convergence in  $\lambda$ -measure. To show (k) it suffices to show that

$$(l) \quad g_n(y) \rightarrow g(y) \quad \text{for all } y$$

and

$$(m) \quad g'_n(y) \rightarrow_\lambda g'(y).$$

To prove (m), note that from (g)

$$\begin{aligned} & \int_{-\infty}^{\infty} |g'_n(y) - g'(y)| dy \\ &= \int_{-\infty}^{\infty} \left| \int_{-\infty}^{\infty} (g'(y - b_n z) - g'(y)) w(z) dz \right| dy \\ &\leq \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |g'(y - b_n z) - g'(y)| dy w(z) dz \\ &\rightarrow 0. \end{aligned}$$

This follows by the  $L_1$ -continuity theorem A.1.1 applied to the inner integral and the dominated convergence theorem applied to the outer integral since  $I(G) < \infty$  implies  $\int |g'(y)| dy < \infty$ . Note that

$$g(x) = \int_{-\infty}^x g' \leq I^{1/2}(G).$$

Since  $\rightarrow_1$  implies  $\rightarrow_\lambda$ , (m) follows directly from (f) and continuity and boundedness of  $g$ . Thus (k) holds, and the proof of lemma 3 is complete.  $\square$

**Proof of lemma 1.** Let  $g_n(x, b)$  be defined by (f) with  $b_n = b$ . By Cauchy's form of Taylor's theorem and the properties of  $w$ ,

$$\begin{aligned} \sqrt{g_n(y)} - \sqrt{g(y)} &= \frac{1}{2} \int_0^1 \left\{ \frac{\partial}{\partial \lambda} g_n(y; \lambda b_n) \right\} g_n^{-1/2}(y; \lambda b_n) d\lambda \\ &= -\frac{1}{2} b_n \int_0^1 g_n^{-1/2}(y; \lambda b_n) \int_{-\infty}^{\infty} z g'(y - \lambda b_n z) w(z) dz d\lambda. \end{aligned}$$

From the inequality

$$(n) \quad (EU)^2 = (E(UV^{-1/2} V^{1/2}))^2 \leq E \frac{U^2}{V} EV,$$

it follows that

$$\begin{aligned} &[\sqrt{g_n(y)} - \sqrt{g(y)}]^2 \\ &\leq \frac{1}{4} b_n^2 \int_0^1 g_n^{-1}(y; \lambda b_n) \left\{ \int_{-\infty}^{\infty} z g'(y - \lambda b_n z) w(z) dz \right\}^2 d\lambda \\ &= \frac{1}{4} b_n^2 \int_0^1 \frac{\left\{ \int_{-\infty}^{\infty} z g'(y - \lambda b_n z) w(z) dz \right\}^2}{\int_{-\infty}^{\infty} g(y - \lambda b_n z) w(z) dz} d\lambda \\ (o) \quad &\leq \frac{1}{4} b_n^2 \int_0^1 \int_{-\infty}^{\infty} \frac{z^2 [g'(y - \lambda b_n z)]^2}{g(y - \lambda b_n z)} w(z) dz d\lambda. \end{aligned}$$

Substitution of (o) into the definition of  $I$  and using  $|\hat{h}_0| \leq c_n$  yields

$$\begin{aligned} I &\leq \frac{1}{4} c_n^2 b_n^2 \int_0^1 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{z^2 [g'(y - \lambda b_n z)]^2}{g(y - \lambda b_n z)} w(z) dz dy d\lambda \\ &= \frac{1}{4} c_n^2 b_n^2 \int_{-\infty}^{\infty} \frac{[g']^2}{g}(t) dt \int_{-\infty}^{\infty} z^2 w(z) dz \\ &\rightarrow 0 \end{aligned}$$

since the integrals are both finite and  $b_n c_n \rightarrow 0$ .  $\square$

**Proof of lemma 2.** Since the logistic density satisfies  $|w'(x)| \leq w(x) \leq 1$  for all  $x$ , it follows from (a), (b), and (f) that for  $i = 0, 1$ ,

$$(p) \quad \text{Var} [\hat{g}^{(i)}(x)] \leq n^{-1} b_n^{-(2i+1)} g_n(x) \leq n^{-1} b_n^{-2(i+1)}.$$

Let  $A$  and  $B$  denote the conditions in the definition of  $\hat{h}_0$ . Then the term  $II$  is bounded by  $II_1 + II_2$  where

$$(q) \quad II_1 = \int_{AB} \left\{ \frac{\hat{g}'(y)}{\hat{g}} - \frac{g_n'(y)}{g_n} \right\}^2 g_n(y) dy,$$

$$(r) \quad II_2 \equiv \int_{[AB]^c} \frac{[g_n']^2}{g_n}(y) dy.$$

But by (p) the expectation  $E(II_1)$  of (q) is bounded by

$$\begin{aligned} & E \int_{AB} \left| \frac{\hat{g}'(y)}{\hat{g}} - \frac{g_n'(y)}{g_n} \right|^2 g_n(y) dy \\ & \leq 2E \int_{AB} \left\{ \left| \frac{\hat{g}'(y)}{\hat{g}} \right| \frac{|\hat{g}(y) - g_n(y)|}{g_n(y)} \right. \\ & \quad \left. + \left| \frac{\hat{g}'(y) - g_n'(y)}{g_n(y)} \right|^2 \right\} g_n(y) dy \\ & \leq 4c_n^2 n^{-1} b_n^{-1} d_n + 4n^{-1} b_n^{-3} d_n \\ & = 4\{b_n^2 c_n^2 + 1\} n^{-1} b_n^{-3} d_n \end{aligned}$$

$$(s) \quad = o(1) \quad \text{by (d) and (e).}$$

Finally,

$$(t) \quad E(II_2) \leq \int_{-\infty}^{\infty} \frac{[g_n']^2}{g_n}(y) \{P(|\hat{g}'(y)| > c_n \hat{g}(y)) + 1_{\{|y| > d_n\}}\} dy,$$

where  $[g_n']^2 / g_n$  is uniformly integrable by the proof of lemma 3. Thus it suffices to show that the probability in (t) converges to 0 in (Lebesgue) measure.

But by (l) and (p)

$$(u) \quad \hat{g}(y) \rightarrow_p g(y) \quad \text{if } nb_n \rightarrow \infty, \quad \text{for all } y,$$

and by (m) and (p)

$$(v) \quad \hat{g}'(y) \rightarrow_p g'(y) \quad \text{if } nb_n^3 \rightarrow \infty$$

except for  $y$  in sets of arbitrarily small Lebesgue measure. Thus if  $c_n \rightarrow \infty$  the probability in (t) converges to 0 for all  $y$  with the exception of sets of small Lebesgue measure, and hence  $II_2 \rightarrow_p 0$ .  $\square$

To complete the proof of proposition 1, it remains only to prove the convergence of estimated information claimed in (38).

**Proof of (38).** Again we only give the proof for  $k = 0$ . First write

$$(w) \quad \hat{I}_{0n} - I_j = \int \hat{h}_0^2(x) (\hat{g}(x) - g(x)) dx \\ + \int [\hat{h}_0^2(x) - h_0^2(x)] g(x) dx$$

where the second term is  $o_p(1)$  by (37). To see this use  $(a^2 - b^2) = (a - b)^2 + 2b(a - b)$  to estimate

$$\left| \int [\hat{h}_0^2(x) - h_0^2(x)] g(x) dx \right|$$

$$\begin{aligned}
&\leq \int [\hat{h}_0(x) - h_0(x)]^2 g(x) dx \\
&\quad + 2 \int |h_0(x)| |\hat{h}_0(x) - h_0(x)| g(x) dx \\
&\leq o_p(1) + 2 \left\{ \int h_0^2(x) g(x) dx \int [\hat{h}_0(x) - h_0(x)]^2 g(x) dx \right\}^{1/2} \\
&= o_p(1).
\end{aligned}$$

Thus it remains only to show that the first term in (w) is  $o_p(1)$ . Since  $|\hat{h}_0(x)| \leq c_n$ , the first term in (w) is bounded by

$$\begin{aligned}
&\int \hat{h}_0^2(x) \frac{|\hat{g}(x) - g_n(x)|}{g_n(x)} g_n(x) dx + \int \hat{h}_0^2(x) |g_n(x) - g(x)| dx \\
&\leq \left\{ \int \hat{h}_0^2(x) g_n(x) dx \int \hat{h}_0^2(x) \frac{(\hat{g}(x) - g_n(x))^2}{g_n(x)} dx \right\}^{1/2} \\
&\quad + c_n^2 \int |g_n(x) - g(x)| dx \\
&\equiv I + II \equiv \sqrt{I_1 I_2} + II.
\end{aligned}$$

Now  $I_1 = O_p(1)$  by lemmas 2 and 3 and finiteness of  $I_0$ , while  $E(I_2) \leq 2c_n^2 n^{-1} b_n^{-1} d_n$  by (p). Hence  $I = o_p(1)$ .

To show that  $II \rightarrow 0$ , note that

$$\begin{aligned}
\int |g_n - g| &= \int |g_n^{1/2} + g^{1/2}| |g_n^{1/2} - g^{1/2}| \\
&\leq 2 \left\{ \int [g_n^{1/2} - g^{1/2}]^2 \right\}^{1/2} \\
&\leq b_n \left\{ I(G) \int z^2 w(z) dz \right\}^{1/2}
\end{aligned}$$

by the argument of the proof of lemma 1. Hence

$$II \leq c_n^2 b_n \left\{ I(G) \int z^2 w(z) dz \right\}^{1/2} \rightarrow 0$$

since  $b_n c_n^2 \rightarrow 0$  by (d). □

**Proof of theorem 3.** As noted in the paragraphs above theorem 3, it suffices to prove (42). (In this proof expectations will be under  $P_n = P_{(\theta_n, G)}$ .) We will repeatedly use applications of the Cauchy-Schwarz inequality, as for example in

$$(a) \quad \left| \frac{1}{n} \sum_{i=1}^n a_i \right|^2 \leq \frac{1}{n} \sum_{i=1}^n |a_i|^2, \quad a_i \in R^k.$$

Indeed,

$$E \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \int [\hat{\Psi}_n(x; \theta_n; \underline{X}_i) - \hat{\Psi}_n(x; \theta_n; \underline{X}_n)] dP_n(x) \right|^2$$

$$\begin{aligned}
(b) \quad &\leq \sum_{i=1}^n E \int |\hat{\Psi}_n(x; \theta_n; \underline{X}_i) - \hat{\Psi}_n(x; \theta_n; \underline{X}_n)|^2 dP_n(x) \\
&= o(1),
\end{aligned}$$



by (45) and the permutation invariance. Combining (b) and (43) with (42), we see that it suffices to prove

$$(c) \quad \frac{1}{\sqrt{n}} \sum_{i=1}^n D_i = o_{P_{(\theta_n, G)}}(1)$$

with

$$D_i = \hat{\psi}_n(X_i; \theta_n; \underline{X}_i) - \int \hat{\psi}_n(x; \theta_n; \underline{X}_i) dP_n(x) - \psi(X_i; \theta_n; G).$$

Thus it suffices to show

$$(d) \quad E \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n D_i \right|^2 = \frac{1}{n} \sum_{i=1}^n E |D_i|^2 + \frac{1}{n} \sum_{j \neq k} E(D_j^T D_k) = o(1).$$

But

$$(e) \quad \begin{aligned} \frac{1}{n} \sum_{i=1}^n E |D_i|^2 &= E |D_n|^2 \\ &\leq E \int |\hat{\psi}_n(x; \theta_n; \underline{X}_n) - \psi(x; \theta_n; G)|^2 dP_n(x) \\ &= o(1) \quad \text{by (44)}. \end{aligned}$$

To show that the second term in (d) is  $o(1)$ , we define, as in the proof of lemma 3.1 of Schick (1987),

$$(f) \quad D_{jk} = E(D_j | \underline{X}_k).$$

Note that  $E(D_k | \underline{X}_k) = 0$  and  $E(D_{kj} | \underline{X}_k) = 0$  and that hence  $E(D_{jk}^T D_k) = 0$  and  $E(D_{jk}^T D_{kj}) = 0$ , for  $j \neq k$ . Thus the absolute value of the second term in (d) equals

$$(g) \quad \begin{aligned} &\left| \frac{1}{n} \sum_{j \neq k} E(D_j - D_{jk})^T (D_k - D_{kj}) \right| \\ &\leq (n-1) |E(D_1 - D_{12})^T (D_2 - D_{21})| \\ &\leq (n-1) \{E |D_1 - D_{12}|^2 E |D_2 - D_{21}|^2\}^{1/2} \\ &= (n-1) E |D_1 - D_{12}|^2. \end{aligned}$$

But, if  $Y = \hat{\psi}_n(X_1; \theta_n; \underline{X}_1)$ ,  $E(E(Y | \underline{X}_1) | \underline{X}_2) = E(E(Y | \underline{X}_2) | \underline{X}_1)$  so that

$$\begin{aligned} E |D_1 - D_{12}|^2 &= E |Y - E(Y | \underline{X}_1) - E(Y | \underline{X}_2) + E(E(Y | \underline{X}_1) | \underline{X}_2)|^2 \\ &= E |[Y - E(Y | \underline{X}_2)] - E([Y - E(Y | \underline{X}_2)] | \underline{X}_1)|^2 \\ &\leq E |Y - E(Y | \underline{X}_2)|^2 \end{aligned}$$

$$\begin{aligned}
 \text{(h)} \quad &= E \int [\hat{\psi}_n(X_1; \theta_n; \underline{X}_1) - \hat{\psi}_n(X_1; \theta_n; x_2, X_3, \dots, X_n)] dP_n(x_2) \Big|^2 \\
 &\leq E \int |\hat{\psi}_n(X_1; \theta_n; \underline{X}_1) - \hat{\psi}_n(X_1; \theta_n; x_2, X_3, \dots, X_n)|^2 dP_n(x_2) \\
 &= E \int \int |\hat{\psi}_n(x; \theta_n; \underline{X}_1) - \hat{\psi}_n(x; \theta_n; x_1, X_3, \dots, X_n)|^2 dP_n(x) dP_n(x_1) \\
 &= E \int |\hat{\psi}_n(x; \theta_n; \underline{X}_1) - \hat{\psi}_n(x; \theta_n; \underline{X}_2)|^2 dP_n(x).
 \end{aligned}$$

Thus (g) converges to zero by (45). Together with (e) this yields (d) and the proof of theorem 3 is complete.  $\square$

**Proof of proposition 2.** Let  $\hat{g}(x; \underline{X}_n)$  and  $\hat{g}'(x; \underline{X}_n)$  be as defined in (a) and (b) of the proof of proposition 1 and

$$\text{(a)} \quad \hat{h}(x; \underline{X}_n) = \begin{cases} \frac{\hat{g}'(x; \underline{X}_n)}{\hat{g}(x; \underline{X}_n) + c_n} & \text{if } |x| \leq d_n, \\ 0 & \text{otherwise,} \end{cases}$$

where  $b_n$ ,  $c_n$ , and  $d_n$  are positive and  $b_n \rightarrow 0$ ,  $c_n \rightarrow 0$ , and  $d_n \rightarrow \infty$  satisfy

$$\text{(b)} \quad b_n^{-3} d_n n^{-1} \rightarrow 0$$

and

$$\text{(c)} \quad b_n^4 c_n^2 n \rightarrow \infty,$$

e.g.,  $c_n = b_n$ ,  $d_n = b_n^{-1}$ , and  $b_n^6 n \rightarrow \infty$ . We note that this construction has been used in paragraph 3.7 of Schick (1987). Furthermore,

$$b_n |\hat{g}'(x; \underline{X}_n)| \leq \hat{g}(x; \underline{X}_n)$$

and hence

$$\text{(d)} \quad |\hat{h}(x; \underline{X}_n)| \leq b_n^{-1}.$$

To prove that  $\hat{h}(\cdot; \underline{X}_n)$  satisfies (46), we first adapt the last step in the proof of lemma 1 as follows:

$$\text{(e)} \quad E(I) \leq$$

$$\begin{aligned}
 &\frac{1}{4} b_n \int_0^1 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{z^2 [g'(y - \lambda b_n z)]^2}{g(y - \lambda b_n z)} P(|\hat{h}(y; \underline{X}_n)| \leq b_n^{-1/2}) w(z) dz dy d\lambda \\
 &+ \frac{1}{4} \int_0^1 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{z^2 [g'(y - \lambda b_n z)]^2}{g(y - \lambda b_n z)} P(|\hat{h}(y; \underline{X}_n)| > b_n^{-1/2}) w(z) dz dy d\lambda.
 \end{aligned}$$

As in the argument at the end of the proof of lemma 2 we obtain

$$P(|\hat{h}(y; \underline{X}_n)| > b_n^{-1/2}) = o(1)$$

for all  $y$  with the exception of sets of arbitrarily small Lebesgue measure. Since  $b_n \rightarrow 0$ , (e) yields  $E(I) = o(1)$ . The proof of lemma 2 with  $c_n = b_n^{-1}$  yields  $E(II) = o(1)$ . Together with lemma 3 and the first part of the proof of proposition 1 we obtain

$$(f) \quad E \int_{|x| \leq d_n} \left\{ \frac{\hat{g}'(x; \underline{X}_n)}{\hat{g}(x; \underline{X}_n)} - \frac{g'(x)}{g(x)} \right\}^2 g(x) dx = o(1).$$

Again by the last part of the proof of lemma 2 we have

$$1 \geq E \left( \frac{c_n}{\hat{g}(x; \underline{X}_n) + c_n} \right)^2 = o(1)$$

for all  $x$ . Together with (f) this implies

$$\begin{aligned} & E \int \{\hat{h}(x; \underline{X}_n) - h_0(x)\}^2 g(x) dx \\ & \leq \int_{|x| > d_n} h_0^2(x) g(x) dx + 2E \int_{|x| \leq d_n} \left\{ \frac{\hat{g}'(x; \underline{X}_n)}{\hat{g}(x; \underline{X}_n)} - \frac{g'(x)}{g(x)} \right\}^2 \left\{ \frac{\hat{g}(x; \underline{X}_n)}{\hat{g}(x; \underline{X}_n) + c_n} \right\}^2 g(x) dx \\ & \quad + 2E \int_{|x| \leq d_n} \left\{ \frac{g'(x)}{g(x)} \right\}^2 \left\{ \frac{c_n}{\hat{g}(x; \underline{X}_n) + c_n} \right\}^2 g(x) dx \\ & = o(1), \end{aligned}$$

and (46) has been proved. Finally, in view of  $|w'(x)| \leq w(x) \leq 1/4$ ,

$$\begin{aligned} & \left| \frac{\hat{g}'(x; \underline{X}_1)}{\hat{g}(x; \underline{X}_1) + c_n} - \frac{\hat{g}'(x; \underline{X}_2)}{\hat{g}(x; \underline{X}_2) + c_n} \right| \\ & = \left| \frac{\hat{g}'(x; \underline{X}_1)}{\hat{g}(x; \underline{X}_1) + c_n} \right. \\ & \quad \left. - \frac{\hat{g}'(x; \underline{X}_1) + b_n^{-2} n^{-1} [w'((x - X_1)/b_n) - w'((x - X_2)/b_n)]}{\hat{g}(x; \underline{X}_1) + c_n + b_n^{-1} n^{-1} [w((x - X_1)/b_n) - w((x - X_2)/b_n)]} \right| \\ & \leq \frac{|\hat{h}(x; \underline{X}_1)|}{4n b_n c_n} + \frac{1}{2n b_n^2 c_n} \leq (b_n^2 c_n n)^{-1}, \end{aligned}$$

and consequently, (c) yields

$$nE \int \{\hat{h}(x; \underline{X}_1) - \hat{h}(x; \underline{X}_2)\}^2 g(x) dx \leq (b_n^4 c_n^2 n)^{-1} = o(1),$$

which is (47). □

# Appendix

## A.1 VECTOR SPACES; LINEAR FUNCTIONALS AND DUAL SPACES

### *Banach Spaces and Hilbert Spaces*

A *normed linear space*  $\mathbf{X}$  is a vector space with a norm  $\|\cdot\| : \mathbf{X} \rightarrow R$  satisfying

- (i)  $\|x\| \geq 0$  for all  $x \in \mathbf{X}$  and  $\|x\| = 0$  iff  $x = 0$ ,
- (1) (ii)  $\|x+y\| \leq \|x\| + \|y\|$  for all  $x, y \in \mathbf{X}$ ,
- (iii)  $\|cx\| = |c| \|x\|$  for all scalars  $c$  and  $x \in \mathbf{X}$ .

If  $\mathbf{X}$  is a normed linear space which is *complete*, so that every Cauchy sequence in  $\mathbf{X}$  has a limit point which is also in  $\mathbf{X}$ , then  $\mathbf{X}$  is a *Banach space*.

**Example 1.**  $\mathbf{X} = R^d$  with  $\|x\|$  the Euclidean norm  $|x| \equiv (\sum_{i=1}^d x_i^2)^{1/2}$ . In this case, we use the notation  $|x|$  throughout the text. We also occasionally use the norm  $|x|_\infty \equiv \max\{|x_i| : i = 1, \dots, d\}$ . Since

$$(2) \quad d^{-1/2} |x| \leq |x|_\infty \leq |x|,$$

these are equivalent norms, and  $\mathbf{X} = R^d$  is a Banach space with either norm.  $\square$

**Example 2.**  $\mathbf{X} = C[0,1] \equiv \{x : [0,1] \rightarrow R : x \text{ is continuous}\}$  with  $\|x\|$  the supremum norm  $\|x\|_\infty \equiv \sup_{0 \leq t \leq 1} |x(t)|$ .  $\square$

**Example 3.** Let  $\mu$  be a positive measure on a measurable space  $(\mathbf{X}, \mathcal{B})$ . For  $f : \mathbf{X} \rightarrow R$  and  $0 < p < \infty$ , let  $\|f\|_p \equiv \left\{ \int_{\mathbf{X}} |f|^p d\mu \right\}^{1/p}$ , and set  $L_p(\mathbf{X}, \mu) \equiv L_p(\mu) \equiv \{f : \|f\|_p < \infty\}$ . If  $1 \leq p < \infty$ , then  $L_p(\mu)$  is a Banach space for every positive measure  $\mu$ .  $\square$

A vector space  $\mathbf{X}$  is a (real) *inner-product space* or *pre-Hilbert space* if there is an inner product  $\langle \cdot, \cdot \rangle : \mathbf{X} \times \mathbf{X} \rightarrow R$  which satisfies

- (i)  $\langle x, y \rangle = \langle y, x \rangle$ ,
- (3) (ii)  $\langle x+y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$ ,
- (iii)  $\langle x, x \rangle \geq 0$  for all  $x \in \mathbf{X}$  with  $\langle x, x \rangle = 0$  iff  $x = 0$ .

Thus  $\|x\| \equiv \sqrt{\langle x, x \rangle}$  is a norm, and if the resulting normed linear space is complete,  $\mathbf{X}$  is a *Hilbert space*. It follows that a Hilbert space is a Banach space with an inner product which defines the norm. Following traditional notation, we will frequently use  $\mathbf{H}$  to denote Hilbert spaces.

**Example 4.**  $\mathbf{H} = R^d$  with  $\langle x, y \rangle \equiv \sum_{i=1}^d x_i y_i$  is a Hilbert space. □

**Example 5.**  $\mathbf{H} = L_2(\mu)$  with  $\langle f, g \rangle \equiv \int_{\mathbf{X}} fg \, d\mu$  is a Hilbert space. □

*Linear Functionals and Dual Spaces*

If  $\mathbf{X}, \mathbf{Y}$  are normed linear spaces and  $T$  is a function from  $\mathbf{X}$  to  $\mathbf{Y}$ , then  $T$  is *linear* if  $T(ax + by) = aT(x) + bT(y)$  for all scalars  $a, b$  and  $x, y \in \mathbf{X}$ .  $T$  is called a *linear operator*. When the range space  $\mathbf{Y}$  is  $R$ ,  $T$  is called a *linear functional*.

**Proposition 1.** If  $T$  is a linear operator from a normed linear space  $\mathbf{X}$  to a normed linear space  $\mathbf{Y}$  which is continuous at  $x_0 \in \mathbf{X}$ , then  $T$  is continuous at every  $x \in \mathbf{X}$ .

**Definition 1.** A linear operator  $T$  is *bounded* if

$$(4) \quad \|T\| \equiv \sup \{ \|T(x)\| : x \in \mathbf{X}, \|x\| \leq 1 \} < \infty .$$

**Proposition 2.** A linear operator  $T$  is bounded if and only if  $T$  is continuous.

**Definition 2.** If  $\mathbf{X}$  is a normed linear space, then the *dual space*  $\mathbf{X}^*$  of  $\mathbf{X}$  is the space of all bounded linear functionals on  $\mathbf{X}$ .

If  $x^* \in \mathbf{X}^*$ , then we write  $\langle x, x^* \rangle_{\mathbf{X}} \equiv \langle x, x^* \rangle \equiv x^*(x)$  for the value of  $x^*$  at  $x \in \mathbf{X}$ . As usual, if there is a one-to-one correspondence between  $\mathbf{X}^*$  and some other set  $\mathbf{Y}$ , we will identify  $\mathbf{X}^*$  with  $\mathbf{Y}$ . According to definition 1,  $\|x^*\| = \sup \{ |x^*(x)| : x \in \mathbf{X}, \|x\| \leq 1 \}$  for every  $x^* \in \mathbf{X}^*$ .

**Proposition 3.**

A.  $\mathbf{X}^*$  is a Banach space.

B.  $\|x\| = \sup \{ |x^*(x)| : x^* \in \mathbf{X}^*, \|x^*\| \leq 1 \}$ .

**Example 6.** If  $\mathbf{X} = R^d$ , then  $\mathbf{X}^* = R^d$  and  $x^*(x) = \sum_{i=1}^d x_i^* x_i = \langle x, x^* \rangle$ . □

**Example 7.** If  $\mathbf{X} = L_p(\mu)$  with  $1 \leq p < \infty$  and  $1/p + 1/q = 1$ , then  $\mathbf{X}^* = L_q(\mu)$ , and  $x^*(x) = \int x^*(t)x(t) \, d\mu(t)$ . □

**Example 8.** If  $\mathbf{H}$  is a Hilbert space with inner product  $\langle \cdot, \cdot \rangle$ , then  $\mathbf{H}^*$  and  $\mathbf{H}$  can be identified since for every  $h^* \in \mathbf{H}^*$  there exists a unique element in  $\mathbf{H}$ , denoted by  $h^*$  too, such that  $h^*(h) = \langle h, h^* \rangle$ . Both example 6 and example 7 with  $p = 2$  are special cases of this; see, e.g., Rudin (1966, theorem 4.12, page 80). □

**Example 9.** If  $\mathbf{X} = C[0, 1]$ , then

$$\mathbf{X}^* = NBV[0, 1]$$

$$\equiv \{ x^* : [0, 1] \rightarrow R : \int_0^1 d|x^*(t)| \equiv |x^*|[0, 1] < \infty \} .$$

This is the Riesz representation theorem; see Rudin (1966, theorem 6.19, page 131) or Hewitt and Stromberg (1965, theorem 20.48, page 364).  $\square$

If  $X = Y = L_p(R, \lambda)$  where  $p \geq 1$ , and  $\lambda$  denotes Lebesgue measure, then for  $h \in R$ , the translation operator  $T_h$  is defined by

$$T_h f(x) \equiv f(x-h) \equiv f_h(x) \quad \text{for } x \in R.$$

**Theorem 1.** ( $L_p$ -continuity theorem)

$$\|T_h f - f\|_p = \|f_h - f\|_p \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

**Proof.** This is a consequence of Lusin's theorem and Vitali's theorem; see, e.g., Hewitt and Stromberg (1965, theorem 13.24, page 199).  $\square$

### Adjoins or Transpose Operators

Let  $X$  and  $Y$  be normed linear spaces with dual spaces  $X^*$  and  $Y^*$  respectively, and suppose that  $A$  is a bounded linear operator from  $X$  to  $Y$ ,  $A \in B(X, Y)$ . Recall that

$$(5) \quad \|A\| \equiv \sup\{\|Ax\| : \|x\| \leq 1\}$$

and that  $B(X, Y)$  is thus a normed space.

**Example 1, continued.** If  $X = R^c$ ,  $Y = R^d$ , and  $A$  is a  $d \times c$ -matrix, then  $A$  is a bounded linear operator from  $X$  to  $Y$  with norm

$$(6) \quad \|A\| \equiv \sup\{|Ax| : |x| \leq 1\}.$$

If, instead of the Euclidean norm  $|\cdot|$ , we use the uniform norm  $|\cdot|_\infty$  on both  $X$  and  $Y$ , then

$$(7) \quad \|A\|_\infty \equiv \sup\{|Ax|_\infty : |x|_\infty \leq 1\}.$$

Note that

$$(8) \quad d^{-1/2} \|A\| \leq \|A\|_\infty \leq c^{1/2} \|A\|.$$

$\square$

If  $Y$  is a Banach space, so is  $B(X, Y)$ . Note that proposition 3.A is the special case of this with  $Y = R$ . The *adjoint* or *transpose*  $A^T$  of  $A$  is defined to be the linear operator from  $Y^*$  to  $X^*$  satisfying

$$(9) \quad \langle x, A^T y^* \rangle_X = \langle Ax, y^* \rangle_Y$$

for all  $x \in X$  and  $y^* \in Y^*$ . Here  $\langle \cdot, \cdot \rangle_X$  on the left side is defined by duality between  $X$  and  $X^*$  and this notation was introduced immediately after definition 2, so  $\langle x, A^T y^* \rangle_X \equiv (A^T y^*)(x)$ . Similarly,  $\langle \cdot, \cdot \rangle_Y$  on the right side is defined by duality of  $Y$  and  $Y^*$ :  $\langle Ax, y^* \rangle_Y = y^*(Ax)$ . Thus  $(A^T y^*)(x) = (y^* A)(x)$  for all  $x \in X$ , or  $A^T y^* = y^* A$ .

The transpose operation satisfies the usual properties associated with the transpose of a matrix; see Rudin (1973, theorem 4.10, page 93 and formula (12.9.6), page 297).

**Proposition 4.**

- A. If  $I : X \rightarrow X$  is the identity,  $I^T = I : X^* \rightarrow X^*$ .
- B. If  $A, B \in B(X, Y)$ , then  $(A + B)^T = A^T + B^T$ .
- C. If  $A \in B(X, Y)$  and  $c \in R$ , then  $(cA)^T = cA^T$ .
- D. If  $A \in B(X, Y)$  and  $B \in B(Y, Z)$ , then  $(BA)^T = A^T B^T$ .
- E. If  $A \in B(X, Y)$  has bounded inverse  $A^{-1}$ , then  $(A^{-1})^T = (A^T)^{-1}$ .
- F. If  $A \in B(X, Y)$ , then  $\|A^T\| = \|A\|$ .
- G. If  $A \in B(X, X)$  with  $X$  a Hilbert space, then  $(A^T)^T = A$ .

Here is a natural and useful formula for computing transposes.

**Proposition 5.** Suppose that the domain  $D(A)$  of  $A$  satisfies  $D(A) \subset L_2(m_1)$ , and that for all  $a \in D(A)$ ,

$$Aa(\cdot) = \int K(s, \cdot) a(s) dm_1(s) \in L_2(m_2).$$

If  $\iint |K(s, t)| |a(s)b(t)| dm_1(s) dm_2(t) < \infty$  for all  $a \in D(A)$ ,  $b \in D(A^T) \subset L_2(m_2)$ , then

$$(10) \quad A^T b(\cdot) = \int K(\cdot, t) b(t) dm_2(t).$$

**Proof.** Straightforward computation using Fubini. □

If  $A \in B(X, Y)$ , the null space and range of  $A$  are denoted by  $N(A)$  and  $R(A)$  respectively:

$$(11) \quad N(A) \equiv \{x \in X : Ax = 0\},$$

$$(12) \quad R(A) \equiv \{y \in Y : Ax = y \text{ for some } x \in X\}.$$

Suppose that  $X$  is a Banach space,  $X_0$  is a subset of  $X$ , and  $X_0^*$  is a subset of  $X^*$ . Then the *annihilators*  $X_0^\perp$  and  ${}^\perp X_0^*$  of  $X_0$  and  $X_0^*$  respectively are defined by

$$(13) \quad X_0^\perp \equiv \{x^* \in X^* : \langle x, x^* \rangle = 0 \text{ for all } x \in X_0\},$$

$$(14) \quad {}^\perp X_0^* \equiv \{x \in X : \langle x, x^* \rangle = 0 \text{ for all } x^* \in X_0^*\}.$$

If  $X$  is a Hilbert space so that  $H^* = H$ , then the annihilators are just orthogonal complements of subsets of  $H$ ; see section A.2. The following theorem relates the null spaces and ranges of an operator  $A \in B(X, Y)$  and its adjoint  $A^T \in B(Y^*, X^*)$ .

**Theorem 2.** If  $X$  and  $Y$  are Banach spaces, and  $A \in B(X, Y)$ , then

$$(15) \quad N(A^T) = R(A)^\perp,$$

and

$$(16) \quad N(A) = {}^\perp R(A^T).$$

**Proof.** See Rudin (1973, theorem 4.12, page 94). □

**Corollary 1.** If  $X$  and  $Y$  are Hilbert spaces, and  $A \in B(X, Y)$ , then

$$(17) \quad N(A^T) = R(A)^\perp$$

and

$$(18) \quad \mathbf{N}(A) = \mathbf{R}(A^T)^\perp,$$

where  $\perp$  denotes orthogonal complement (see section A.2).

**Proposition 6.** (Polar decomposition). Suppose that  $A$  is a continuous linear map from a Hilbert space  $\mathbf{H}$  to another Hilbert space  $\mathbf{K}$ . Then:

- A. There exists a unique self-adjoint positive definite operator  $(A^T A)^{1/2} : \mathbf{H} \rightarrow \mathbf{H}$  such that  $A^T A = (A^T A)^{1/2} (A^T A)^{1/2}$  (so  $(A^T A)^{1/2}$  is the square root of  $A^T A$ ).
- B.  $\mathbf{R}((A^T A)^{1/2}) = \mathbf{R}(A^T)$ .
- C. If  $k \in \overline{\mathbf{R}(A)}$  and  $h \in \overline{\mathbf{R}(A^T)}$  satisfy  $A^T k = (A^T A)^{1/2} h$ , then  $\|k\|_{\mathbf{K}} = \|h\|_{\mathbf{H}}$ .

**Proof.** (Van der Vaart (1991)) Part A is standard; see, e.g., Rudin (1973, 12.33) or Reed and Simon (1972, theorem V.1.9, page 196). To prove B, note that

$$\begin{aligned} \|Ah\|_{\mathbf{K}}^2 &= \langle Ah, Ah \rangle_{\mathbf{K}} = \langle A^T A h, h \rangle_{\mathbf{H}} \\ &= \langle (A^T A)^{1/2} h, (A^T A)^{1/2} h \rangle_{\mathbf{H}} = \|(A^T A)^{1/2} h\|_{\mathbf{H}}^2. \end{aligned}$$

Thus there is an isometry  $U$  from  $\mathbf{R}((A^T A)^{1/2})$  to  $\mathbf{R}(A)$  defined by  $U((A^T A)^{1/2} h) = Ah$ . We can then extend  $U$  to an isometry of  $\overline{\mathbf{R}((A^T A)^{1/2})}$  onto  $\overline{\mathbf{R}(A)}$ , and then to a partial isometry on  $\mathbf{H}$  by defining it to be zero on  $\mathbf{R}((A^T A)^{1/2})^\perp$ . Then  $A^T = (A^T A)^{1/2} U^T$  where  $U^T$  acts as the inverse of  $U$  on  $\overline{\mathbf{R}(A)}$  and  $\mathbf{R}((A^T A)^{1/2})^\perp = \mathbf{N}((A^T A)^{1/2})$ . Thus,

$$A^T \mathbf{K} = (A^T A)^{1/2} U^T \mathbf{K} = (A^T A)^{1/2} \overline{\mathbf{R}((A^T A)^{1/2})} = (A^T A)^{1/2} \mathbf{H},$$

proving B.

Finally, if  $A^T k = (A^T A)^{1/2} h$ , then  $U^T k - h \in \mathbf{N}((A^T A)^{1/2}) = \mathbf{R}((A^T A)^{1/2})^\perp$ . But since  $k \in \overline{\mathbf{R}(A)}$  and  $h \in \overline{\mathbf{R}(A^T)}$ ,  $U^T k - h \in \overline{\mathbf{R}(A^T)} = \overline{\mathbf{R}((A^T A)^{1/2})}$  by B, and taken together this yields  $U^T k = h$ . The conclusion follows since  $U^T$  is an isometry on  $\overline{\mathbf{R}(A)}$ .  $\square$

### Inverses and Ranges

**Proposition 7.** Suppose that  $\mathbf{X}$  and  $\mathbf{Y}$  are normed linear spaces and  $A : \mathbf{X} \rightarrow \mathbf{Y}$  is linear. Then:

- A.  $A$  has a continuous inverse  $A^{-1}$  from  $\mathbf{R}(A)$  to  $\mathbf{X}$  if and only if  $\|Ax\| \geq c\|x\|$  for all  $x \in \mathbf{X}$  where  $c > 0$ .
- B. (Banach's theorem) If  $\mathbf{X}$  and  $\mathbf{Y}$  are Banach spaces and  $A$  is continuous and one-to-one ( $\mathbf{N}(A) = \{0\}$ ), then  $A^{-1}$  is continuous (i.e. bounded) if and only if  $\mathbf{R}(A)$  is closed.
- C. If  $\mathbf{X}$  and  $\mathbf{Y}$  are Banach spaces and  $A$  is bounded, then  $A$  is onto  $\mathbf{Y}$  if and only if  $\|A^T y^*\| \geq d\|y^*\|$  for all  $y^* \in \mathbf{Y}^*$  for some  $d > 0$ .



D. If  $X$  and  $Y$  are Banach spaces and  $A$  is bounded, then  $R(A)$  is closed if and only if  $R(A^T)$  is closed.

**Proof.** For  $A$ , see Yosida (1974, corollary I.6.3, page 43). For  $B$ , see Jörgens (1982, theorem 5.1, page 81), or Rudin (1973, corollary 2.12, page 48), and corollary 3 below. For  $C$  and  $D$ , see Rudin (1973, theorem 4.15, page 97, and theorem 4.14, page 96).  $\square$

**Corollary 2.** If  $X = Y$  is a Hilbert space and  $A = A^T$ , then  $A$  has a bounded inverse  $A^{-1} : X \rightarrow X$  if and only if  $\|Ax\| \geq c\|x\|$  for all  $x \in X$  for some  $c > 0$ . Then  $\|A^{-1}\| \leq 1/c$ .

Suppose that  $A$  is a bounded operator from  $X$  to  $Y$ . It is not true, in general that  $R(A)$  is closed in  $Y$ . Here are two simple examples illustrating the difficulty; both examples are from Halmos (1982, page 29).

**Example 10.** Let  $X = Y = l_2$ , and define  $A : l_2 \rightarrow l_2$  by  $Ax = (x_1, \frac{1}{2}x_2, \frac{1}{3}x_3, \dots)$  for  $x \in l_2 = \{x \in R^\infty : \sum_{i=1}^\infty x_i^2 < \infty\}$ . Then  $A$  is bounded with  $\|A\| = 1$  and  $R(A) = \{y \in l_2 : \sum_{i=1}^\infty i^2 y_i^2 < \infty\}$ . Since this set contains all  $y$  of the form  $y = (y_1, y_2, \dots, y_k, 0, 0, \dots)$  with all zeros after a finite number of coordinates,  $R(A)$  is dense in  $l_2$ . On the other hand,  $y = (1, \frac{1}{2}, \frac{1}{3}, \dots) \in l_2$ , but is not in  $R(A)$ . Thus  $R(A)$  is not closed.  $\square$

**Example 11.** Let  $X = Y = L_2(0, 1)$ , and define  $A : L_2(0, 1) \rightarrow L_2(0, 1)$  by  $Ax(u) = ux(u)$  for  $x \in L_2(0, 1)$  and  $0 \leq u \leq 1$ . Then  $A$  is bounded and  $R(A) = \{y \in L_2(0, 1) : \int_0^1 u^{-2} y^2(u) du < \infty\}$ . This set is dense in  $L_2(0, 1)$  since it contains all functions of the form  $x(u)1_{[e, 1]}(u)$  for  $x \in L_2(0, 1)$  and  $0 < \epsilon < 1$ . But  $y_1(u) \equiv 1$  and  $y_2(u) = \sqrt{u}$  are in  $L_2(0, 1)$  but not in  $R(A)$ . Thus  $R(A)$  is not closed.  $\square$

Of course, these examples are connected with the failure of the inverse operator  $A^{-1}$  to be bounded: in example 10,  $A^{-1}y = (y_1, 2y_2, 3y_3, \dots)$  for  $y \in R(A)$  is not bounded as an operator on  $l_2$ ; and in example 11,  $A^{-1}y(u) = u^{-1}y(u)$  for  $y \in R(A)$  is not bounded as an operator on  $L_2(0, 1)$ . The following corollary gives a simple condition for closedness of  $R(A)$  for a bounded operator  $A$  in terms of boundedness of  $A^{-1}$ .

**Corollary 3.** Suppose that  $X$  and  $Y$  are normed linear spaces,  $A : X \rightarrow Y$  is a bounded linear operator,  $A^{-1}$  exists on  $R(A)$  and is bounded, and  $X$  is complete. Then  $R(A)$  is closed.

**Proof.** Suppose that  $Ax_n \rightarrow y \in Y$ . We want to show that  $y = Ax$  for some  $x$ . Now

$$\|A^{-1}(Ax_n) - A^{-1}(Ax_m)\| \leq \|A^{-1}\| \|Ax_n - Ax_m\| \rightarrow 0$$

since  $Ax_n \rightarrow y$  and  $\|A^{-1}\| < \infty$  by the boundedness of  $A^{-1}$ . Thus  $\{x_n\} = \{A^{-1}(Ax_n)\}$  is Cauchy in  $X$ . Hence, by completeness of  $X$ ,  $x_n \rightarrow$  some  $x \in X$ . Hence  $Ax_n \rightarrow Ax$  by continuity (boundedness) of  $A$ . But  $Ax_n \rightarrow y$ , so  $y = Ax$ .  $\square$

The following lemma is closely related to the Banach-Steinhaus theorem;

see, e.g., Rudin (1973, pages 43–46). The present formulation is from Van der Vaart (1991).

**Lemma 1.** Suppose that  $X$  is a normed linear space and  $Y$  is a Banach space. If  $f : X \rightarrow Y$  satisfies  $y^* \circ f \in X^*$  for every  $y^* \in Y^*$ , then  $f$  is continuous and linear.

**Proof.** For all  $y^* \in Y^*$ ,  $\alpha_1, \alpha_2 \in R$  and  $x_1, x_2 \in X$  we have

$$(a) \quad y^*(f(\alpha_1 x_1 + \alpha_2 x_2) - \alpha_1 f(x_1) - \alpha_2 f(x_2)) = 0$$

since  $y^* \in Y^*$  and  $y^* \circ f \in X^*$ . Hence  $f$  is linear since  $Y^*$  separates points of  $Y$ ; see, e.g., corollary 3.4 of Rudin (1973).

Consider the embedding  $\phi$  of  $Y$  into  $Y^{**}$  with, for every  $y \in Y$ ,

$$(b) \quad \phi(y)(y^*) = y^*(y), \quad y^* \in Y^*.$$

Let  $U$  be the unit ball in  $X$ . We have to show that  $f(U) \subset Y$  is bounded. Since  $y^* \circ f \in X^*$  is bounded,

$$(c) \quad \{\phi(f(u))(y^*) : u \in U\} = \{y^* \circ f(u) : u \in U\} = y^* \circ f(U)$$

is bounded for all  $y^* \in Y^*$ . It follows by the Banach-Steinhaus theorem (see sections 2.6 and 2.4 of Rudin (1973)) that the collection  $\phi(f(U))$  of linear functionals from  $Y^*$  to  $R$  is equicontinuous and hence uniformly bounded. Consequently  $|\phi(f(u))(y^*)|$  is bounded uniformly in  $u \in U$  and  $y^* \in Y^*$  with  $\|y^*\| \leq 1$ . Since  $\phi : Y \rightarrow Y^{**}$  is an isometry (section 4.5 of Rudin (1973)), it follows that

$$\begin{aligned} \sup_{u \in U} \|f(u)\| &= \sup_{u \in U} \|\phi(f(u))\| \\ &= \sup_{u \in U} \sup\{|\phi(f(u))(y^*)| : \|y^*\| \leq 1\} < \infty \end{aligned}$$

and hence that  $f(U)$  is bounded. □

### The $R$ and $L$ Operators

Now we introduce two simple operators, and use them to illustrate some of the above definitions and propositions. Proposition 8 below is used in our treatments of the Cox proportional hazards model in examples 3.4.2 and 5.5.2, of the censored regression model in examples 4.6.4 and 6.6.2, and of the random truncation model in example 6.4.2.

Let  $F$  be a continuous df on  $R$ , let  $H \equiv L_2(F)$ , and let  $H_0 \equiv \{a \in H : \int a dF = 0\}$ . We define the operators  $R$  and  $L$  by

$$\begin{aligned} (19) \quad R a(t) &\equiv a(t) - \frac{\int_t^\infty a dF}{1 - F(t)} \\ &= -E\{a(X) - a(t) \mid X > t\} \end{aligned}$$

and

$$(20) \quad Lb(t) \equiv b(t) - \int_{-\infty}^t b \, d\Lambda = b(t) - \int_{-\infty}^t b \frac{dF}{1-F} \\ = b(t) - \int_{-\infty}^{\infty} 1_{[t \geq s]} b(s) \, d\Lambda(s).$$

Note that  $L$  is a ‘‘martingale operator’’: if  $X \sim F$  and

$$(21) \quad M(t) = 1_{[X \leq t]} - \int_{-\infty}^t 1_{[X \geq s]} \, d\Lambda(s)$$

is the corresponding counting process martingale, then

$$(22) \quad Lb(X) = \int_{-\infty}^{\infty} b \, dM.$$

The operators  $R$  and  $L$  arise as the (logarithmic) derivatives of the maps  $Hf \equiv f/(1-F) \equiv \lambda$  and  $D\lambda \equiv \lambda \exp[-\int_0^{\cdot} \lambda(s) \, ds] \equiv f$  which map a density function  $f$  to its hazard function  $\lambda$ , and vice-versa. As a result the  $R$  and  $L$  operators are closely related. For use and motivation of the  $R$  and  $L$  operators, see sections 3.4 and 4.6. The following proposition records some of the important properties of  $R$  and  $L$ .

**Proposition 8.** Suppose that the df  $F$  is continuous. The operators  $R$  and  $L$  mapping  $\mathbf{H} = L_2(F)$  to itself, and defined by (19) and (20) respectively, satisfy:

- A.  $R$  and  $L$  are bounded with  $\|R\| = \|L\| = 1$ .
- B. For  $a \in \mathbf{H}$ ,  $L \circ Ra = a - E a(X)$  and  $R \circ La = a$ ; thus  $R^{-1} = L$  on  $\mathbf{H}_0$ .
- C.  $L = R^T$  and  $R = L^T$ ; hence  $L$  and  $R$  are isometries of  $\mathbf{H}_0$  (or unitary transformations):  $\|La\| = \|a\| = \|Ra\|$  for all  $a \in \mathbf{H}_0$  and  $L^T L = R^T R = \text{identity on } \mathbf{H}_0$ .
- D.  $\text{Var}[a(X)] = E[Ra(X)]^2$ , or, with  $A$  defined by  $A(t) \equiv E[a(X) | X > t]$ ,  
 $\text{Var}[a(X)] = E[a(X) - A(X)]^2$ .
- E.  $\mathbf{N}(R) = \{\text{constants}\}$ ,  $\mathbf{R}(R) = \mathbf{H}$ ,  $\mathbf{N}(L) = \{0\}$ , and  $\mathbf{R}(L) = \mathbf{H}_0$ .

**Proof.** We first prove B–E assuming that  $R$  and  $L$  are bounded, and then return to the proof of A.

B follows by direct calculation via Fubini. It can be easily understood in the case  $a \in \mathbf{H}_0$  in view of the fact that both  $D \circ Hf = f$  and  $H \circ D\lambda = \lambda$ ; (logarithmic) differentiation of these identities yields the claim when  $a \in \mathbf{H}_0$ .

To prove C, let  $\langle \cdot, \cdot \rangle$  denote the inner product in  $L_2(F)$ , and let  $a, b \in \mathbf{H}$ . Then  $Lb \in \mathbf{H}$  by the boundedness of  $L$ , and Fubini’s theorem yields

$$\langle a, Lb \rangle = \int_{-\infty}^{\infty} \left\{ b(t) - \int_{-\infty}^t b \, d\Lambda \right\} a(t) \, dF(t) \\ = \int_{-\infty}^{\infty} b(t) \left\{ a(t) - \frac{\int_{-\infty}^{\infty} a \, dF}{1-F(t)} \right\} \, dF(t)$$

(a)  $\qquad\qquad\qquad = \langle Ra, b \rangle,$

and hence  $L^T = R$ . Then  $R^T = (L^T)^T = L$  by proposition 4.G. Now it follows that  $R$  and  $L$  are isometries of  $H_0$ : for  $a \in H_0$  we have

(b)  $\qquad\|La\|^2 = \langle La, La \rangle$   
 $\qquad\qquad\qquad = \langle a, L^T La \rangle \quad \text{by definition of } L^T$   
 $\qquad\qquad\qquad = \langle a, RL a \rangle \quad \text{by } L^T = R$   
 $\qquad\qquad\qquad = \langle a, a \rangle = \|a\|^2 \quad \text{by the second identity of part B.}$

A similar calculation shows that  $\|Ra\|^2 = \|a\|^2$  for  $a \in H_0$ , but more generally it follows from B that

$$\begin{aligned} \text{Var}[a(X)] &= \langle a, a - E a(X) \rangle \\ &= \langle a, L R a \rangle \quad \text{by B} \\ &= \langle R a, R a \rangle \quad \text{by C} \\ &= E[a(X) - A(X)]^2 \quad \text{by (19),} \end{aligned}$$

and this proves D.

Now for E.  $R(R) = H$  follows from A and the second identity of B,  $R \circ L a = a$ . Next,  $R(L) = H_0$  follows from A,  $R(R) = H$ , and the first identity of B. Then  $N(R) = R(L)^\perp = H_0^\perp = \{\text{constants}\}$ , and  $N(L) = R(R)^\perp = H^\perp = \{0\}$  by corollary 1.

It remains only to prove A. Since  $L^T = R$ , by the theory of adjoints it suffices to prove that one of  $R$  and  $L$  is bounded. We will give two different proofs. The first shows that  $L$  is bounded by a martingale argument, while the second uses Hardy's inequality to show that  $R$  is bounded. Yet another proof can be based on use of the Cauchy-Schwarz inequality and Fubini's theorem.

First proof of A: We use martingale theory as presented, e.g., in appendix B of Shorack and Wellner (1986). For  $a \in L_2(F)$  set

(c)  $\qquad\mathbb{Z}(t) \equiv \int_{-\infty}^t a \, dM$

where  $M$  is the counting process martingale of (21) with predictable variation process  $\langle M \rangle(t) = \int_{-\infty}^t 1_{[X \geq s]} \, d\Lambda(s)$ . Then  $\mathbb{Z}$  is a square integrable martingale with

(d)  $\qquad\mathbb{Z}(t) \xrightarrow{\text{a.s.}} La(X) \quad \text{as } t \rightarrow \infty,$

and with predictable variation process

$$\langle \mathbb{Z} \rangle(t) = \int_{-\infty}^t a^2 \, d\langle M \rangle;$$

see, e.g., Shorack and Wellner (1986, formula (B.3.2), page 891, and example

B.2.3.) Thus,

$$\begin{aligned}
 E[\mathbb{Z}^2(t)] &= E[\langle \mathbb{Z} \rangle(t)] \\
 &= E\left[\int_{-\infty}^t a^2 d\langle M \rangle\right] \\
 &= E\left[\int_{-\infty}^t a^2(s) 1_{[X \geq s]} d\Lambda(s)\right] \\
 \text{(e)} \quad &= \int_{-\infty}^t a^2(s) dF(s).
 \end{aligned}$$

Hence by Fatou's lemma and (e),

$$\begin{aligned}
 \text{(f)} \quad E\{[La(X)]^2\} &= E\left\{\lim_{t \rightarrow \infty} \mathbb{Z}^2(t)\right\} \\
 &\leq \liminf_{t \rightarrow \infty} \int_{-\infty}^t a^2 dF \\
 &= E[a^2(X)] < \infty.
 \end{aligned}$$

It follows that  $L$  is bounded (as an operator from  $L_2(F)$  to  $L_2(F)$ ) and  $\|L\| \leq 1$ . Then, by part C,  $\|L\| = \|R\| = 1$ .

The second proof of A uses Hardy's inequality. Let  $T : L_2(0,1) \rightarrow L_2(0,1)$  be defined by

$$\text{(g)} \quad Th(x) = \frac{1}{x} \int_0^x h(y) dy \quad \text{for } h \in L_2(0,1).$$

Then  $T$  is bounded with  $\|T\| = 2$ : i.e.

$$\text{(23)} \quad \int_0^1 \left\{ \frac{1}{x} \int_0^x h(y) dy \right\}^2 dx \leq 4 \int_0^1 h^2(y) dy.$$

This is Hardy's inequality; see Hardy, Littlewood, and Polya (1952, page 240), Dunford and Schwartz (1958, page 582), or Rudin (1973, page 107); it can be proved by use of the Cauchy-Schwarz inequality and Fubini's theorem.

Boundedness of  $Ra = a - \int_{-\infty}^{\infty} a dF / (1 - F)$  follows from boundedness of  $T$  by the probability integral transformation: since the first term of  $R$  is clearly bounded, it suffices to show that

$$\text{(h)} \quad \int_{-\infty}^{\infty} \left\{ \frac{1}{1 - F(t)} \int_t^{\infty} a dF \right\}^2 dF(t) \leq 4 \int_{-\infty}^{\infty} a^2 dF.$$

But by letting  $u = 1 - F(t) = \bar{F}(t)$ , the left side of (h) equals

$$\begin{aligned}
 &\int_0^1 \left\{ \frac{1}{u} \int_0^u a(F^{-1}(1-s)) ds \right\}^2 du \\
 &\leq 4 \int_0^1 [a(F^{-1}(1-s))]^2 ds \quad \text{by (23)} \\
 &= 4 \int_{-\infty}^{\infty} a^2 dF.
 \end{aligned}$$

Thus (h) holds and  $R$  is bounded. □

We conclude this section with some remarks on proposition 8.

**Remark 1.** The equality  $\|La\| = \|a\|$  in  $C$  also follows from the martingale representation (22) of  $L$  and part A.

**Remark 2.** The variance identity D generalizes the formula

$$(24) \quad \text{Var}[X] = E[e^2(X)]$$

with  $e(x) = E(X - x | X > x)$  for  $X \sim F$  continuous noted by Pyke (1965), and extended to the case of discontinuous  $F$  by Hall and Wellner (1981): they show that

$$(25) \quad \text{Var}[X] = E[e(X)e(X-)]$$

always holds. These formulas were extended to any  $a \in L_2(F)$  and arbitrary  $F$  by Shorack and Wellner (1986, formula (6.4.7), page 283): from part B and (22) it follows that

$$(26) \quad a(X) - E a(X) = \int_{-\infty}^{\infty} Ra \, dM,$$

and hence a martingale calculation yields

$$(27) \quad \begin{aligned} \text{Var}[a(X)] &= E[(Ra)^2(X)(1 - \Delta\Lambda(X))] \\ &= E[e_a^2(X)(1 - \Delta\Lambda(X))]. \end{aligned}$$

(Note that  $1 - \Delta\Lambda(x) = (1 - F(x))/(1 - F(x-))$ .)

**Remark 3.** These and other properties of  $R$  and  $L$  have been independently discovered by Efron and Johnstone (1990). (Their  $A$  is our  $R$  and their  $B$  is our  $L$ .) They point out the following consequence of proposition 6.D. For a smooth family of functions  $\{h_\theta : \theta \in \Theta \subset R\}$ , write

$$\dot{h}_\theta(t) \equiv \frac{\partial}{\partial \theta} h_\theta(t).$$

Suppose that  $\{f_\theta : \theta \in \Theta\}$  is a smooth family of functions with  $\dot{f}_\theta/f_\theta \in L_2(F_\theta) \equiv L_2(F)$ . Then it follows immediately from the interpretation of  $R$  as the logarithmic derivative of the map  $H$  from densities to hazards that

$$\frac{\dot{\lambda}_\theta}{\lambda_\theta}(t) = \frac{\partial}{\partial \theta} \log \lambda_\theta(t) = R \left( \frac{\dot{f}_\theta}{f_\theta} \right) (t),$$

and hence D yields an identity for the Fisher information  $I_\theta$  for  $\theta$ :

$$(28) \quad I(\theta) = E_\theta \left[ \left( \frac{\dot{f}_\theta}{f_\theta} \right)^2 \right] = E_\theta \left[ \left( \frac{\dot{\lambda}_\theta}{\lambda_\theta} \right)^2 \right].$$

A.2 ORTHOGONALITY AND PROJECTION FORMULAS

If  $\mathbf{H}$  is a Hilbert space and  $x, y \in \mathbf{H}$  satisfy  $\langle x, y \rangle = 0$ , we say that  $x$  and  $y$  are *orthogonal* and write  $x \perp y$ . If  $x \perp y$  for all  $y \in \mathbf{H}_0 \subset \mathbf{H}$ , write  $x \perp \mathbf{H}_0$ . If  $x, y \in \mathbf{H}$  are orthogonal, then the ‘‘Pythagorean theorem’’ holds:

$$(1) \quad \|x + y\|^2 = \|x\|^2 + \|y\|^2 \quad \text{if } x \perp y.$$

For  $x \in \mathbf{H}$  define the *orthogonal complement* of  $x$  in  $\mathbf{H}$ ,  $x^\perp$ , by

$$(2) \quad x^\perp \equiv \{y \in \mathbf{H} : \langle x, y \rangle = 0\},$$

and if  $\mathbf{H}_0$  is a subspace of  $\mathbf{H}$ , the orthogonal complement of  $\mathbf{H}_0$  is the closed subspace

$$(3) \quad \mathbf{H}_0^\perp \equiv \{y \in \mathbf{H} : \langle x, y \rangle = 0 \quad \text{for all } x \in \mathbf{H}_0\}.$$

**Notation.** If  $\mathbf{H}_0$  is any set of elements of  $\mathbf{H}$ , write  $\overline{\mathbf{H}_0}$  for the closure of  $\mathbf{H}_0$ .

Note that

$$(\mathbf{H}_0^\perp)^\perp = \overline{\mathbf{H}_0}, \quad (\mathbf{H}_1 + \mathbf{H}_2)^\perp = \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp,$$

and hence

$$(4) \quad \overline{\mathbf{H}_1 + \mathbf{H}_2} = (\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)^\perp.$$

**Proposition 1.** If  $\mathbf{H}_0$  is a closed subspace of a Hilbert space  $\mathbf{H}$ , then every  $y \in \mathbf{H}$  has a unique decomposition as  $y = x_1 + x_2$  where  $x_1 \in \mathbf{H}_0$ ,  $x_2 \in \mathbf{H}_0^\perp$ .

**Proof.** See Rudin (1966, theorem 4.11, page 79). □

**Definition 1.** If  $\mathbf{H}_0$  is a closed subspace of a Hilbert space  $\mathbf{H}$ , and  $x = x_1 + x_2$  where  $x_1 \in \mathbf{H}_0$  and  $x_2 \in \mathbf{H}_0^\perp$ , then  $x_1$  is called the *orthogonal projection of  $x$  on  $\mathbf{H}_0$* , denoted by  $\Pi(x | \mathbf{H}_0)$ . By proposition 1,  $x_1 = \Pi(x | \mathbf{H}_0)$  is unique.

Here is a useful way of restating proposition 1 and definition 1.

**Theorem 1.** (Projection theorem). Suppose that  $\mathbf{H}_0$  is a closed subspace of a Hilbert space  $\mathbf{H}$ . Then for every  $x \in \mathbf{H}$  there is a unique vector  $\Pi(x) \equiv \Pi(x | \mathbf{H}_0) \in \mathbf{H}_0$  satisfying

$$\|x - \Pi(x | \mathbf{H}_0)\| = \min_{y \in \mathbf{H}_0} \|x - y\|.$$

Moreover,  $\Pi(x | \mathbf{H}_0)$  is determined uniquely by

$$(5) \quad x - \Pi(x | \mathbf{H}_0) \perp \mathbf{H}_0.$$

The following proposition gives the key properties of the *projection operator*  $\Pi$ :

**Proposition 2.**

A. If  $\mathbf{H}_0$  is a closed subspace of  $\mathbf{H}$ , then the projection operator  $\Pi(\cdot | \mathbf{H}_0) \equiv \Pi$  satisfies:

- (i)  $\Pi$  is linear:  $\Pi(x + y) = \Pi(x) + \Pi(y)$ .
- (ii)  $\Pi$  is idempotent:  $\Pi^2 = \Pi$ .
- (iii)  $\Pi$  is self-adjoint:  $\Pi^T = \Pi$ .

B. Properties (ii) and (iii) are equivalent to

$$(iv) \quad \langle \Pi x, \Pi y \rangle = \langle \Pi x, y \rangle \text{ for all } x, y \in \mathbf{H}.$$

C. If an operator  $\Pi$  satisfies (i), (ii), and (iii) (or equivalently (i) and (iv)), then  $\Pi$  is the (orthogonal) projection operator on  $\mathbf{H}_0 \equiv \mathbf{R}(\Pi)$  which is closed.

**Proof.**

A. To prove (i), note that  $x + y - \Pi(x) - \Pi(y)$  is  $\perp \mathbf{H}_0$ . By taking  $x$  to be  $\Pi(x | \mathbf{H}_0)$ , (ii) follows immediately from (5). To prove (iii), note that by (A.1.9), (5), and (i) and (ii) it suffices to show that

$$(a) \quad \langle x, \Pi y^* \rangle = \langle \Pi x, y^* \rangle \quad \text{for all } x, y^* \in \mathbf{H}_0.$$

But (a) is easy, since, for  $x, y^* \in \mathbf{H}_0$ ,

$$\langle x, \Pi y^* \rangle = \langle x, y^* \rangle = \langle \Pi x, y^* \rangle.$$

B. (ii) and (iii) clearly imply (iv). To show that (iv) implies (iii), note that

$$\begin{aligned} \langle \Pi x, y \rangle &= \langle \Pi x, \Pi y \rangle && \text{by (iv)} \\ &= \langle \Pi y, \Pi x \rangle && \text{by (A.1.3)} \\ &= \langle \Pi y, x \rangle && \text{by (iv)} \\ &= \langle x, \Pi y \rangle && \text{by (A.1.3),} \end{aligned}$$

so (iii) holds. To prove (ii), note that

$$\begin{aligned} \langle \Pi^2 x, y \rangle &= \langle \Pi x, \Pi y \rangle && \text{by (iii)} \\ &= \langle \Pi x, y \rangle && \text{by (iv),} \end{aligned}$$

and hence

$$(b) \quad \langle \Pi^2 x - \Pi x, y \rangle = 0 \quad \text{for all } y.$$

By taking  $y = \Pi^2 x - \Pi x$ , (b) implies that  $\Pi^2 x - \Pi x = 0$  for all  $x$ , or  $\Pi^2 = \Pi$ .

C. We first show that  $\mathbf{H}_0 \equiv \mathbf{R}(\Pi)$  is closed. Note that  $\mathbf{R}(\Pi) = \{x \in \mathbf{H} : \Pi x = x\}$ ; the right side is contained in  $\mathbf{R}(\Pi)$  trivially, and if  $y \in \mathbf{R}(\Pi)$ ,  $y = \Pi x$  for some  $x \in \mathbf{H}$ , so by (ii),  $y = \Pi x = \Pi^2 x = \Pi y$ .

Now (iv) implies that  $\|\Pi\| \leq 1$ : for  $x \in \mathbf{H}$ ,

$$\|\Pi x\|^2 = \langle \Pi x, \Pi x \rangle = \langle \Pi x, x \rangle \leq \|\Pi x\| \|x\|,$$

so



$$\|\Pi x\| \leq \|x\|.$$

Hence  $\Pi$  is bounded, and by proposition A.1.2,  $\Pi$  is continuous.

Finally, suppose  $\{x_n\} \subset \mathbf{R}(\Pi)$  and  $x_n \rightarrow x$ . Then, on the one hand,

$$\Pi x_n \rightarrow \Pi x \quad \text{by continuity of } \Pi,$$

while, on the other hand, since  $x_n \in \mathbf{R}(\Pi)$ ,

$$\Pi x_n = x_n \rightarrow x.$$

Hence  $\Pi x = x \in \mathbf{R}(\Pi)$ .

Now  $\Pi$  is the projection onto  $\mathbf{H}_0$  since, for  $y \in \mathbf{H}$  and  $x \in \mathbf{H}_0$ ,

$$\begin{aligned} \langle y - \Pi y, x \rangle &= \langle y, x \rangle - \langle \Pi y, x \rangle \\ &= \langle y, x \rangle - \langle y, \Pi x \rangle \quad \text{by (iii)} \\ &= \langle y, x \rangle - \langle y, x \rangle \quad \text{since } x \in \mathbf{H}_0 \\ &= 0. \end{aligned}$$

□

**Proposition 3.**

A. If  $\Pi$  is a projection operator,  $\Pi = \Pi(\cdot | \mathbf{H}_0)$  for some subspace  $\mathbf{H}_0$  of  $\mathbf{H}$ , then  $\Pi$  is bounded and  $\|\Pi\| = 1$  if  $\mathbf{H}_0 \neq \{0\}$ .

B. If  $\mathbf{H}_1, \mathbf{H}_2$  are subspaces of  $\mathbf{H}$  with  $\mathbf{H}_1 \subset \mathbf{H}_2$ , then

$$(7) \quad \Pi(\cdot | \mathbf{H}_1) = \Pi(\Pi(\cdot | \mathbf{H}_2) | \mathbf{H}_1)$$

or, with  $\Pi_i \equiv \Pi(\cdot | \mathbf{H}_i), i = 1, 2$ ,

$$(8) \quad \Pi_1 = \Pi_1 \circ \Pi_2.$$

C. If  $\mathbf{H}_1, \mathbf{H}_2$  are any subspaces of  $\mathbf{H}$  such that  $\mathbf{H}_1 \perp \mathbf{H}_2$ , then

$$(9) \quad \Pi(\cdot | \mathbf{H}_1 + \mathbf{H}_2) = \Pi(\cdot | \mathbf{H}_1) + \Pi(\cdot | \mathbf{H}_2).$$

**Proof.** A.  $\|\Pi\| \leq 1$  was established in the proof of proposition 2.C; to get  $\|\Pi\| \geq 1$  when  $\dim(\mathbf{H}) \geq 1$ , take  $x \in \mathbf{H}_0$  so that  $\Pi x = x$  and  $\|\Pi x\| = \|x\|$ .

To prove B, write

$$\begin{aligned} (a) \quad x - \Pi(\Pi(x | \mathbf{H}_2) | \mathbf{H}_1) &= x - \Pi(x | \mathbf{H}_1) + \Pi(x | \mathbf{H}_1) - \Pi(\Pi(x | \mathbf{H}_2) | \mathbf{H}_1) \\ &= x - \Pi(x | \mathbf{H}_1) + \Pi(x - \Pi(x | \mathbf{H}_2) | \mathbf{H}_1) \\ &= x - \Pi(x | \mathbf{H}_1) \end{aligned}$$

since  $x - \Pi(x | \mathbf{H}_2) \perp \mathbf{H}_2 \supset \mathbf{H}_1$ .

Now for C: If  $\mathbf{H}_1 \perp \mathbf{H}_2$ , then we have

$$(b) \quad \Pi(x | \mathbf{H}_2) \perp \mathbf{H}_1.$$

Then, by (5),  $x - \Pi(x | \mathbf{H}_1) - \Pi(x | \mathbf{H}_2) \perp \mathbf{H}_1$ . Interchanging  $\mathbf{H}_1$  and  $\mathbf{H}_2$  in this argument together with (5) yields (9).  $\square$

**Proposition 4.** Suppose that  $\mathbf{H}_1, \mathbf{H}_2$  are closed subspaces of  $\mathbf{H}$ . Then

$$(10) \quad \Pi(\cdot | \mathbf{H}_1) = \Pi(\cdot | \mathbf{H}_1 \cap \mathbf{H}_2^\perp) + \Pi(\cdot | \overline{\Pi(\mathbf{H}_2 | \mathbf{H}_1)}).$$

An important special case occurs when  $\mathbf{H}_2 \subset \mathbf{H}_1$ . Then (10) becomes

$$(11) \quad \Pi(\cdot | \mathbf{H}_1) = \Pi(\cdot | \mathbf{H}_1 \cap \mathbf{H}_2^\perp) + \Pi(\cdot | \mathbf{H}_2).$$

**Proof.** Note that  $\Pi(\mathbf{H}_2 | \mathbf{H}_1) = \mathbf{H}_1 \cap (\mathbf{H}_1^\perp + \mathbf{H}_2)$  is the set of those elements of  $\mathbf{H}_1$  that can be written as the difference of an element of  $\mathbf{H}_2$  and its projection onto  $\mathbf{H}_1^\perp$ . Hence

$$\mathbf{H}_1 \cap (\mathbf{H}_1 \cap \mathbf{H}_2^\perp)^\perp = \mathbf{H}_1 \cap \overline{(\mathbf{H}_1^\perp + \mathbf{H}_2)} = \overline{\Pi(\mathbf{H}_2 | \mathbf{H}_1)}$$

holds and

$$\mathbf{H}_1 = \mathbf{H}_1 \cap \mathbf{H}_2^\perp + \overline{\Pi(\mathbf{H}_2 | \mathbf{H}_1)}$$

with

$$\mathbf{H}_1 \cap \mathbf{H}_2^\perp \perp \overline{\Pi(\mathbf{H}_2 | \mathbf{H}_1)}.$$

The result follows from proposition 3.C.  $\square$

**Theorem 2.** Suppose  $A \in B(\mathbf{G}, \mathbf{H})$  and that  $y \in \mathbf{H}$  is fixed. Then  $x \in \mathbf{G}$  minimizes  $\|y - Ax\|$  if and only if  $A^T Ax = A^T y$ . Thus, if  $(A^T A)^{-1}$  exists,

$$(12) \quad x = (A^T A)^{-1} A^T y$$

and the projection operator  $\Pi = \Pi(\cdot | \mathbf{H}_0)$  onto  $\mathbf{H}_0 = \mathbf{R}(A)$  is given by

$$(13) \quad \Pi(\cdot | \mathbf{H}_0) = A(A^T A)^{-1} A^T.$$

**Notation:** If  $H_0$  is any set of elements of  $\mathbf{H}$ , write  $[H_0]$  for the closed linear span of  $H_0$ .

An important special case of theorem 2 is:

**Example 1.** Let  $\mathbf{G} = R^p$ , let  $\mathbf{H}$  be a Hilbert space, and suppose that  $x_1, \dots, x_p$  are elements of  $\mathbf{H}$ . Write  $x = (x_1, \dots, x_p)^T$ . Then if  $A : R^p \rightarrow \mathbf{H}$  is defined by  $A\lambda = x^T \lambda = \sum_{i=1}^p \lambda_i x_i$  for  $\lambda \in R^p$ ,  $\mathbf{R}(A) = \mathbf{H}_0 = [x_1, \dots, x_p]$  is a  $p$ -dimensional subspace of  $\mathbf{H}$  if and only if the Gram matrix

$$(14) \quad G_x \equiv \| \langle x_i, x_j \rangle \| = \| \langle x, x^T \rangle \|^2$$

is positive definite. It is easily verified that  $A^T A : R^p \rightarrow R^p$  is represented by  $G_x$  and, if  $\mathbf{H}_0$  is  $p$ -dimensional so that  $G_x^{-1} = \| \langle x_i, x_j \rangle \|^{-1}$  exists, then for  $y \in \mathbf{H}$

$$(15) \quad \Pi(y | \mathbf{H}_0) = x^T \| \langle x, x^T \rangle \|^{-1} \langle x, y \rangle$$

and

$$(16) \quad \Pi(y | H_0^\perp) = y - x^T \| \langle x, x^T \rangle \|^{-1} \langle x, y \rangle .$$

These formulae are easily checked by verifying (5): in vector notation

$$\begin{aligned} \langle x, (y - x^T \| \langle x, x^T \rangle \|^{-1} \langle x, y \rangle) \rangle \\ = \langle x, y \rangle - \| \langle x, x^T \rangle \| \| \langle x, x^T \rangle \|^{-1} \langle x, y \rangle \\ = 0, \end{aligned}$$

where 0 denotes a  $p$ -vector of zeros. □

**Proposition 5.** Suppose  $H = H_1 + H_2$ . Let  $a_1, a_2$  denote elements of  $H_1, H_2$  respectively, and suppose that

- (i)  $h \perp H_2,$
- (ii)  $h - a_1 \perp H_1,$
- (iii)  $a_1 - a_2 \perp H_2.$

Then

$$h = \frac{\|a_1\|^2(a_1 - a_2)}{\|a_1 - a_2\|^2} .$$

**Proof.** Let  $H_0 = [a_1, a_2]$  and write  $h = \alpha_1 a_1 + \alpha_2 a_2 + b, b \perp H_0$ .  
By (iii)

$$(a) \quad 0 = \langle a_1 - a_2, a_2 \rangle$$

and by (i)

$$(b) \quad 0 = \langle h, a_2 \rangle = \alpha_1 \langle a_1, a_2 \rangle + \alpha_2 \|a_2\|^2$$

hold. From (a) and (b) we obtain

$$(c) \quad h = \alpha_1(a_1 - a_2) + b, \quad b \perp H_0,$$

and hence by (i) and (iii) also

$$(d) \quad b \perp H_2 .$$

Furthermore, by (ii)

$$(\alpha_1 - 1)a_1 - \alpha_1 a_2 + b \in H_1^\perp$$

or

$$(e) \quad b \in H_1^\perp + H_0 .$$

Together, (c), (d), and (e) yield

$$b \perp H_0 + H_2 + (H_1 \cap H_0^\perp) = H,$$

and

$$(f) \quad b = 0 .$$

Finally, (ii) implies

$$0 = \langle h - a_1, a_1 \rangle$$

and hence by (a), (c), and (f)

$$(g) \quad \alpha_1 = \frac{\|a_1\|^2}{\|a_1\|^2 - \|a_2\|^2} = \frac{\|a_1\|^2}{\|a_1 - a_2\|^2}.$$

Combining (c), (f), and (g) we obtain the proposition.  $\square$

### A.3 CONDITIONAL EXPECTATION FORMULAS

Suppose that  $(X, \mathcal{B}, P)$  is a probability space, and let  $X, Y, Z, \dots$  denote random variables defined thereon. Let  $\mathcal{B}_0$  denote a sub- $\sigma$ -field of  $\mathcal{B}$  and let  $E|X| < \infty$ . The conditional expectation of  $X$  given  $\mathcal{B}_0$ ,  $E(X | \mathcal{B}_0)$ , is any  $\mathcal{B}_0$ -measurable random variable satisfying

$$E[1_A E(X | \mathcal{B}_0)] = E[1_A X] \quad \text{for all } A \in \mathcal{B}_0.$$

In this appendix, we collect a variety of useful formulas concerning conditional expectation, with emphasis on conditional expectations for random variables in  $L_2(P)$ . The basic result is:

**Proposition 1.** Suppose that  $\mathcal{B}_0$  is a sub- $\sigma$ -field of  $\mathcal{B}$  and let  $\mathbf{H}_0 \equiv \{h \in L_2(P) : h \text{ is } \mathcal{B}_0 \text{ measurable}\}$ . Then the projection in  $L_2(P)$  onto  $\mathbf{H}_0$  is given by

$$(1) \quad \Pi_0(h | \mathbf{H}_0) = E[h(X) | \mathcal{B}_0].$$

**Proof.** Let  $h \in L_2(P)$ . Then for  $h_0 \in \mathbf{H}_0$

$$\begin{aligned} E[(h - E(h | \mathcal{B}_0))h_0] &= E\{E[(h - E(h | \mathcal{B}_0))h_0 | \mathcal{B}_0]\} \\ &= E\{h_0 E[(h - E(h | \mathcal{B}_0)) | \mathcal{B}_0]\} \\ &= 0 \end{aligned}$$

by the properties of conditional expectation. Hence the orthogonality condition (A.2.5) holds.  $\square$

The sub- $\sigma$ -field  $\mathcal{B}_0$  of proposition 1 is often determined by invariance considerations. Suppose that  $\mathbf{P}$  is a (dominated) family of probability measures which are invariant with respect to a group  $\mathbf{T}$  of transformations  $T$  on  $\mathbf{X}$ :

$$(2) \quad \mathbf{P} = \{P \in \mathbf{M} : PT^{-1} = P \text{ for all } T \in \mathbf{T}\}.$$

Then

$$\dot{\mathbf{P}} = \{h \in L_2^0(P) : h \text{ is almost invariant under } \mathbf{T}\}$$

where  $h$  is *almost invariant* if and only if for each  $T \in \mathbf{T}$

$$(3) \quad h(Tx) = h(x) \quad \text{a.s. } P;$$

the null set in (3) may depend on  $T$ . Recall  $L_2^0(P) = \{h \in L_2(P) :$

$\int h dP = 0$ ). Under mild conditions every almost invariant function  $h$  is a.s. equal to an *invariant function*; see Lehmann (1986, theorem 6.4, page 297), or theorem 5 of Berk and Bickel (1968). Under the conditions of these theorems it follows that we may take

$$(4) \quad \dot{P} = \{ h \in L_2^0(P) : h \text{ is invariant under } T \} \\ = \{ h \in L_2^0(P) : h \text{ is } \mathcal{B}_T \text{-measurable} \}$$

where the invariant  $\sigma$ -field  $\mathcal{B}_T$  is defined by

$$(5) \quad \mathcal{B}_T \equiv \{ A \in \mathcal{B} : T(A) = A \text{ for all } T \in T \}.$$

**Proposition 2.** Suppose that  $\dot{P}$  is as in (4). Then

$$(6) \quad \Pi_0(h | \dot{P}) = E(h | \mathcal{B}_T) - E(h) \quad \text{for } h \in L_2(P).$$

**Proof.** This follows from (4) and proposition 1. □

**Proposition 3.** Suppose that  $T$  is a compact group of measurable transformations on  $X$  such that  $P$  is invariant under  $T$ :

$$PT^{-1} = P \quad \text{for all } T \in T.$$

Let  $\mathcal{B}_T$  be as in (5). Let  $m$  be right Haar probability measure on  $T$ ,  $m(Tt) = m(T)$  for all  $t \in T$  and  $T$  an open subset of  $T$ ; see, e.g., Cohn (1980, chapter 9). Then, if  $E | h(X) | < \infty$ ,

$$(7) \quad E(h | \mathcal{B}_T)(x) = \int_T h(tx) dm(t) \quad \text{a.s.}$$

**Proof.** Suppose  $w$  is bounded and invariant ( $\mathcal{B}_T$ -measurable). Then

$$(a) \quad E w(X)h(X) = \int \int w(x)h(x) dP(x) dm(t) \\ = \int \int w(tx)h(tx) dP(x) dm(t) \quad (\text{invariance of } P) \\ = \int \int w(x)h(tx) dP(x) dm(t) \quad (\text{invariance of } w) \\ = \int w(x) \left( \int h(tx) dm(t) \right) dP(x) \quad (\text{Fubini}).$$

Fubini and a.s. existence of  $\int_T h(tx) dm(t)$  can be justified by applying (a) with  $w = 1$  and  $|h|$ . The invariance of  $m$  makes  $\int_T h(tx) dm(t)$  invariant. □

Here are two simple examples of the formulas in propositions 1 and 3.

**Example 1. Symmetry.**

Suppose that  $X = R$  with its Borel  $\sigma$ -field  $\mathcal{B}$ , and let

$$\mathcal{B}_0 = \{ h^{-1}(A) : A \in \mathcal{B}, h : R \rightarrow R \text{ symmetric about } 0 \},$$

$H_0 = \{ h \in L_2(P) : h \text{ symmetric} \}$ . Suppose that  $X \sim P$  on  $(R, \mathcal{B})$ . Define  $Q$  by

$$Q(A) \equiv \frac{1}{2} \{ P(X \in A) + P(X \in -A) \} \quad \text{for } A \in \mathcal{B},$$

so that  $P \ll Q$ . Let  $dP/dQ$  denote the Radon-Nikodym derivative. Then for

$h \in L_2(P)$ ,

$$\begin{aligned} \Pi_0(h | \mathbf{H}_0)(x) &= E(h(X) | \mathcal{B}_0)(x) \\ (8) \qquad \qquad \qquad &= \frac{1}{2} \left\{ h(x) \frac{dP}{dQ}(x) + h(-x) \frac{dP}{dQ}(-x) \right\} \quad \text{a.s.} \end{aligned}$$

If  $P$  is symmetric about 0,  $Q = P$ ,  $dP/dQ = 1$ , and (8) simplifies to

$$\begin{aligned} \Pi_0(h | \mathbf{H}_0)(x) &= E(h(X) | \mathcal{B}_0)(x) \\ (9) \qquad \qquad \qquad &= \frac{1}{2} \{h(x) + h(-x)\}. \end{aligned}$$

This second formula also follows from proposition 3 with  $\mathbf{T} = \{T, e\}$  where  $T(x) = -x$  and  $m\{T\} = m\{e\} = 1/2$ . Note that

$$\mathbf{H}_0 = \{h \in L_2(P) : hT = h\} \equiv \mathbf{H}_s$$

and

$$\mathbf{H}_a \equiv \{h \in L_2(P) : hT = -h\}$$

are orthogonal subspaces of  $L_2(P)$  when  $P$  is symmetric about 0. □

**Example 2. Exchangeable random variables.**

Suppose that  $\mathbf{X} = \mathbf{Y} \times \mathbf{Y}$ ,  $\mathcal{B} = \mathcal{A} \times \mathcal{A}$ , and let

$$Tx = T(y_1, y_2) = (y_2, y_1),$$

$$\mathcal{B}_0 \equiv \{B \in \mathcal{B} : TB = B\},$$

$$\mathbf{H}_0 = \{h \in L_2(P) : hT = h\}.$$

Suppose that  $X \sim P$  on  $(\mathbf{X}, \mathcal{B})$ . Define  $Q$  by

$$Q(A) \equiv \frac{1}{2} \{P(X \in A) + P(TX \in A)\} \quad \text{for } A \in \mathcal{B},$$

so that  $P \ll Q$ . Let  $dP/dQ$  denote the Radon-Nikodym derivative. Then for  $h \in L_2^0(P)$

$$\begin{aligned} \Pi_0(h | \mathbf{H}_0) &= E(h | \mathcal{B}_0)(x) \\ (10) \qquad \qquad \qquad &= \frac{1}{2} \left\{ h(x) \frac{dP}{dQ}(x) + h(Tx) \frac{dP}{dQ}(Tx) \right\} \quad \text{a.s.} \end{aligned}$$

If  $P$  is exchangeable, i.e.,  $PT^{-1} = P$ , then  $Q = P$ ,  $dP/dQ = 1$ , and (10) simplifies to

$$\begin{aligned} (11) \qquad \Pi_0(h | \mathbf{H}_0) &= \frac{1}{2} \{h(x) + h(Tx)\} \\ &= \frac{1}{2} \{h(y_1, y_2) + h(y_2, y_1)\}. \end{aligned}$$

This formula also follows from proposition 3 with  $\mathbf{T} = \{T, e\}$  and  $m\{T\} = m\{e\} = \frac{1}{2}$ . Note that

$$\mathbf{H}_0 = \{h \in L_2(P) : hT = h\} \equiv \mathbf{H}_s$$

and

$$\mathbf{H}_a \equiv \{h \in L_2(P) : hT = -h\},$$

are orthogonal subspaces of  $L_2(P)$  and

$$(12) \quad L_2(P) = \mathbf{H}_0 + \mathbf{H}_a. \quad \square$$

**Proposition 4.** Suppose that  $X \sim P \ll \mu$  with  $\mu$   $\sigma$ -finite, that  $T : \mathbf{X} \rightarrow \mathbf{X}'$  is measurable from  $(\mathbf{X}, \mathcal{B})$  to  $(\mathbf{X}', \mathcal{B}')$ , and that  $\mu T^{-1}$  is  $\sigma$ -finite. Then

$$(13) \quad \frac{dPT^{-1}}{d\mu T^{-1}}(t) = E_\mu \left( \frac{dP}{d\mu} \mid T = t \right) \quad \text{a.e. } \mu T^{-1}.$$

**Proof.** Note that the standard theory of conditional expectations can be extended from probability measures to  $\sigma$ -finite measures  $\mu$ ; cf. Loève (1978, page 10). In particular, for  $f$  a  $\mathcal{B}$ -measurable function with  $\int |f| d\mu < \infty$ , the conditional expectation  $E_\mu(f \mid \mathcal{B}_0)$  under  $\mu$  of  $f$  with respect to the sub- $\sigma$ -field  $\mathcal{B}_0$  of  $\mathcal{B}$  is an essentially unique  $\mathcal{B}_0$ -measurable function  $\tilde{f}$  satisfying  $\int_A \tilde{f} d\mu = \int_A f d\mu$  for all  $A \in \mathcal{B}_0$ .

Let  $A \in \mathcal{B}'$ . Then

$$\begin{aligned} PT^{-1}(A) &= P(T(X) \in A) \\ &= \int_{[T(x) \in A]} \frac{dP}{d\mu}(x) d\mu(x) \\ &= E_\mu \left[ 1_{[T \in A]} E_\mu \left( \frac{dP}{d\mu} \mid T \right) \right] \\ &= \int_{[T(x) \in A]} E_\mu \left( \frac{dP}{d\mu} \mid T = T(x) \right) d\mu(x) \\ &= \int_A E_\mu \left( \frac{dP}{d\mu} \mid T = t \right) d\mu T^{-1}(t). \end{aligned} \quad \square$$

**Proposition 5.** Let  $\mathcal{B}_0$  be a sub  $\sigma$ -field of  $\mathcal{B}$ , let  $b$  be a fixed  $r$ -vector of  $L_2(P)$  functions with  $E[bb^T \mid \mathcal{B}_0]$  nonsingular a.s., and let

$$(14) \quad \mathbf{H}_0 \equiv \{[ab] : a \text{ is } \mathcal{B}_0\text{-measurable}\}.$$

Then, for  $h \in L_2(P)$ ,

$$(15) \quad \Pi_0(h \mid \mathbf{H}_0) = E(hb^T \mid \mathcal{B}_0)[E(bb^T \mid \mathcal{B}_0)]^{-1} b.$$

**Proof.** Let  $h \in L_2(P)$  and  $\lambda^T a b \in \mathbf{H}_0$  for some  $\mathcal{B}_0$ -measurable function  $a$  and  $\lambda \in R^r$ . Then

$$\begin{aligned} &E \left[ \left( h - E(hb^T \mid \mathcal{B}_0)[E(bb^T \mid \mathcal{B}_0)]^{-1} b \right) b^T \lambda a \right] \\ &= EE \left[ \left( h - E(hb^T \mid \mathcal{B}_0)[E(bb^T \mid \mathcal{B}_0)]^{-1} b \right) b^T \lambda a \mid \mathcal{B}_0 \right] \end{aligned}$$

$$\begin{aligned}
 &= E \left[ a \left( E(hb^T | \mathcal{B}_0) - E(hb^T | \mathcal{B}_0) [E(bb^T | \mathcal{B}_0)]^{-1} E(bb^T | \mathcal{B}_0) \right) \lambda \right] \\
 &= 0.
 \end{aligned}$$

Hence the orthogonality condition (A.2.5) holds. □

*Conditional Expectations and Martingales*

Now suppose that the basic probability space  $(X, \mathcal{B}, P)$  has an increasing family of sub- $\sigma$ -fields  $\{\mathcal{B}_t\}_{t \geq 0}$  of  $\mathcal{B}$ . It is well known that if  $Y \in L_1(P)$ , then the process  $\{Y_t\}_{t \geq 0}$  defined by

$$(16) \quad Y_t \equiv E\{X | \mathcal{B}_t\} \quad \text{for } t \geq 0$$

is a uniformly integrable martingale with respect to the  $\sigma$ -fields  $\{\mathcal{B}_t\}$ . The process  $\{Y_t\}$  is sometimes called *Doob's martingale*; see, e.g., Karlin and Taylor (1975, page 246) for the discrete time version of this, and Elliott (1982, page 36) for the continuous time result. Conversely, every uniformly integrable martingale can be written as the conditional expectation of some integrable function  $a$  as in (16)); see, e.g., Elliott (1982, page 36), or Liptser and Shirayev (1977, theorem 2.7, page 45).

This structure is often present in problems involving censoring for a particular choice of the  $\sigma$ -fields  $\{\mathcal{B}_t\}$ , as follows. Let  $X \sim F$  be defined on a probability space  $(X, \mathcal{B}, P)$ , and consider the filtration  $\{\mathcal{B}_t\}_{t \geq 0}$  defined by

$$(17) \quad \mathcal{B}_t \equiv \sigma\{1_{[X \leq s]} : s \leq t\} = \sigma\{X \wedge t, 1_{[X \leq t]}\}$$

for  $t \geq 0$ . Since we are most interested in score functions, i.e., elements of  $L_2(F)$ , we start with a fixed function  $a \in L_2(F)$ , and then define, for  $t \geq 0$ ,

$$(18) \quad \mathbf{Y}(t) \equiv E\{a(X) | \mathcal{B}_t\}.$$

Thus  $\{\mathbf{Y}(t)\}$  is the Doob martingale corresponding to  $a(X)$  and the particular family  $\{\mathcal{B}_t\}$  of (17). For our present particular filtration  $\{\mathcal{B}_t\}$  it is easily shown that

$$\begin{aligned}
 (19) \quad \mathbf{Y}(t) &= 1_{[X \leq t]} a(X) + 1_{[X > t]} E(a(X) | X > t) \\
 &\equiv 1_{[X \leq t]} a(X) + 1_{[X > t]} A(t).
 \end{aligned}$$

Now we want to relate the martingale  $\mathbf{Y}$  to the counting process martingale

$$\begin{aligned}
 (20) \quad \mathbf{M}(t) &\equiv N(t) - \mathbf{A}(t) \\
 &= 1_{[X \leq t]} - \int_{-\infty}^t 1_{[X \geq s]} d\Lambda(s)
 \end{aligned}$$

defined in formula (A.1.21) and in example 3.4.2 (method 2). The connection involves the  $R$  operator introduced in section A.1. The following proposition is apparently due to Chou and Meyer (1974), (1975): see theorem 2 and formulas (9) and (10) of Chou and Meyer (1974, page 1563); and see proposition 2 and formula (13) of Chou and Meyer (1975, pages 231–232). (The  $(q, M, H, h)$  of Chou and Meyer (1975) correspond to our  $(M, \mathbf{Y}, a, Ra)$ . Their result is



somewhat more general in that they only require  $a \in L_1(F)$ .) For more on related martingale representation theorems see Liptser and Shirayev (1978, chapter 19); or see Shorack and Wellner (1986, appendix B and chapter 6). We will use this proposition in connection with proposition A.5.5 to compute score functions in models involving censoring.

**Proposition 6.** Suppose that  $a \in L_2(F)$  and  $E a(X) = 0$ . Then the martingale  $\mathbf{Y}(t)$  in (18) and (19) is related to the counting process martingale  $\mathbf{M}$  of (20) by

$$(21) \quad \mathbf{Y}(t) = \int_{-\infty}^t Ra(s) d\mathbf{M}(s),$$

where, as in (A.1.19),

$$(22) \quad Ra(t) \equiv a(t) - \frac{\int_{-\infty}^t a dF}{1 - F(t)}.$$

**Proof.** First write  $\mathbf{M} = \mathbf{N} - \mathbf{A}$  and note that since  $Ra = a - A$ ,

$$(a) \quad \int_{-\infty}^t Ra d\mathbf{N} = a(X)1_{[X \leq t]} - A(X)1_{[X \leq t]}.$$

But then

$$(b) \quad \begin{aligned} \int_{-\infty}^t Ra d\mathbf{A} &= \int_{-\infty}^{t \wedge X} Ra d\Lambda \\ &= -A(X)1_{[X \leq t]} - A(t)1_{[X > t]} \end{aligned}$$

since, by a calculation using Fubini's theorem (justified in view of the boundedness of  $R$  established in proposition A.1.8.A) and  $\int a dF = 0$  twice,

$$\begin{aligned} \int_{-\infty}^u Ra d\Lambda &= \int_{-\infty}^u a d\Lambda \\ &\quad - \int_{-\infty}^u \left\{ \frac{1}{1 - F(s)} \int_s^\infty a dF \right\} d\Lambda(s) \\ &= -A(u). \end{aligned}$$

Subtracting (a) and (b) yields (21). □

**Corollary 1.** The predictable variation process of the martingale transform  $\mathbf{Y}$  is

$$(23) \quad \langle \mathbf{Y} \rangle(t) = \int_{-\infty}^t (Ra)^2(s) 1_{[X \geq s]} d\Lambda(s)$$

so that both

$$(24) \quad E[\mathbf{Y}^2(t)] = E\langle \mathbf{Y} \rangle(t) = \int_{-\infty}^t (Ra)^2 dF$$

from (23) and

$$(25) \quad E[\mathbf{Y}^2(t)] = \int_{-\infty}^t a^2 dF + A^2(t)(1 - F(t))$$

directly from (19).

**Proof.** (23) follows immediately from (21) and standard martingale theory; see, e.g., Shorack and Wellner (1986, formula (B.3.2), page 891).  $\square$

### A.4 PROJECTION ON SUMSPACES AND ACE

#### *Projection on the Sum of Two Subspaces $\mathbf{H}_1 + \mathbf{H}_2$*

Suppose that  $\mathbf{H}_1, \mathbf{H}_2$  are closed subspaces of a Hilbert space  $\mathbf{H}$  with inner product  $\langle \cdot, \cdot \rangle$  and norm  $\|\cdot\|$ . Throughout this section we write  $h_n \rightarrow h$  to mean convergence in the norm of  $\mathbf{H}$ :  $\|h_n - h\| \rightarrow 0$ .

If  $\mathbf{H}_1 \perp \mathbf{H}_2$ , then  $\mathbf{H}_1 \cap \mathbf{H}_2 = \{0\}$ ,  $\mathbf{H}_1 + \mathbf{H}_2$  is closed, and projection onto the sumspace  $\mathbf{H}_1 + \mathbf{H}_2$  is simply the sum of the projections onto  $\mathbf{H}_1$  and  $\mathbf{H}_2$ ; recall (A.2.9). We now consider projection onto  $\overline{\mathbf{H}_1 + \mathbf{H}_2}$ , the closure of  $\mathbf{H}_1 + \mathbf{H}_2$  without requiring that  $\mathbf{H}_1 \perp \mathbf{H}_2$ ,  $\mathbf{H}_1 \cap \mathbf{H}_2 = \{0\}$ , or that  $\mathbf{H}_1 + \mathbf{H}_2$  is closed. The basic theorem is due to Von Neumann. His lecture notes from 1933 are published in Von Neumann (1950).

Let  $P_i \equiv \Pi(\cdot | \mathbf{H}_i)$ , for  $i = 1, 2$ , be the orthogonal projection operators onto  $\mathbf{H}_i$ ,  $i = 1, 2$ , and write

$$(1) \quad Q_i \equiv I - P_i = \Pi(\cdot | \mathbf{H}_i^\perp), \quad i = 1, 2,$$

for the projection operator onto  $\mathbf{H}_i^\perp$ ,  $i = 1, 2$ . Also write  $Q \equiv \Pi(\cdot | \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)$ .

**Theorem 1.** (Von Neumann). For any  $h \in \mathbf{H}$

$$(2) \quad \|(Q_1 Q_2)^m h - Q h\| \rightarrow 0 \quad \text{as } m \rightarrow \infty$$

and

$$(3) \quad \|(Q_2 Q_1)^m h - Q h\| \rightarrow 0 \quad \text{as } m \rightarrow \infty.$$

All the proofs are collected at the end of this section.

Figure 1 is the geometric picture behind theorem 1.

Since

$$(4) \quad (\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)^\perp = \overline{\mathbf{H}_1 + \mathbf{H}_2},$$

Von Neumann's theorem immediately yields the following corollary.

**Corollary 1.** For any  $h \in \mathbf{H}$ ,

$$(5) \quad [I - (Q_1 Q_2)^m] h \rightarrow \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2}) \quad \text{as } m \rightarrow \infty,$$

and similarly for (3).

The convergence in (5) entails two complications: The first complication results from the fact that  $\mathbf{H}_1 + \mathbf{H}_2 \subset \overline{\mathbf{H}_1 + \mathbf{H}_2}$  with the possibility of strict containment; see, e.g., Halmos (1982, page 29), or Kober (1939) for examples. We give yet another set of simple examples below. The second complication is that  $\mathbf{H}_1 \cap \mathbf{H}_2 \neq \{0\}$  is possible, so that (even if  $\mathbf{H}_1 + \mathbf{H}_2$  is closed) a decomposition of  $h \in \mathbf{H}_1 + \mathbf{H}_2$  as  $h_1 + h_2$  with  $h_i \in \mathbf{H}_i$ ,  $i = 1, 2$ , is not unique. To

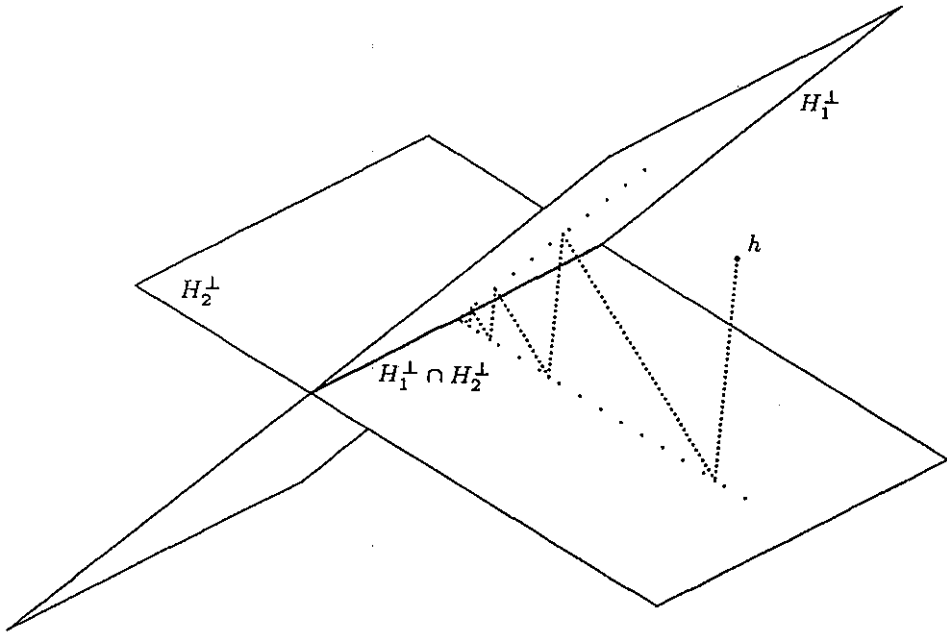


FIGURE 1. Alternating projections.

address these complications, and to give another perspective on (5), it is convenient to introduce the following “backfitting” procedure: set

$$h_2^{(1)} \equiv P_2 h ,$$

and proceed inductively: for  $m \geq 1$ , set

$$h_1^{(m)} \equiv P_1(h - h_2^{(m)})$$

(6) 
$$h_2^{(m+1)} \equiv P_2(h - h_1^{(m)}) .$$

When the projection operators  $P_i$  are conditional expectation operators, this is just the “inner loop” of the (population version of the) *alternating conditional expectation* or ACE algorithm of Breiman and Friedman (1985).

It follows immediately from (6), by using

$$(I - P_1)P_1 = 0 = (I - P_2)P_2 ,$$

that, for  $m \geq 2$ ,

$$\begin{aligned} h - h_1^{(m)} - h_2^{(m)} &= (I - P_1)(h - h_2^{(m)}) \\ &= (I - P_1)(h - h_1^{(m-1)} - h_2^{(m)}) \\ &= (I - P_1)(I - P_2)(h - h_1^{(m-1)}) \\ &= (I - P_1)(I - P_2)(h - h_1^{(m-1)} - h_2^{(m-1)}) , \end{aligned}$$

and that

$$h - h_1^{(1)} - h_2^{(1)} = (I - P_1)(I - P_2)h .$$

Hence, by induction and (1),

$$(7) \quad h - h_1^{(m)} - h_2^{(m)} = (Q_1 Q_2)^m h, \quad m \geq 1.$$

Comparison of (5) and (7) suggests that if  $H_1 + H_2$  is closed and  $H_1 \cap H_2 = \{0\}$ , so that  $\Pi(h | \overline{H_1 + H_2}) = h_1^* + h_2^*$  uniquely, then  $h_i^{(m)} \rightarrow h_i^*$  as  $m \rightarrow \infty$ ,  $i = 1, 2$ . This leads to the following reformulation and strengthening of theorem 1.

**Theorem 2.**

- A.  $\|h_1^{(m)} + h_2^{(m)} - \Pi(h | \overline{H_1 + H_2})\| \rightarrow 0$  as  $m \rightarrow \infty$ .
- B.  $\|h_1^{(m)} + h_2^{(m)} - \Pi(h | \overline{H_1 + H_2})\| \leq \rho^{2(m-1)} \|\Pi(h | \overline{H_1 + H_2})\|$   
 where  $\rho$  is the cosine of the minimum angle  $\tau$  between  $H_1$  and  $H_2$  (relative to  $(H_1 \cap H_2)^\perp$ ):

$$(8) \quad \rho \equiv \rho(H_1, H_2) \\ \equiv \sup \{ \langle h_1, h_2 \rangle : h_i \in H_i \cap (H_1 \cap H_2)^\perp, \|h_i\| \leq 1 \} \\ (9) \quad = \sup \{ \langle h_1, h_2 \rangle : h_i \in H_i^\perp \cap (H_1^\perp \cap H_2^\perp)^\perp, \|h_i\| \leq 1 \} \\ \equiv \rho(H_1^\perp, H_2^\perp) \\ < 1 \quad \text{if and only if } H_1 + H_2 \text{ is closed.}$$

- C. Suppose that  $\Pi(h | \overline{H_1 + H_2}) = h_1^* + h_2^*$  with  $h_1^* \in H_1 \cap (H_1 \cap H_2)^\perp$  and  $h_2^* \in H_2$ . Then

$$(10) \quad \|h_i^{(m)} - h_i^*\| \rightarrow 0, \quad i = 1, 2.$$

If  $H_1 + H_2$  is closed, then the convergence in (10) is geometric with the same rate as in B.

$$D. \quad \|(Q_2 Q_1)^m - Q\| = \rho^{2m-1}.$$

**Remark 1.** Part A of theorem 2 is, by (7), simply a restatement of corollary 1 of Von Neumann's theorem 1. Theorem 1 (and part A of theorem 2) was first proved by Von Neumann in his lectures on operator theory at Princeton in 1933-1934, and then later independently rediscovered by others, including Aronszajn (1950), Nakano (1953), and Wiener (1955).

**Remark 2.** The measure  $\rho \equiv \cos \tau$  of the angle  $\tau$  between  $H_1$  and  $H_2$  was apparently first introduced by Friedrichs (1937). Part B of theorem 2 was proved by Aronszajn (1950, section 12, pages 375-380). He also stated the equality of angles (9), but did not give a proof. Kato (1976, page 221) gives a proof of (9) in a more general setting.

**Remark 3.** Aronszajn (1950) proved  $\|(Q_2 Q_1)^m - Q\| \leq \rho^{2m-1}$ ; this was refined to give the equality in part D of theorem 2 by Kayalar and Weinert (1988). Other less sharp bounds were proved by Franchetti and Light (1986) ( $\rho^m$  rather than  $\rho^{2(m-1)}$ ) and Deutsch (1983) ( $(1 - \rho^2)^{-1} \rho^{2m-1} \|h - P_1 P_2 h\|$ ). Franchetti and Light (1986) show that the convergence in A can be arbitrarily slow if  $H_1 + H_2$  is not closed; and if  $\Pi(h | \overline{H_1 + H_2}) \notin H_1 + H_2$ , then

$\|h_1^{(m)}\| \rightarrow \infty$  as  $m \rightarrow \infty$ .

**Remark 4.** If the roles of  $P_1$  and  $P_2$  are reversed in the definition of the backfitting algorithm (6), then the sequences  $h_1^{(m)}$  and  $h_2^{(m)}$  in  $\mathbf{C}$  converge to  $h_1^* \in \mathbf{H}_1$  and  $h_2^* \in \mathbf{H}_2 \cap (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp$  respectively.

**Remark 5.** When the  $\mathbf{H}_i$  are subspaces of  $L_2(P)$  for a probability measure  $P$  (see (16)–(19) below), the convergence in norm (in  $L_2(P)$ ) can be replaced by almost sure convergence: see, e.g., Burkholder and Chow (1961), Burkholder (1962), and Rota (1962).

**Remark 6.** Suppose that  $\mathbf{A}_i, i = 1, 2$ , are closed subspaces of a Banach space  $\mathbf{B}$ . For a general theorem concerning existence of a continuous projection (satisfying  $\Pi_1^2 = \Pi_1$ )  $\Pi_1$  with range  $\mathbf{A}_1$  and null space  $\mathbf{A}_2$ , see Rudin (1973, page 126, theorem 5.16). (In theorem 5.16, note that the hypothesis that  $\mathbf{X} = \mathbf{A}_1 + \mathbf{A}_2$  is an  $F$ -space implies that  $\mathbf{X}$  is closed by Rudin (1973, page 20, theorem 1.20).)

Suppose that  $\mathbf{H}_1 + \mathbf{H}_2$  is closed. It is often useful to think of the projection  $h^* \equiv \Pi(h \mid \mathbf{H}_1 + \mathbf{H}_2)$  in terms of a pair of equations satisfied by the components of any decomposition  $h^* = h_1^* + h_2^* \in \mathbf{H}_1 + \mathbf{H}_2$  with  $h_i^* \in \mathbf{H}_i$ .

**Proposition 1.** Suppose that  $\mathbf{H}_1 + \mathbf{H}_2$  is closed, and let  $h \in \mathbf{H}$ ,  $h_i^* \in \mathbf{H}_i, i = 1, 2$ . Then

$$(11) \quad \Pi(h \mid \mathbf{H}_1 + \mathbf{H}_2) = h_1^* + h_2^*$$

if and only if  $h - h_1^* - h_2^* \in \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp$ , if and only if  $h_1^*$  and  $h_2^*$  satisfy the two equations

$$(12) \quad h_1^* = P_1(h - h_2^*)$$

and

$$(13) \quad h_2^* = P_2(h - h_1^*).$$

Substitution of (12) into (13), and vice versa, shows that the coupled equations (12) and (13) can be rewritten in terms of the uncoupled equation

$$(14) \quad h_2^* = P_2 Q_1 h + P_2 P_1 h_2^*,$$

and then  $h_2^*$  is determined by (12). Alternatively, solve the uncoupled equation

$$(15) \quad h_1^* = P_1 Q_2 h + P_1 P_2 h_1^*$$

for  $h_1^*$ , and then  $h_2^*$  is determined by (13).

Frequently  $\mathbf{X} = (X_1, X_2), \mathbf{H}_i = \{b(X_i) \in L_2(P)\}$ , and

$$P_i \equiv \Pi(\cdot \mid \mathbf{H}_i) = E(\cdot \mid X_i).$$

Then, with  $h \equiv h(X_1, X_2)$  having  $E h^2(X_1, X_2) < \infty$ , (12) and (13) can be rewritten as

$$(16) \quad h_1^*(X_1) = E\{h(X_1, X_2) - h_2^*(X_2) \mid X_1\}$$

and

$$(17) \quad h_2^*(X_2) = E\{h(X_1, X_2) - h_1^*(X_1) \mid X_2\}.$$

Moreover, the iterative algorithm (6) becomes

$$(18) \quad h_1^{(m)}(X_1) = E(h(X_1, X_2) - h_2^{(m)}(X_2) \mid X_1),$$

$$h_2^{(m+1)}(X_2) = E(h(X_1, X_2) - h_1^{(m)}(X_1) \mid X_2),$$

and this is what has been called the *alternating conditional expectation* or ACE, algorithm by Breiman and Friedman (1985). For this reason, we introduce and use the notation

$$(19) \quad \begin{aligned} ACE(h) &\equiv ACE(h \mid X_1, X_2) \equiv h_1^*(X_1) + h_2^*(X_2) \\ &= \Pi(h \mid \mathbf{H}_1 + \mathbf{H}_2) \end{aligned}$$

where  $h_1^*$  and  $h_2^*$  are not well-defined unless  $\mathbf{H}_1 + \mathbf{H}_2$  is closed, and not unique unless  $\mathbf{H}_1 \cap \mathbf{H}_2 = \{0\}$  so that  $\mathbf{H}_1 + \mathbf{H}_2 = \mathbf{H}_1 \oplus \mathbf{H}_2$ .

If  $P$  is a Hilbert-Schmidt operator, then it is compact; see, e.g., Reed and Simon (1972, theorem VI.22, page 210). If  $P_1, P_2$  are conditional expectations and if  $X_1, X_2$  have joint density  $f(x_1, x_2)$  and marginal densities  $f_1(x_1), f_2(x_2)$ , then for a function  $h(X_2)$

$$\begin{aligned} (P_1 h)(x_1) &= E(h(X_2) \mid X_1 = x_1) \\ &= \int h(x_2) \frac{f(x_1, x_2)}{f_1(x_1)} dx_2 \\ &= \int h(x_2) \frac{f(x_1, x_2)}{f_1(x_1)f_2(x_2)} f_2(x_2) dx_2 \\ &\equiv E(h(X_2)K(x_1, X_2)) \end{aligned}$$

with

$$K(x_1, x_2) \equiv \frac{f(x_1, x_2)}{f_1(x_1)f_2(x_2)}.$$

Thus  $P_1 : \mathbf{H}_2 \rightarrow \mathbf{H}_1$  is Hilbert-Schmidt if

$$\iint K^2(x_1, x_2) f_1(x_1)f_2(x_2) dx_1 dx_2 = \iint \frac{f^2(x_1, x_2)}{f_1(x_1)f_2(x_2)} dx_1 dx_2 < \infty;$$

see, e.g., theorem VI.23, of Reed and Simon (1972).

The following proposition gives several useful facts concerning closure of  $\mathbf{H}_1 + \mathbf{H}_2$ . For any operator  $P$  and subspace  $\mathbf{H}_0$  of  $\mathbf{H}$ , let  $P \mid \mathbf{H}_0$  denote the restriction of  $P$  to  $\mathbf{H}_0$ .

**Proposition 2.**

- A.  $\mathbf{H}_1 + \mathbf{H}_2$  is closed if and only if there is a constant  $c > 0$  such that for every  $h \in \mathbf{H}_1 + \mathbf{H}_2$ , there is a decomposition  $h = h_1 + h_2$  with  $h_i \in \mathbf{H}_i$ ,  $i = 1, 2$ , such that

$$\|h_1 + h_2\| \geq c \max \{ \|h_1\|, \|h_2\| \} .$$

- B. If either  $P_1 | \mathbf{H}_2$  or  $P_2 | \mathbf{H}_1$  is compact, then  $\mathbf{H}_1 + \mathbf{H}_2$  is closed.
- C. A sufficient condition for compactness of  $P_1 | \mathbf{H}_2$  is that  $P_1 | \mathbf{H}_2$  is Hilbert-Schmidt.
- D.  $\rho$  defined in (8) is  $< 1$  if and only if  $\mathbf{H}_1 + \mathbf{H}_2$  is closed.

Here is a simple example showing that  $\mathbf{H}_1 + \mathbf{H}_2$  is not necessarily closed.

**Example 1.** Suppose that  $U, W \sim \text{Uniform}(0, 1)$  are independent. Let  $\delta$  be a Bernoulli random variable with

$$P(\delta = 1 | U = u) = u, \quad 0 \leq u \leq 1 .$$

Set

$$V = U + \delta W .$$

Let  $P$  be the joint distribution of  $(U, V)$  on  $[0, 1] \times [0, 2]$ , and let

$$\mathbf{H}_1 \equiv \{a(U) : E a(U) = 0, E a^2(U) < \infty\} ,$$

$$\mathbf{H}_2 \equiv \{b(V) : E b^2(V) < \infty\} .$$

It is easily seen that  $\mathbf{H}_1 \cap \mathbf{H}_2 = \{0\}$ . To show that  $\mathbf{H}_1 + \mathbf{H}_2$  is not closed, we exhibit sequences  $\{h_i^{(n)}\}$ ,  $i = 1, 2$ , such that  $\|h_1^{(n)} + h_2^{(n)}\| \rightarrow 0$ , but  $\|h_1^{(n)}\| \rightarrow 1$  as  $n \rightarrow \infty$ . Thus part A of proposition 2 fails, and  $\mathbf{H}_1 + \mathbf{H}_2$  is not closed.

For  $n \geq 1$ , let

$$h_n(u) \equiv \sqrt{n} \left( 1_{[0, 1/n]}(u) - \frac{1}{n} \right), \quad 0 \leq u \leq 1 ,$$

and set

$$h_1^{(n)}(U) \equiv h_n(U) ,$$

$$h_2^{(n)}(V) \equiv -h_n(V) .$$

Then

$$E h_1^{(n)}(U) = 0, \quad E [h_1^{(n)}(U)]^2 = 1 - \frac{1}{n} ,$$

and

$$\begin{aligned} E [h_1^{(n)}(U) + h_2^{(n)}(V)]^2 &= E \{ \delta [h_n(U) - h_n(U+W)]^2 \} \\ &\leq E \{ \delta n 1_{[U \leq n^{-1} < U+W]} \} \leq E \{ U n 1_{[U \leq n^{-1}]} \} \\ &\leq P \left( U \leq \frac{1}{n} \right) = \frac{1}{n} \rightarrow 0 \quad \text{as } n \rightarrow \infty . \end{aligned}$$

Here is a function  $h(U, V) \in \overline{\mathbf{H}_1 + \mathbf{H}_2} - (\mathbf{H}_1 + \mathbf{H}_2)$ : let

$$h(U, V) = (U^{-1/2} - 2) - (V^{-1/2} - 2) .$$

Then  $E(U^{-1/2} - 2) = 0, E(U^{-1/2} - 2)^2 = \infty$ , so  $h(U, V) \neq h_1(U) + h_2(V)$

with  $h_i \in \mathbf{H}_i, i = 1, 2$ . But

$$\begin{aligned} E[h^2(U, V)] &= E(U^{-1/2} - V^{-1/2})^2 \\ &= E\{\delta(U^{-1/2} - (U + W)^{-1/2})^2\} \\ &\leq E\{\delta U^{-1}\} = 1 < \infty. \end{aligned}$$

While the joint distribution of  $(U, V)$  in this example is not absolutely continuous with respect to Lebesgue measure on the plane, it is easy to construct similar examples in which  $(U, V)$  does have a density. For example, let  $U$  and  $W$  be as before, and set  $V \equiv U + U^2 W$ . Then with  $P$  the distribution of  $(U, V)$ , the sum space  $\mathbf{H}_1 + \mathbf{H}_2$  again fails to be closed by a similar calculation.  $\square$

Another simple example is given by Rényi (1959): take  $P$  to be the uniform distribution on the region shown in Figure 2 (where the two curves are both tangent to the diagonal at 0).

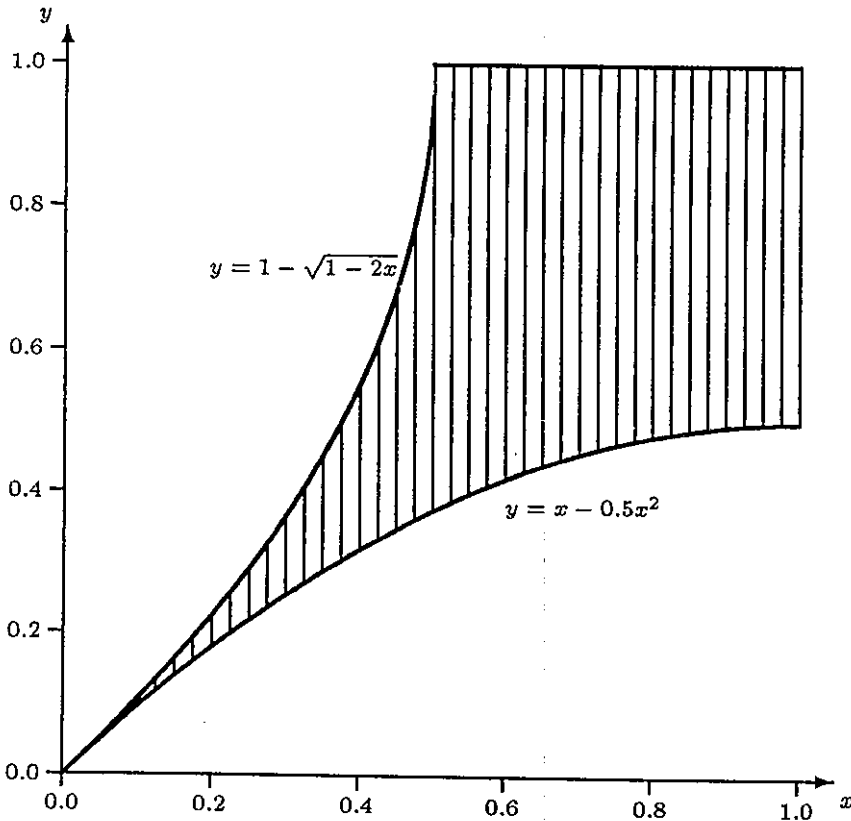


FIGURE 2. Rényi's distribution with  $\rho = 1$ .

*Projection on the Sum of  $r \geq 2$  subspaces  $\mathbf{H}_1 + \dots + \mathbf{H}_r$*

Now suppose that  $\mathbf{H}_1, \dots, \mathbf{H}_r$  are closed subspaces of  $\mathbf{H}$ , and let  $P_i \equiv \Pi(\cdot | \mathbf{H}_i), i = 1, \dots, r$  be the orthogonal projection operators onto  $\mathbf{H}_i$ , and  $Q_i \equiv I - P_i = \Pi(\cdot | \mathbf{H}_i^\perp), i = 1, \dots, r$ .



**Theorem 3.** (Halperin). For any  $h \in \mathbf{H}$ ,

$$(20) \quad \|(Q_1 \dots Q_r)^m h - \Pi(h | \mathbf{H}_1^\perp \cap \dots \cap \mathbf{H}_r^\perp)\| \rightarrow 0 \quad \text{as } m \rightarrow \infty.$$

Or, with  $T \equiv Q_1 \dots Q_r$  and  $\mathbf{H}_+^\perp \equiv \mathbf{H}_1^\perp \cap \dots \cap \mathbf{H}_r^\perp$ ,

$$\|T^m h - \Pi(h | \mathbf{H}_+^\perp)\| \rightarrow 0 \quad \text{as } m \rightarrow \infty.$$

Since

$$(21) \quad (\mathbf{H}_1^\perp \cap \dots \cap \mathbf{H}_r^\perp)^\perp = \overline{\mathbf{H}_1 + \dots + \mathbf{H}_r} \equiv \mathbf{H}_+,$$

Halperin's theorem immediately yields the following corollary:

**Corollary 2.** For any  $h \in \mathbf{H}$ ,

$$(22) \quad [I - T^m]h \rightarrow \Pi(h | \mathbf{H}_+) \quad \text{as } m \rightarrow \infty.$$

It is easy to extend the backfitting algorithm (6) to  $r > 2$ . Set  $h_i^{(0)} = 0$ ,  $i = 1, \dots, r$ , and proceed inductively as follows: for  $m \geq 1$ , set

$$(23) \quad h_i^{(m)} \equiv P_i(h - \sum_{j=1}^{i-1} h_j^{(m-1)} - \sum_{j=i+1}^r h_j^{(m)}), \quad i = r, \dots, 1,$$

(where  $\sum_{j=r}^s a_j \equiv 0$  if  $r > s$ ). Then, just as in the case  $r = 2$ , we obtain

$$(24) \quad h - h_1^{(m)} - \dots - h_r^{(m)} = (Q_1 \dots Q_r)^m h, \quad m \geq 1.$$

**Theorem 4.** Let  $h \in \mathbf{H}$ . Then:

A.  $\|\sum_{i=1}^r h_i^{(m)} - \Pi(h | \mathbf{H}_+)\| \rightarrow 0 \quad \text{as } m \rightarrow \infty.$

B.  $\|\sum_{i=1}^r h_i^{(m)} - \Pi(h | \mathbf{H}_+)\| \leq \rho_+^m \|\Pi(h | \mathbf{H}_+)\|$

where

$$\rho_+^2 \leq 1 - \prod_{j=1}^{r-1} \sin^2 \tau_j$$

and

$$\cos \tau_j \equiv \rho(\mathbf{H}_j^\perp, \mathbf{H}_{j+1}^\perp \cap \dots \cap \mathbf{H}_r^\perp), \quad j = 1, \dots, r-1.$$

C. If  $\mathbf{H}_1 + \dots + \mathbf{H}_r$  is closed and  $\|\sum_{i=1}^r h_i\| = 0$  with  $h_i \in \mathbf{H}_i$ ,  $i = 1, \dots, r$  implies  $\|h_i\| = 0$ ,  $i = 1, \dots, r$ , then

$$(25) \quad \|h_i^{(m)} - h_i^*\| \rightarrow 0 \quad \text{as } m \rightarrow \infty,$$

where  $h_i^*$ ,  $i = 1, \dots, r$  are the components of  $h^* \equiv h_1^* + \dots + h_r^* \equiv \Pi(h | \mathbf{H}_+)$ .

**Remark 7.** Part A of theorem 4 is, by (24), simply a restatement of corollary 2 to Halperin's theorem 3.

**Remark 8.** Part B of theorem 4 is due to Smith, Solomon, and Wagner (1977); see Kayalar and Weinert (1988) for further bounds.

*Finite-dimensional  $\mathbf{H}_1$  or  $\mathbf{H}_2$* 

We now specialize to the case when one of the subspaces  $\mathbf{H}_1, \mathbf{H}_2$  is finite-dimensional. Let  $\mathbf{H}_1, \mathbf{H}_2$  be subspaces of  $\mathbf{H}$  such that

$$\mathbf{H}_2 = [b_1, \dots, b_k]$$

and, as before, let  $P_1 \equiv \Pi(\cdot | \mathbf{H}_1)$ .

**Theorem 5.** Suppose, without loss of generality, that

$$[b_1 - P_1 b_1, \dots, b_k - P_1 b_k] = [b_1 - P_1 b_1, \dots, b_d - P_1 b_d]$$

and

$$\dim [b_1 - P_1 b_1, \dots, b_k - P_1 b_k] = d \leq k.$$

Then

$$\begin{aligned} (26) \quad \Pi(h | \mathbf{H}_1 + \mathbf{H}_2) &= P_1 h + \sum_{j=1}^d \lambda_j(h) (b_j - P_1 b_j) \\ &= (P_1 h - \sum_{j=1}^d \lambda_j(h) P_1 b_j) + \sum_{j=1}^d \lambda_j(h) b_j, \end{aligned}$$

where  $\lambda(h) = (\lambda_1(h), \dots, \lambda_d(h))^T$  is given by

$$(27) \quad \lambda(h) = B^{-1} (\langle b_1 - P_1 b_1, h \rangle, \dots, \langle b_d - P_1 b_d, h \rangle)^T$$

with

$$(28) \quad B = \left( \langle b_i - P_1 b_i, b_j - P_1 b_j \rangle \right)_{d \times d}.$$

**Remark 9.** Note that  $d = k$  if and only if  $\mathbf{H}_1 \cap \mathbf{H}_2 = \{0\}$  and  $\dim [b_1, \dots, b_k] = k$ . Suppose that  $d = k$  and define

$$\begin{aligned} B_1 &\equiv \left( \langle P_1 b_i, P_1 b_j \rangle \right)_{k \times k}, \\ B &\equiv \left( \langle b_i - P_1 b_i, b_j - P_1 b_j \rangle \right)_{k \times k}, \end{aligned}$$

and

$$B_0 \equiv \left( \langle b_i, b_j \rangle \right)_{k \times k}.$$

Then

$$B = B_0 - B_1 = B_0(I - B_0^{-1} B_1)$$

and

$$B^{-1} = (I - B_0^{-1} B_1)^{-1} B_0^{-1}.$$

Now suppose that

$$\mathbf{H} \equiv \{h(Y, Z) : E h^2(Y, Z) < \infty\},$$

(with the usual inner product  $\langle h_1, h_2 \rangle = E h_1(Y, Z) h_2(Y, Z)$ ),

$$\mathbf{H}_1 \equiv \{a(Y) : E a^2(Y) < \infty, E a(Y) = 0\},$$

$$\mathbf{H}_2 \equiv \{b(Z) : E b^2(Z) < \infty, E b(Z) = 0\},$$

where  $Z$  takes values in  $\{z_1, \dots, z_k\} = \{1, \dots, k\}$  without loss of generality. Then  $\mathbf{H}_2$  is finite-dimensional with basis

$$b_j \equiv 1_{[Z=j]} - P(Z=j), \quad \text{for } j = 1, \dots, k,$$

and  $P_1 = E(\cdot | Y)$ . Furthermore

$$\begin{aligned} d &= \dim [b_1 - P_1 b_1, \dots, b_k - P_1 b_k] \\ &= \dim [\{1_{[Z=j]} - P(Z=j|Y)\}, j = 1, \dots, k] \\ &\leq k-1 \end{aligned}$$

since

$$(29) \quad \sum_{j=1}^k \left\{ 1_{[Z=j]} - P(Z=j|Y) \right\} = 0.$$

**Proposition 3.**  $d = k-1$  if and only if there exists no nonconstant function  $q : \{1, \dots, k\} \rightarrow R$  such that

$$(30) \quad P(q(Z) = a(Y)) = 1 \quad \text{for some } a,$$

or equivalently

$$(31) \quad q(Z) = E(q(Z) | Y) \quad \text{a.s.}$$

If the condition of proposition 3 holds, then, using the notation of (19), define the right side of (26) to be

$$\begin{aligned} (32) \quad ACE(h | Y, Z) &\equiv E(h | Y) - Eh + \sum_{j=1}^{k-1} \lambda_j(h) \left\{ 1_{[Z=j]} - P(Z=j|Y) \right\} \\ &\equiv ACE_1(h | Y, Z) + ACE_2(h | Y, Z), \end{aligned}$$

where

$$(33) \quad ACE_1(h | Y, Z) \equiv E(h | Y) - Eh - \sum_{j=1}^{k-1} \lambda_j(h) E(b_j | Y),$$

$$(34) \quad ACE_2(h | Y, Z) \equiv \sum_{j=1}^{k-1} \lambda_j(h) b_j,$$

and

$$(35) \quad \lambda = B^{-1} c,$$

with  $B = (B_{ij})_{(k-1) \times (k-1)}$  where

$$(36) \quad B_{ij} = E\{P(Z=i|Y)(\delta_{ij} - P(Z=j|Y))\},$$

and  $c$  is a  $(k - 1)$  vector with

$$(37) \quad c_j = E\{h(Y, Z)(1_{[Z=j]} - P(Z=j|Y))\}, \quad j = 1, \dots, k-1.$$

The above generalizes easily to the case of  $r > 1$  infinite-dimensional spaces  $\mathbf{H}_1, \dots, \mathbf{H}_r$  with closed sumspace  $\sum_{j=1}^r \mathbf{H}_j$ . Consider projections onto  $\mathbf{H}_1 + \dots + \mathbf{H}_r + \mathbf{H}_Z \equiv \mathbf{H}_+$  where

$$\mathbf{H}_j \equiv \{a(Y_j): E a^2(Y_j) < \infty, E a(Y_j) = 0\}, \quad j = 1, \dots, r$$

with  $\mathbf{H}_1 \cap \dots \cap \mathbf{H}_r = \{0\}$ , and  $\mathbf{H}_Z \equiv \{b(Z): E b^2(Z) < \infty, E b(Z) = 0\}$ . In analogy with our notation in the case  $r=1$  we write  $\Pi(h | \mathbf{H}_+) = ACE(h | Y_1, \dots, Y_r, Z)$ . Here are some calculations relating this projection, again in the case when  $Z$  takes values in  $\{1, \dots, k\}$  to  $ACE(h | Y_1, \dots, Y_r)$ .

Since  $\mathbf{H}_1 \cap \dots \cap \mathbf{H}_r = \{0\}$ ,

$$(38) \quad ACE(h | Y_1, \dots, Y_r) = a_1(Y_1, h) + \dots + a_r(Y_r, h),$$

where  $a_1, \dots, a_r$  are uniquely defined. Let  $ACE_j(h | Y_1, \dots, Y_r) \equiv a_j(Y_j, h)$ , for  $j = 1, \dots, r$ . Suppose that  $Z$  takes values in  $\{1, \dots, k\}$ , and there is no nonconstant function  $q$  on  $\{1, \dots, k\}$  such that

$$P(q(Z) = \sum_{j=1}^r a_j(Y_j)) = 1 \quad \text{for some } a_1, \dots, a_r.$$

Let  $b_j = 1_{[Z=j]} - P(Z=j)$  as before. Then we can define

$$\begin{aligned} \Pi(h | [b_j - ACE(b_j | Y_1, \dots, Y_r), j = 1, \dots, k]) \\ = \sum_{j=1}^{k-1} \lambda_j(h) (b_j - ACE(b_j | Y_1, \dots, Y_r)) \end{aligned}$$

computable as before by (A.2.15). Then

$$(39) \quad \begin{aligned} ACE(h | Y_1, \dots, Y_r, Z) \\ = \sum_{j=1}^r ACE_j(h | Y_1, \dots, Y_r, Z) + ACE_{r+1}(h | Y_1, \dots, Y_r, Z), \end{aligned}$$

where, for  $j = 1, \dots, r$ ,

$$(40) \quad \begin{aligned} ACE_j(h | Y_1, \dots, Y_r, Z) \\ = a_j(Y_j, h) - \sum_{i=1}^{k-1} \lambda_i(h) ACE_j(b_i | Y_1, \dots, Y_r) \end{aligned}$$

and

$$(41) \quad ACE_{r+1}(h | Y_1, \dots, Y_r, Z) = \sum_{j=1}^k \lambda_j(h) (1_{[Z=j]} - P(Z=j)).$$

*Proofs*

We give complete proofs for the case  $r = 2$  in theorems 1, 2, and 5; for  $r > 2$  we refer to Halperin (1962) and Smith, Solomon, and Wagner (1977).

**Proof of theorem 1.** We prove (3), and then (2) follows by symmetry. Consider the sequence of operators

$$(a) \quad Q_1, Q_2Q_1, Q_1Q_2Q_1, \dots$$

Let  $A_m$  denote the  $m$ th operator in this sequence. For  $g, h \in \mathbf{H}$  we have, by proposition A.2.2.A,

$$\langle A_m g, A_n h \rangle = \langle A_{m+n-\varepsilon} g, h \rangle,$$

where

$$\varepsilon = \begin{cases} 1 & \text{if } m \text{ and } n \text{ have the same parity,} \\ 0 & \text{if } m \text{ and } n \text{ have opposite parity.} \end{cases}$$

Therefore, for  $h \in \mathbf{H}$ ,

$$\begin{aligned} \|A_m h - A_n h\|^2 &= \langle A_m h - A_n h, A_m h - A_n h \rangle \\ &= \langle A_m h, A_m h \rangle - \langle A_m h, A_n h \rangle \\ &\quad + \langle A_n h, A_n h \rangle - \langle A_n h, A_m h \rangle \\ &= \langle A_{2m-1} h, h \rangle + \langle A_{2n-1} h, h \rangle - 2 \langle A_{m+n-\varepsilon} h, h \rangle \\ (b) \quad &= \langle A_{2m-1} h, h \rangle + \langle A_{2n-1} h, h \rangle - 2 \langle A_{2k-1} h, h \rangle \end{aligned}$$

for some  $k$  since  $m + n - \varepsilon$  is always odd. Thus if  $\lim_{k \rightarrow \infty} \langle A_{2k-1} h, h \rangle$  exists,  $\{A_m h\}$  is Cauchy in  $\mathbf{H}$ , and hence converges. Now

$$\langle A_{2k-1} h, h \rangle = \langle A_k h, A_k h \rangle = \|A_k h\|^2 \geq 0,$$

and similarly  $\langle A_{2k+1} h, h \rangle \geq 0$ . But

$$A_{k+1} h = Q_1 A_k h \quad \text{or} \quad Q_2 A_k h$$

so  $\|A_{k+1} h\| \leq \|A_k h\|$  since  $\|Q_i h\| \leq \|h\|$  because the  $Q_i$  are projections. Hence  $\langle A_{2k-1} h, h \rangle \geq \langle A_{2k+1} h, h \rangle$ , and we have

$$\langle A_1 h, h \rangle \geq \langle A_3 h, h \rangle \geq \dots \geq 0.$$

Thus

$$(c) \quad \lim_{k \rightarrow \infty} \langle A_{2k-1} h, h \rangle \quad \text{exists,}$$

and by (b), (c), and completeness of  $\mathbf{H}$  it follows that  $\lim_{n \rightarrow \infty} A_n h$  exists; call it  $A h$ .

Then  $A$  is a linear operator on  $\mathbf{H}$ , and it remains only to show that  $A = \Pi(\cdot | \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)$ . Now for any  $g, h \in \mathbf{H}$ ,

$$\langle A g, A h \rangle = \lim_{m \rightarrow \infty, n \rightarrow \infty} \langle A_m g, A_n h \rangle$$

$$\begin{aligned}
 &= \lim_{m \rightarrow \infty, n \rightarrow \infty} \langle A_{m+n-\varepsilon} g, h \rangle \\
 &= \langle Ag, h \rangle,
 \end{aligned}$$

and it follows from proposition A.2.2.B that  $A$  is a projection operator,  $A = \Pi(\cdot | \mathbf{H}_0)$  for some closed subspace  $\mathbf{H}_0$  of  $\mathbf{H}$ .

Let  $h \in \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp$ . Then  $Q_1 h = Q_2 h = h$ ,  $A_n h = h$ , and hence  $A h = h$ . Thus  $h \in \mathbf{H}_0$ , and it follows that  $\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp \subset \mathbf{H}_0$ .

Since  $Q_1 A_{2i} = A_{2i+1}$  and  $Q_2 A_{2i-1} = A_{2i}$ , it follows from continuity of  $Q_1$  and  $Q_2$  that  $Q_1 A = A$  and  $Q_2 A = A$ . Letting  $h \in \mathbf{H}$  and setting  $g \equiv Ah \in \mathbf{H}_0$  yields

$$Q_1 g = Q_1 Ah = Ah = g,$$

and similarly,  $Q_2 g = g$ , so that  $\text{range}(A) \equiv \mathbf{H}_0 \subset \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp$ .

Thus  $\mathbf{H}_0 = \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp$  and  $A = \Pi(\cdot | \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)$ . □

**Proof of theorem 2.A.** We give a second proof in the spirit of the arguments of Friedrichs (1937), Aronszajn (1950), and Wiener (1955). From the definitions of  $Q_1$ ,  $Q_2$  and (4), it follows that

$$\begin{aligned}
 \text{(a)} \quad Q_1 Q_2 (h - \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2})) &= Q_1 Q_2 (\Pi(h | \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)) \\
 &= \Pi(h | (\overline{\mathbf{H}_1 + \mathbf{H}_2})^\perp) = h - \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2}).
 \end{aligned}$$

Hence by induction, for any  $m \geq 1$ ,

$$\text{(b)} \quad h - \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2}) = (Q_1 Q_2)^m (h - \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2})).$$

Subtract (7) from (b) to obtain, using  $Q_2^2 = Q_2$ ,

$$\begin{aligned}
 \text{(c)} \quad h_1^{(m)} + h_2^{(m)} - \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2}) & \\
 &= - (Q_1 Q_2)^m \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2}) \\
 &= - Q_1 (Q_2 Q_1 Q_2)^{m-1} \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2}).
 \end{aligned}$$

Now  $T \equiv Q_2 Q_1 Q_2$  is a self-adjoint positive operator with norm bounded by 1. Since  $T$  is the identity on  $\mathbf{H}_+^\perp \equiv \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp$ ,  $\overline{\mathbf{H}_+} \equiv \overline{\mathbf{H}_1 + \mathbf{H}_2}$  is an invariant subspace of  $T$ : for  $x \in \mathbf{H}_+^\perp$  and  $y \in \overline{\mathbf{H}_+}$ ,  $\langle Ty, x \rangle = \langle y, Tx \rangle = \langle y, x \rangle = 0$ , hence  $Ty \in \overline{\mathbf{H}_+}$ . Moreover, 1 is not an eigenvalue of  $T | \overline{\mathbf{H}_1 + \mathbf{H}_2}$ . Thus standard theory implies that the right side of (c) converges to 0. Here is a detailed argument.

First, note that  $U \equiv T | \overline{\mathbf{H}_1 + \mathbf{H}_2}$  has a spectral resolution representation; see, e.g., Kantorovich and Akilov (1982, sections IX.5.1–IX.5.5, pages 258–270). In particular, there is an increasing family of projections  $\{I_t : 0 \leq t \leq 1\}$  such that  $I_{t+} = I_t$ ,  $I_0 = 0$  (since  $U$  is positive),  $I_1 = 1$  (since  $\|U\| \leq 1$ ), and  $I_{1-} = I_1$  (since 1 is not an eigenvalue of  $U$ ), and

$$\text{(d)} \quad \|U^{m-1} f\|^2 = \int_0^1 t^{2(m-1)} d\langle f, I_t f \rangle, \quad f \in \overline{\mathbf{H}_1 + \mathbf{H}_2}.$$

The dominated convergence theorem, (d), and  $I_{1-} = I_1$  therefore imply that

$$\|(Q_2 Q_1 Q_2)^{m-1} \Pi(h | \overline{H_1 + H_2})\| \rightarrow 0 \quad \text{as } m \rightarrow \infty .$$

This and (c) complete the proof of A. □

**Proof of theorem 2.C.** First note that, by proposition 1,

$$\begin{aligned} \text{(a)} \quad & (P_1 - P_1 P_2)h \\ &= (P_1 - P_1 P_2)(h - h_1^* - h_2^*) + (P_1 - P_1 P_2)h_2^* \\ &\quad + (P_1 - P_1 P_2)h_1^* \\ &= 0 + 0 + (I - P_1 P_2)h_1^* \end{aligned}$$

since  $P_i h_i^* = h_i^*$ . We now show that

$$\text{(b)} \quad h_1^{(m)} = (I - (P_1 P_2)^m)h_1^*, \quad m \geq 1 .$$

Clearly

$$\begin{aligned} h_1^{(1)} &= P_1(h - h_2^{(1)}) = P_1(h - P_2 h) \\ &= (P_1 - P_1 P_2)h \\ &= (I - P_1 P_2)h_1^* \quad \text{by (a)} . \end{aligned}$$

Suppose (b) holds for some  $m \geq 1$ . Then

$$\begin{aligned} h_1^{(m+1)} &= P_1(h - h_2^{(m+1)}) \\ &= P_1(h - P_2(h - h_1^{(m)})) \\ &= (P_1 - P_1 P_2)h + P_1 P_2 h_1^{(m)} \\ &= (I - (P_1 P_2)^{m+1})h_1^* \quad \text{by (a) and (b)} . \end{aligned}$$

Hence by induction (b) holds for all  $m \geq 1$ . Therefore

$$h_1^* - h_1^{(m)} = (P_1 P_2)^m h_1^* ,$$

where, since  $h_1^* \in (H_1 \cap H_2)^\perp$ ,

$$\begin{aligned} \|(P_1 P_2)^m h_1^*\| &= \|(P_1 P_2)^m h_1^* - \Pi(h_1^* | H_1 \cap H_2)\| \\ &\rightarrow 0 \quad \text{as } m \rightarrow \infty \end{aligned}$$

by theorem 1 applied to  $(P_1 P_2)^m$ .

Since

$$\|h_1^{(m)} + h_2^{(m)} - h_1^* - h_2^*\| \rightarrow 0$$

by part A, it follows that  $\|h_2^{(m)} - h_2^*\| \rightarrow 0$  as well. □

The proof of theorem 2.B requires several preparatory lemmas.

**Lemma 1.** Let  $A : H \rightarrow H$  be a self-adjoint positive operator. Then  $\sup\{\langle x, Ax \rangle : x \in H, \|x\| = 1\} = \|A\|$ .

**Proof.** Let  $x \in H, \|x\| = 1$ . Then

$$\langle x, Ax \rangle \leq \|Ax\| ,$$

and hence

$$\sup\{\langle x, Ax \rangle : x \in \mathbf{H}, \|x\| = 1\} \leq \|A\|.$$

To prove that equality holds, let  $\alpha = \|A\|$ . Then, since the norm of a self-adjoint positive operator is a spectral point, there are  $x_n \in \mathbf{H}$  with  $\|x_n\| = 1$  such that  $\|Ax_n - \alpha x_n\| \rightarrow 0$ . Therefore

$$\begin{aligned} \langle x_n, Ax_n \rangle &= \langle x_n, \alpha x_n \rangle + \langle x_n, Ax_n - \alpha x_n \rangle \\ &\rightarrow \alpha = \|A\|. \end{aligned}$$

□

**Lemma 2.**

$$\begin{aligned} \|P_2 P_1 | (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp \|^2 &= \|P_1 P_2 P_1 | (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp \|^2, \\ \|Q_1 Q_2 | (\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)^\perp \|^2 &= \|Q_2 Q_1 Q_2 | (\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)^\perp \|^2. \end{aligned}$$

**Proof.** For  $h \in \mathbf{H}$  we have, by self-adjointness and idempotency of  $P_1, P_2$ ,

$$\begin{aligned} (a) \quad \|P_2 P_1 h\|^2 &= \langle P_2 P_1 h, P_2 P_1 h \rangle \\ &= \langle P_2 P_1 h, P_1 h \rangle = \langle P_1 P_2 P_1 h, h \rangle. \end{aligned}$$

Since  $(\mathbf{H}_1 \cap \mathbf{H}_2)^\perp$  is invariant under  $P_1 P_2 P_1$  as in the proof of theorem 2.A, (a) and lemma 1 imply that

$$\begin{aligned} \sup_{h \in (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp} \|P_2 P_1 h\|^2 &= \sup_{h \in (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp} \langle P_1 P_2 P_1 h, h \rangle \\ &= \|P_1 P_2 P_1 | (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp \|^2. \end{aligned}$$

□

**Lemma 3.** Let  $\rho(\mathbf{H}_1, \mathbf{H}_2)$  be as defined in theorem 2.B. Then

$$\begin{aligned} \|P_1 P_2 P_1 | (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp \|^2 &= \rho(\mathbf{H}_1, \mathbf{H}_2) = \rho(\mathbf{H}_1^\perp, \mathbf{H}_2^\perp) \\ &= \|Q_2 Q_1 Q_2 | (\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)^\perp \|^2. \end{aligned}$$

**Proof.** We first prove the first and last equalities. Now

$$\begin{aligned} \rho(\mathbf{H}_1, \mathbf{H}_2) &\equiv \sup\{\langle h_1, h_2 \rangle : h_i \in \mathbf{H}_i \cap (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp, \\ &\quad \|h_i\| \leq 1, i = 1, 2\} \\ &= \sup\{\langle P_1 h_1, P_2 h_2 \rangle : h_i \in (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp, \\ &\quad \|h_i\| \leq 1, i = 1, 2\} \\ &= \sup\{\langle P_2 P_1 h_1, h_2 \rangle : h_i \in (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp, \\ &\quad \|h_i\| \leq 1, i = 1, 2\} \\ &= \sup\{\|P_2 P_1 h_1\| : h_1 \in (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp, \|h_1\| \leq 1\} \end{aligned}$$

by Schwarz



(a) 
$$= \|P_2 P_1 | (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp\| .$$

Thus the first equality follows from (a) and lemma 2, and the third equality follows from the first by symmetry.

The second equality follows from the fact that  $\gamma(\mathbf{H}_1, \mathbf{H}_2) = \gamma(\mathbf{H}_1^\perp, \mathbf{H}_2^\perp)$  by Kato (1976, theorem 4.8, page 221), where  $\gamma^2(\mathbf{H}_1, \mathbf{H}_2)$  is defined in Kato (1976, page 219) as (see problem 6.27, page 55, of Kato (1976)),

$$\begin{aligned} \gamma^2(\mathbf{H}_1, \mathbf{H}_2) &= \inf_{h_1 \in \mathbf{H}_1, h_1 \perp \mathbf{H}_2} \frac{\|h_1 - P_2 h_1\|^2}{\|h_1 - \Pi(h_1 | \mathbf{H}_1 \cap \mathbf{H}_2)\|^2} \\ &= \inf_{h_1 \in \mathbf{H}_1^*, \|h_1\|=1} \|h_1 - P_2 h_1\|^2, \end{aligned}$$

with  $\mathbf{H}_1^* \equiv \mathbf{H}_1 \cap (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp$

$$= 1 - \sup_{h_1 \in \mathbf{H}_1^*, \|h_1\|=1} \langle P_2 h_1, P_2 h_1 \rangle$$

$$= 1 - \sup_{h_1 \in \mathbf{H}_1^*, \|h_1\|=1} \langle P_1 P_2 P_1 h_1, h_1 \rangle$$

(b) 
$$= 1 - \rho^2(\mathbf{H}_1, \mathbf{H}_2)$$

by lemma 1 and the first equality. □

**Proof of theorem 2.B.** From (c) of the proof of theorem 2.A,

$$\begin{aligned} \|h_1^{(m)} + h_2^{(m)} - \Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2})\| \\ \leq \|Q_2 Q_1 Q_2 | (\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)^\perp\|^{m-1} \|\Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2})\| \end{aligned}$$

(a) 
$$= \rho^{2(m-1)} \|\Pi(h | \overline{\mathbf{H}_1 + \mathbf{H}_2})\| \quad \text{by lemma 3.}$$

That  $\rho < 1$  if and only if  $\mathbf{H}_1 + \mathbf{H}_2$  is closed follows from

(b) 
$$1 - \rho^2(\mathbf{H}_1, \mathbf{H}_2) = \gamma^2(\mathbf{H}_1, \mathbf{H}_2)$$

as in (b) of the proof of lemma 3, and from Kato (1976, theorem 4.2, page 219), which asserts that  $\mathbf{H}_1 + \mathbf{H}_2$  is closed if and only if  $\gamma(\mathbf{H}_1, \mathbf{H}_2) > 0$ . □

**Proof of theorem 2.D.** We first show that

(a) 
$$(Q_2 Q_1)^m - Q = (Q_2 Q_1)^m P$$

where  $P \equiv I - Q = \Pi(\cdot | \overline{\mathbf{H}_1 + \mathbf{H}_2})$ . Now

(b) 
$$\begin{aligned} (Q_2 Q_1)^m &= (Q_2 Q_1)^m (Q + P) = (Q_2 Q_1)^m Q + (Q_2 Q_1)^m P \\ &= Q + (Q_2 Q_1)^m P \end{aligned}$$

since  $\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp \subset \mathbf{H}_i^\perp$  implies  $Q_i Q = Q, i = 1, 2$ . Thus (a) holds.

As in the proof of theorem 2.A, note that  $\overline{\mathbf{H}_1 + \mathbf{H}_2}$  is an invariant subspace of  $Q_i$ . Consequently,

$$(c) \quad Q_i | \overline{\mathbf{H}_1 + \mathbf{H}_2} = \Pi(\cdot | \mathbf{H}_i^\perp \cap \overline{\mathbf{H}_1 + \mathbf{H}_2}) | \overline{\mathbf{H}_1 + \mathbf{H}_2}$$

and, by proposition A.2.3.B,

$$(d) \quad Q_i P = \Pi(\cdot | \mathbf{H}_i^\perp \cap \overline{\mathbf{H}_1 + \mathbf{H}_2}).$$

With

$$(e) \quad R_i \equiv Q_i P,$$

Equation (c) also yields

$$Q_2 Q_1 P = Q_2 P Q_1 P = R_2 R_1$$

and even, by (a),

$$(f) \quad (Q_2 Q_1)^m - Q = (R_2 R_1)^m;$$

hence

$$(g) \quad \|(Q_2 Q_1)^m - Q\| = \|(R_2 R_1)^m\|.$$

Since  $(R_2 R_1)^m$  has adjoint  $(R_1 R_2)^m$ , it follows that

$$(h) \quad \|(R_2 R_1)^m\|^2 = \|(R_1 R_2)^m (R_2 R_1)^m\| = \|(R_1 R_2 R_1)^{2m-1}\|,$$

where  $R_1 R_2 R_1$  is self-adjoint, and hence normal. Thus by Kato (1976, page 56),

$$(i) \quad \|(R_1 R_2 R_1)^{2m-1}\| = \|R_1 R_2 R_1\|^{2m-1}.$$

The claim follows from (g)–(i),

$$\|R_1 R_2 R_1\| = \|Q_1 Q_2 Q_1 P\| = \|Q_1 Q_2 Q_1 | \overline{\mathbf{H}_1 + \mathbf{H}_2}\|$$

and lemma 3. □

**Proof of proposition 1.** For all  $h \in \mathbf{H}$ ,  $h_i^* \in \mathbf{H}_i$ ,  $i = 1, 2$ , we have (11) if and only if

$$h - h_1^* - h_2^* \perp \mathbf{H}_1 + \mathbf{H}_2 = (\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp)^\perp,$$

if and only if

$$h - h_1^* - h_2^* \in \overline{\mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp} = \mathbf{H}_1^\perp \cap \mathbf{H}_2^\perp,$$

if and only if

$$(a) \quad P_i(h - h_1^* - h_2^*) = 0, \quad i = 1, 2.$$

Clearly (a) is equivalent to (12) and (13). □

**Proof of proposition 2.** Part A is proved in Kober (1939), and part D was established in the proof of theorem 2.B. Now for B. If either  $P_1 | \mathbf{H}_2$  or  $P_2 | \mathbf{H}_1$  is compact, so is  $P_1 P_2 P_1$  since the composition of two operators  $U \circ V$  is compact if either  $U$  or  $V$  is compact. Now  $P_1 P_2 P_1$  is a positive compact operator,

and therefore has at most a finite number of eigenvalues greater than  $\varepsilon > 0$ ; see, e.g., Rudin (1973, theorem 4.24, page 101). Since  $P_1 P_2 P_1 \upharpoonright (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp$  has norm bounded by 1 and 1 is not an eigenvalue,  $\|P_1 P_2 P_1 \upharpoonright (\mathbf{H}_1 \cap \mathbf{H}_2)^\perp\| < 1$ . Then B follows from lemma 3 and part D. Finally, C follows from standard theory; see, e.g., Reed and Simon (1972, theorem VI.22, page 210).  $\square$

**Proof of theorem 5.**

$$\mathbf{H}_1 + \mathbf{H}_2 = \mathbf{H}_1 + [b_1 - P_1 b_1, \dots, b_d - P_1 b_d].$$

Hence, by (A.2.9), since  $b_j - P_1 b_j \perp \mathbf{H}_1, j = 1, \dots, d$ ,

$$(a) \quad \Pi(h \mid \mathbf{H}_1 + \mathbf{H}_2) = \Pi(h \mid \mathbf{H}_1) + \Pi(h \mid [b_1 - P_1 b_1, \dots, b_d - P_1 b_d]).$$

An application of (A.2.15) to the second term on the right side in (a) yields (26).  $\square$

**Proof of proposition 3.** Suppose first that  $d < k - 1$ . Then there are constants  $\alpha_1, \dots, \alpha_{k-1}$  not all zero such that

$$(a) \quad \sum_{j=1}^{k-1} \alpha_j \left( 1_{[Z=j]} - P(Z=j \mid Y) \right) = 0 \quad \text{a.s.}$$

Let

$$q(Z) \equiv \sum_{j=1}^{k-1} \alpha_j 1_{[Z=j]}$$

and

$$a(Y) \equiv \sum_{j=1}^{k-1} \alpha_j P(Z=j \mid Y).$$

Then  $q$  is nonconstant and (30) holds in view of (a).

Conversely if a  $q$  satisfying (31) exists, then

$$(b) \quad \sum_{j=1}^k q(j) \left( 1_{[Z=j]} - P(Z=j \mid Y) \right) = 0 \quad \text{a.s.}$$

or, by subtracting  $q(k)$  times (29) from (b),

$$\sum_{j=1}^{k-1} (q(j) - q(k)) \left( 1_{[Z=j]} - P(Z=j \mid Y) \right) = 0 \quad \text{a.s.}$$

Since  $q$  is nonconstant, not all  $q(j) - q(k)$  vanish and hence (a) holds.  $\square$

## A.5 DERIVATIVES

Let  $T$  be a function from a normed linear space  $\mathbf{X}$  to another normed linear space  $\mathbf{Y}$ . A variety of possible definitions of a "derivative" of  $T$  are possible in this setting. One of the simplest and most useful is:

**Definition 1.** Let  $x \in \mathbf{X}$ . Then  $T$  is said to be *Gâteaux differentiable* at  $x$  if

$$(1) \quad \dot{T}_x(h) \equiv \lim_{\lambda \rightarrow 0} \frac{T(x + \lambda h) - T(x)}{\lambda}$$

exists for every  $h \in X$ .

Note that this definition does not use the fact that  $X$  is normed, but only that  $X$  is a linear space. Thus the existence of a Gâteaux derivative is a very weak requirement, and in fact  $T$  may fail to be continuous and still have a Gâteaux derivative. A stronger notion of a derivative of  $T$  is:

**Definition 2.** Let  $x \in X$ . Then  $T$  is *Fréchet differentiable* at  $x$  if there exists a function  $\dot{T}_x : X \rightarrow Y$  which is linear and continuous such that

$$(2) \quad \lim_{\|h\| \rightarrow 0} \frac{\|T(x+h) - T(x) - \dot{T}_x(h)\|}{\|h\|} = 0.$$

We will also write  $\dot{T}(x)$  or sometimes just  $\dot{T}$  for  $\dot{T}_x$ . For problems in statistics, the  $X$  space is often either a parameter space, frequently denoted by  $\Theta$ , or a collection of probability distributions, which we denote by  $P$  (or  $S$ ; see section 3.2).

Note that (2) is equivalent to

$$(3) \quad \frac{\|T(x + h_n) - T(x) - \dot{T}_x(h_n)\|}{\|h_n\|} \rightarrow 0 \quad \text{for any } \|h_n\| \rightarrow 0,$$

and to

$$(4) \quad \left\| \frac{T(x + \varepsilon_n h_n) - T(x)}{\varepsilon_n} - \dot{T}_x(h_n) \right\| \rightarrow 0$$

for any  $\varepsilon_n \rightarrow 0$  in  $R$  and  $\|h_n\|$  bounded. This second equivalent formulation of Fréchet differentiability leads naturally to another notion of differentiability intermediate to Gâteaux and Fréchet differentiability.

**Definition 3.** Let  $x \in X$ . Then  $T$  is *Hadamard (or compactly) differentiable* at  $x$  if there exists a continuous linear function  $\dot{T}_x : X \rightarrow Y$  such that

$$(5) \quad \left\| \frac{T(x + \varepsilon_n h_n) - T(x)}{\varepsilon_n} - \dot{T}_x(h_n) \right\| \rightarrow 0$$

for any  $\varepsilon_n \rightarrow 0$  and any  $\{h_n\}_{n \geq 0} \subset X$  with  $\|h_n - h_0\| \rightarrow 0$ . Here is another way of looking at the three kinds of derivative defined so far:

**Definition 4.** Let  $S$  be a collection of subsets of  $X$ , and for  $x, h \in X$ , set

$$(6) \quad \text{Rem}(T, x, h) \equiv T(x + h) - T(x) - \dot{T}_x(h).$$

Then  $T$  is  $S$ -differentiable at  $x$  with derivative  $\dot{T}_x$  if, for all  $S \in S$ ,

$$(7) \quad \frac{\text{Rem}(T, x, \varepsilon h)}{\varepsilon} \rightarrow 0 \quad \text{as } \varepsilon \rightarrow 0 \quad \text{uniformly in } h \in S.$$

Now definitions 1–3 can be reformulated as follows:

G: If  $S = \{\text{all singeltons of } X\}$ ,  $S$ -differentiability is Gâteaux differentiability.

- H: If  $\mathcal{S} = \{\text{all compact subsets of } \mathbf{X}\}$ ,  $\mathcal{S}$ -differentiability is Hadamard (or compact) differentiability.
- F: If  $\mathcal{S} = \{\text{all bounded subsets of } \mathbf{X}\}$ ,  $\mathcal{S}$ -differentiability is Fréchet differentiability.

Written this way, it is clear that Fréchet differentiability (at  $x$ ) implies Hadamard differentiability (at  $x$ ) which in turn implies Gâteaux differentiability (at  $x$ ).

**Remark 1.** When  $\mathbf{X} = R$ , all three definitions of the derivative (at a point  $x$ ), Gâteaux, Hadamard, and Fréchet, are equivalent. When  $\mathbf{X} = R^k$ , the Hadamard and Fréchet derivatives still coincide, but the Gâteaux derivative is strictly weaker (as in the classical calculus; coordinatewise differentiability does not imply “joint differentiability”).

Remarkably however, if we insist on “continuity of the derivative,” as  $x$  varies in  $\mathbf{X}$ , then the three theories again coincide if  $\mathbf{X}$  is finite dimensional, and continuous Gâteaux differentiability is equivalent to continuous Hadamard differentiability if  $\mathbf{X}$  is a Banach space.

As Reeds (1976) and Gill (1988) argue with some justice, Hadamard differentiability is well suited for applications in statistics.

**Proposition 1.** Suppose that  $T : \mathbf{X} \rightarrow \mathbf{Y}$  is Fréchet (or Hadamard) differentiable at  $x \in \mathbf{X}$ . Then:

- A.  $\dot{T}_x$  is unique
- B.  $T$  is continuous at  $x$
- C. The Hadamard and Gâteaux derivatives exist and equal the Fréchet (or Hadamard) derivative.
- D. If, in addition,  $U : \mathbf{Y} \rightarrow \mathbf{Z}$ ,  $\dot{U}$  exists as a Fréchet (Hadamard) derivative at  $T(x)$ , and  $V \equiv U \circ T : \mathbf{X} \rightarrow \mathbf{Z}$ , then  $V$  is Fréchet (Hadamard) differentiable at  $x$  and  $\dot{V} = \dot{U}_T \circ \dot{T}$ ; i.e.,  $\dot{V}_x = \dot{U}_{T(x)} \circ \dot{T}_x : \mathbf{X} \rightarrow \mathbf{Z}$ .
- E. Suppose  $T : \mathbf{X} \rightarrow \mathbf{Y}$  is continuously Gâteaux differentiable in a neighborhood of  $x_0 \in \mathbf{X}$ , but not necessarily Fréchet or Hadamard differentiable. Then the Gâteaux derivative  $\dot{T}_{x_0}$  at  $x_0$ ,

$$\frac{\partial T}{\partial \lambda}(x_0 + \lambda b) |_{\lambda=0} = \dot{T}_{x_0}(b),$$

is a linear operator. If, further, for some  $\varepsilon > 0$ ,

$$\sup \left\{ \left\| \frac{\partial T}{\partial \lambda}(x_0 + \lambda b) \right\|_{\mathbf{Y}} : \|b\|_{\mathbf{X}} \leq 1, \quad |\lambda| \leq \varepsilon \right\} = M_1 < \infty,$$

then  $T$  is Fréchet differentiable with derivative  $\dot{T}_x(\cdot)$  for all  $x$  in an open neighborhood  $N$  of  $x_0$  and  $N$  can be taken such that

$$\sup \{ \|\dot{T}_x(\cdot)\| : x \in N \} = M_1.$$

If, further,

$$\sup \{ \left\| \frac{\partial^2 T}{\partial \lambda^2}(x_0 + \lambda b) \right\|_{\mathbf{Y}} : \|b\|_{\mathbf{X}} \leq 1, \quad |\lambda| \leq \varepsilon \} = M_2 < \infty,$$

then  $N$  can be chosen such that  $x \rightarrow \dot{T}_x(\cdot)$  is continuous in the operator topology on  $N$  and,

$$\| \dot{T}_{x_2}(\cdot) - \dot{T}_{x_1}(\cdot) \|_Y \leq M_2 \| x_2 - x_1 \|_X$$

for all  $x_1, x_2 \in N$ .

**Proof.** For A–D see, e.g., Luenberger (1969, sections 7.2 and 7.3, pages 171–177), and Flett (1980, sections 4.1 and 4.2, pages 251–273). For E, see Flett (1980, proposition 4.1.7, page 257), and Kantorovich and Akilov (1982, pages 499–501).  $\square$

In fact, Hadamard differentiability is the weakest kind of differentiation that satisfies the chain rule, D. An equivalent definition of Hadamard differentiability is easily formulated in terms of the chain rule:  $T$  is Hadamard differentiable at  $x$  if and only if for all  $U : R \rightarrow X$  which are differentiable at 0 (in the ordinary sense; see remark 1) and satisfy  $U(0) = x, T \circ U : R \rightarrow Y$  is differentiable at 0 (in the ordinary sense) with derivative  $\dot{T}_x \circ \dot{U}_0$ ; see Averbukh and Smolyanov (1967, pages 201 and 203).

The following variants of Hadamard differentiability will be useful in many particular cases.

**Definition 5.** Suppose that  $X_0 \subset X$  is a subspace of  $X$ . Then  $T : X \rightarrow Y$  is *Hadamard differentiable (tangentially to  $X_0$ )* (at  $x$ ) if (5) holds for all  $\varepsilon_n \rightarrow 0$  and  $\{h_n\} \subset X$  with  $\|h_n - h_0\| \rightarrow 0$  and  $h_0 \in X_0$ . If (5) holds for all  $\varepsilon_n \rightarrow 0$  and  $\{h_n\} \subset X_0$  with  $\|h_n - h_0\| \rightarrow 0$ , then  $T$  is *Hadamard differentiable along  $X_0$* .

**Remark 2.** We will also use definition 5 when  $T$  is defined on a subset  $E \subset X$  which is not a subspace itself, but  $E$  has a *tangent space  $\dot{E}$* : for all  $h_0 \in \dot{E}$  there are  $h_n \rightarrow h_0, \varepsilon_n \rightarrow 0$ , so that  $x + \varepsilon_n h_n \in E$  for all  $n \geq 1$ . In fact, with this modification, Hadamard differentiability of  $T$  along  $\dot{E}$  is just what we call *pathwise differentiability* in sections 3.3 and 5.3.

Part C of Proposition 1 justifies the following procedure: we first calculate a Gâteaux derivative of a functional  $T$ , and then proceed to verify that the Gâteaux derivative is, in fact, a Fréchet derivative.

**Example 1.** Suppose that  $T : R^d \rightarrow R$  has continuous partial derivatives  $\dot{T}_i, i = 1, \dots, d$ , at  $x \in R^d$ . Then  $T$  is Fréchet differentiable at  $x$  and

$$(8) \quad \dot{T}_x(h) = \dot{T}h = \sum_{i=1}^d \dot{T}_i h_i \equiv (\nabla T)^T h \quad \text{for } h \in R^d.$$

(See Luenberger (1969, example 7.2.4, page 173).)  $\square$

**Example 2.** Let  $\mu$  denote a fixed positive measure on  $(X, \mathcal{B})$  and suppose that  $T : \Theta \subset R^k \rightarrow L_2(\mu)$ . (We now use  $x$  to denote the argument of the  $L_2(\mu)$  functions; note that this is the situation in section 2.1 where  $T$  is taken to be  $s$  or  $r$ .) Then, if  $\dot{T}$  exists, it is a continuous linear function from  $\Theta$  into  $L_2(\mu)$  such that

$$(11) \quad \| T(\theta+h) - T(\theta) - \dot{T}(\theta)h \| = o(|h|)$$

where  $\| \cdot \|$  denotes the norm in  $L_2(\mu)$  and  $|\cdot|$  denotes the Euclidean norm in  $R^k$ . By a Gâteaux derivative calculation, it is easily seen that if (4) holds, then, by the linearity of  $\dot{T}$ , there is a vector of  $L_2(\mu)$  functions  $\dot{T}(\theta) = (\dot{T}_1(\theta), \dots, \dot{T}_k(\theta))^T$  such that

$$(12) \quad \dot{T}_\theta h = \sum_{i=1}^k \dot{T}_i(\theta) h_i = \dot{T}(\theta)^T h .$$

(Note that the argument  $x$  of the  $L_2$  functions is suppressed throughout.) When such a vector  $\dot{T}(\theta)$  exists, then  $T$  is also called *differentiable in quadratic mean*; see, e.g., Le Cam (1970). A sufficient condition for the existence of such a vector  $\dot{T}(\theta)$  is that  $T(\theta) = T(\theta)(x)$  be continuously differentiable with respect to  $\theta$  for  $\mu$  a.e.  $x$  with gradient  $\dot{T}(\theta) = \nabla T(\theta)$  with  $\dot{T}_i(\theta) \in L_2(\mu)$  for  $i = 1, \dots, k$ , and that  $\Sigma(\theta) \equiv \int \dot{T}(\theta) \dot{T}(\theta)^T d\mu$  be continuous in  $\theta$ : a proof is contained in the proof of Proposition 2.1.1.  $\square$

**Example 3.** Suppose that  $\mu$  is a fixed positive measure on  $(X, \mathcal{B})$  as in example 2, but now suppose that (changing notation to match with section 3.2)  $v \equiv T : \mathcal{S} \subset L^2(\mu) \rightarrow R$ . A variety of functions of interest in statistics can be expressed in this form. For example, if  $p = dP/d\mu$  and  $s = \sqrt{p} \in \mathcal{S} \subset L^2(\mu)$  as in chapter 3, then  $v(s) = \int x s^2(x) d\mu(x) = \int x dP(x)$  is the mean, assuming  $\int |x| dP(x) < \infty$ . If such a function  $v$  has a Fréchet derivative, or a Hadamard derivative in the sense of definition 5 and remark 1, then  $\dot{v} : \dot{\mathcal{S}} \rightarrow R$ , is linear and continuous, and hence by the Riesz representation theorem  $\dot{v}$  is of the form

$$(9) \quad \dot{v}(t) = \langle \dot{v}, t \rangle \quad \text{for } t \in \dot{\mathcal{S}}$$

for some  $L_2(\mu)$  function  $\dot{v} = \dot{v}_s$ ; see example A.1.8. Although the mean functional

$$v(s) = \int x s^2(x) d\mu(x)$$

is not Fréchet differentiable at any  $s \in \mathcal{S} = \{s : \int |x| s^2(x) d\mu(x) < \infty\}$ ,  $v$  is Hadamard differentiable along (or pathwise differentiable on) appropriate subsets, e.g., for all  $\mathcal{S}_M \equiv \{s : P \in \mathcal{P} \text{ with } P[-M, M] = 1\}$  (see section 3.3), or  $\mathcal{S}_M \equiv \{s : \int x^2 dP(x) = \int x^2 s^2(x) d\mu(x) \leq M\}$ . The following useful proposition (due to Van der Vaart (1988c)) establishes this last claim for general linear functionals  $v(P)$ . See Beran (1977b) for construction of a family of location functionals  $v$  which are Fréchet differentiable at every  $s \in \mathcal{M}^{1/2}$ .  $\square$

**Proposition 2.** Suppose that  $v : \mathcal{P} \rightarrow R$  is of the form  $v(P) = \int g dP$  where  $g$  is a fixed function. If

$$(10) \quad \sup_{P \in \mathcal{P}} E_P g^2(X) < \infty ,$$

then  $v$  is pathwise differentiable at every  $P_0 \in \mathcal{P}$  with derivative  $\dot{v} \in L_2(P_0)$  given by  $\dot{v} \equiv \dot{v}(P_0) = g - E_0 g(X)$ .

**Proof.** Let  $\{s(\eta)\}$  be a path in  $\mathbf{S}$  with tangent  $t \in \dot{\mathbf{S}} \subset L_2(\mu)$ . Then, with  $\langle \dot{v}(s_0), t \rangle = \langle \dot{v}(P_0), 2t/s_0 \rangle_0$  (see definition 3.3.1), we have

$$\begin{aligned} v(s(\eta)) - v(s(0)) - \eta \langle \dot{v}(s_0), t \rangle &= \int g \{s^2(\eta) - s^2(0) - 2\eta t s_0\} d\mu \\ &= \int g \{s(\eta) - s(0) - \eta t\} \{s(\eta) + s(0)\} d\mu \\ &\quad + \eta \int g \{s(\eta) - s(0)\} t d\mu \\ (a) \qquad \qquad \qquad &\equiv \eta \{A(\eta) + B(\eta)\}. \end{aligned}$$

But

$$\begin{aligned} |A(\eta)| &\leq \left[ \int g^2 \{s(\eta) + s(0)\}^2 d\mu \right]^{1/2} \left[ \int \{\eta^{-1}(s(\eta) - s(0)) - t\}^2 d\mu \right]^{1/2} \\ &\leq \left[ 2 \{E_\eta(g^2) + E_0(g^2)\} \right]^{1/2} o(1) \end{aligned}$$

$$(b) \qquad = o(1) \quad \text{by (10),}$$

and

$$\begin{aligned} |B(\eta)| &\leq \left| \int_{\{|g| \leq 1/\sqrt{\eta}\}} g t (s(\eta) - s(0)) d\mu \right| \\ &\quad + \left| \int_{\{|g| > 1/\sqrt{\eta}\}} g t (s(\eta) - s(0)) d\mu \right| \\ &\leq \left\{ \frac{1}{\eta} \int t^2 d\mu \int (s(\eta) - s(0))^2 d\mu \right\}^{1/2} \\ &\quad + \left\{ \int_{\{|g| > 1/\sqrt{\eta}\}} t^2 d\mu (E_\eta(g^2) + E_0(g^2)) \right\}^{1/2} \end{aligned}$$

$$(c) \qquad = o(1) + o(1) \quad \text{as } \eta \rightarrow 0$$

using (10),  $t \in L_2(\mu)$ , and  $\eta^{-2} \int (s(\eta) - s(0))^2 d\mu \rightarrow \int t^2 d\mu$ . Combining (b) and (c) with (a) completes the proof.  $\square$

For our applications,  $T: \Theta \subset R^k \rightarrow L_2(\mu)$  in example 2 is a parametric family of square-root density functions with respect to  $\mu$  on  $(X, \mathcal{B})$ :

$$s(\theta) = \sqrt{p(\cdot, \theta)} \quad \text{where} \quad \int p(x, \theta) d\mu(x) = 1$$

for all  $\theta \in \Theta$ . We therefore take  $T = s$  temporarily. Recall the definition of a regular parametrization given in section 2.1.

**Proposition 3.** Suppose that  $\mathbf{P} = \{P_\theta: \theta \in \Theta\}$  is regular. Then:

- A.  $I(\theta) = 4 \int \dot{s}(\theta) \dot{s}(\theta)^T d\mu$  is continuous on  $\Theta$ .
- B.  $s(\theta + h) - s(\theta) = \int_0^1 h^T \dot{s}(\theta + uh) du \mu$  a.e., if  $h$  is sufficiently small.
- C. The density  $p(\cdot, \theta)$  is Fréchet differentiable in  $L_1(\mu)$  uniformly in  $\theta \in K$  compact with  $\dot{p}(\theta) = 2s(\theta) \dot{s}(\theta)$ :



$$\int |p(x, \theta + h) - p(x, \theta) - h^T \dot{p}(x, \theta)| d\mu(x) = o(|h|).$$

D.  $p(\theta + h) - p(\theta) = \int_0^1 h^T \dot{p}(\theta + uh) du \quad \mu \text{ a.e.}$

E.  $\int_{[s(\theta) = 0]} p(\theta + h) d\mu = o(|h|^2).$

F.  $\int_{[s(\theta) = 0]} |\dot{s}(\theta)|^2 d\mu = 0.$

**Proof.** To prove A, note that for any  $1 \leq i, j \leq k$ ,

$$\begin{aligned} & |I_{ij}(\theta + h) - I_{ij}(\theta)| \\ &= 4 \left| \int \dot{s}_i(\theta + h) \dot{s}_j(\theta + h) d\mu - \int \dot{s}_i(\theta) \dot{s}_j(\theta) d\mu \right| \\ \text{(a)} \quad &\leq 4 \int |\dot{s}_i(\theta + h)| |\dot{s}_j(\theta + h) - \dot{s}_j(\theta)| d\mu \\ &\quad + 4 \int |\dot{s}_j(\theta)| |\dot{s}_i(\theta + h) - \dot{s}_i(\theta)| d\mu \\ &\leq 4 \|\dot{s}_i(\theta + h)\| \|\dot{s}_j(\theta + h) - \dot{s}_j(\theta)\| \\ &\quad + 4 \|\dot{s}_j(\theta)\| \|\dot{s}_i(\theta + h) - \dot{s}_i(\theta)\| \\ &\rightarrow 0 \quad \text{as } h \rightarrow 0 \end{aligned}$$

by continuity of  $\theta \rightarrow \dot{s}_i(\theta), i = 1, \dots, k$ , in  $L_2(\mu)$ .

Part B is a generalized version of the Newton-Leibniz fundamental theorem of calculus for functions with values in a Banach space; see, e.g., Kantorovich and Akilov (1982, chapter XVII, section 1.7, pages 503–504).

To prove C, write the integrand of the left-hand side as

$$\begin{aligned} & s^2(\theta + h) - s^2(\theta) - 2s(\theta)h^T \dot{s}(\theta) \\ \text{(b)} \quad &= 2s(\theta)[s(\theta + h) - s(\theta) - h^T \dot{s}(\theta)] + [s(\theta + h) - s(\theta)]^2. \end{aligned}$$

This identity yields  $L_1(\mu)$ -differentiability of  $\theta \rightarrow p(\theta) = s^2(\theta)$ . Assertion D follows in the same way as B. Finally, by the nonnegativity and Fréchet differentiability of  $s$  we have

$$\begin{aligned} & \int_{[s(\theta) = 0]} p(\theta + h) d\mu \leq \int_{[s(\theta) = 0]} (s(\theta + h) + s(\theta - h))^2 d\mu \\ \text{(c)} \quad &\leq 2 \int_{[s(\theta) = 0]} (s(\theta + h) - \dot{s}^T(\theta)h)^2 d\mu \\ &\quad + 2 \int_{[s(\theta) = 0]} (s(\theta - h) + \dot{s}^T(\theta)h)^2 d\mu \\ &= o(|h|^2), \end{aligned}$$

and consequently

$$\begin{aligned} \text{(d)} \quad & \int_{[s(\theta) = 0]} (\dot{s}^T(\theta)h)^2 d\mu \leq 2 \int_{[s(\theta) = 0]} s^2(\theta + h) d\mu \\ &\quad + 2 \int_{[s(\theta) = 0]} (s(\theta + h) - \dot{s}^T(\theta)h)^2 d\mu \end{aligned}$$

$$= o(|h|^2).$$

□

Hájek (1972) provided the following useful set of conditions for Fréchet differentiability of  $s: \Theta \subset R \rightarrow L_2(\mu)$  where  $\mu$  is Lebesgue measure. Suppose that:

- (i) In some neighborhood of  $\theta_0 \in \Theta \subset R$  the functions  $p: R \times \Theta \rightarrow R$ , denoted  $p(x, \theta)$ , are absolutely continuous in  $\theta$  for all  $x \in R$ .
- (ii) For every  $\theta$  in some neighborhood of  $\theta_0$  the  $\theta$ -derivative  $\dot{p}(x, \theta) = (\partial/\partial\theta)p(x, \theta)$  exists for almost all (Lebesgue measure  $\mu$ )  $x \in R$ .
- (iii) The Fisher information

$$I(\theta) \equiv \int_{-\infty}^{\infty} \frac{[\dot{p}(x, \theta)]^2}{p(x, \theta)} 1_{[p(x, \theta) > 0]} d\mu(x)$$

exists, is continuous at  $\theta_0$ , and  $I(\theta_0) > 0$ .

**Proposition 4.** (Hájek, 1972). If (i)–(iii) hold, then:

A.  $s \equiv \sqrt{p}: R \times \Theta \rightarrow R$  is absolutely continuous in  $\theta$  in a neighborhood of  $\theta_0$  for almost all  $x \in R$ .

B.  $\theta_0$  is a regular point of the model  $\mathbf{P} \equiv \{P_\theta \ll \mu \text{ with density } p(\cdot, \theta), \theta \in \Theta\}$ . In particular,

$$\|s(\theta) - s(\theta_0) - (\theta - \theta_0)\dot{s}(\theta_0)\| = o(|\theta - \theta_0|),$$

where

$$\dot{s}(\theta_0) \equiv \frac{\dot{p}(\theta_0)}{2p^{1/2}(\theta_0)} 1_{[p(\theta_0) > 0]}.$$

**Proof.** See the appendix of Hájek (1972) or Ibragimov and Has'minskii (1981, page 121). Part A of the conclusion also follows from standard real variable theory; see, e.g., Graves (1946, theorem 7, page 222) with  $\xi = p$  and  $\lambda(x) = 1/(2\sqrt{x})$ . □

**Corollary 1.** Suppose that  $g$  is a density which is absolutely continuous with respect to Lebesgue measure  $\mu$  on  $(R, \mathcal{B})$ ,

$$g(x) = \int_{-\infty}^x g'(y) dy, \quad \text{with} \quad I(g) \equiv \int \frac{[g']^2}{g} < \infty.$$

Then:

- A.  $s \equiv g^{1/2}$  is absolutely continuous with  $s'(x) = g'(x)/2g^{1/2}(x)$ ,
- B.  $\| [s(\cdot + \eta) - s]/\eta \|^2 \rightarrow \int [s']^2 d\mu$  as  $\eta \rightarrow 0$ ,
- C.  $\|s(\cdot + \eta) - s - \eta s'\| = o(\eta)$ ,
- D. The distributions with density  $g(\cdot - \theta)$  with respect to Lebesgue measure on  $(R, \mathcal{B})$ ,  $\theta \in R$ , form a regular parametric model.

**Proof.** With  $p(x, \theta) = g(x - \theta)$  the conditions of proposition 4 hold by the fundamental theorem of calculus. Note that the Fisher information  $I(g)$  is independent of  $\theta$  and that it has to be positive, since  $g$  cannot be constant. Consequently, A, B, and C follow from proposition 4. Furthermore, D holds in view of the  $L_2$ -continuity in  $\theta$  of  $s'(\cdot - \theta)$ ; see theorem A.1.1.

A direct proof of these results is given in Hájek and Šidák (1967, lemma VI.2.1.a, pages 211, 212). □

A regular parametric model  $\mathbf{P}$  is preserved by an arbitrary measurable mapping  $T \equiv T(X)$  of the original data  $X \sim P_\theta \in \mathbf{P}$ . This is stated by Ibragimov and Has'minskii (1981, theorem 7.2, page 70), and generalizations have been proved by Van der Vaart (1988a, appendix A.3), and by Le Cam and Yang (1988). With a slightly different definition of regularity this phenomenon has been proved already by Pitman (1979); see his Theorem 3.8, page 19. The following proposition gives a simple version of this type of result.

**Proposition 5.** Suppose that  $\mathbf{P} = \{P_\theta : \theta \in \Theta \subset R^k\}$  is regular. Let  $T : \mathbf{X} \rightarrow \mathbf{T}$  be a statistic, and let

$$\mathbf{P}_T \equiv \{P_\theta T^{-1} : \theta \in \Theta\}.$$

Take  $\mu$  such that  $\mu T^{-1}$  is  $\sigma$ -finite and let

$$p_T(t, \theta) = \frac{dP_\theta T^{-1}}{d\mu T^{-1}}(t)$$

denote the density of  $P_\theta T^{-1}$ ,  $s_T = p_T^{1/2}$ ,  $l_T = \log p_T$ , etc. Then:

A.  $\theta \rightarrow s_T(\theta)$  is continuously Fréchet differentiable and

$$(13) \quad \dot{s}_T(\theta) = \frac{1}{2} s_T(\theta) E[\dot{l}(X, \theta) | T], \quad \dot{l}_T(\theta) = E[\dot{l}(X, \theta) | T].$$

B. The information for  $\theta$  in the family  $\mathbf{P}_T$  is no larger than that in the family  $\mathbf{P}$ :

$$(14) \quad I(\theta, \mathbf{P}_T) = E\{E[\dot{l}(X, \theta) | T] E[\dot{l}^T(X, \theta) | T]\} \\ \leq E[\dot{l}^T(X, \theta)] = I(\theta, \mathbf{P}).$$

C. Moreover, if  $I(\theta, \mathbf{P}_T)$  is nonsingular for all  $\theta$ , then  $\mathbf{P}_T$  is regular.

**Remark 3.** If  $T : R^2 \rightarrow R \times \{0, 1\}$ ,  $T(x, y) = (y, 1[x \leq y])$  and if  $\mu$  is Lebesgue measure, then  $\mu T^{-1}$  is *not* a  $\sigma$ -finite measure; see, e.g., example 6.6.1. However, we may always change  $\mu$  into a dominating probability measure  $\nu$  by taking  $\nu(A) = \sum_{i=1}^\infty 2^{-i} \mu(A \cap A_i) / \mu(A_i)$  where  $\{A_i\}$  is a partition of  $\mathbf{X}$  with  $0 < \mu(A_i) < \infty$ . Note that the corresponding measure  $\nu T^{-1}$  is a  $\sigma$ -finite dominating measure for  $\mathbf{P}_T$ , and in fact a probability measure also. Consequently, since the regularity of a family does not depend on the choice of the dominating measure, the assumption on  $\sigma$ -finiteness of proposition 5 is not a real restriction.

For the proof of this proposition we need the following inequality.

**Lemma 1.** Let  $(X, Y, Z)^T$  be a random vector in  $R^3$  with  $X, Y \geq 0$  a.s. and  $EX, EY$ , and  $E|Z|$  finite. Then

$$(15) \quad \left( \sqrt{EX} - \sqrt{EY} - \frac{EZ}{\sqrt{EX} + \sqrt{EY}} \right)^2 \\ \leq E \left( \sqrt{X} - \sqrt{Y} - \frac{Z}{\sqrt{X} + \sqrt{Y}} \right)^2.$$

The same inequality holds with  $EX$  replaced by  $\int f d\mu$ , etc., where  $\mu$  is a  $\sigma$ -finite measure (not necessarily a probability measure).

**Proof.** Let  $A$  and  $B$  be random variables with  $E|A| < \infty$ . By Cauchy-Schwarz

$$(a) \quad (EA 1_{[B \neq 0]})^2 = \left( E \frac{A}{B} B \right)^2 \leq E \left( \frac{A}{B} \right)^2 EB^2.$$

Since  $E|A| 1_{[B=0]} > 0$  implies  $E(A/B)^2 = \infty$ , we conclude that  $E(A/B)^2 \geq (EA)^2 / EB^2$ . In particular,

$$(b) \quad E \left( \sqrt{X} - \sqrt{Y} - \frac{Z}{\sqrt{X} + \sqrt{Y}} \right)^2 = E \left( \frac{X - Y - Z}{\sqrt{X} + \sqrt{Y}} \right)^2 \geq \frac{\{E(X - Y - Z)\}^2}{E(\sqrt{X} + \sqrt{Y})^2}.$$

Furthermore,

$$(c) \quad E(\sqrt{X} + \sqrt{Y})^2 = EX + EY + 2E\sqrt{X}\sqrt{Y} \\ \leq EX + EY + 2\sqrt{EXEY} = (\sqrt{EX} + \sqrt{EY})^2.$$

Combining (b) and (c) we obtain

$$(d) \quad E \left( \sqrt{X} - \sqrt{Y} - \frac{Z}{\sqrt{X} + \sqrt{Y}} \right)^2 \geq \left( \frac{EX - EY - EZ}{\sqrt{EX} + \sqrt{EY}} \right)^2$$

and hence the lemma. □

**Proof of proposition 5.** In view of

$$(s(\theta + h) + s(\theta))^{-1} |s(\theta + h) - s(\theta)| \leq 1,$$

we have, for  $i = 1, \dots, k$  and  $K > 0$ ,

$$(a) \quad \|\dot{s}_i(\theta) - 2(s(\theta + h) + s(\theta))^{-1} \dot{s}_i(\theta) s(\theta)\|^2 \\ = \|\dot{s}_i(\theta) (s(\theta + h) + s(\theta))^{-1} (s(\theta + h) - s(\theta))\|^2 \\ \leq \int_{|\dot{s}_i(\theta)| > K s(\theta)} \dot{s}_i^2(\theta) s^{-2}(\theta) dP_\theta + K^2 \int (s(\theta + h) - s(\theta))^2 d\mu.$$

It follows from the regularity of  $\mathbf{P}$  that  $s(\theta)$  is continuous in  $\theta$ , that  $\dot{s}_i(\theta) s^{-1}(\theta)$  belongs to  $L_2(P_\theta)$ , and hence that the left side of (a) is  $o(1)$  as  $|h| \rightarrow 0$ . Together with (2.1.2) of the definition of regularity, this yields

$$(b) \quad \|s(\theta + h) - s(\theta) - 2(s(\theta + h) + s(\theta))^{-1} s^T(\theta) h s(\theta)\| = o(|h|).$$

By proposition A.3.4 it follows that

$$(c) \quad s_T^2(\theta) = E_\mu(s^2(\theta) | T) \quad \text{a.e. } \mu T^{-1}.$$

It follows from (b), (c), and lemma 1 that

$$(d) \quad \|s_T(\theta + h) - s_T(\theta) - (s_T(\theta + h) + s_T(\theta))^{-1} E_\mu(\dot{s}^T(\theta) h | T)\| = o(|h|).$$

Denoting

$$(e) \quad \dot{s}_T(\theta) = (2s_T(\theta))^{-1} E_\mu(\dot{p}(\theta) | T)$$

and noting that

$$(f) \quad \dot{s}_{Ti}^2(\theta) \leq E_\mu(\dot{s}_i^2(\theta) | T) \quad \text{a.e.} \quad \mu T^{-1},$$

we see that  $\dot{s}_{Ti}(\theta) s_T^{-1}(\theta)$  belongs to  $L_2(P_\theta T^{-1})$  and that hence (d) implies continuity of  $s_T(\theta)$  in  $\theta$ . Consequently by the same argument as in (a) and using (f), we get

$$\begin{aligned} (g) \quad & \| \dot{s}_{Ti}(\theta) - 2(s_T(\theta + h) + s_T(\theta))^{-1} \dot{s}_{Ti}(\theta) s_T(\theta) \|^2 \\ &= \| \dot{s}_{Ti}(\theta) (s_T(\theta + h) + s_T(\theta))^{-1} (s_T(\theta + h) - s_T(\theta)) \|^2 \\ &\leq \int E_\mu(\dot{s}_i^2(\theta) | T) (s_T(\theta + h) + s_T(\theta))^{-2} (s_T(\theta + h) - s_T(\theta))^2 d\mu T^{-1} \\ &\leq \int_{|\dot{s}_i(\theta)| > K s(\theta)} \dot{s}_i^2(\theta) s^{-2}(\theta) dP_\theta + K^2 \int (s_T(\theta + h) - s_T(\theta))^2 d\mu T^{-1}. \end{aligned}$$

Together with (d) this implies

$$(h) \quad \| s_T(\theta + h) - s_T(\theta) - \dot{s}_T^T(\theta) h \| = o(|h|).$$

Furthermore, we have

$$\begin{aligned} (i) \quad & \| \dot{s}_{Ti}(\theta + h) - \dot{s}_{Ti}(\theta) \| \\ &\leq \| (s_T(\theta + h) + s_T(\theta))^{-1} E_\mu(\dot{p}_i(\theta + h) - \dot{p}_i(\theta) | T) \| \\ &\quad + \| \dot{s}_{Ti}(\theta + h) (s_T(\theta + h) + s_T(\theta))^{-1} (s_T(\theta + h) - s_T(\theta)) \| \\ &\quad + \| \dot{s}_{Ti}(\theta) (s_T(\theta + h) + s_T(\theta))^{-1} (s_T(\theta + h) - s_T(\theta)) \|. \end{aligned}$$

In view of lemma A.7.2.B, continuity of  $\dot{s}(\theta)$  in  $\theta$  implies uniform square integrability of  $\dot{s}_i(\theta + h) s^{-1}(\theta + h)$ . Arguing as in (g), we see by lemma A.7.3 that this uniform integrability implies that the second term on the right side of (i) is  $o(1)$  as  $|h| \rightarrow 0$ , while (g) itself shows that the third term is  $o(1)$ . Applying Cauchy-Schwarz twice (see (c) of the proof of lemma 1), we obtain

$$\begin{aligned} (j) \quad & (E_\mu(\dot{p}_i(\theta + h) - \dot{p}_i(\theta) | T))^2 \\ &\leq E_\mu(\{(s(\theta + h) + s(\theta))^{-1} (\dot{p}_i(\theta + h) - \dot{p}_i(\theta))\}^2 | T) \\ &\quad \cdot E_\mu((s(\theta + h) + s(\theta))^2 | T) \\ &\leq E_\mu(\{2(s(\theta + h) + s(\theta))^{-1} (\dot{s}_i(\theta + h) s(\theta + h) - \dot{s}_i(\theta) s(\theta))\}^2 | T) \\ &\quad \cdot (s_T(\theta + h) + s_T(\theta))^2, \end{aligned}$$

and by (a) with  $\theta$  and  $\theta + h$  interchanged, by the continuity of  $\dot{s}_i(\theta)$  in  $\theta$ , and by (a) itself, we obtain the convergence to 0 of the first term in the right side of (i). Consequently,  $\dot{s}_T(\theta)$  is continuous in  $\theta$ , and the proof of A is complete.

Part B follows easily from

$$a^T I(\theta, \mathbf{P})a - a^T I(\theta, \mathbf{P}_T)a = E((\dot{I}(X, \theta) - E(\dot{I}(X, \theta) | T))^T a)^2 \geq 0,$$

while C follows from A and the definition of a regular parametric model.  $\square$

### A.6 METRICS ON CLASSES OF PROBABILITY MEASURES AND PROBABILITY INEQUALITIES

Let  $(X, |\cdot|)$  be a metric space with Borel  $\sigma$ -field  $\mathcal{B}$  and let  $\{P_\theta : \theta \in \Theta\} \equiv \mathbf{P} \subset \mathbf{M}$  be a collection of probability measures on  $(X, \mathcal{B})$ . Inference about  $\theta$  is determined to a large degree by smoothness properties of the map  $\theta \rightarrow P_\theta$ . To talk about smoothness, we need topologies or a metric on  $\mathbf{M}$ . We list several metrics and their properties.

The *variational distance*  $d_v(P, Q)$  is

$$\begin{aligned} (1) \quad d_v(P, Q) &\equiv 2 \sup \{ |P(A) - Q(A)| : A \in \mathcal{B} \} \\ &= \int \left| \frac{dP}{d\mu} - \frac{dQ}{d\mu} \right| d\mu \\ &= 2 \left\{ 1 - \int \left( \frac{dP}{d\mu} \wedge \frac{dQ}{d\mu} \right) d\mu \right\} \end{aligned}$$

for any measure  $\mu$  dominating both  $P$  and  $Q$ . With this distance,  $\mathbf{P}$  can be viewed as a subset of the Banach space of finite signed measures on  $X$  with the variational norm.

The *Hellinger distance*  $d_H(P, Q)$  is the square root of

$$\begin{aligned} (2) \quad d_H^2(P, Q) &\equiv \int \left| \left( \frac{dP}{d\mu} \right)^{1/2} - \left( \frac{dQ}{d\mu} \right)^{1/2} \right|^2 d\mu \\ &= 2 \left( 1 - \int \sqrt{\frac{dP}{d\mu} \frac{dQ}{d\mu}} d\mu \right) \end{aligned}$$

for  $\mu$  dominating both  $P$  and  $Q$ . It is well-known that  $d_H$  and  $d_v$  induce the same topologies and that

$$(3) \quad d_H^2 \leq d_v \leq d_H(4 - d_H^2)^{1/2} \leq 2d_H.$$

The Hellinger metric is of particular importance for  $\mathbf{P}$  dominated by some fixed measure  $\mu$ , since the correspondence between  $P$  and  $\sqrt{dP/d\mu} \equiv s \equiv s(\cdot, P)$  then yields a representation of  $\mathbf{P}$  as a subset of  $L_2(\mu)$ , the Hilbert space of all  $\mu$  square integrable functions on  $X$ . The metric induced on  $\mathbf{P}$  by this correspondence is just  $d_H$ .

If  $X_1, \dots, X_n$  are independent with distributions  $P_1, \dots, P_n$  or  $Q_1, \dots, Q_n$  and  $P^{(n)}$  and  $Q^{(n)}$  denote the corresponding product laws, then it is easily checked that

$$(4) \quad d_H^2(P^{(n)}, Q^{(n)}) = 2 - 2 \prod_{i=1}^n \left[ 1 - \frac{1}{2} d_H^2(P_i, Q_i) \right].$$

Thus in the i.i.d. case when  $P_i = P$  and  $Q_i = Q$  for  $i = 1, \dots, n$ , it follows that

$$(5) \quad d_H^2(P^{(n)}, Q^{(n)}) = 2 \{ 1 - [1 - \frac{1}{2}d_H^2(P, Q)]^n \}$$

$$\rightarrow 0 \quad \text{as } d_H(P, Q) \rightarrow 0.$$

$$\rightarrow 2 \quad \text{as } n \rightarrow \infty \text{ if } d_H(P, Q) > 0.$$

Now let  $\delta_x$  denote the probability measure with mass one at  $x$  and let

$$(6) \quad IP_n \equiv n^{-1} \sum_{i=1}^n \delta_{X_i}$$

denote the empirical measure of the sample. It is useful to have available metrics on  $\mathbf{M}$  which are suited to  $IP_n$  in the following sense. We say that a metric  $d$  on  $\mathbf{M}$  is *compatible* with  $IP_n$  if  $d(IP_n, P)$  is a (measurable) random variable for every  $P \in \mathbf{M}$  and  $n^{1/2} d(IP_n, P) = O_P(1)$  uniformly on  $\mathbf{M}$ . That is,

$$(7) \quad \limsup_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} \sup \{ P(d(IP_n, P) \geq M n^{-1/2}) : P \in \mathbf{M} \} = 0.$$

When  $\mathbf{X} = R^d$ , simple examples of such metrics are the *Kolmogorov-Smirnov* metric

$$(8) \quad d_K(P, Q) \equiv \sup_x |F_P(x) - F_Q(x)|$$

where  $F_P$  is the distribution function of  $P$ , or the related Hilbert metrics such as

$$(9) \quad d_\mu(P, Q) \equiv \{ \int [F_P(x) - F_Q(x)]^2 d\mu(x) \}^{1/2}$$

for a finite measure  $\mu$ . The compatibility of  $d_K$  and  $d_\mu$  given by (8) and (9) follows from Kiefer (1961).

The *Prohorov* metric  $d_{pr}$  is defined by

$$(10) \quad d_{pr}(P, Q) \equiv \inf \left\{ \varepsilon > 0 : \begin{array}{l} P(A) \leq Q(A^\varepsilon) + \varepsilon, \\ Q(A) \leq P(A^\varepsilon) + \varepsilon \end{array} \quad \text{for all } A \in \mathcal{B} \right\},$$

where  $A^\varepsilon \equiv \{x : |x - y| < \varepsilon \text{ for some } y \in A\}$ . The Prohorov metric is of importance primarily because it metrizes weak convergence:  $d_{pr}(P_n, P) \rightarrow 0$  if and only if  $\int g(x) dP_n(x) \rightarrow \int g(x) dP(x)$  for all bounded and continuous functions  $g$  on  $\mathbf{X}$ . The results of Kersting (1978) imply that  $d_{pr}$  is *not* a compatible metric.

The following theorem gives an extremely useful way of interpreting the Prohorov metric  $d_{pr}$ .

**Theorem 1.** (Strassen). Let  $\varepsilon > 0$ , and let  $P, Q \in \mathbf{M}$ . Then the following are equivalent:

- (i)  $d_{pr}(P, Q) \leq \varepsilon$ .
- (ii)  $X \sim P, Y \sim Q$  can be defined on a common probability space  $(\mathbf{X} \times \mathbf{X}, \mathcal{B} \times \mathcal{B}, Pr)$  so that

$$(11) \quad Pr(|X - Y| \leq \varepsilon) \geq 1 - \varepsilon.$$

**Proof.** See Strassen (1965), Dudley (1968), (1989), and Schay (1974).  $\square$

Now let  $1 \leq r < \infty$ , suppose that  $X$  is a separable Banach space with norm  $\| \cdot \|$ , and let

$$(12) \quad P_r \equiv \{ P \in \mathcal{P}: E \|X\|^r < \infty \text{ where } X \sim P \}.$$

The *Mallows metric*  $d_r$  is defined for  $P, Q \in P_r$  by

$$(13) \quad d_r(P, Q) \equiv \inf\{ E \|X - Y\|^r : X \sim P, Y \sim Q \}.$$

It is easily shown that the infimum in (13) is obtained; see Bickel and Freedman (1981). When  $X = R$  and  $\|x\| = |x|$ , then

$$(14) \quad d_r(P, Q) = \int_0^1 |F_P^{-1}(t) - F_Q^{-1}(t)|^r dt,$$

and furthermore

$$(15) \quad d_1(P, Q) = \int_0^1 |F_P^{-1}(t) - F_Q^{-1}(t)| dt = \int_{-\infty}^{\infty} |F_P(x) - F_Q(x)| dx.$$

**Proposition 1.** If  $P_n, P \in P_r, r \geq 1$ , then the following are equivalent:

- (i)  $d_r(P_n, P) \rightarrow 0$ .
- (ii)  $P_n \rightarrow P$  weakly, and  $\int \|x\|^r dP_n(x) \rightarrow \int \|x\|^r dP(x)$ .
- (iii)  $P_n \rightarrow P$  weakly, and  $\|x\|^r$  is  $P_n$ -uniformly integrable.

**Proof.** See Bickel and Freedman (1981).  $\square$

### Probability Inequalities

In the conclusion of this section we state several useful probability inequalities.

An inequality which is used repeatedly to relate our convolution theorems in chapters 2, 3, and 5 to comparisons of probabilities of convex sets is the following basic inequality due to Anderson (1955), which depends in turn on the Brunn-Minkowski inequality:

**Lemma 1.** (Anderson (1955)). Suppose that  $X$  has a distribution on  $R^d$  which is symmetric about zero and unimodal, and that  $Y$  on  $R^d$  is independent of  $X$  (but otherwise arbitrary). Then

$$(16) \quad P(X + Y \in C) \leq P(X \in C)$$

for every measurable set  $C$  in  $R^d$  which is convex and symmetric about 0.

**Proof.** See Anderson (1955, page 172, theorem 2); Ibragimov and Has'minskii (1981, lemma II.10.1, page 155); or Das Gupta (1980) for an interesting discussion of this inequality and related results.  $\square$

**Corollary 1.** Anderson's inequality holds for  $X \sim N_d(0, \Sigma)$  for any  $\Sigma$ .

Recall that  $l : R^d \rightarrow R^+$  is called *subconvex* (or bowl-shaped) if

$$l(0) = 0 \leq l(x) \text{ for every } x \in R^d$$

$$l(x) = l(-x).$$

$$\{x : l(x) \leq c\} \text{ is convex and closed for every } c \in R.$$



**Corollary 2.** Suppose that  $l : R^d \rightarrow R^+$  is subconvex. Suppose that  $X$  satisfies the hypotheses of lemma 1. Then

$$(17) \quad E l(X + y) \geq E l(X) \quad \text{for all } y \in R^d.$$

**Proof.** See Ibragimov and Has'minskii (1981, lemma II.10.2, page 157).  $\square$

Now we state a useful inequality for the empirical process due to Alexander (1984), (1987). Let  $\mathbf{F}$  be a collection of functions on  $R^d$  such that  $\|f\|_\infty \leq 1$  for all  $f \in \mathbf{F}$ . Define the metric entropy  $H(u, \mathbf{F})$  of  $\mathbf{F}$  by

$$(18) \quad H(u, \mathbf{F}) = \log N(u, \mathbf{F})$$

where

$$N(u, \mathbf{F}) = \min\{k : \text{there exist } f_1, \dots, f_k \in \mathbf{F} \text{ such that} \\ \min_{1 \leq i \leq k} \|f - f_i\|_\infty < u \text{ for all } f \in \mathbf{F}\}.$$

Let  $P$  be a probability measure on  $R^d$ . Suppose that  $X_1, X_2, \dots, X_n$  are i.i.d. with distribution  $P$ , let  $\mathbb{P}_n$  denote the empirical measure of the (first  $n$ )  $X_i$ 's, and let  $\mathbf{X}_n \equiv \sqrt{n}(\mathbb{P}_n - P)$  denote the empirical process. The following inequality is from Alexander (1984, corollary 2.2), as corrected in Alexander (1987).

Let  $\alpha \equiv \sup\{\text{Var}(f(X_1)) : f \in \mathbf{F}\}$ ,  $0 < r < 2$ , and suppose that for a universal constant  $A$ ,

$$(19) \quad H(u, \mathbf{F}) \leq Au^{-r}, \quad \text{for all } u > 0.$$

**Lemma 2.** (Alexander). Suppose that (19) holds,  $\mathbf{F}$  is as above, and  $\varepsilon > 0$ ,  $n > 2$ . Then there exists a constant  $K = K(r, \varepsilon, A)$  so that, if

$$(20) \quad \lambda \geq K\{\alpha^{1/2-r/4} \vee n^{(r-2)/2(r+2)}\},$$

then

$$(21) \quad Pr^*\{\|\mathbf{X}_n\|_{\mathbf{F}} \geq \lambda\} = Pr^*\{\sup\{|n^{-1/2} \sum_{i=1}^n (f(X_i) - \int f dP)| : \\ f \in \mathbf{F}\} \geq \lambda\} \\ \leq 5 \exp\left(- (1-\varepsilon) \frac{\lambda^2}{2\alpha} \left(1 + \frac{\lambda}{3n^{1/2}\alpha}\right)^{-1}\right).$$

The following corollary is immediate:

**Corollary 3.** Suppose that  $\alpha = \alpha_n$ ,  $\lambda = \mu_n$  so that  $\mathbf{F} = \mathbf{F}_n = \mathbf{F}_n(\alpha_n, r)$  and, in addition to (19) and (20), we have

$$(22) \quad \alpha_n \rightarrow 0, \quad \text{and} \quad \frac{\mu_n}{n^{1/2}\alpha_n} \rightarrow 0.$$

Then, for  $n \geq$  some  $N_\varepsilon$ ,

$$(23) \quad Pr^*\{\|\mathbf{X}_n\|_{\mathbf{F}_n} \geq \mu_n\} \leq 5 \exp\left(- (1-2\varepsilon) \frac{\mu_n^2}{2\alpha_n}\right).$$

**Proof.** By the second part of (22),  $(1 + \mu_n/3n^{1/2}\alpha_n) \leq 1 + \varepsilon$  for  $n \geq N_\varepsilon$ . Then (23) follows since  $(1 - \varepsilon)/(1 + \varepsilon) > 1 - 2\varepsilon$ . Note that it follows from (20) that  $\mu_n^2/\alpha_n \geq K^2\alpha_n^{-r/2} \rightarrow \infty$  by the first part of (22) and  $r > 0$ .  $\square$

**Remark 1.** As in example 1, page 1045, of Alexander (1984), let  $F = \{f: [0,1]^d \rightarrow R : \|f\|_\infty \leq C, \|D^p f\|_\infty \leq C, 1 \leq |p| \leq m\}$  where  $D^p f$  is a partial derivative of order  $|p|$  and  $p = (p_1, \dots, p_d) \in (\mathbb{Z}^+)^d$ ,  $|p| \equiv \sum_{j=1}^d p_j$ , and  $m > d/2$ . Then (19) holds with  $r = d/m$ .

## A.7 LIMIT THEOREMS, WEAK CONVERGENCE, AND TIGHTNESS

We begin with several results connected with uniform integrability and convergence in distribution of random variables and vectors  $X_n$  to  $X$  with distributions  $P_n$  and  $P$  respectively. Whenever necessary we will assume that  $X_n$  and  $X$  are defined on the same probability space.

### Definition 1.

A.  $X_n \rightarrow_d X$  or  $L(X_n) \rightarrow L(X)$ ; i.e.,  $X_n$  converges in distribution (law) to  $X$ , if and only if for all  $g$  bounded and continuous

$$E g(X_n) \rightarrow E g(X)$$

as  $n \rightarrow \infty$ . If the  $X_n$  are random variables, taking values in  $R$ , this is equivalent to

$$F_n(x) \equiv P(X_n \leq x) \rightarrow P(X \leq x) = F(x)$$

at all continuity points  $x$  of the distribution function  $F$ .

B.  $X_n \rightarrow_p X$ ; i.e.,  $X_n$  converges to  $X$  in probability, if and only if for every  $\varepsilon > 0$

$$P(|X_n - X| > \varepsilon) \rightarrow 0$$

as  $n \rightarrow \infty$ .

C.  $X_n \rightarrow_r X$ ; i.e.,  $X_n$  converges in the  $r$ th mean to  $X$ , if and only if  $E|X|^r < \infty$  and

$$E|X_n - X|^r \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

**Definition 2.** A sequence of random vectors  $\{X_n\}$  is *uniformly integrable* if and only if

$$(1) \quad \lim_{\lambda \rightarrow \infty} \limsup_{n \rightarrow \infty} E|X_n| 1_{\{|X_n| \geq \lambda\}} = 0.$$

**Lemma 1.** ( $L_r$ -convergence theorem) Let  $r > 0$ . Then  $X_n \rightarrow_r X$  if and only if

$$(i) \quad X_n \rightarrow_p X,$$

$$(ii) \quad E|X_n|^r \rightarrow E|X|^r.$$

**Proof.** See Loève (1977, section 9.4.C, page 165). □

**Lemma 2.** Suppose  $X_n \rightarrow_d X$ .

- A. If  $\{X_n\}$  is uniformly integrable, then  $EX_n \rightarrow EX$ .
- B. If  $X_n \geq 0, X \geq 0$ , with  $E|X| < \infty$  and  $EX_n \rightarrow EX$ , then  $\{X_n\}$  is uniformly integrable.

**Proof.** See theorem 5.4, page 32 of Billingsley (1968). □

**Lemma 3.** Suppose  $X_n \rightarrow_d X$  and  $|X_n| \leq V_n$  where  $\{V_n\}$  is uniformly integrable. Then  $EX_n \rightarrow EX$ .

**Proof.** Since  $|x|1_{|x| \geq \lambda}$  is increasing in  $|x|$ , uniform integrability of  $V_n$  implies it for  $X_n$  and lemma 2.A applies. □

**Lemma 4.** Suppose that:

- (i)  $X_n \rightarrow_d X$  and  $X_n - Y_n \rightarrow_p 0$  (or equivalently  $(X_n, Y_n) \rightarrow_d (X, X)$ ).
- (ii) For some  $r \geq 1$ ,  $E|X|^r < \infty$  and both  $E|X_n|^r \rightarrow E|X|^r$  and  $E|Y_n|^r \rightarrow E|X|^r$ .

Then  $X_n - Y_n \rightarrow_r 0$ .

**Proof.** By (i), (ii), and lemma 2.B,  $\{|X_n|^r + |Y_n|^r\}$  are uniformly integrable. But  $|X_n - Y_n|^r \leq 2^{r-1} \{|X_n|^r + |Y_n|^r\}$ , and hence the conclusion follows from lemma 3. □

Combining lemmas 1 through 3 we arrive at:

**Theorem 1.** If  $X_n \rightarrow_p X$ ,  $|X_n|^r \leq Y_n$  where  $Y_n \rightarrow_p Y$  and  $EY_n \rightarrow EY < \infty$ , then  $X_n \rightarrow_r X$ . Conversely, if  $X_n \rightarrow_r X$ , then  $Y_n = |X_n|^r$  satisfies the hypotheses of the direct part of the theorem.

**Proof.** By lemma 2.B  $\{Y_n\}$  is uniformly integrable and hence by lemma 3  $E|X_n|^r \rightarrow E|X|^r$ . Lemma 1 yields the result. □

These results can be generalized by considering vector-valued measurable functions  $\{f_n\}$  on some measure space  $(X, \mathcal{A}, \mu)$  where  $\mu$  is a  $\sigma$ -finite measure on  $\mathcal{A}$ .

**Definition 3.**

- A.  $f_n \rightarrow_\mu f$  if and only if  $\mu\{|f_n - f| \geq \varepsilon\} \rightarrow 0$ .
- B.  $f_n \rightarrow_r f$  if and only if  $\int |f_n - f|^r d\mu \rightarrow 0$ .

Note that evidently  $f_n \rightarrow_r f$  implies  $f_n \rightarrow_\mu f$ . Furthermore this definition of convergence in  $\mu$ -measure agrees with section 6.3, page 116, of Loève (1977), but a weaker version of it is given in definition 2.11.2, page 93, of Bauer (1972).

**Theorem 2.** Suppose that  $f_n \rightarrow_\mu f$  and  $|f_n|^r \leq g_n$  where  $g_n \rightarrow_\mu g$  and  $\int g_n d\mu \rightarrow \int g d\mu < \infty$ . Then  $f_n \rightarrow_r f$ .

**Proof.** See sections 7.2.C and 9.4.B of Loève (1977, pages 126, 164, 165, respectively). □

The next lemma, which is related to lemma 1, is sometimes called Vitali's theorem.

**Lemma 5.** If  $f_n \rightarrow f$  a.e.  $\mu$ , and

$$\limsup_{n \rightarrow \infty} \int |f_n|^r d\mu \leq \int |f|^r d\mu < \infty$$

for  $r > 0$ , then

$$\int |f_n - f|^r d\mu \rightarrow 0$$

as  $n \rightarrow \infty$ .

**Proof.** See Novinger (1972) and theorems II.4.2 and V.1.3 of Hájek and Šidák (1967).  $\square$

**Lemma 6.** If  $n^{-1/2} \sum_{i=1}^n Z_{ni} \rightarrow_p 0$  with  $Z_{n1}, \dots, Z_{nn}$  i.i.d. random variables, then  $Z_{n1} \rightarrow_p 0$ .

**Proof.** Since the  $Z_{ni}$ 's are i.i.d. the convergence of the characteristic function of  $n^{-1/2} \sum_{i=1}^n Z_{ni}$  to 1 implies

$$n(Ee^{itn^{-1/2}Z_{n1}} - 1) \rightarrow 0,$$

and hence the (uniform) asymptotic negligibility condition

$$P(n^{-1/2}|Z_{n1}| \geq \tau) \rightarrow 0.$$

By the degenerate convergence theorem (Loève (1977, section 23.5, page 329)) we obtain

$$nE(n^{-1/2}Z_{n1})^2 1_{[|n^{-1/2}Z_{n1}| < \tau]} \rightarrow 0.$$

Consequently

$$\begin{aligned} P(|Z_{n1}| \geq \varepsilon) &\leq P(|Z_{n1}| \geq \tau\sqrt{n}) \\ &\quad + \varepsilon^{-2} E(Z_{n1}^2 1_{[|Z_{n1}| < \tau\sqrt{n}]} \rightarrow 0. \end{aligned} \quad \square$$

### Some Uniform in $P$ Limit Theorems

Suppose that  $X_1, \dots, X_n, \dots$  are i.i.d. real-valued random variables with probability distribution  $P \in \mathbf{M} = \{\text{all probability distributions on } R\}$ . Let  $\mu(P) = \int x dP(x)$  and  $\sigma^2(P) \equiv \int (x - \mu(P))^2 dP(x)$ , when they exist.

**Theorem 3.** (Chung's uniform strong law of large numbers). Let  $\mathbf{P} \subset \mathbf{M}$  be any collection of probability distributions such that

$$(2) \quad \lim_{\lambda \rightarrow \infty} \sup \{E_P |X| 1_{[|X| \geq \lambda]} : P \in \mathbf{P}\} = 0.$$

Then for every  $\varepsilon > 0$

$$(3) \quad \sup \{P(\sup_{m \geq n} |\bar{X}_m - \mu(P)| > \varepsilon) : P \in \mathbf{P}\} \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

**Proof.** See Chung (1951).  $\square$

This uniform law of large numbers should be compared with the definition of uniform consistency given in section 2.2.

**Theorem 4.** (Petrov). Suppose  $\mu(P) = 0$  and  $\sigma^2(P) = 1$ , and let  $S_n \equiv X_1 + \dots + X_n$ . Then there is an absolute constant  $c > 0$  such that for all  $\varepsilon > 0, t \in R, n = 1, 2, \dots$ ,

$$(4) \quad |P\left(\frac{S_n}{\sqrt{n}} \leq t\right) - \Phi(t)| \leq c \{E_P(X^2 1_{[|X| \geq \varepsilon\sqrt{n}]}) + \varepsilon\}.$$

Hence if  $\mathbf{P}$  is a family of probability measures with  $\mu(P) = 0$  and  $\sigma^2(P) = 1$  for all  $P \in \mathbf{P}$ , and

$$(5) \quad \limsup_{\lambda \rightarrow \infty} \{E_P(X^2 1_{\{|X| \geq \lambda\}}) : P \in \mathbf{P}\} = 0,$$

then

$$(6) \quad \sup \left\{ \sup_{t \in \mathcal{K}} \left| P(n^{-1/2} S_n \leq t) - \Phi(t) \right| : P \in \mathbf{P} \right\} \rightarrow 0$$

as  $n \rightarrow \infty$ .

**Proof.** This follows from theorem 8, page 118, of Petrov (1975). □

Theorem 4 should be compared with definition 2.2.2 of a uniformly asymptotically normal estimate.

### Weak Convergence and Tightness Criteria

Suppose that  $(M, d)$  is a metric space equipped with its Borel  $\sigma$ -field  $\mathcal{M}_B$ . Let  $C_b(M)$  denote the collection of real-valued, bounded, continuous functions on  $M$ .

**Definition 4.** If  $\{P_n\}, P$  are probability measures on  $(M, \mathcal{M}_B)$  satisfying

$$\int f dP_n \rightarrow \int f dP \quad \text{as } n \rightarrow \infty \quad \text{for all } f \in C_b(M),$$

then we say that  $P_n$  converges in distribution (or law) to  $P$ , or that  $P_n$  converges weakly to  $P$ , and we write  $P_n \rightarrow_d P$ . Similarly, if  $\mathbf{X}_n, \mathbf{X}$  are random elements in  $M$  with

$$E f(\mathbf{X}_n) \rightarrow E f(\mathbf{X}) \quad \text{for all } f \in C_b(M),$$

then we write  $\mathbf{X}_n \rightarrow_d \mathbf{X}$ .

On the real line there are only two ways for mass to “escape”: it can only leave the space by going off to  $\pm\infty$ . For measures on a general metric space, however, there are more ways to leave the space, and hence the following notions of tightness and relative compactness play key roles.

**Definition 5.** Let  $\mathbf{P}$  be a collection of probability measures on  $(M, \mathcal{M}_B)$ . Then  $\mathbf{P}$  is *tight* (or *uniformly tight*) if and only if for every  $\varepsilon > 0$  there is a compact set  $K \equiv K_\varepsilon \subset M$  with

$$P(K) > 1 - \varepsilon \quad \text{for all } P \in \mathbf{P}.$$

**Definition 6.** Let  $\mathbf{P}$  be a family of probability measures on  $(M, \mathcal{M}_B)$ . We call  $\mathbf{P}$  *relatively compact* if every sequence  $\{P_n\} \subset \mathbf{P}$  contains a weakly convergent subsequence. Thus every  $\{P_n\} \subset \mathbf{P}$  contains a subsequence  $\{P_{n'}\}$  with  $P_{n'} \rightarrow$  some  $Q$  (not necessarily in  $\mathbf{P}$ ).

**Proposition 1.** Let  $(M, d)$  be a separable metric space.

- A. (Le Cam). If  $P_n \rightarrow_d P$ , then  $\{P_n\}$  is tight.
- B. If  $P_n \rightarrow_d P$ , then  $\{P_n\}$  is relatively compact.

- C. If  $\{P_n\}$  is relatively compact and the set of limit points is just the one-point set  $\{P\}$ , then  $P_n \rightarrow_d P$ .

The key theorem relating tightness and relative compactness is due to Prohorov (1956) in the case of  $(M, d)$  complete and separable.

**Theorem 5.** (Prohorov). Let  $\mathbf{P}$  be a collection of probability measures on  $(M, \mathcal{M}_P)$ .

- A. If  $\mathbf{P}$  is tight, it is relatively compact.  
 B. Suppose  $(M, d)$  is separable and complete. If  $\mathbf{P}$  is relatively compact, it is tight.

**Proof.** See Billingsley (1968, section 6, pages 35–40) for the proof given by Varadarajan (1958), (1961). The proof progresses from  $M = R^k$  to  $R^\infty$  to  $M$   $\sigma$ -compact, and finally to general  $M$ , at each step using the next simpler case.  $\square$

Once weak convergence of a sequence of measures or processes has been established, it is often convenient to work in terms of almost surely convergent versions of those processes, the existence of which is guaranteed by the following theorem.

**Theorem 6.** (Skorokhod (1956)). Suppose that  $(M, d)$  is complete and separable and the measures  $\{P_n: n \geq 0\}$  thereon satisfy  $P_n \rightarrow_d P_0$ . Then there exist random elements  $\{X_n: n \geq 0\}$  with values in  $M$  ( $X_n(\omega) \in M$  for all  $\omega \in \Omega$ ) all defined on the common probability space  $(\Omega, \mathcal{A}, P) \equiv ([0, 1], \mathcal{B}([0, 1]), \text{Lebesgue})$  with  $X_n \sim P_n$  and satisfying  $d(X_n(\omega), X_0(\omega)) \rightarrow 0$  for each  $\omega \in \Omega$ .

Theorem 6 has been extended in several ways: Dudley (1968) removed the hypothesis of completeness of  $M$ , and Wichura (1970) proved an extension to nonseparable  $M$ ; see, e.g., Shorack and Wellner (1986, theorem 2.3.4, page 47). In appendix section 8 we define “weak convergence” still more generally, for possibly nonmeasurable random functions, and give a further extension of theorem 6 in the extended theory due to Dudley (1985).

The following useful theorem is an easy consequence of Skorokhod’s theorem 6.

**Theorem 7.** (Continuous mapping theorem). Let  $M$  be separable. Suppose  $X_n \rightarrow_d X$  in  $M$  and suppose  $g: M \rightarrow M'$  (where  $(M', d')$  is another metric space) is continuous a.e.  $P = P_X$ ; i.e.,  $P(X \in C(g)) = 1$  where  $C(g) = \{x \in M: g \text{ is continuous at } x\}$ . Then  $g(X_n) \rightarrow_d g(X)$ .

The continuous mapping theorem has several useful corollaries:

**Corollary 1.** (Slutsky’s theorem). Suppose that multiplication by reals and addition in  $M$  are well defined, and that  $A_n \rightarrow_p a \in R$ ,  $B_n \rightarrow_p b \in M$ , and  $X_n \rightarrow_d X$  in  $(M, d)$ . Then  $A_n X_n + B_n \rightarrow_d aX + b$ .

**Corollary 2.** Suppose that  $(M, d)$  and  $(M', d')$  are vector spaces, that  $T: M \rightarrow M'$  is a (fixed) bounded linear operator, and that  $X_n \rightarrow_d X$ . Then  $TX_n \rightarrow_d TX$ .

**Corollary 3.** Suppose that  $(M, d)$  and  $(M', d')$  are vector spaces, that  $\{T_n\}$  is a sequence of bounded operators from  $M$  to  $M'$  satisfying  $\|T_n - T\| \rightarrow 0$

(where  $\|\cdot\|$  denotes the operator norm),  $\limsup_{n \rightarrow \infty} \|T_n\| < \infty$ , and that  $X_n \rightarrow_d X$  in  $(M, d)$ . Then  $T_n X_n \rightarrow_d TX$ .

Now we will establish a useful uniform convergence result connected with convexity. Let  $\mathbf{W}$  be the class of all functions  $w : N \rightarrow R^m$  such that

- (i) For all  $u \in R^m$ ,  $t \in N$ , the map  $\lambda \rightarrow u^T w(t + \lambda u)$  is monotone nondecreasing.
- (ii)  $w(t) = 0$  for a unique  $t \in N$ .

**Theorem 8.** Suppose that  $\{U_n(\cdot)\}$  is a sequence of random vector fields from  $R^m$  to  $R^m$  such that

- (i)  $U_n(\cdot) \in \mathbf{W}$  with probability 1.
- (ii)  $U_n(t) \rightarrow_p U(t)$  for all  $t \in R^m$ .
- (iii)  $U$  is continuous with probability 1.

Then, for all  $M > 0$ ,

$$\sup\{|U_n(t) - U(t)| : |t| \leq M\} \rightarrow_p 0 \quad \text{as } n \rightarrow \infty.$$

The proof is based on the following lemma:

**Lemma 7.** Suppose that  $w \in \mathbf{W}$ . Let  $e_1, \dots, e_m$  be an orthonormal system in  $R^m$  and let  $K$  be a cube in  $R^m$  with center  $t_0$  and vertices

$$t_i \equiv t_0 + a \sum_{j=1}^m d_{ij} e_j, \quad i = 1, \dots, 2^m,$$

with  $d_{ij} \in \{-1, +1\}$ ,  $j = 1, \dots, m$ . Let  $K'$  be a similar cube with center  $t_0$  and vertices  $t_0 + b \sum_{i=1}^m d_{ij} e_j$  where  $0 < b < a$ . Then, for all  $t \in K'$  and all  $t^* \in \text{domain}(w)$ ,

$$\begin{aligned} & |w(t) - w(t^*)| \\ & \leq \frac{a+b}{a-b} \sqrt{m} \max\{|w(t_i) - w(t^*)| : i = 1, \dots, 2^m\}. \end{aligned}$$

**Proof.** Write  $t = (t_1, \dots, t_m)^T$ ,  $t_i = (t_{i1}, \dots, t_{im})^T$ , and  $w = (w_1, \dots, w_m)^T$ . By monotonicity of  $(t - t_i)^T w(t + \lambda(t - t_i))$ , it follows that

$$(t - t_i)^T (w(t) - w(t_i)) \geq 0 \quad \text{for } i = 1, \dots, 2^m;$$

or equivalently

$$(a) \quad - \sum_{j=1}^m (t_j - t_{ij})(w_j(t) - w_j(t^*)) \leq \sum_{j=1}^m (t_j - t_{ij})(w_j(t^*) - w_j(t_i)).$$

By the definition of vertices, there exists an  $i$  such that

$$(b) \quad (t_j - t_{ij})(w_j(t) - w_j(t^*)) \leq 0 \quad \text{for } j = 1, \dots, m.$$

Fix  $i$  at this value. Combining (a) and (b) we obtain

$$(c) \quad \sum_{j=1}^m |t_j - t_{ij}| |w_j(t) - w_j(t^*)| \leq \sum_{j=1}^m |t_j - t_{ij}| |w_j(t^*) - w_j(t_i)|.$$

Since  $t \in K'$ ,  $a - b \leq |t_j - t_{ij}| \leq a + b$ . Therefore (c) implies

$$\begin{aligned} |w(t) - w(t^*)| &\leq \sum_{j=1}^m |w_j(t) - w_j(t^*)| \\ &\leq (a - b)^{-1} \sum_{j=1}^m |t_j - t_{ij}| |w_j(t) - w_j(t^*)| \\ &\leq (a - b)^{-1} \sum_{j=1}^m |t_j - t_{ij}| |w_j(t_i) - w_j(t^*)| \\ &\leq \frac{a + b}{a - b} \sum_{j=1}^m |w_j(t_i) - w_j(t^*)| \\ &\leq \frac{a + b}{a - b} \sqrt{m} |w(t_i) - w(t^*)|. \end{aligned}$$

□

**Proof of theorem 8.** Let  $K$  be a cube of side  $A$ . Given  $\varepsilon_1, \varepsilon_2 > 0$ , decompose  $K$  into  $k = [A/\delta]^m$  cubes  $K_j$  of side  $\delta < \varepsilon_1$  such that

$$(a) \quad P(B) \geq 1 - \varepsilon_2,$$

where

$$B \equiv \{ \sup\{ |U(t) - U(s)| : s, t \in K, |t - s| \leq 3\sqrt{m}\delta \} \leq \varepsilon_1 \}.$$

This is possible by (iii). For each cube  $K_j$ , let  $\bar{K}_j$  be the cube generated by  $K_j$  and its  $3^m - 1$  neighbors. Let  $t_1$  be a vertex of  $\bar{K}_j$ . By lemma 7 with  $a = 3\delta$ ,  $b = \delta$ , it follows that, on  $B$ ,

$$(b) \quad \begin{aligned} \sup\{ |U_n(t) - U_n(t_1)| : t \in K_j \} \\ \leq 2\sqrt{m} \max\{ |U_n(t_i) - U_n(t_1)| : t_i \text{ a vertex of } \bar{K}_j \} \end{aligned}$$

and hence, on  $B$ ,

$$\begin{aligned} \sup\{ |U_n(t) - U(t)| : t \in K_j \} \\ \leq \sup\{ |U_n(t) - U_n(t_1)| : t \in K_j \} \\ \quad + |U_n(t_1) - U(t_1)| + \varepsilon_1 \\ (c) \quad \leq 2\sqrt{m} \max\{ |U_n(t_i) - U(t_i)| + |U(t_i) - U(t_1)| \\ \quad + |U(t_1) - U_n(t_1)| : t_i \text{ a vertex of } \bar{K}_j \} \\ \quad + |U_n(t_1) - U(t_1)| + \varepsilon_1 \\ \leq (2\sqrt{m} + 1)\varepsilon_1 \\ \quad + (4\sqrt{m} + 1) \max\{ |U_n(t_i) - U(t_i)| : t_i \text{ a vertex of } \bar{K}_j \}. \end{aligned}$$

Apply (ii) to conclude from (b) that



$$\begin{aligned}
 &P(\sup\{|U_n(t) - U(t)| : t \in K\} \geq (6\sqrt{m} + 2)\varepsilon_1) \\
 &\leq \varepsilon_2 \\
 &\quad + \sum_{j=1}^k P(\max\{|U_n(t_j) - U(t_j)| : t_j \text{ is a vertex of } \bar{K}_j\} \geq \varepsilon_1) \\
 &\rightarrow \varepsilon_2 \quad \text{as } n \rightarrow \infty.
 \end{aligned}$$

Since  $\varepsilon_1$ ,  $\varepsilon_2$ , and  $A$  are arbitrary, the proposition follows. □

### A.8 HOFFMANN-JØRGENSEN-DUDLEY WEAK CONVERGENCE THEORY

#### *Definitions and Basic Results*

Often the random functions with which we work are, unfortunately, *not* measurable with respect to the Borel  $\sigma$ -field  $\mathcal{M}_B$  of the metric space  $M$ , and hence do not induce probability measures thereon. This typically occurs when  $M$  is nonseparable; for example  $D[0, 1]$  with the supremum (or uniform) metric  $\|\cdot\|_\infty$ , or  $l^\infty(\mathbb{F})$  with the supremum metric  $\|\cdot\|_F$ . The theory we outline below, due to Hoffmann-Jørgensen (1984) and Dudley (1985) following an evolution from Dudley (1966), gives up the goal of inducing distributions on  $M$  equipped with some  $\sigma$ -field of subsets. It gives a theory of “weak convergence of laws without laws being defined”—except asymptotically. This section is adapted from Van der Vaart and Wellner (1990).

Suppose that  $(\Omega, \mathcal{A}, P)$  is a probability space and  $f : \Omega \rightarrow R$  is a (completely arbitrary) function.

**Definition 1.**  $f^*$  denotes any measurable function from  $(\Omega, \mathcal{A})$  to  $(\bar{R}, \mathcal{B})$  such that:

- (i)  $f^* \geq f$  a.s.
- (ii) If  $h \geq f$  and  $h$  is measurable, then  $h \geq f^*$  a.s.

We will show below that the *measurable covering function*  $f^*$  of  $f$  exists and is unique. Define

$$(1) \quad f_* \equiv -((-f)^*).$$

Then  $f_* : (\Omega, \mathcal{A}) \rightarrow \bar{R}$  is measurable with:

- (iii)  $f_* \leq f$  a.s.
- (iv) If  $h \leq f$  and  $h$  is measurable, then  $h \leq f_*$  a.s.

We summarize (i) and (ii), and (iii) and (iv), respectively, by the notation

$$(2) \quad \begin{aligned} f^* &\equiv \text{essinf}\{h : h \geq f, h \text{ measurable}\}, \\ f_* &\equiv \text{esssup}\{h : h \leq f, h \text{ measurable}\}. \end{aligned}$$

This is essentially the same definition of essential supremum as used in the sequential analysis literature; see, e.g., Chow, Robbins, and Siegmund (1971, page 8). Let

$$(3) \quad \begin{aligned} E_P^* f &\equiv \inf \{ E_P h : f \leq h \text{ and } E_P h \text{ is defined} \}, \\ E_{*P} f &\equiv \sup \{ E_P h : f \geq h \text{ and } E_P h \text{ is defined} \}; \end{aligned}$$

here “ $E_P h$  is defined” means that  $h$  is measurable and at least one of  $h^+$  and  $h^-$  is integrable.  $E_P^* f$  and  $E_{*P} f$  are of interest because the definition of weak convergence which we will develop here is stated in terms of  $E_P^* f$ . The measurable covering function is of interest because  $E^* f = E f^*$ . In fact, it is shown in lemma 2 below that

$$(4) \quad E_P^* f = E_P f^* \quad \text{and} \quad E_{*P} f = E_P f_*$$

whenever the right sides are defined. We also define

$$\begin{aligned} P^*(A) &\equiv \inf \{ P(B) : B \supset A, B \in \mathcal{A} \}, \\ P_*(A) &\equiv \sup \{ P(B) : B \subset A, B \in \mathcal{A} \}. \end{aligned}$$

Furthermore, let  $A^*$  denote any  $\mathcal{A}$ -measurable set such that

- (i)  $A^* \supset A$ .
- (ii) If  $B \supset A$ , and  $B$  is measurable, then  $B \supset A^*$ .

We then define  $A_* \equiv ((A^c)^*)^c$ .

Our first two lemmas summarize the basic properties of  $f^*$ ,  $f_*$ ,  $A^*$ , and  $A_*$ . All the proofs for this section are given at the end of the section. The following lemma 1 and A, B, and H of the second lemma are given in Dudley and Philipp (1983), while K of lemma 2 is in Dudley (1984, lemma 3.1.6).

**Lemma 1.**  $f^*$  exists; moreover we can choose  $f^* \geq f$  everywhere. If  $f$  is measurable, then  $f^* = f$ .

**Lemma 2.** (Some facts about  $f^*$  and  $f_*$ )

- A.  $f_* + g_* \leq (f + g)_* \leq f^* + g_* \leq (f + g)^* \leq f^* + g^*$  a.s.; here the second and third inequality hold only on the set where  $f^* + g_*$  is well defined.
- B.  $(f \pm g)^* = f^* \pm g$  a.s., and  $(f \pm g)_* = f_* \pm g$  a.s. if  $g$  is measurable.
- C. If  $f + g$  is measurable, then  $f + g = f^* + g_*$  a.s.
- D.  $|f^* - g^*| \leq |f - g|^*$  a.s. whenever both sides are defined a.s.
- E.  $|f_* - g_*| \leq |f - g|^*$  a.s. whenever both sides are defined a.s.
- F. If  $g : \Omega \rightarrow R$  is measurable, then  $(fg)^* = f^* g 1_{[g > 0]} + f_* g 1_{[g < 0]}$  a.s.
- G. For any  $A \subset \Omega$ ,  $1_A^* = 1_{A^*}$ ,  $1_{A_*} = (1_A)_*$  a.s.
- H.  $(A^*)^c = (A^c)_*$ .
- I.  $E_P^* f = E_P f^*$  if  $E_P f^*$  is defined; otherwise  $E_P^* f = \infty$ .  $E_{*P} f = E_P f_*$  if  $E_P f_*$  is defined; otherwise  $E_{*P} f = -\infty$ .

- J.  $P^*(A) = E1_A^* = P(A^*); P_*(A) = E(1_A)_* = P(A_*)$ .
- K.  $[f > c]^* = [f^* > c]$  a.s. and  $P^*(f > c) = P(f^* > c), c \in R$ .
- L.  $[f \geq c]_* = [f_* \geq c]$  a.s. and  $P_*(f \geq c) = P(f_* \geq c), c \in R$ .
- M. For any index set  $I, \sup_{i \in I} f_i^* \leq (\sup_{i \in I} f_i)^*$  a.s. with equality if  $I$  is countable.

Suppose that  $(M, d)$  is a metric space (nonseparable in general),  $\{(\mathbf{X}_n, \mathcal{A}_n, P_n)\}_{n \geq 0}$  is a sequence of probability spaces, and

$$(5) \quad \mathbf{X}_n : \mathbf{X}_n \rightarrow M, \quad \text{for } n = 0, 1, 2, \dots$$

are arbitrary maps. Let  $C_b(M)$  be the collection of bounded, continuous functions  $h$  from  $M$  to  $R$ . In the following we write  $E^*h(\mathbf{X}_n)$  for  $E_{P_n}^*h(\mathbf{X}_n)$  where  $h$  is a function from  $M$  to  $R$ ; and, similarly, we write  $P^*(\mathbf{X}_n \in K)$  for  $P_n^*(\mathbf{X}_n \in K)$  where  $K \subset M$ .

**Definition 2.** We say that  $\mathbf{X}_n$  converges weakly to a random element  $\mathbf{X}_0$  in  $(M, \mathcal{M}_B)$ , and write  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$ , if for every  $h \in C_b(M)$ ,

$$(6) \quad E^*h(\mathbf{X}_n) \rightarrow Eh(\mathbf{X}_0), \quad \text{as } n \rightarrow \infty.$$

Call  $\mathbf{X}_0$  separable if there exists a separable Borel set  $M_0 \subset M$  with  $P(\mathbf{X}_0 \in M_0) = 1$ . With the possible exception of set-theoretic pathological cases, it is no loss of generality to assume that a random element in the Borel  $\sigma$ -field is separable. See the discussion in Dudley (1985, pages 148, 149).

Since (6) holds for  $-h$ , it follows from definition 2 that, for  $h \in C_b(M)$ , also

$$(7) \quad E_*h(\mathbf{X}_n) \rightarrow Eh(\mathbf{X}_0) \quad \text{as } n \rightarrow \infty.$$

Note that to show  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$ , it suffices to prove

$$(8) \quad \limsup_{n \rightarrow \infty} E^*h(\mathbf{X}_n) \leq Eh(\mathbf{X}_0)$$

for all  $h \in C_b(M)$ : (8) implies that

$$\limsup_{n \rightarrow \infty} E^*(-h(\mathbf{X}_n)) \leq E(-h(\mathbf{X}_0)),$$

or

$$\liminf_{n \rightarrow \infty} E_*h(\mathbf{X}_n) \geq Eh(\mathbf{X}_0),$$

and hence

$$(9) \quad \begin{aligned} Eh(\mathbf{X}_0) &\leq \liminf_{n \rightarrow \infty} E_*h(\mathbf{X}_n) \leq \liminf_{n \rightarrow \infty} E^*h(\mathbf{X}_n) \\ &\leq \limsup_{n \rightarrow \infty} E^*h(\mathbf{X}_n) \leq Eh(\mathbf{X}_0), \end{aligned}$$

so that (6) holds.

As in the usual weak convergence several alternative formulations of the definition are possible; these are usually given in the form of a ‘‘portmanteau’’

theorem. It is stated by Andersen and Dobrić (1987) (see their remark 2.13, page 168) and has been also used by Giné and Zinn (1986, page 61).

**Lemma 3.** (Portmanteau theorem). The following statements are equivalent:

- (i)  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  as  $n \rightarrow \infty$ .
- (ii)  $\lim_{n \rightarrow \infty} E^*h(\mathbf{X}_n) = Eh(\mathbf{X}_0)$  for every (real) bounded Lipschitz continuous function  $h$ .
- (iii)  $\liminf_{n \rightarrow \infty} P_*(\mathbf{X}_n \in G) \geq P(\mathbf{X}_0 \in G)$  for every open set  $G \subset M$ .
- (iv)  $\limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \in F) \leq P(\mathbf{X}_0 \in F)$  for every closed set  $F \subset M$ .
- (v)  $\lim_{n \rightarrow \infty} P^*(\mathbf{X}_n \in A) = P(\mathbf{X}_0 \in A)$  for every set  $A$  with  $[\mathbf{X}_0 \in A]$  measurable and  $P(\mathbf{X}_0 \in \partial A) = 0$ .
- (vi)  $\liminf_{n \rightarrow \infty} E_*h(\mathbf{X}_n) \geq Eh(\mathbf{X}_0)$  for every lower semicontinuous function  $h$  bounded from below.
- (vii)  $\limsup_{n \rightarrow \infty} E^*h(\mathbf{X}_n) \leq Eh(\mathbf{X}_0)$  for every upper semicontinuous function  $h$  bounded from above.
- (viii)  $\lim_{n \rightarrow \infty} E^*h(\mathbf{X}_n) = Eh(\mathbf{X}_0)$  for every (real) bounded  $P\mathbf{X}_0^{-1}$ -continuous function  $h$  such that  $h(\mathbf{X}_0)$  is measurable.

An appropriate generalization of the tightness definition for this theory is as follows:

**Definition 3.** The sequence  $\{\mathbf{X}_n\}_{n \geq 1}$  is *asymptotically tight* if and only if for every  $\varepsilon > 0$  there exists a compact set  $K \subset M$  such that for all  $\delta > 0$

$$(10) \quad \liminf_{n \rightarrow \infty} P_*(\mathbf{X}_n \in K^\delta) \geq 1 - \varepsilon;$$

where  $K^\delta \equiv \{x \in M : d(x, K) < \delta\}$ , or, equivalently

$$(11) \quad \limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \notin K^\delta) < \varepsilon.$$

This definition is due to Hoffmann-Jørgensen (1984) in a development from Dudley (1966); see Andersen and Dobrić (1987, page 167).

Call a random element  $\mathbf{X}_0$  in  $(M, \mathcal{M}_B)$  *tight* if for every  $\varepsilon > 0$  there exists a compact set  $K \subset M$  such that  $P(\mathbf{X}_0 \in K) \geq 1 - \varepsilon$ . We note that, by theorem 3.2 of Parthasarathy (1967) or theorem 1.4 of Billingsley (1968), any separable  $\mathbf{X}_0$  in a *complete* metric space  $(M, d)$  is tight.

With these definitions, most of the usual results in the theory of weak convergence of measures on metric spaces as outlined in Billingsley (1968), and elsewhere, have analogues or extensions to the nonmeasurable Hoffmann-Jørgensen and Dudley theory. One notable exception, however, is the direct part of Prohorov's theorem (theorem 6.1, page 37, Billingsley (1968)), as will be seen in example 1 below. We begin with analogues of several results from the standard theory.

**Lemma 4.** (Weak convergence to a tight limit implies asymptotic tightness). Suppose that  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  where  $\mathbf{X}_0$  is tight. Then  $\{\mathbf{X}_n\}$  is asymptotically tight.

For  $n = 1, 2, \dots$ , let  $\mathbf{X}_n : \mathbf{X}_n \rightarrow M_1$  and  $\mathbf{Y}_n : \mathbf{X}_n \rightarrow M_2$  be maps into metric spaces  $(M_1, d_1)$  and  $(M_2, d_2)$ . Equip  $M_1 \times M_2$  with the metric

$$d_m((x_1, x_2), (y_1, y_2)) = d_1(x_1, y_1) \vee d_2(x_2, y_2).$$

**Lemma 5.** (Marginal asymptotic tightness implies joint asymptotic tightness). If both  $\{\mathbf{X}_n\}$  and  $\{\mathbf{Y}_n\}$  are asymptotically tight, then  $\{(\mathbf{X}_n, \mathbf{Y}_n)\}$  is asymptotically tight too.

**Proposition 1.** (Continuous mapping theorem). Suppose that  $(M, d)$  and  $(M', d')$  are metric spaces and that:

- (i)  $g : M \rightarrow M'$  is continuous on a Borel set  $M_0 \subset M$ ,
- (ii)  $\mathbf{X}_0$  has all its values in  $M_0$ .

Then  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  implies that  $g(\mathbf{X}_n) \Rightarrow g(\mathbf{X}_0)$ .

**Remark 1.** Note that for real-valued  $\mathbf{X}_n$ 's,  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  if and only if both  $\mathbf{X}_n^* \Rightarrow \mathbf{X}_0$  and  $\mathbf{X}_{n^*} \Rightarrow \mathbf{X}_0$ . This is proved straightforwardly using the portmanteau theorem and the inequalities  $P_*(a < \mathbf{X}_n < b) \geq P(\mathbf{X}_n^* < b) - P(\mathbf{X}_{n^*} \leq a)$  and  $P^*(\mathbf{X}_n \leq b) \geq P(\mathbf{X}_n^* \leq b) \geq 1 - P^*(\mathbf{X}_n \geq b)$ , which can be obtained from lemma 2. It allows various results for the Hoffmann-Jørgensen and Dudley theory to be deduced from known results in the classical theory. The following proposition is an example of this.

**Proposition 2.** (Pratt's theorem). Suppose that  $L_n, \mathbf{X}_n, U_n, n = 0, 1, 2, \dots$  are real-valued with:

- (i)  $L_n \leq \mathbf{X}_n \leq U_n$ ,
- (ii)  $L_n \Rightarrow L_0, \mathbf{X}_n \Rightarrow \mathbf{X}_0, U_n \Rightarrow U_0$ ,
- (iii)  $\lim_{n \rightarrow \infty} E_*(L_n) = E(L_0), \lim_{n \rightarrow \infty} E^*(U_n) = E(U_0)$ ,

where all the indicated expectations are finite. Then

$$\lim_{n \rightarrow \infty} E_*(\mathbf{X}_n) = E(\mathbf{X}_0) = \lim_{n \rightarrow \infty} E^*(\mathbf{X}_n).$$

*Le Cam's Third Lemma*

For  $n = 1, 2, \dots$  let  $P_n$  and  $Q_n$  be probability measures on measurable spaces  $(\mathbf{X}_n, \mathcal{A}_n)$ . Let  $\Lambda_n$  be the log-likelihood ratio of  $Q_n$  with respect to  $P_n$ ; i.e., given densities  $q_n$  and  $p_n$  with respect to a  $\sigma$ -finite dominating measure  $\mu_n$  (e.g.,  $\mu_n = P_n + Q_n$ ), let

$$\Lambda_n = \log \left( \frac{q_n}{p_n} \right),$$

where  $\log(a/b) = -\infty$  if  $a = 0 < b$ ,  $+\infty$  if  $b = 0 < a$ , and 0 if  $a = b$ .

Let  $\mathbf{X}_n : \mathbf{X}_n \rightarrow M$  as before. Equip  $M \times \bar{R}$  with the metric

$$d_m((x, r), (y, s)) = d(x, y) \vee \arctan |r - s|.$$

**Lemma 6.** (An extension of Le Cam's third lemma). Let  $\{Q_n\}$  be contiguous to  $\{P_n\}$ , and suppose that

$$(12) \quad (\mathbf{X}_n, \Lambda_n) \Rightarrow (\mathbf{X}_0, \Lambda_0) \quad \text{under } P_n,$$

where  $(\mathbf{X}_0, \Lambda_0) : \mathbf{X}_0 \rightarrow M \times R$  is Borel measurable. Then

$$(13) \quad \mathbf{X}_n \Rightarrow \mathbb{Z} \quad \text{under } Q_n,$$

where  $\mathbb{Z} : \mathbf{X}_0 \rightarrow M$  is Borel measurable and

$$(14) \quad P(\mathbb{Z} \in B) = E 1_B(\mathbf{X}_0) e^{\Lambda_0}.$$

Furthermore, if  $\mathbf{X}_0$  is separable (or tight), then  $\mathbb{Z}$  may be taken to be separable (or tight) too.

**Remark 2.** If  $\mathbf{X}_0$  is Borel measurable in  $M$  and has separable range, then  $(\mathbf{X}_0, \Lambda_0)$  is Borel measurable in  $(M, \bar{R})$ . *Proof:* Let  $M_0$  be the range of  $\mathbf{X}$ . Trivially  $(\mathbf{X}_0, \Lambda_0)$  is measurable as a map in  $M_0 \times \bar{R}$  with the product  $\sigma$ -field of the Borel  $\sigma$ -fields of  $M_0$  and  $\bar{R}$ . The latter  $\sigma$ -field equals the Borel  $\sigma$ -field of  $M_0 \times \bar{R}$ , by separability of  $M_0$  and  $R$ . Now for any Borel set  $B$  in  $M \times \bar{R}$ ,  $\{(\mathbf{X}_0, \Lambda_0) \in B\} = \{(\mathbf{X}_0, \Lambda_0) \in B \cap (M_0 \times \bar{R})\} \in \mathcal{A}_0$  since  $B \cap (M_0 \times \bar{R})$  is a Borel set in  $M_0 \times \bar{R}$ .

### Prohorov's Theorem

Here is an example to show that the direct half of Prohorov's theorem (asymptotic tightness implies relative compactness) fails for  $\Rightarrow$  without additional hypotheses.

**Example 1.** Let  $(\mathbf{X}_n, \mathcal{A}_n, P_n) = (\mathbf{X}, \mathcal{A}, P) = ([0, 1], \{\emptyset, [0, 1]\}, \lambda)$  for all  $n \geq 1$  where  $\lambda$  denotes Lebesgue measure. Let  $M = [0, 1]$ , and define  $\mathbf{X}(\omega) = \omega$  for  $\omega \in \mathbf{X} \equiv [0, 1]$ . Consider the sequence  $\{\mathbf{X}_n\}$  defined by  $\mathbf{X}_n \equiv \mathbf{X}$ ,  $n \geq 1$ . Then  $P_*(\mathbf{X}_n \in [0, 1]) = \lambda_*(\mathbf{X} \in [0, 1]) = 1$  for all  $n = 1, 2, \dots$ , so  $\{\mathbf{X}_n\}$  is tight. But  $\mathbf{X}^* = 1$ ,  $\mathbf{X}_* = 0$ ; and hence for the bounded, continuous function  $h(x) = x$  on  $M$  it follows that

$$E^* h(\mathbf{X}_n) = E \mathbf{X}_n^* = 1, \quad \text{and}$$

$$E_* h(\mathbf{X}_n) = E \mathbf{X}_n^* = 0 \quad \text{for all } n = 1, 2, \dots$$

Thus  $\{\mathbf{X}_n\}$  has no subsequence for which  $\Rightarrow$  holds. □

Of course, the key difficulty here is that  $\Rightarrow$  implies some "measurability in the limit" since (6) and (7) imply that

$$(15) \quad E^* h(\mathbf{X}_n) - E_* h(\mathbf{X}_n) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

for every  $h \in C_b(M)$ . In the example there is no measurability to start with, so there is no measurability in the limit, and weak convergence in the sense of definition 2 fails; thus it is not true in general that  $\mathbf{X} \Rightarrow \mathbf{X}$ . In fact  $\mathbf{X} \Rightarrow \mathbf{X}$

is true if and only if  $\mathbf{X}$  is Borel measurable. This leads to the following definition:

**Definition 4.** The sequence  $\{\mathbf{X}_n\}$  is *asymptotically measurable* if and only if

$$E^* h(\mathbf{X}_n) - E_* h(\mathbf{X}_n) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

for every  $h \in C_b(M)$ .

Here is a modification of Prohorov's theorem appropriate for the present (Hoffmann-Jørgensen and Dudley) definition of weak convergence.

**Theorem 1.** (Extension of Prohorov's theorem). Suppose that:

- (i)  $\{\mathbf{X}_n\}$  is asymptotically tight,
- (ii)  $\{\mathbf{X}_n\}$  is asymptotically measurable.

Then there exists a subsequence  $\{\mathbf{X}_{n'}\} \subset \{\mathbf{X}_n\}$  such that  $\mathbf{X}_{n'} \Rightarrow$  some tight  $\mathbf{X}_0$ .

Note that the converse half of Prohorov's theorem for sequences (relative compactness implies asymptotic tightness) is trivial if relative compactness of  $\{\mathbf{X}_n\}$  is understood to mean, as usual, that every subsequence  $\{n'\}$  has a further subsequence  $\{n''\}$  such that  $\mathbf{X}_{n''} \Rightarrow \mathbf{X}_0$  for some tight  $\mathbf{X}_0$ .

Asymptotic measurability is often hard to establish directly, so it is of some importance to have useful criteria which can be checked in special cases. The main object of these criteria is to only have to check that (15) holds for a smaller class of functions than all of  $C_b(M)$ . Here is a rather abstract version of this which can be applied in many special cases.

**Lemma 7.** Suppose that  $\{\mathbf{X}_n\}$  is asymptotically tight and that (15) holds for every  $h$  in a subalgebra  $\mathbf{H}$  of  $C_b(M)$  that separates points of  $M$  (i.e., for all  $x, y \in M, x \neq y$ , there exists an  $h \in \mathbf{H}$  with  $h(x) \neq h(y)$ ). Then  $\{\mathbf{X}_n\}$  is asymptotically measurable.

Theorem 1 in combination with lemma 7 has a number of useful corollaries:

**Corollary 1.** For  $n = 1, 2, \dots$  suppose that  $\mathbf{X}_n : \mathbf{X}_n \rightarrow M_1$  and  $\mathbf{Y}_n : \mathbf{X}_n \rightarrow M_2$  are maps with

$$\mathbf{X}_n \Rightarrow \mathbf{X}_0 \quad \text{and} \quad \mathbf{Y}_n \Rightarrow \mathbf{Y}_0,$$

where  $\mathbf{X}_0$  and  $\mathbf{Y}_0$  are tight Borel measurable maps into  $M_1$  and  $M_2$  respectively. Then there exists a subsequence  $\{n'\} \subset \{n\}$  such that

$$(\mathbf{X}_{n'}, \mathbf{Y}_{n'}) \Rightarrow (\mathbf{X}_0, \mathbf{Y}_0) \quad \text{as } n' \rightarrow \infty$$

for some tight joint law  $\mathbf{L}(\mathbf{X}_0, \mathbf{Y}_0)$  on the product space  $(M_1 \times M_2, \mathcal{M}_{B1} \times \mathcal{M}_{B2})$ .

As in classical weak convergence theory for separable metric spaces, joint convergence for the full (joint) sequence in corollary 1 fails in general: consider  $\{(\mathbf{X}_n, \mathbf{Y}_n); n \geq 1\}$  in  $[0, 1]^2$  with  $(\mathbf{X}_{2n}, \mathbf{Y}_{2n})$  uniformly distributed on the line  $y = x$  and  $(\mathbf{X}_{2n+1}, \mathbf{Y}_{2n+1})$  uniformly distributed on the line  $y = 1 - x$ . Then  $\mathbf{X}_n \sim \text{Uniform}(0, 1)$  and  $\mathbf{Y}_n \sim \text{Uniform}(0, 1)$  for all  $n$  so that  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  and  $\mathbf{Y}_n \Rightarrow \mathbf{Y}_0$  with  $\mathbf{X}_0$  and  $\mathbf{Y}_0 \sim \text{Uniform}(0, 1)$ . But the full sequence

$\{(\mathbf{X}_n, \mathbf{Y}_n); n \geq 1\}$  does not converge weakly. Joint convergence for the full sequence can hold under additional hypotheses; see, e.g., Billingsley (1968, page 27, theorems 4.4 and 4.5).

For an arbitrary set  $F$ , let  $l^\infty(F)$  denote the collection of all bounded real functions on  $F$ . It will be equipped with the supremum norm  $\|z\|_\infty = \sup_{f \in F} |z(f)|$ .

**Corollary 2.** Suppose that  $M = l^\infty(F)$  and that:

- (i)  $\{\mathbf{X}_n\}$  is asymptotically tight,
- (ii)  $\mathbf{X}_n(f)$  is  $\mathcal{A}_n - \mathcal{B}$ -measurable for every  $f \in F$  and  $n = 1, 2, \dots$

Then there exist  $\{n'\} \subset \{n\}$  such that  $\mathbf{X}_{n'} \Rightarrow$  some tight  $\mathbf{X}_0$ .

When  $M = l^\infty(F)$ , asymptotic tightness of  $\mathbf{X}_n$  in the sense of definition 3 is equivalent to (asymptotic)  $\rho$ -equicontinuity of  $\mathbf{X}_n$  for some pseudometric  $\rho$  on  $F$  together with (asymptotic) boundedness; see, e.g., Andersen and Dobrić (1987, page 167). This parallels the classical results in Billingsley (1968), theorem 8.2. Also, a tight law  $L(\mathbf{X}_0)$  on  $l^\infty(F)$  is completely determined by its finite-dimensional laws  $L(\mathbf{X}_0(f_1), \dots, \mathbf{X}_0(f_m))$ . Hence if (i) and (ii) of corollary 2 hold and, moreover,

$$L(\mathbf{X}_n(f_1), \dots, \mathbf{X}_n(f_m)) \rightarrow L(\mathbf{X}_0(f_1), \dots, \mathbf{X}_0(f_m)),$$

then  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$ .

**Corollary 3.** Let  $\mathcal{A}$  be a  $\sigma$ -field of subsets of  $M$  such that  $\{h \in C_b(M) : h \text{ is } \mathcal{A}\text{-measurable}\}$  separates points of  $M$ . Suppose that:

- (i)  $\{\mathbf{X}_n\}$  is asymptotically tight.
- (ii)  $\mathbf{X}_n$  is  $(\mathcal{A}_n - \mathcal{A})$  measurable for every  $n = 1, 2, \dots$

Then there exists  $\{n'\} \subset \{n\}$  such that  $\mathbf{X}_{n'} \Rightarrow$  some tight  $\mathbf{X}_0$ .

An example is obtained by letting  $\mathcal{A}$  be the  $\sigma$ -field generated by the closed balls (cf. Gaenssler (1983), Dudley (1966)). Then  $h(x) = (1 - m d(x, x_1))^+$  satisfies  $h(x_1) = 1$  and  $h(x_2) = 0$  for sufficiently large  $m$ . Actually, corollary 2 is a special case of corollary 3, too.

### Measurability in the Ball $\sigma$ -Field

An alternative theory of weak convergence due to Dudley (1966) (see also Gaenssler (1983)) is applicable when the  $\mathbf{X}_n$ 's are measurable in the  $\sigma$ -field  $\mathcal{M}_b$  of subsets of  $M$  generated by the closed balls. In this theory,  $\mathbf{X}_n \Rightarrow_{\text{Dudley}} \mathbf{X}_0$  means that

$$(16) \quad Eh(\mathbf{X}_n) \rightarrow Eh(\mathbf{X}_0) \quad \text{for all } h \in C_b(M, \mathcal{M}_b)$$

where  $C_b(M, \mathcal{M}_b)$  is the collection of all bounded continuous functions defined on  $M$  which are  $\mathcal{M}_b$ -measurable. In view of results of Dudley (1966) (cf. Gaenssler (1983, theorem 28, page 47)), (16) is equivalent to

$$(17) \quad \int^* h dP_n \circ \mathbf{X}_n^{-1} \rightarrow Eh(\mathbf{X}_0) \quad \text{for all } h \in C_b(M)$$

provided  $\mathbf{X}_0$  has separable range.



When  $\mathbf{X}_n$  is measurable in the ball  $\sigma$ -field  $\mathcal{M}_b$  and  $\mathbf{X}_0$  has separable range the following proposition says that  $\Rightarrow_{\text{Dudley}}$  is equivalent to  $\Rightarrow$ .

**Proposition 3.** (Equivalence of  $\Rightarrow_{\text{Dudley}}$  and  $\Rightarrow$  under  $\mathcal{M}_b$ -measurability). Suppose that  $\mathbf{X}_n : \mathbf{X}_n \rightarrow M$  is  $\mathcal{M}_b$ -measurable for  $n = 0, 1, \dots$  and  $P(\mathbf{X}_0 \in M_0) = 1$  for a separable subset  $M_0$  of  $M$ . Then  $\mathbf{X}_n \Rightarrow_{\text{Dudley}} \mathbf{X}_0$  if and only if  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$ .

On the other hand, Dudley (1985) gives an example to show that  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  is *not* equivalent to

$$\int^* h dP_n \circ \mathbf{X}_n^{-1} \rightarrow Eh(\mathbf{X}_0) \quad \text{for all } h \in C_b(M),$$

(where each  $P_n \circ \mathbf{X}_n^{-1}$  is defined on  $\{B \subset M : \mathbf{X}_n^{-1}(B) \in \mathcal{A}_n\}$ ), which would be a natural extension of (17) in the case that the  $\mathbf{X}_n$ 's are not  $\mathcal{M}_b$ -measurable.

### Convergence in Outer Probability, Convergence Almost Uniformly, and Continuous Mapping Theorems

Now suppose that every  $\mathbf{X}_n$  is defined on the same probability space  $(\mathbf{X}, \mathcal{A}, P)$ ,  $n = 0, 1, \dots$ . We say that  $\mathbf{X}_n$  converges *almost uniformly* to  $\mathbf{X}_0$  if

$$d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_{a.s.} 0 \quad \text{as } n \rightarrow \infty.$$

We say that  $\mathbf{X}_n$  converges in *outer probability* to  $\mathbf{X}_0$  if

$$d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_p 0 \quad \text{as } n \rightarrow \infty.$$

It is clear that convergence almost uniformly implies convergence in outer probability. By lemma 2.K, the latter is equivalent to, for every  $\varepsilon > 0$

$$(18) \quad P^*(d(\mathbf{X}_n, \mathbf{X}_0) > \varepsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

The definition of almost uniform convergence has a similar translation:

**Proposition 4.** (Dudley). Let  $(\mathbf{X}, \mathcal{A}, P)$  be a probability space,  $(M, d)$  a metric space, and  $\mathbf{X}_n$  functions from  $\mathbf{X}$  into  $M$ ,  $n = 0, 1, \dots$ . Then the following are equivalent:

- (i)  $d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_{a.s.} 0$ .
- (ii) For every  $\varepsilon > 0$ ,  $P^*(\sup_{m \geq n} d(\mathbf{X}_m, \mathbf{X}_0) > \varepsilon) \rightarrow 0$  as  $n \rightarrow \infty$ .
- (iii) For each  $\delta > 0$ , there is some  $B \in \mathcal{A}$  with  $P(B) > 1 - \delta$  such that  $\mathbf{X}_n \rightarrow \mathbf{X}_0$  uniformly on  $B$ .

**Proof.** See proposition 1.1 of Dudley (1985). □

However, it is *not true* that (i)–(iii) are equivalent to

$$P^*\left(\bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} \{d(\mathbf{X}_m, \mathbf{X}_0) > \varepsilon\}\right) = 0.$$

This is due to the fact that

$$\left(\prod_{i=1}^{\infty} f_i\right)^* \leq \prod_{i=1}^{\infty} f_i^*$$

is true in general, but the reverse inequality fails without additional hypotheses

as shown by the following example given by Dudley (1985). (Also see lemma 2.3 in Dudley and Philipp (1983).) Thus almost uniform convergence is strictly stronger than convergence for (almost) all  $\omega$ .

**Example 2.** There exist nonmeasurable subsets  $A_1 \supset A_2 \supset \dots$  of  $[0, 1]$  with Lebesgue measure  $\lambda$  with  $\bigcap_{i=1}^{\infty} A_i = \emptyset$  and  $\lambda^*(A_i) = 1$  for all  $i$ ; see, e.g., Cohn (1980, page 35). Hence  $(\prod_{i=1}^{\infty} 1_{A_i})^* = 0^* = 0$ , but  $\prod_{i=1}^{\infty} 1_{A_i}^* = 1$ .

Both convergence in outer probability and almost uniformly are preserved by a continuous map if the limit is Borel measurable.

**Proposition 5.** (Continuous mapping theorem, continued). Suppose that  $(M, d)$  and  $(M', d')$  are metric spaces and that:

- (i)  $g : M \rightarrow M'$  is continuous on a Borel set  $M_0 \subset M$ .
- (ii)  $\mathbf{X}_0$  is Borel measurable and  $P(\mathbf{X}_0 \in M_0) = 1$ .

Then:

- A.  $d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_p 0$  implies that  $d'(g(\mathbf{X}_n), g(\mathbf{X}_0))^* \rightarrow_p 0$ .
- B.  $d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_{\text{a.s.}} 0$  implies that  $d'(g(\mathbf{X}_n), g(\mathbf{X}_0))^* \rightarrow_{\text{a.s.}} 0$ .

The continuous mapping (or Mann-Wald) theorems given in propositions 1 and 5 have a further useful extension to a sequence of functions  $g_n$ . The following proposition generalizes Billingsley (1968, theorem 5.5, pages 33–34), and Gaenssler (1983, theorem 5, page 56). The basic idea is apparently due to H. Rubin; see Topsoe (1967) for a discussion and a refined version of the theorem in the standard (separable  $M$ ) weak convergence theory. It is a very useful result in connection with Hadamard-differentiable functions and the delta method; see Wellner (1989) and Sheehy and Wellner (1988).

**Proposition 6.** (Extended continuous mapping theorem). Suppose that  $g, g_n : M \rightarrow M'$  are functions which satisfy:

- (i)  $\mathbf{X}_0$  takes all its values in some Borel subset  $M_0$  of  $M$ .
- (ii) For every  $x \in M_0$  and every sequence  $\{x_n\}$  with  $x_n \rightarrow x, g_n(x_n) \rightarrow g(x)$ .

Then the restriction of  $g$  to  $M_0$  is continuous and:

- A.  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  implies  $g_n(\mathbf{X}_n) \Rightarrow g(\mathbf{X}_0)$ .
- B.  $d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_p 0$  implies  $d'(g_n(\mathbf{X}_n), g(\mathbf{X}_0))^* \rightarrow_p 0$ .
- C.  $d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_{\text{a.s.}} 0$  implies  $d'(g_n(\mathbf{X}_n), g(\mathbf{X}_0))^* \rightarrow_{\text{a.s.}} 0$ .

Convergence in outer probability is stronger than convergence  $\Rightarrow$ :

**Lemma 8.** Let every  $\mathbf{X}_n$  be defined on the same probability space  $(\mathbf{X}, \mathcal{A}, P)$  and assume  $d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_p 0$  where  $\mathbf{X}_0$  is Borel measurable. Then  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$ .

This lemma is actually a special case of the following lemma:

**Lemma 9.** For  $n = 0, 1, \dots$  let  $\mathbf{X}_n : \mathbf{X}_n \rightarrow M$  and  $\mathbf{Y}_n : \mathbf{X}_n \rightarrow M$  be arbitrary maps. Suppose that  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  where  $\mathbf{X}_0$  is Borel measurable and  $d(\mathbf{X}_n, \mathbf{Y}_n)^* \rightarrow_p 0$ . Then  $\mathbf{Y}_n \Rightarrow \mathbf{X}_0$ .

The almost surely convergent construction of Skorokhod (1956) (see, e.g., Billingsley (1971)) has a counterpart in the Hoffmann-Jørgensen theory of weak convergence; the counterpart is due to Dudley (1985). To state Dudley's theorem, we first introduce the notion of a perfect function.

**Definition 5.** Let  $(\tilde{\mathbf{X}}, \tilde{\mathcal{A}}, \tilde{P})$  be a probability space, let  $(\mathbf{X}, \mathcal{A})$  be a measurable space, and let  $\phi : \tilde{\mathbf{X}} \rightarrow \mathbf{X}$  be measurable. Let  $P = \tilde{P} \circ \phi^{-1}$ . Then  $\phi$  is *perfect* if and only if for any bounded real-valued function  $g$  on  $\mathbf{X}$

$$E_P^* g = E_{\tilde{P}}^* g \circ \phi.$$

As shown by Dudley (1985, theorem 2, page 146), this is equivalent to  $(g \circ \phi)^* = g^* \circ \phi$  a.s.  $\tilde{P}$  and to  $P^* = \tilde{P}^* \circ \phi^{-1}$ .

The key property of perfect functions that makes them useful and important is the following: suppose that  $\mathbf{X} : (\mathbf{X}, \mathcal{A}, P) \rightarrow (M, d)$  where  $(M, d)$  is a metric space. Suppose there exists a perfect measurable function  $\phi : (\tilde{\mathbf{X}}, \tilde{\mathcal{A}}, \tilde{P}) \rightarrow (\mathbf{X}, \mathcal{A}, P)$  with  $P = \tilde{P} \circ \phi^{-1}$ . Define  $\tilde{\mathbf{X}} : (\tilde{\mathbf{X}}, \tilde{\mathcal{A}}) \rightarrow (M, d)$  by  $\tilde{\mathbf{X}} = \mathbf{X} \circ \phi$ . Then for any set  $B \subset M$

$$(19) \quad \tilde{P}^*(\tilde{\mathbf{X}} \in B) = P^*(\mathbf{X} \in B);$$

this follows from the definition since, with  $g(\omega) \equiv 1_{[\mathbf{X}(\omega) \in B]}$  for  $\omega \in \mathbf{X}$ ,

$$P^*(\mathbf{X} \in B) = \tilde{P}^*(\mathbf{X} \circ \phi \in B) = \tilde{P}^*(\tilde{\mathbf{X}} \in B).$$

In this sense we can say that " $\tilde{\mathbf{X}} =_d \mathbf{X}$ ".

Here is Dudley's (1985) fourth-generation Skorokhod theorem.

**Theorem 2.** (Skorokhod-Dudley-Wichura-Dudley). Let  $(M, d)$  be any metric space,  $(\mathbf{X}_n, \mathcal{A}_n, P_n)$  any probability spaces, and  $\mathbf{X}_n$  a function from  $\mathbf{X}_n$  into  $M$  for each  $n = 0, 1, \dots$ . Suppose that  $\mathbf{X}_0$  has separable range  $M_0 \subset M$  and is measurable. Then  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  if and only if there exists a probability space  $(\tilde{\mathbf{X}}, \tilde{\mathcal{A}}, \tilde{P})$  and perfect measurable functions  $\phi_n$  from  $(\tilde{\mathbf{X}}, \tilde{\mathcal{A}})$  to  $(\mathbf{X}_n, \mathcal{A}_n)$  for each  $n = 0, 1, \dots$  such that:

- (i)  $\tilde{P} \circ \phi_n^{-1} = P_n$  on  $\mathcal{A}_n$  for each  $n = 0, 1, \dots$ ,
- (ii)  $\tilde{\mathbf{X}}_n \equiv \mathbf{X}_n \circ \phi_n \rightarrow \mathbf{X}_0 \circ \phi_0 \equiv \tilde{\mathbf{X}}_0$  almost uniformly; i.e.,

$$(20) \quad d(\tilde{\mathbf{X}}_n, \tilde{\mathbf{X}}_0)^* \rightarrow_{\text{a.s.}} 0 \quad \text{as } n \rightarrow \infty.$$

**Proof.** For a complete discussion of theorem 2 and the notion of perfect functions, see Dudley (1985). For another discussion and exposition, see Pollard (1990). □

*Proofs*

**Proof of lemma 1.** Set  $L = \inf\{\int \arctan h dP : h \geq f, h \text{ measurable}\}$ . Choose  $h_m \geq f$  measurable with  $\int \arctan h_m dP \downarrow L$ . Set  $f_m \equiv \inf_{k \leq m} h_k$ .

Define  $f^*(\omega) \equiv \lim_{m \rightarrow \infty} f_m(\omega)$  (finite or  $-\infty$ ). Then  $f^*(\omega) \geq f(\omega)$  everywhere,  $f^*$  is measurable, and  $\int \arctan f^* dP = L$ . Moreover, for every  $h \geq f$  that is measurable, it follows that

$$(a) \quad \int \arctan(f^* \wedge h) dP = \lim_{m \rightarrow \infty} \int \arctan(f_m \wedge h) dP \geq L = \int \arctan f^* dP.$$

Since  $\arctan f^* - \arctan(f^* \wedge h) \geq 0$ , and

$$(b) \quad \int (\arctan f^* - \arctan(f^* \wedge h)) dP \leq 0, \quad \text{by (a),}$$

it follows that

$$(c) \quad \arctan f^* = \arctan(f^* \wedge h) \quad \text{a.s.,}$$

implying that  $f^* = f^* \wedge h$  a.s., or  $f^* \leq h$  a.s. □

**Proof of lemma 2.** *Proof of A:* The first and last inequality of A are trivial. By the first inequality  $(f + g)^* + (-f)^* \leq g^*$ , which is the second inequality. Similarly, by the last inequality  $f^* \leq (f + g)^* + (-g)^*$ , which proves the third inequality.

*Proof of B:* The first equality in B follows directly from the third and fourth inequality of A with  $g^* = g^* = g$ . Likewise, the second equality of B is implied by the first two inequalities of A.

*Proof of C:* C follows immediately from the second and third inequality of A.

*Proof of D:* For D, note that by the third inequality of A,  $f^* - g^* \leq (f - g)^*$ , which is obviously smaller than  $|f - g|^*$ . E may be proved analogously.

*Proof of F:* To prove F, first note that

$$(fg)^* \leq (fg 1_{[g > 0]})^* + (fg 1_{[g < 0]})^* \leq f^* g 1_{[g > 0]} + f^* g 1_{[g < 0]}.$$

Let  $fg \leq h$  with  $h$  measurable. Then  $0 \leq h 1_{[g = 0]}$ ,  $f 1_{[g > 0]} \leq (h/g) 1_{[g > 0]}$ , and  $f 1_{[g < 0]} \geq (h/g) 1_{[g < 0]}$ , and hence

$$f^* 1_{[g > 0]} \leq (h/g) 1_{[g > 0]}, \quad f^* 1_{[g < 0]} \geq (h/g) 1_{[g < 0]}.$$

Consequently,

$$f^* g 1_{[g > 0]} + f^* g 1_{[g < 0]} \leq h 1_{[g > 0]} + h 1_{[g < 0]} \leq h$$

and the proof of F is complete.

*Proof of G:* let  $h = 1_A^*$ . Then  $A^* = [h \geq 1]$  satisfies  $1_A \leq 1_{A^*}$  and  $1_{A^*} \leq h$ . Note that H is trivial.

*Proof of I:* First,  $E_P f \leq \int f^* dP \equiv E_P f^*$  is trivially true. On the other hand, if  $h \geq f$ ,  $E_P h$  is well defined, and  $E_P f^*$  is well defined, then  $E_P h \geq E_P f^*$ , and hence  $E_P f \geq E_P f^*$ .

*Proof of J:*

$$\begin{aligned} P^*(A) &\equiv \inf\{P(B) : B \supset A, B \text{ measurable}\} \\ &\geq \inf\{E h : h \geq 1_A, \text{ measurable}\} \end{aligned}$$

$$\begin{aligned}
 &= E 1_A^* \quad \text{by I} \\
 &= E 1_A \quad \text{by G} \\
 &= P(A^*) \geq P^*(A).
 \end{aligned}$$

*Proof of K:* We show that  $B \supset \{f > c\}$  with  $B$  measurable implies that  $B \supset \{f^* > c\}$ . Thus  $\{f^* > c\} \subset \{f > c\}^*$ .

Define  $g \equiv f^* 1_A + (f^* \wedge c) 1_{A^c}$ . Then  $g \geq f$  a.s. and is measurable, and hence  $g \geq f^*$ . But this implies  $f^* \wedge c \geq f^*$  on  $A^c$ , or  $A^c \subset \{f^* \leq c\}$ , and hence  $A \supset \{f^* > c\}$ .

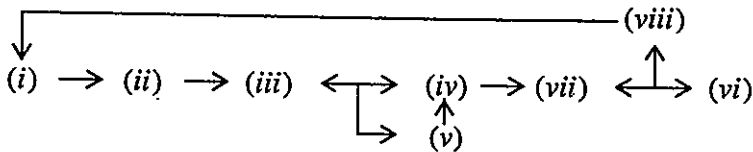
Note that this implies the first assertion:  $\{f > c\}^* \subset \{f^* > c\}$  is easy since  $f^* \geq f$ . The second assertion follows by J and G.

*Proof of L:* Follows easily from K and H.

*Proof of M:* First,  $f_i \leq \sup_{i \in I} f_i \leq (\sup_{i \in I} f_i)^*$  for all  $i \in I$ , so  $f_i^* \leq (\sup_{i \in I} f_i)^*$  for all  $i \in I$ , and hence  $\sup_{i \in I} f_i^* \leq (\sup_{i \in I} f_i)^*$ .

If  $I$  is countable, then  $\sup_{i \in I} f_i^*$  is measurable and  $\sup_{i \in I} f_i^* \geq f_i^* \geq f_i$  for all  $i \in I$ , and therefore  $\sup_{i \in I} f_i^* \geq \sup_{i \in I} f_i$ . This implies  $\sup_{i \in I} f_i^* \geq (\sup_{i \in I} f_i)^*$ .  $\square$

**Proof of lemma 3.** We will establish the following implications:



The equivalence of (iii) and (iv) is trivial; similarly, the equivalence of (vi) and (vii) is trivial. Furthermore, (i) implies (ii) trivially since every Lipschitz continuous function is continuous. That (viii) implies (i) is easy since a continuous function is  $P_{\mathbf{X}_0}^{-1}$ -continuous.

Now we show that (ii) implies (iii). First, there exists a sequence of Lipschitz continuous functions  $\{h_m\}$  with  $0 \leq h_m \leq 1_G$  and  $h_m \uparrow 1_G$  as  $m \rightarrow \infty$ ; for example take  $h_m(x) = 1 - (1 - m d(x, G^c))^+$ . Now for every  $m = 1, 2, \dots$

$$\liminf_{n \rightarrow \infty} P_*(\mathbf{X}_n \in G) \geq \liminf_{n \rightarrow \infty} E_* h_m(\mathbf{X}_n) = E h_m(\mathbf{X}_0).$$

By monotone convergence  $E h_m(\mathbf{X}_0) \uparrow P(\mathbf{X}_0 \in G)$  as  $m \rightarrow \infty$ .

To show that (iv) implies (vii), assume without loss of generality that  $0 \leq h \leq 1$ . Fix an integer  $r$ . For  $j = 0, \dots, r$ , set  $F_j = \{x : h(x) \geq j/r\}$  (a closed set by upper semicontinuity of  $h$ ) and  $h_r(x) = r^{-1} \sum_{j=0}^r 1_{F_j}(x)$ . Then  $h_r \geq h$  and  $\|h_r - h\|_\infty \leq 1/r$ . Now

$$\begin{aligned}
 \limsup_{n \rightarrow \infty} E_* h(\mathbf{X}_n) &\leq \limsup_{n \rightarrow \infty} E_* h_r(\mathbf{X}_n) \leq \frac{1}{r} \sum_{j=0}^r \limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \in F_j) \\
 &\leq \frac{1}{r} \sum_{j=0}^r P(\mathbf{X}_0 \in F_j) \quad \text{by (iv)}
 \end{aligned}$$

$$= Eh_r(\mathbf{X}_0).$$

Let  $r \rightarrow \infty$  to get the conclusion.

*Proof that (v) implies (iv):* define  $\bar{F}^\varepsilon = \{x : d(x, F) \leq \varepsilon\}$  and  $\partial F^\varepsilon \subset \{x : d(x, F) = \varepsilon\}$ . Thus we have  $\partial F^{\varepsilon_1} \cap \partial F^{\varepsilon_2} = \emptyset$  if  $\varepsilon_1 \neq \varepsilon_2$ , so that  $P(\mathbf{X}_0 \in \partial F^\varepsilon)$  can be nonzero for at most countably many  $\varepsilon > 0$ . Choose  $\varepsilon_m \downarrow 0$  such that  $P(\mathbf{X}_0 \in \partial F^{\varepsilon_m}) = 0, m = 1, 2, \dots$ . Then

$$\limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \in F) \leq \limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \in \bar{F}^{\varepsilon_m}) = P(\mathbf{X}_0 \in \bar{F}^{\varepsilon_m}).$$

Letting  $m \rightarrow \infty$  yields (iv).

*Proof that (vi) and (vii) imply (viii):* Suppose that  $h$  is bounded and  $P\mathbf{X}_0^{-1}$ -continuous. Then

$$h_u \equiv \inf\{g : g \geq h, g \text{ upper semicontinuous}\},$$

$$h_l \equiv \sup\{g : g \leq h, g \text{ lower semicontinuous}\},$$

are upper and lower semicontinuous respectively,  $h_l \leq h \leq h_u$ , and

$$P(h_l(\mathbf{X}_0) < h_u(\mathbf{X}_0)) = 0.$$

Thus we find that

$$Eh_l(\mathbf{X}_0) = Eh(\mathbf{X}_0) = Eh_u(\mathbf{X}_0).$$

Therefore, by (vi) and (vii) we obtain

$$\begin{aligned} Eh(\mathbf{X}_0) = Eh_l(\mathbf{X}_0) &\leq \liminf_{n \rightarrow \infty} E_* h_l(\mathbf{X}_n) \leq \liminf_{n \rightarrow \infty} E_* h(\mathbf{X}_n) \\ &\leq \liminf_{n \rightarrow \infty} E^* h(\mathbf{X}_n) \leq \limsup_{n \rightarrow \infty} E^* h(\mathbf{X}_n) \\ &\leq \limsup_{n \rightarrow \infty} E^* h_u(\mathbf{X}_n) \leq Eh_u(\mathbf{X}_0) = Eh(\mathbf{X}_0). \end{aligned}$$

Since the extreme terms are equal, this yields (viii).

To complete the proof, we show that (iii) and (iv) together imply (v): For  $A \subset M$  with  $P(\mathbf{X}_0 \in \partial A) = 0$ ,

$$\begin{aligned} P(\mathbf{X}_0 \in \text{Int} A) &\leq \liminf_{n \rightarrow \infty} P_*(\mathbf{X}_n \in \text{Int} A) \\ &\leq \liminf_{n \rightarrow \infty} P^*(\mathbf{X}_n \in A) \leq \limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \in A) \\ &\leq \limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \in \bar{A}) \leq P(\mathbf{X}_0 \in \bar{A}). \end{aligned}$$

But the left and right sides are equal since  $P(\mathbf{X}_0 \in \partial A) = 0$ , and hence (v) follows.  $\square$

**Proof of lemma 4.** For every  $\varepsilon > 0$  there exists a compact set  $K$  with  $P(\mathbf{X}_0 \in K) \geq 1 - \varepsilon$ . By the implication (i) implies (iii) of lemma 3,

$$\liminf_{n \rightarrow \infty} P_*(\mathbf{X}_n \in K^\delta) \geq P(\mathbf{X}_0 \in K^\delta) \geq P(\mathbf{X}_0 \in K) \geq 1 - \varepsilon. \quad \square$$

**Proof of lemma 5.** Let  $\varepsilon > 0$ , and let  $K_1$  and  $K_2$  be compact sets in  $M_1$  and  $M_2$  respectively with  $\liminf_{n \rightarrow \infty} P^*(\mathbf{X}_n \in K_1^\delta) \geq 1 - \varepsilon$  and  $\liminf_{n \rightarrow \infty} P^*(\mathbf{Y}_n \in K_2^\delta) \geq 1 - \varepsilon$  for every  $\delta > 0$ . Then  $(K_1 \times K_2)^\delta = K_1^\delta \times K_2^\delta$  and

$$\begin{aligned} & \limsup_{n \rightarrow \infty} P^*((\mathbf{X}_n, \mathbf{Y}_n) \in (K_1^\delta \times K_2^\delta)^c) \\ & \leq \limsup_{n \rightarrow \infty} P^*((\mathbf{X}_n, \mathbf{Y}_n) \in (K_1^\delta)^c \times M_2 \cup M_1 \times (K_2^\delta)^c) \\ & \leq \limsup_{n \rightarrow \infty} \{P^*(\mathbf{X}_n \in (K_1^\delta)^c) + P^*(\mathbf{Y}_n \in (K_2^\delta)^c)\} \\ & \leq 2\varepsilon. \end{aligned}$$

□

**Proof of proposition 1.** This is almost as in Billingsley (1968): Let  $F \subset M'$  be closed and let  $C(g) \subset M$  denote the continuity set of  $g$ . Then, since  $\overline{g^{-1}(F)} \subset C^c(g) \cup g^{-1}(F)$  and  $P(\mathbf{X}_0 \in C^c(g)) = 0$ ,

$$\begin{aligned} \limsup_{n \rightarrow \infty} P^*(g(\mathbf{X}_n) \in F) &= \limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \in g^{-1}(F)) \\ &\leq \limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \in \overline{g^{-1}(F)}) \\ &\leq P(\mathbf{X}_0 \in \overline{g^{-1}(F)}) \end{aligned}$$

by (i) implies (iv) of lemma 3

$$\begin{aligned} & \leq P(\mathbf{X}_0 \in C(g)^c \cup g^{-1}(F)) \\ & = P(\mathbf{X}_0 \in g^{-1}(F)) \\ & = P(g(\mathbf{X}_0) \in F) \end{aligned}$$

which yields the conclusion by (iv) implies (i) of lemma 3. □

**Proof of proposition 2.** Note that

$$(a) \quad L_{n^*} \leq \mathbf{X}_{n^*} \leq \mathbf{X}_n^* \leq U_n^*$$

where  $L_{n^*} \Rightarrow L_0$ ,  $\mathbf{X}_{n^*} \Rightarrow \mathbf{X}_0$ ,  $\mathbf{X}_n^* \Rightarrow \mathbf{X}_0$ ,  $U_n^* \Rightarrow U_0$  by hypothesis (ii) and remark 1. Furthermore, by lemma 2.I and hypothesis (iii)

$$E(L_{n^*}) = E^*(L_n) \rightarrow E(L_0), \quad E(U_n^*) = E^*(U_n) \rightarrow E(U_0).$$

The conclusion follows from two applications of the classical weak convergence version of Pratt's theorem. See Pratt (1960) or, e.g., Shirayev (1984, exercise 18, page 209). □

**Proof of lemma 6.** Let  $h \in C_b(M)$ . Define  $h(\mathbf{X}_n)^*$  as in definition 1, but with the a.s. with respect to  $\mu_n \equiv P_n + Q_n$ . Then  $h(\mathbf{X}_n)^*$  works for both  $P_n$  and  $Q_n$ . But, since  $h$  is bounded and  $Q_n(\{dP_n/d\mu_n = 0\}) \rightarrow 0$  by contiguity (cf. lemma A.9.4),

$$(a) \quad E_{Q_n}^* h(\mathbf{X}_n) = E_{P_n}^*(h(\mathbf{X}_n)e^{\Lambda_n}) + o(1).$$

But by the continuous mapping theorem (proposition 1)

$$(b) \quad h(\mathbf{X}_n)e^{\Lambda_n} \Rightarrow h(\mathbf{X}_0)e^{\Lambda_0}.$$

Furthermore,

$$(c) \quad 0 \leq |h(\mathbf{X}_n)e^{\Lambda_n}| \leq \|h\|_\infty e^{\Lambda_n} \quad \text{for } n \geq 0,$$

and, by contiguity

$$(d) \quad E_{P_n}^*(e^{\Lambda_n}) = E_{P_n}(e^{\Lambda_n}) \rightarrow 1 = E(e^{\Lambda_0}).$$

Hence by Pratt's theorem (proposition 2), the right side of (a) converges to  $Eh(\mathbf{X}_0)e^{\Lambda_0}$ , and hence (13) holds.

The second assertion of the lemma is obvious.  $\square$

**Proof of theorem 1.** For  $m = 1, 2, \dots$ , let  $K_m$  be a compact such that  $\liminf_{n \rightarrow \infty} P_*(\mathbf{X}_n \in K_m^\delta) \geq 1 - 1/m$  for every  $\delta > 0$ .

**Step 1.** There exists  $\{n'\}$  such that  $\lim_{n' \rightarrow \infty} E^*h(\mathbf{X}_{n'})$  exists for every  $h \in C_b(M)$ ; say

$$(a) \quad \lim_{n' \rightarrow \infty} E^*h(\mathbf{X}_{n'}) = T(h).$$

**Proof of step 1.** Since  $K_m$  is compact,  $C(K_m)$  is separable. The same is true for  $\{h \in C(K_m) : \|h\|_\infty \leq 1\}$ . By the Tietze extension theorem (Jame-son (1974, theorem 12.4, page 113)), every  $h \in C(K_m)$  with  $\sup_{y \in K_m} |h(y)| \leq 1$  has an extension to an  $h \in C_b(M)$  with  $\|h\|_\infty \leq 1$ . Thus, there exists a countable subset of  $\{h \in C_b(M) : \|h\|_\infty \leq 1\}$ , of which the restrictions to  $K_m$  are dense in  $\{h \in C(K_m) : \|h\|_\infty \leq 1\}$ .

Let  $\{h_j\}_{j=1}^\infty \subset \{h \in C_b(M) : \|h\|_\infty \leq 1\}$  have this property for every  $m = 1, 2, \dots$ . By a diagonalization argument one can find  $\{n'\} \subset \{n\}$  such that

$$E^*h_j(\mathbf{X}_{n'}) \rightarrow T(h_j), \quad \text{as } n \rightarrow \infty$$

for every  $j = 1, 2, \dots$  and numbers  $T(h_j) \in [-1, 1]$ .

Fix  $h \in C_b(M)$  with  $\|h\|_\infty \leq 1$ . Given  $\varepsilon > 0$  and  $m$  there exists  $h_j$  with

$$(b) \quad \sup_{y \in K_m} |h(y) - h_j(y)| < \varepsilon.$$

As a consequence, there exists  $\delta > 0$  such that

$$(c) \quad \sup_{y \in K_m^\delta} |h(y) - h_j(y)| < 2\varepsilon.$$

(Indeed, suppose that there exists  $y_p \in K_m^{1/p}$  with  $|h(y_p) - h_j(y_p)| \geq \varepsilon$  for every  $p = 1, 2, \dots$ . Then  $d(y_p, K_m) \leq 1/p \rightarrow 0$  as  $p \rightarrow \infty$ , so that there exists  $x_p \in K_m$  with  $d(y_p, x_p) \rightarrow 0$ . Extract a subsequence  $\{p'\}$  with  $x_{p'} \rightarrow x \in K_m$ . Then  $y_{p'} \rightarrow x$  too, and by continuity of  $h - h_j$  it would follow that  $|h(x) - h_j(x)| \geq \varepsilon$ .)



Now since  $1_{K_n^\delta}(\mathbf{X}_n)^* + 1_{(K_n^\delta)^c}(\mathbf{X}_n)^* \equiv 1$  by lemma 2.C,

$$|E^* h(\mathbf{X}_n) - E^* h_j(\mathbf{X}_n)|$$

$$\begin{aligned} (d) \quad &\leq |E(h(\mathbf{X}_n)^* - h_j(\mathbf{X}_n)^*) 1_{K_n^\delta}(\mathbf{X}_n)^*| \\ &\quad + |E h(\mathbf{X}_n)^* 1_{(K_n^\delta)^c}(\mathbf{X}_n)^*| + |E h_j(\mathbf{X}_n)^* 1_{(K_n^\delta)^c}(\mathbf{X}_n)^*| \\ &\leq E |h(\mathbf{X}_n) - h_j(\mathbf{X}_n)|^* 1_{A_n} + 2P_n^*(\mathbf{X}_n \in (K_m^\delta)^c), \end{aligned}$$

for some measurable  $A_n \in \mathbf{A}_n$  with  $A_n \subset [\mathbf{X}_n \in K_m^\delta]$ ; this follows from lemma 2.D and 2.G. The last expression is smaller than  $2\varepsilon + 2/m$  for sufficiently large  $n$ .

We conclude that  $\{E^* h(\mathbf{X}_{n'})\}$  has the property that for every  $\eta > 0$  there is a converging sequence  $\{c_{n'}\}$  of numbers such that  $|E^* h(\mathbf{X}_{n'}) - c_{n'}|$  is eventually less than  $\eta$ . Thus  $E^* h(\mathbf{X}_{n'})$  converges, and (a) is proved.

**Step 2.** The map  $T : C_b(M) \rightarrow R$  is an abstract integral; i.e.,  $T(h)$  is linear, positive, and continuous at 0 on  $C_b(M)$ .

**Proof of step 2.** A. First linearity: For  $h_1, h_2 \in C_b(M)$

$$\begin{aligned} T(h_1 + h_2) &= \lim_{n' \rightarrow \infty} E^*(h_1 + h_2)(\mathbf{X}_{n'}) \\ &\leq \lim_{n' \rightarrow \infty} (E^* h_1(\mathbf{X}_{n'}) + E^* h_2(\mathbf{X}_{n'})) \\ &= T(h_1) + T(h_2) = \lim_{n' \rightarrow \infty} (E_* h_1(\mathbf{X}_{n'}) + E_* h_2(\mathbf{X}_{n'})) \\ &\leq \lim_{n' \rightarrow \infty} E_*(h_1(\mathbf{X}_{n'}) + h_2(\mathbf{X}_{n'})) = T(h_1 + h_2) \end{aligned}$$

by (a) and (ii). By a similar argument  $T(ch) = cT(h)$  for  $c \in R$  and  $h \in C_b(M)$ .

B.  $T$  is positive. This is trivial.

C. If  $\{h_p\} \subset C_b(M)$  with  $h_p \downarrow 0$  pointwise, then  $T(h_p) \downarrow 0$ . Indeed, by Dini's theorem  $h_p \downarrow 0$  uniformly on compacts. Fix  $\varepsilon > 0$  and  $m$ . For sufficiently large  $p$ ,

$$\sup_{y \in K_m} |h_p(y)| < \varepsilon.$$

Then, by the argument leading from (b) to (c), there exists  $\delta \equiv \delta_p > 0$  such that

$$\sup_{y \in K_m^\delta} |h_p(y)| < 2\varepsilon.$$

Thus for sufficiently large  $p$

$$\begin{aligned} T(h_p) &= \lim_{n' \rightarrow \infty} (E h_p(\mathbf{X}_{n'})^* 1_{K_m^\delta}(\mathbf{X}_{n'})^* + E h_p(\mathbf{X}_{n'})^* 1_{(K_m^\delta)^c}(\mathbf{X}_{n'})^*) \\ &\leq 2\varepsilon + \|h_1\|_\infty / m, \end{aligned}$$

proving that  $T$  is continuous from above at 0. Consequently,  $T$  is continuous at 0 and step 2 is proved.

**Step 3.** Let  $\mathcal{U}(C_b(M))$  be the smallest  $\sigma$ -field of subsets of  $M$  with respect

to which all  $h \in C_b(M)$  are measurable. As is well-known,  $\mathcal{U}(C_b(M))$  equals the Borel  $\sigma$ -field of  $M$ ; see, e.g., Bauer (1972, theorem 7.2.4, page 206). By the Daniell-Stone theorem (Bauer (1972, theorem 7.1.4, page 197)), there exists a probability measure  $L$  on  $\mathcal{U}(C_b(M))$  with

$$T(h) = \int h dL, \quad h \in C_b(M).$$

**Step 4.**  $L(K_m) \geq 1 - \frac{1}{m}$ .

**Proof of step 4.** Indeed there exist  $\{h_p\} \subset C_b(M)$  with

$$1_{K_m} \leq h_p \leq 1 \quad \text{and} \quad h_p \downarrow 1_{K_m} \quad \text{as} \quad p \rightarrow \infty;$$

for example, take  $h_p(x) = (1 - p d(x, K_m))^+$ . Thus

$$L(K_m) = \lim_{p \rightarrow \infty} \int h_p dL = \lim_{p \rightarrow \infty} \lim_{n' \rightarrow \infty} E_* h_p(\mathbf{X}_{n'}).$$

Let  $0 < r < 1$ . Since  $h_p = 1$  on  $K_m$ , there exists a  $\delta \equiv \delta(p, r) > 0$  such that

$$\inf_{y \in K_m^{\delta}} |h_p(y)| > r.$$

Thus

$$\lim_{n' \rightarrow \infty} E_* h_p(\mathbf{X}_{n'}) \geq \liminf_{n' \rightarrow \infty} r P_*(\mathbf{X}_{n'} \in K_m^{\delta}) \geq r(1 - \frac{1}{m}),$$

where the last inequality follows by the choice of  $K_m$  at the beginning of the proof, and the claim follows.  $\square$

**Proof of lemma 7.** Fix  $\varepsilon > 0$  and a compact set  $K$  such that  $\limsup_{n \rightarrow \infty} P^*(\mathbf{X}_n \notin K^{\delta}) \leq \varepsilon$  for every  $\delta > 0$ . By the Stone-Weierstrass theorem the restrictions of the functions in  $\mathbf{H}$  to  $K$  are uniformly dense in  $C_b(K)$ . Hence given  $h \in C_b(M)$  there exists  $g \in \mathbf{H}$  with  $|h(x) - g(x)| \leq \varepsilon/4$  for all  $x \in K$ . Using the compactness of  $K$ , there is a  $\delta > 0$  such that  $|h(x) - g(x)| \leq \varepsilon/3$  for every  $x \in K^{\delta}$ ; see the argument leading from (b) to (c) in the proof of theorem 1. Then

$$\begin{aligned} P(h(\mathbf{X}_n)^* - h(\mathbf{X}_n)_* > \varepsilon) &\leq P(\{h(\mathbf{X}_n)^* - h(\mathbf{X}_n)_* > \varepsilon\} \cap \{\mathbf{X}_n \in K^{\delta}\}_*) + \varepsilon \\ &\leq P\left(g(\mathbf{X}_n)^* - g(\mathbf{X}_n)_* > \frac{\varepsilon}{3}\right) + \varepsilon \\ &\leq \frac{3E\{g(\mathbf{X}_n)^* - g(\mathbf{X}_n)_*\}}{\varepsilon} + \varepsilon \\ &\rightarrow \varepsilon \quad \text{as } n \rightarrow \infty \end{aligned}$$

by hypothesis and lemma 2.I. Hence  $h(\mathbf{X}_n)^* - h(\mathbf{X}_n)_* \rightarrow_p 0$ . By dominated convergence this implies that  $E\{h(\mathbf{X}_n)^* - h(\mathbf{X}_n)_*\} \rightarrow 0$ .  $\square$

**Proof of corollary 1.** By lemmas 4 and 5, hypothesis (i) of theorem 1 is satisfied. Thus we only need to check (ii) of theorem 1.

For  $f \geq 0$  and  $g \geq 0$

$$(a) \quad \begin{aligned} f(\mathbf{X}_n)_* g(\mathbf{Y}_n)_* &\leq (f(\mathbf{X}_n)g(\mathbf{Y}_n))_* \leq f(\mathbf{X}_n)g(\mathbf{Y}_n) \\ &\leq (f(\mathbf{X}_n)g(\mathbf{Y}_n))^* \leq f(\mathbf{X}_n)^* g(\mathbf{Y}_n)^* . \end{aligned}$$

Thus

$$\begin{aligned} E^* f(\mathbf{X}_n)g(\mathbf{Y}_n) - E_* f(\mathbf{X}_n)g(\mathbf{Y}_n) &\leq E f(\mathbf{X}_n)^* g(\mathbf{Y}_n)^* - E f(\mathbf{X}_n)_* g(\mathbf{Y}_n)_* \\ &\leq E |f(\mathbf{X}_n)^* - f(\mathbf{X}_n)_*| |g(\mathbf{Y}_n)^*| \\ &\quad + E |f(\mathbf{X}_n)_*| |g(\mathbf{Y}_n)^* - g(\mathbf{Y}_n)_*| \\ &\leq \|g\|_\infty E(f(\mathbf{X}_n)^* - f(\mathbf{X}_n)_*) \\ &\quad + \|f\|_\infty E(g(\mathbf{Y}_n)^* - g(\mathbf{Y}_n)_*) \end{aligned}$$

$$(b) \quad \rightarrow 0 \quad \text{by (4) and (15) .}$$

For  $f \geq 0, g \geq 0$  and arbitrary  $a, b \in R$ ,

$$\begin{aligned} E^*(a+f)(\mathbf{X}_n)(b+g)(\mathbf{Y}_n) - E_*(a+f)(\mathbf{X}_n)(b+g)(\mathbf{Y}_n) &\leq ab + E^*(bf(\mathbf{X}_n)) + E^*(ag(\mathbf{Y}_n)) + E^*(f(\mathbf{X}_n)g(\mathbf{Y}_n)) \\ &\quad - \{ab + E_*(bf(\mathbf{X}_n)) + E_*(ag(\mathbf{Y}_n)) + E_*(f(\mathbf{X}_n)g(\mathbf{Y}_n))\} \\ &\quad \text{using lemma 2.A} \end{aligned}$$

$$(c) \quad \rightarrow 0 \quad \text{using lemma 2.F and the } \geq 0 \quad \text{case .}$$

Finally, for linear combinations,

$$\begin{aligned} E^* \sum_i f_i(\mathbf{X}_n)g_i(\mathbf{Y}_n) - E_* \sum_i f_i(\mathbf{X}_n)g_i(\mathbf{Y}_n) &\leq \sum_i (E^* f_i(\mathbf{X}_n)g_i(\mathbf{Y}_n) - E_* f_i(\mathbf{X}_n)g_i(\mathbf{Y}_n)) \quad \text{by lemma 2.A} \\ &\rightarrow 0 \quad \text{using (c) .} \end{aligned}$$

Now note that the collection  $\mathbf{H}$  of all functions of the form  $h(x,y) = \sum_i f_i(x)g_i(y)$  is an algebra and separates points of  $M = M_1 \times M_2$ . Thus  $\{(\mathbf{X}_n, \mathbf{Y}_n)\}$  is asymptotically measurable by lemma 7. The conclusion follows from theorem 1.  $\square$

**Proof of corollary 2.** We again verify the hypothesis (ii) of theorem 1 via lemma 7. If we take

$$(a) \quad \mathbf{H} \equiv \{h : M \rightarrow R : h(x) = g(x(f_1), \dots, x(f_m)) \text{ for some } f_1, \dots, f_m \in \mathbf{F} \text{ and } g \in C_b(R^m) \text{ for some } m \geq 1 \},$$

then (ii) implies that (15) holds for all  $h \in \mathbf{H}$ . Moreover,  $\mathbf{H}$  is an algebra and separates points of  $M = l^\infty(\mathbf{F})$ . Thus  $\{\mathbf{X}_n\}$  is asymptotically measurable by lemma 7. The conclusion follows from theorem 1.  $\square$

**Proof of corollary 3.** Use lemma 7 applied with  $\mathbf{H}$  equal to the set of measurable  $f \in C_b(M)$  to deduce asymptotic measurability of  $\{\mathbf{X}_n\}$ . The conclusion follows from theorem 1.  $\square$

**Proof of proposition 3.** Let  $h : M \rightarrow R$  be bounded and continuous, and suppose  $\mathbf{X}_n \Rightarrow_{\text{Dudley}} \mathbf{X}_0$ . Then

$$\begin{aligned} \text{(a)} \quad E^*h(\mathbf{X}_n) &\equiv \inf\left\{ \int g \, dP_n : g : \mathbf{X}_n \rightarrow R, g \geq h \circ \mathbf{X}_n, g \text{ measurable} \right\} \\ &\leq \inf\left\{ \int g \, dP_n : g = f \circ \mathbf{X}_n, f \geq h, f \mathcal{M}_b\text{-measurable} \right\} \\ &= \inf\left\{ \int f \, dP_n \circ \mathbf{X}_n^{-1} : f \geq h, f \mathcal{M}_b\text{-measurable} \right\} \\ &= \int^* h \, dP_n \circ \mathbf{X}_n^{-1}. \end{aligned}$$

Hence

$$\text{(b)} \quad \limsup_{n \rightarrow \infty} E^*h(\mathbf{X}_n) \leq \lim_{n \rightarrow \infty} \int^* h \, dP_n \circ \mathbf{X}_n^{-1} = \int h \, dP \circ \mathbf{X}_0^{-1} = Eh(\mathbf{X}_0),$$

and this implies  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  by (8).

On the other hand, if  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$ , then

$$E^*h(\mathbf{X}_n) \rightarrow Eh(\mathbf{X}_0) \quad \text{for all } h \in C_b(M),$$

so, in particular, for  $h \in C_b(M, \mathcal{M}_b) \subset C_b(M)$ ,  $(\mathcal{A}_n - \mathcal{M}_b)$ -measurability of  $\mathbf{X}_n$  implies that

$$Eh(\mathbf{X}_n) = E^*h(\mathbf{X}_n) \rightarrow Eh(\mathbf{X}_0);$$

i.e.,  $\mathbf{X}_n \Rightarrow_{\text{Dudley}} \mathbf{X}_0$ .  $\square$

**Proof of proposition 5.** A: Let  $\varepsilon, \delta > 0$ ; for  $k = 1, 2, \dots$  set

$$B_k \equiv \{x \in M : \text{for some } y \in M \text{ with } d(x, y) < 1/k, \\ d'(g(x), g(y)) > \varepsilon\}.$$

Then  $C_k \equiv B_k \cap M_0$  is relatively open in  $M_0$ , and hence  $C_k \in \mathcal{M}_b$ . Moreover, the sets  $C_k \downarrow \emptyset$  as  $k \rightarrow \infty$  by continuity of  $g$  on  $M_0$ . Choose  $k$  so large that  $P(\mathbf{X}_0 \in C_k) < \delta$ . Then

$$\begin{aligned} &[d'(g(\mathbf{X}_n), g(\mathbf{X}_0)) > \varepsilon] \\ &= [d'(g(\mathbf{X}_n), g(\mathbf{X}_0)) > \varepsilon] \cap ([\mathbf{X}_0 \in C_k^c] \cap [d(\mathbf{X}_n, \mathbf{X}_0)^* < \frac{1}{k}])^c \\ &\quad \cup [d'(g(\mathbf{X}_n), g(\mathbf{X}_0)) > \varepsilon] \cap [\mathbf{X}_0 \in C_k^c] \cap [d(\mathbf{X}_n, \mathbf{X}_0)^* < \frac{1}{k}] \end{aligned}$$

$$\text{(a)} \quad \subset [\mathbf{X}_0 \in C_k] \cup [d(\mathbf{X}_n, \mathbf{X}_0)^* \geq \frac{1}{k}] \cup [\mathbf{X}_0 \in M_0^c]$$

by the definition of  $C_k$ , and hence

$$P^*(d'(g(\mathbf{X}_n), g(\mathbf{X}_0)) > \varepsilon) \leq P(\mathbf{X}_0 \in C_k) + P(d(\mathbf{X}_n, \mathbf{X}_0)^* \geq \frac{1}{k})$$

$$(b) \qquad \qquad \qquad \leq \delta + \delta = 2\delta$$

for  $n$  sufficiently large.

B: Let  $C_k, k = 1, 2, \dots$  be the sets in the proof of A. Then, by (a),

$$(c) \quad [ \sup_{m \geq n} d'(g(\mathbf{X}_m), g(\mathbf{X}_0)) > \varepsilon ] = \bigcup_{m=n}^{\infty} [d'(g(\mathbf{X}_m), g(\mathbf{X}_0)) > \varepsilon]$$

$$\subset [ \mathbf{X}_0 \in C_k ] \cup \bigcap_{m=n}^{\infty} [d(\mathbf{X}_m, \mathbf{X}_0)^* \geq \frac{1}{k}] \cup [ \mathbf{X}_0 \in M_0^c ]$$

$$\subset [ \mathbf{X}_0 \in C_k ] \cup [ \sup_{m \geq n} d(\mathbf{X}_m, \mathbf{X}_0)^* \geq \frac{1}{k} ] \cup [ \mathbf{X}_0 \in M_0^c ],$$

and hence

$$(d) \quad P^* \left( \sup_{m \geq n} d'(g(\mathbf{X}_m), g(\mathbf{X}_0)) > \varepsilon \right)$$

$$\leq P(\mathbf{X}_0 \in C_k) + P(\sup_{m \geq n} d(\mathbf{X}_m, \mathbf{X}_0)^* \geq \frac{1}{k})$$

$$\leq \delta + \delta = 2\delta$$

for  $n$  sufficiently large. □

**Proof of proposition 6.** First we show that  $g|_{M_0}$  is continuous. Note that

(ii) is equivalent to:

$$(a) \quad \text{for every } x \in M_0 \quad \text{both } g_n(x) \rightarrow g(x)$$

$$\text{and } \lim_{\delta \rightarrow 0} \limsup_{n \rightarrow \infty} \omega_{g_n}(S(x, \delta)) = 0,$$

where  $S(x, \delta) \subset M$  is the open sphere of radius  $\delta$  centered at  $x \in M_0$  and  $\omega_g(S)$  denotes the oscillation of  $g$  on the set  $S$ . Hence for  $x, y \in M_0$ ,

$$d'(g(y), g(x)) \leq d'(g(y), g_n(y)) + d'(g_n(y), g_n(x)) + d'(g_n(x), g(x))$$

$$\leq d'(g(y), g_n(y)) + \omega_{g_n}(S(x, 2d(x, y))) + d'(g_n(x), g(x))$$

$$\rightarrow 0$$

as  $n \rightarrow \infty$  and then  $d(x, y) \downarrow 0$ .

**First proof of A.** Our first proof is based on the representation theorem 2.

Since  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$ , by theorem 2 there exist  $\tilde{\mathbf{X}}_n =_d \mathbf{X}_n, n = 0, 1, \dots$  such that  $\tilde{\mathbf{X}}_n \rightarrow \tilde{\mathbf{X}}_0$  a.u. (almost uniformly). Then

$$\tilde{\mathbf{Y}}_n \equiv g_n(\tilde{\mathbf{X}}_n) \rightarrow_{\text{a.u.}} g(\tilde{\mathbf{X}}_0) \equiv \tilde{\mathbf{Y}}_0 \quad \text{by C,}$$

so  $\tilde{\mathbf{Y}}_n \Rightarrow \tilde{\mathbf{Y}}_0$  by lemma 8. Thus for  $f \in C_b(M')$

$$E^* f(g_n \mathbf{X}_n) = E^* f(g_n \mathbf{X}_n \phi_n)$$

$$= E^* f(\tilde{\mathbf{Y}}_n) \rightarrow Ef(\tilde{\mathbf{Y}}_0) = Ef(g\mathbf{X}_0).$$

□

**Second proof of A.** This proof is an adaptation of the classical proof given in Billingsley (1968). Let  $G \subset M'$  be open. By lemma 3, it suffices to show that

$$(b) \quad \liminf_{n \rightarrow \infty} P_*(g_n(\mathbf{X}_n) \in G) \geq P(g(\mathbf{X}_0) \in G).$$

Now (ii) implies that for every  $x \in M_0$  and every  $\varepsilon > 0$  there exist  $k, \delta > 0$  such that  $i \geq k$  and  $d(x, y) < \delta$  imply  $d'(g(x), g_i(y)) < \varepsilon$ . Set

$$T_k \equiv \bigcap_{i \geq k} \{x \in M : g_i(x) \in G\} = \bigcap_{i \geq k} g_i^{-1}G.$$

Then, with  $T_k^0 \equiv \text{interior}(T_k)$ ,

$$(c) \quad T_k^0 \subset T_k \subset g_n^{-1}G \quad \text{for } n \geq k.$$

Note that  $x \in M_0 \cap g^{-1}G$  implies that for  $\delta$  sufficiently small and  $k$  sufficiently large, whenever  $d(x, y) < \delta$  and  $i \geq k$  it follows that  $g_i(y) \in G$ ; hence the open sphere of radius  $\delta$  about  $x$ ,  $S(x, \delta) \subset T_k$ , so that  $x \in T_k^0$ , an open set. Thus

$$(d) \quad g^{-1}G = (g^{-1}G \cap M_0) \cup (g^{-1}G \cap M_0^c) \\ \subset \bigcup_k T_k^0 \cup M_0^c.$$

Since  $P(\mathbf{X}_0 \in M_0^c) = 0$  by (i),

$$(e) \quad P(g(\mathbf{X}_0) \in G) \leq P(\mathbf{X}_0 \in \bigcup_k T_k^0) \quad \text{by (d)} \\ \leq P(\mathbf{X}_0 \in T_k^0) + \varepsilon$$

for  $k$  sufficiently large since

$$T_k \uparrow \text{ implies } T_k^0 \uparrow$$

$$\leq \liminf_{n \rightarrow \infty} P_*(\mathbf{X}_n \in T_k^0) + \varepsilon$$

by  $\mathbf{X}_n \Rightarrow \mathbf{X}_0$  and lemma 3(iii)

$$\leq \liminf_{n \rightarrow \infty} P_*(g_n(\mathbf{X}_n) \in G) + \varepsilon \quad \text{by (c).}$$

Since  $\varepsilon$  is arbitrary, (b) holds.

**Proof of B.** For  $x \in M$  and  $\varepsilon > 0$  define

$$(f) \quad k(x, \varepsilon) \equiv \min \{ k : \text{for all } y \text{ with } d(x, y) < 1/k,$$

$$\text{for all } n \geq k, \quad d'(g_n(y), g(x)) \leq \varepsilon \}.$$

Note that this is well defined: if the set of  $k$ 's in (f) is empty, then there exist points  $y_k \in M$  and integers  $n_k$  such that  $n_{k-1} < n_k$ ,  $d(y_k, x) < 1/k$ , and  $d'(g_{n_k}(y_k), g(x)) > \varepsilon$  for  $k \geq 1$ , contradicting the hypothesis (ii).

**Claim.**  $k(\cdot, \varepsilon)$  is a measurable function.

**Proof.** First we show that

$$(g) \quad k(x, \varepsilon) \leq \liminf_{m \rightarrow \infty} k(x_m, \varepsilon)$$

for any sequence  $\{x_m\} \subset M_0$  with  $x_m \rightarrow x \in M_0$ .

Since  $k$  is integer-valued, the liminf is achieved for some subsequence  $\{x_{m'}\}$ , and in fact  $\liminf_m k(x_m, \varepsilon) = k(x_{m'}, \varepsilon) \equiv k'$  for all  $m'$  sufficiently large. Suppose that the right side is finite for some subsequence (if not, then the inequality is trivially true). If  $d(x, y) < 1/k'$ , then there exists an  $m_0$  such that  $d(x_{m'}, y) < 1/k'$  for all  $m' \geq m_0$ . Hence

$$(h) \quad d'(g_n(y), g(x_{m'})) \leq \varepsilon \quad \text{for all } n \geq k'.$$

Since  $g|_{M_0}$  is continuous as we have shown at the beginning of the proof of proposition 6, we can let  $m' \rightarrow \infty$  in (h) to obtain

$$(i) \quad d'(g_n(y), g(x)) \leq \varepsilon \quad \text{for all } n \geq k'.$$

Hence

$$k(x, \varepsilon) \leq k' = \liminf_{m \rightarrow \infty} k(x_m, \varepsilon).$$

and (g) holds. It follows from this that for any fixed integer  $K$  the set  $\{x \in M_0: k(x, \varepsilon) \leq K\}$  is closed, and hence  $k(\cdot, \varepsilon)$  is measurable, proving the claim.

Now we are ready to prove B. Let  $\varepsilon > 0$ . By the claim,  $k(\mathbf{X}_0, \varepsilon)$  is a (proper) random variable, and there exists a  $k_0 = k_0(\varepsilon)$  such that

$$(j) \quad P(k(\mathbf{X}_0, \varepsilon) > k_0) < \frac{\varepsilon}{2}.$$

Since  $d(\mathbf{X}_n, \mathbf{X}_0)^* \rightarrow_p 0$ , for all  $n \geq n_0(\varepsilon)$ ,

$$(k) \quad P(d(\mathbf{X}_n, \mathbf{X}_0)^* \geq \frac{1}{k_0}) < \frac{\varepsilon}{2}.$$

Then, with

$$B_n \equiv \{d'(g_n(\mathbf{X}_n), g(\mathbf{X}_0))^* > \varepsilon\},$$

$$(l) \quad C_n \equiv \{d(\mathbf{X}_n, \mathbf{X}_0)^* \geq \frac{1}{k_0}\},$$

$$D \equiv \{k(\mathbf{X}_0, \varepsilon) > k_0\},$$

for  $n \geq \max\{n_0, k_0\}$ ,

$$\begin{aligned} P(B_n) &\leq P(B_n \cap (C_n^c \cap D^c)) + P(B_n \cap (C_n^c \cap D^c)^c) \\ &\leq P(\emptyset) + P(C_n) + P(D) \quad \text{by (f)} \\ &< \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \quad \text{by (j)-(l),} \end{aligned}$$

and B holds.

**Proof of C.** Let  $\varepsilon > 0$  and write

$$(m) \quad \tilde{B}_n \equiv \left\{ \sup_{m \geq n} d'(g_m(\mathbf{X}_m), g(\mathbf{X}_0))^* > \varepsilon \right\} \\ = \bigcup_{m=n}^{\infty} \left\{ d'(g_m(\mathbf{X}_m), g(\mathbf{X}_0))^* > \varepsilon \right\}$$

and

$$(n) \quad \tilde{C}_n \equiv \left\{ \sup_{m \geq n} d(\mathbf{X}_m, \mathbf{X}_0)^* \geq \frac{1}{k_0} \right\} \supset \bigcup_{m=n}^{\infty} \left\{ d(\mathbf{X}_m, \mathbf{X}_0)^* \geq \frac{1}{k_0} \right\}.$$

By hypothesis, there is an  $n_0 = n_0(\varepsilon)$  such that

$$(o) \quad P(\tilde{C}_n) < \frac{\varepsilon}{2} \quad \text{for all } n \geq n_0.$$

Hence

$$P(\tilde{B}_n) \leq P(\tilde{B}_n \cap (\tilde{C}_n^c \cap D^c)) + P(\tilde{B}_n \cap (\tilde{C}_n^c \cap D)^c) \\ \leq P(\emptyset) + P(\tilde{C}_n) + P(D) \quad \text{by (f)} \\ < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \quad \text{by (j), (n), and (o),}$$

for all  $n \geq \max\{n_0, k_0\}$ , and C holds. □

**Proof of lemma 9.** Let  $F \subset M$  be closed and  $\varepsilon > 0$ . Then

$$P_n^*(\mathbf{Y}_n \in F) \leq P_n^*(\mathbf{Y}_n \in F \quad \text{and} \quad d(\mathbf{X}_n, \mathbf{Y}_n)^* \leq \varepsilon) \\ + P_n(d(\mathbf{X}_n, \mathbf{Y}_n)^* > \varepsilon) \\ \leq P_n^*(\mathbf{X}_n \in \bar{F}^\varepsilon) + o(1).$$

Thus, by lemma 3,  $\limsup_{n \rightarrow \infty} P^*(\mathbf{Y}_n \in F) \leq P(\mathbf{X}_0 \in \bar{F}^\varepsilon)$  for every  $\varepsilon > 0$ . Finally let  $\varepsilon \downarrow 0$ ; the conclusion then follows from (iv) implies (i) of the portmanteau lemma 3. □

**Proof of lemma 8.** Apply lemma 9 with  $\mathbf{X}_n = \mathbf{X}_0$  and  $\mathbf{Y}_n = \mathbf{X}_n$ ,  $n = 1, 2, \dots$ . □

## A.9 CONTIGUITY

Consider a sequence of statistical problems (with only two sequences of probability measures) with:

Measure spaces	$(\mathbf{X}_n, \mathcal{A}_n, \mu_n),$
Probability measures	$P_n \ll \mu_n, \quad Q_n \ll \mu_n,$
Densities	$p_n \equiv \frac{dP_n}{d\mu_n}, \quad q_n \equiv \frac{dQ_n}{d\mu_n},$



$$\text{Likelihood ratios } L_n \equiv \begin{cases} q_n/p_n & \text{if } p_n > 0, \\ 1 & \text{if } q_n = p_n = 0, \\ \infty & \text{if } q_n > 0 = p_n, \end{cases}$$

$$\text{Log-likelihood ratios } \Lambda_n \equiv \log L_n .$$

**Definition 1.** The sequence  $\{Q_n\}$  is *contiguous* to  $\{P_n\}$  if for every sequence  $B_n \in \mathcal{A}_n$  for which  $P_n(B_n) \rightarrow 0$  it follows that  $Q_n(B_n) \rightarrow 0$ .

Thus contiguity of  $\{Q_n\}$  to  $\{P_n\}$  means that  $Q_n$  is “asymptotically absolutely continuous” with respect to  $P_n$  in the sense of domination of measures. We therefore denote contiguity of  $\{Q_n\}$  to  $\{P_n\}$  by  $\{Q_n\} \triangleleft \{P_n\}$ , a notation due to Witting and Nölle (1970). Two sequences are contiguous to each other if both  $\{Q_n\} \triangleleft \{P_n\}$  and  $\{P_n\} \triangleleft \{Q_n\}$  and we then write  $\{P_n\} \triangleleft \triangleright \{Q_n\}$ . Contiguity has been introduced by Le Cam (1960). We will state Le Cam’s lemmas 1, 2, and 3 and give proofs along the lines of Hájek and Šidák (1967).

**Definition 2.** The sequence  $\{Q_n\}$  is *asymptotically orthogonal* to  $\{P_n\}$  if there exists  $B_n \in \mathcal{A}_n$  such that  $Q_n(B_n) \rightarrow 1$  and  $P_n(B_n) \rightarrow 0$ .

**Lemma 1.** (Le Cam’s first lemma). Suppose that  $L(L_n | P_n) \rightarrow L(L)$  and  $E(L) = 1$ . Then  $\{Q_n\} \triangleleft \{P_n\}$ .

**Proof.** Let  $B_n \in \mathcal{A}_n$  with  $P_n(B_n) \rightarrow 0$ . By the Neyman-Pearson lemma there is a critical function  $\phi_n \equiv 1_{[L_n > k_n]} + \gamma_n 1_{[L_n = k_n]}$  such that  $E_{P_n}(\phi_n) = \alpha_n \equiv P_n(B_n) \rightarrow 0$  and

$$Q_n(B_n) \leq E_{Q_n}(\phi_n) .$$

But for any fixed  $0 < y < \infty$

$$\begin{aligned} Q_n(B_n) &\leq E_{Q_n}(\phi_n) = E_{Q_n}(\phi_n 1_{[L_n \leq y]}) + E_{Q_n}(\phi_n 1_{[L_n > y]}) \\ &\leq y E_{P_n}(\phi_n) + E_{Q_n}(1_{[L_n > y]}) \end{aligned}$$

$$(a) \quad \leq y P_n(B_n) + 1 - E_{P_n}(L_n 1_{[L_n \leq y]}) .$$

Let  $\varepsilon > 0$  and choose  $y$  to be a continuity point of  $L(L)$  such that  $1 - E(L 1_{[L \leq y]}) < \varepsilon/2$ ; this is possible since  $E(L) = 1$  by hypothesis. Then  $L(L_n | P_n) \rightarrow L(L)$  implies that  $E_{P_n}(L_n 1_{[L_n \leq y]}) \rightarrow E(L 1_{[L \leq y]})$  and hence  $1 - E_{P_n}(L_n 1_{[L_n \leq y]}) < \varepsilon$  for  $n$  sufficiently large. Since  $P_n(B_n) \rightarrow 0$  we also have  $y P_n(B_n) < \varepsilon$  for  $n$  sufficiently large, and hence it follows from (a) that  $Q_n(B_n) < 2\varepsilon$  for  $n$  sufficiently large.  $\square$

**Corollary 1.** If  $L(\Lambda_n | P_n) \rightarrow L(\Lambda) = N(-\frac{1}{2}\sigma^2, \sigma^2)$ , then  $\{Q_n\} \triangleleft \{P_n\}$ .

**Proof.** Note that  $L(L) = L(e^{\sigma Z - \sigma^2/2})$  where  $L(Z) = N(0,1)$  and hence  $E(L) = 1$ .  $\square$

Now suppose that  $(X_n, \mathcal{A}_n, \mu_n) = (X, \mathcal{A}, \mu)^n$  and  $X_n = (X_{n1}, \dots, X_{nn}) \in X_n$  and that

$$p_n(x_n) = \prod_{i=1}^n f_{ni}(x_{ni}), \quad P_n \equiv \prod_{i=1}^n P_{ni},$$

$$q_n(x_n) = \prod_{i=1}^n g_{ni}(x_{ni}), \quad Q_n \equiv \prod_{i=1}^n Q_{ni},$$

so that

$$(1) \quad \Lambda_n = \sum_{i=1}^n \log \left( \frac{g_{ni}}{f_{ni}}(X_{ni}) \right) \quad \begin{cases} < \infty & \text{a.s. } P_n, \\ > -\infty & \text{a.s. } Q_n. \end{cases}$$

Suppose that the summands in (1) satisfy the *uniform asymptotic negligibility (UAN)* condition

$$(2) \quad \lim_{n \rightarrow \infty} \max_{1 \leq i \leq n} P_n(|\frac{g_{ni}}{f_{ni}}(X_{ni}) - 1| > \epsilon) = 0 \quad \text{for all } \epsilon > 0.$$

To get random variables with finite variance (to which classical central limit theorems may be applied), let

$$(3) \quad W_n \equiv 2 \sum_{i=1}^n \left\{ \frac{g_{ni}^{1/2}}{f_{ni}^{1/2}}(X_{ni}) - 1 \right\}$$

and note that under  $P_n$

$$\text{Var} \left( \frac{g_{ni}^{1/2}}{f_{ni}^{1/2}}(X_{ni}) \right) \leq E \left( \frac{g_{ni}}{f_{ni}}(X_{ni}) \right) = \int 1_{[f_{ni} > 0]} g_{ni} d\mu \leq 1.$$

The following lemma reduces the proof of asymptotic normality of  $\Lambda_n$  to the problem of establishing asymptotic normality of  $W_n$ .

**Lemma 2.** (Le Cam's second lemma). Suppose that the UAN condition (2) holds and  $L(W_n | P_n) \rightarrow N(-\frac{1}{4}\sigma^2, \sigma^2)$ . Then

$$(4) \quad \Lambda_n - (W_n - \frac{1}{4}\sigma^2) = o_p(1)$$

under  $P_n$  and hence

$$(5) \quad L(\Lambda_n | P_n) \rightarrow N(-\frac{1}{2}\sigma^2, \sigma^2).$$

The proof of lemma 2 involves a long truncation argument, and is therefore deferred to the end of the section.

**Corollary 2.** If  $f_n$  is a sequence of densities such that

$$\| \sqrt{n} (f_n^{1/2} - f^{1/2}) - \delta \| \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

where  $\| \cdot \|$  is the  $L_2(\mu)$ -metric and  $\delta \in L_2(\mu)$ , then with  $p_n(x) \equiv \prod_{i=1}^n f(x_i)$  and  $q_n(x) \equiv \prod_{i=1}^n f_n(x_i)$  it follows that, under  $p_n$ ,

$$(6) \quad \Lambda_n - \left( \frac{2}{\sqrt{n}} \sum_{i=1}^n \frac{\delta}{f^{1/2}}(X_i) - 2 \|\delta\|^2 \right) = o_p(1)$$

and hence

$$(7) \quad L(\Lambda_n | P_n) \rightarrow N\left(-\frac{1}{2}\sigma^2, \sigma^2\right)$$

with  $\sigma^2 = 4 \|\delta\|^2$ .

**Proof.** Note that the hypothesis of the corollary implies both

$$(a) \quad n \|f_n^{1/2} - f^{1/2}\|^2 \rightarrow \|\delta\|^2,$$

$$(b) \quad \|f_n^{1/2} - f^{1/2}\|^2 \rightarrow 0,$$

and, via

$$\begin{aligned} 1 &= \int f_n \, d\mu = \int [f^{1/2} + n^{-1/2}\delta + (f_n^{1/2} - f^{1/2} - n^{-1/2}\delta)]^2 \, d\mu \\ &= 1 + n^{-1/2} \int \delta f^{1/2} \, d\mu + o(n^{-1/2}), \end{aligned}$$

the orthogonality relation

$$(c) \quad \int \delta f^{1/2} \, d\mu = 0.$$

Consequently,

$$\begin{aligned} E P_n \left( \left| \frac{f_n}{f}(X_i) - 1 \right| \geq \varepsilon \right) &\leq E \left| \frac{f_n}{f} - 1 \right| \\ &= E \left( \left| \frac{f_n^{1/2}}{f^{1/2}} - 1 \right| \left| \frac{f_n^{1/2}}{f^{1/2}} + 1 \right| \right) \\ &\leq \|f_n^{1/2} - f^{1/2}\| \|f_n^{1/2} + f^{1/2}\| \\ &\rightarrow 0 \quad \text{uniformly in } 1 \leq i \leq n \end{aligned}$$

as  $n \rightarrow \infty$ , and the UAN condition (2) holds. Furthermore, the hypothesis, (c) and (a) imply

$$\begin{aligned} &E \left\{ W_n - \frac{2}{\sqrt{n}} \sum_{i=1}^n \frac{\delta}{f^{1/2}}(X_i) + \|\delta\|^2 \right\}^2 \\ &= 4n E \left\{ \frac{f_n^{1/2}}{f^{1/2}}(X_1) - 1 - \frac{1}{\sqrt{n}} \frac{\delta}{f^{1/2}}(X_1) + \frac{1}{2n} \|\delta\|^2 \right\}^2 \\ &\quad + n(n-1) \left[ E \left\{ 2 \frac{f_n^{1/2}}{f^{1/2}}(X_1) - 2 - \frac{2}{\sqrt{n}} \frac{\delta}{f^{1/2}}(X_1) + \frac{1}{n} \|\delta\|^2 \right\} \right]^2 \\ &= o(1) + \left(1 - \frac{1}{n}\right) \left[ -n \|f_n^{1/2} - f^{1/2}\|^2 + \|\delta\|^2 \right]^2 \\ &\rightarrow 0. \end{aligned}$$

Thus

$$(d) \quad W_n - \frac{2}{\sqrt{n}} \sum_{i=1}^n \frac{\delta}{f^{1/2}}(X_i) + \|\delta\|^2 = o_P(1).$$

Since, by (c),

$$L\left(\frac{2}{\sqrt{n}} \sum_{i=1}^n \frac{\delta}{f^{1/2}}(X_i) \mid P_n\right) \rightarrow N(0, 4\|\delta\|^2)$$

it follows that

$$L(W_n \mid P_n) \rightarrow N(-\|\delta\|^2, 4\|\delta\|^2)$$

and hence, by lemma 2, that

$$\begin{aligned} \Lambda_n &= W_n - \|\delta\|^2 + o_P(1) \\ &= \frac{2}{\sqrt{n}} \sum_{i=1}^n \frac{\delta}{f^{1/2}}(X_i) - 2\|\delta\|^2 + o_P(1). \end{aligned}$$

□

Now suppose that under

$$(8) \quad P_n: X_{n1}, \dots, X_{nn} \text{ are i.i.d. } f,$$

and under

$$(9) \quad Q_n: X_{n1}, \dots, X_{nn} \text{ are independent with densities } f_{n1}, \dots, f_{nn}$$

with respect to  $\mu$ . Assume that  $a_n = (a_{n1}, \dots, a_{nn})$ ,  $n \geq 1$ , is a vector of real constants which satisfy

$$(10) \quad \max_{1 \leq i \leq n} \frac{a_{ni}^2}{|a_n|^2} \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

and suppose there exists  $\delta \in L_2(\mu)$  such that

$$(11) \quad \sum_{i=1}^n \|(f_{ni}^{1/2} - f^{1/2}) - \frac{a_{ni}}{|a_n|} \delta\|^2 \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

**Corollary 3.** Suppose that (8)–(11) hold. Then, under  $P_n$ ,

$$(12) \quad \Lambda_n - (Z_n - 2\|\delta\|^2) = o_P(1),$$

where

$$(13) \quad Z_n \equiv 2 \sum_{i=1}^n \frac{a_{ni}}{|a_n|} \frac{\delta}{f^{1/2}}(X_{ni})$$

and

$$(14) \quad L(Z_n \mid P_n) \rightarrow N(0, 4\|\delta\|^2) \quad \text{as } n \rightarrow \infty.$$

**Proof.** See Shorack and Wellner (1986, pages 154, 163–165). □

Note that Corollary 2 is the special case of Corollary 3 with all  $a_{ni} = 1$  and

$f_{ni} = f_n$  for all  $i = 1, \dots, n$ . Now we return to the general situation as described at the beginning of this appendix.

**Lemma 3.** (Le Cam's third lemma). Suppose that a statistic  $T_n$  satisfies

$$(15) \quad L((T_n, \Lambda_n)^T | P_n) \rightarrow L((T, \Lambda)^T) = N_2 \left( \begin{pmatrix} \mu \\ -\sigma^2/2 \end{pmatrix}, \begin{pmatrix} \tau^2 & c \\ c & \sigma^2 \end{pmatrix} \right).$$

Then

$$(16) \quad L((T_n, \Lambda_n)^T | Q_n) \rightarrow L((T + c, \Lambda + \sigma^2)^T) \\ = N_2 \left( \begin{pmatrix} \mu + c \\ \sigma^2/2 \end{pmatrix}, \begin{pmatrix} \tau^2 & c \\ c & \sigma^2 \end{pmatrix} \right).$$

**Remark 1.** If  $T_n$  is asymptotically linear and  $\Lambda_n$  is asymptotically linear, for example as in (6) or (12), then verification of (15) is straightforward via the multivariate central limit theorem.

**Remark 2.** This is a particular instance of lemma A.8.6.

Recall definition A.7.2 of uniform integrability of random variables  $X_n$  with distribution  $P_n$ .

**Proposition 1.**  $\{X_n\}$  is uniformly integrable if and only if both of the following hold:

- (i)  $\sup_{n \geq 1} E(|X_n|) < \infty$ .
- (ii) For all sequences  $\{B_n\}$  with  $B_n \in \mathcal{A}_n$ , the convergence  $P_n(B_n) \rightarrow 0$  implies  $E(|X_n| 1_{B_n}) \rightarrow 0$ .

**Proof.** See e.g. Billingsley (1968, problem 5.5, page 34) or Chow and Teicher (1978, pages 92–93). □

**Lemma 4.** (Hall and Loynes (1977).)  $\{Q_n\} \triangleleft \{P_n\}$  if and only if  $\{L_n\}$  is uniformly integrable with respect to  $\{P_n\}$  and  $Q_n(p_n = 0) \rightarrow 0$ .

**Proof.** First note that for  $B_n \in \mathcal{A}_n$  we have

$$\begin{aligned} Q_n(B_n) &= \int 1_{B_n} dQ_n \\ &= \int 1_{B_n \cap \{p_n=0\}} dQ_n + \int 1_{B_n \cap \{p_n>0\}} L_n dP_n \\ &= \int 1_{B_n \cap \{p_n=0\}} dQ_n + \int 1_{B_n} L_n dP_n \\ (a) \quad &\leq Q_n(p_n = 0) + \int 1_{B_n} L_n dP_n \\ (b) \quad &\geq \int 1_{B_n} L_n dP_n . \end{aligned}$$

Thus if  $P_n(B_n) \rightarrow 0$ ,  $L_n$  is uniformly integrable and  $Q_n(p_n = 0) \rightarrow 0$ , then  $Q_n(B_n) \rightarrow 0$  by (a) and proposition 1, so  $\{Q_n\} \triangleleft \{P_n\}$ .

Conversely, if  $\{Q_n\} \triangleleft \{P_n\}$  so that  $P_n(B_n) \rightarrow 0$  implies  $Q_n(B_n) \rightarrow 0$ , then (b) implies that  $\int 1_{B_n} L_n dP_n \rightarrow 0$  so (ii) of proposition 1 holds. Part (i) of proposition 1 holds trivially since  $E_{P_n}(L_n) = \int L_n dP_n = \int 1_{\{p_n>0\}} dQ_n \leq 1$ , and therefore  $\{L_n\}$  is uniformly integrable with respect to  $\{P_n\}$  by proposition

1. Since  $P_n(p_n = 0) = 0$ , contiguity implies that  $Q_n(p_n = 0) \rightarrow 0$ . □

**Proof of lemma 3.**  $\{Q_n\} \triangleleft \{P_n\}$  by corollary 1 to Le Cam's first lemma. Hence by lemma 4,  $L_n$  is uniformly integrable and  $Q_n(p_n = 0) \rightarrow 0$  as  $n \rightarrow \infty$ .

Now let  $f: R^2 \rightarrow R$  be bounded and continuous. Then

$$\begin{aligned}
 E_{Q_n} f(T_n, \Lambda_n) &= E_{Q_n} f(T_n, \Lambda_n) \{ 1_{[p_n > 0]} + 1_{[p_n = 0]} \} \\
 \text{(a)} \qquad \qquad \qquad &= E_{P_n} f(T_n, \Lambda_n) L_n + E_{Q_n} f(T_n, \Lambda_n) 1_{[p_n = 0]} \\
 \text{(b)} \qquad \qquad \qquad &\rightarrow E \{ f(T, \Lambda) L \} \\
 \text{(c)} \qquad \qquad \qquad &= E f(T + c, \Lambda + \sigma^2),
 \end{aligned}$$

where (b) holds since  $f(T_n, \Lambda_n) L_n$  is uniformly integrable by uniform integrability of  $L_n$  and boundedness of  $f$ , and since the second term in (a) is bounded by  $\|f\|_\infty Q_n(p_n = 0) \rightarrow 0$ . It remains only to establish (c).

To verify (c), let  $U, V, W$  be i.i.d. and hence  $N(0,1)$ . Note that  $L = \exp(\sigma U - \sigma^2/2)$  can be viewed as the likelihood ratio of a  $N(\sigma, 1)$  and a  $N(0, 1)$  distribution. Consequently, for any bounded measurable function  $g$ ,

$$\text{(d)} \qquad E(g(U,V)L) = E g(W + \sigma, V),$$

since  $W + \sigma$  has a  $N(\sigma, 1)$  distribution. Now  $(T, \Lambda)^T$  has the same distribution as

$$\left( \rho\tau U + (1 - \rho^2)^{1/2}\tau V + \mu, \sigma U - \frac{\sigma^2}{2} \right)^T$$

with  $\rho\sigma\tau = c$ , and we obtain by (d)

$$\begin{aligned}
 E \{ f(T, \Lambda) L \} &= E f \left( \rho\tau W + \rho\sigma\tau + (1 - \rho^2)^{1/2}\tau V + \mu, \sigma W + \frac{\sigma^2}{2} \right) \\
 &= E f(T + c, \Lambda + \sigma^2).
 \end{aligned}$$

Hence (c) and (16) hold. □

*Contiguity of Product Measures*

Throughout the following we assume that the  $P_n$ 's and  $Q_n$ 's are product measures as in (1). By section A.6 we have

$$\text{(17)} \qquad d_H^2(P, Q) \leq d_v(P, Q) \leq 2 d_H(P, Q)$$

and

$$\text{(18)} \qquad d_H^2(P_n, Q_n) = 2 \left\{ 1 - \prod_{i=1}^n \left( -\frac{1}{2} d_H^2(P_{ni}, Q_{ni}) \right) \right\}.$$

**Proposition 2.** (Oosterhoff and Van Zwet (1979))

A. If  $\sum_{i=1}^n d_H^2(P_{ni}, Q_{ni}) = o(1)$ , then  $\{P_n\} \triangleleft \{Q_n\}$ .

B. If  $\{Q_n\} \triangleleft \{P_n\}$  then  $\sum_{i=1}^n d_H^2(P_{ni}, Q_{ni}) = o(1)$ .

**Proof.** The hypothesis in A implies that  $\sum_{i=1}^n \log(1 - d_H^2(P_{ni}, Q_{ni})/2) = o(1)$ , which by (18) yields  $d_H^2(P_n, Q_n) = o(1)$ , and hence  $d_v(P_n, Q_n) = o(1)$  by (17). But this implies the conclusion of A.

To prove B, suppose that

$$\limsup_{n \rightarrow \infty} d_H(P_n, Q_n) = 2^{1/2};$$

then by (17)  $\limsup_{n \rightarrow \infty} d_v(P_n, Q_n) = 1$  which contradicts  $\{Q_n\} \triangleleft \{P_n\}$ .

Thus  $\limsup_{n \rightarrow \infty} d_H^2(P_n, Q_n) < 2$  and hence  $\liminf_{n \rightarrow \infty} \prod_{i=1}^n (1 - d_H^2(P_{ni}, Q_{ni})/2) > 0$  or  $\limsup_{n \rightarrow \infty} \sum_{i=1}^n d_H^2(P_{ni}, Q_{ni}) < \infty$ .  $\square$

**Theorem 1.** (Oosterhoff and Van Zwet (1979)).  $\{Q_n\} \triangleleft \{P_n\}$  if and only if

$$(19) \quad \limsup_{n \rightarrow \infty} \sum_{i=1}^n d_H^2(P_{ni}, Q_{ni}) < \infty$$

and

$$(20) \quad \lim_{n \rightarrow \infty} \sum_{i=1}^n Q_{ni} \left( \frac{g_{ni}}{f_{ni}}(X_{ni}) \geq c_n \right) = 0 \quad \text{whenever } c_n \rightarrow \infty.$$

**Proof.** See Oosterhoff and Van Zwet (1979).  $\square$

**Corollary 4.** If  $P_{ni} = P$  and  $Q_{ni} = Q_{n1}$  (and  $f_{ni} = f_{n1} \equiv f$ ,  $g_{ni} = g_{n1} \equiv g_n$ ) for all  $i = 1, \dots, n$ , then  $\{Q_n\} \triangleleft \{P_n\}$  if and only if

$$\limsup_{n \rightarrow \infty} n d_H^2(Q_{n1}, P) < \infty$$

and

$$n Q_{n1} \left( \frac{g_n}{f}(X_{n1}) \geq c_n \right) \rightarrow 0 \quad \text{whenever } c_n \rightarrow \infty.$$

**Theorem 2.** If  $\Lambda_n$  is as in (1), then

$$L(\Lambda_n | P_n) \rightarrow N\left(-\frac{1}{2}\sigma^2, \sigma^2\right)$$

and

$$\lim_{n \rightarrow \infty} \max_{1 \leq i \leq n} P_{ni}(|\log(\frac{g_{ni}}{f_{ni}}(X_{ni}))| \geq \varepsilon) = 0 \quad \text{for every } \varepsilon > 0$$

if and only if

$$(21) \quad \sum_{i=1}^n d_H^2(P_{ni}, Q_{ni}) \rightarrow \frac{1}{4}\sigma^2,$$

$$(22) \quad \sum_{i=1}^n Q_{ni} \left( \frac{g_{ni}}{f_{ni}}(X_{ni}) \geq 1 + \varepsilon \right) \rightarrow 0 \quad \text{for every } \varepsilon > 0,$$

and

$$(23) \quad \sum_{i=1}^n P_{ni} \left( \frac{f_{ni}}{g_{ni}}(X_{ni}) \geq 1 + \varepsilon \right) \rightarrow 0 \quad \text{for every } \varepsilon > 0,$$

and, if and only if (21) holds and for every  $\varepsilon > 0$

$$\sum_{i=1}^n \int 1_{[|g_{ni} - f_{ni}| \geq \varepsilon g_{ni}]} (g_{ni}^{1/2} - f_{ni}^{1/2})^2 d\mu \rightarrow 0.$$

**Proof.** Oosterhoff and van Zwet (1979). □

**Proof of lemma 2.** As noted before, the following proof is from Hájek and Šidák (1967).

For any function  $h$  with second derivative  $h''$  we have

$$\begin{aligned} h(x) &= h(x_0) + (x - x_0)h'(x_0) \\ &\quad + \frac{1}{2}(x - x_0)^2 \int_0^1 2(1 - \lambda)h''(x_0 + \lambda(x - x_0)) d\lambda \end{aligned}$$

by integration by parts. Thus for  $h(x) = \log(1 + x)$ ,

$$\log(1 + x) = x - \frac{1}{2}x^2 \int_0^1 \frac{2(1 - \lambda)}{(1 + \lambda x)^2} d\lambda.$$

Thus, with  $T_{ni} \equiv 2\{(g_{ni}^{1/2}/f_{ni}^{1/2})(X_{ni}) - 1\}$ ,

$$\begin{aligned} (a) \quad \log \left( \frac{g_{ni}}{f_{ni}}(X_{ni}) \right) &= 2 \log \left( 1 + \frac{1}{2} T_{ni} \right) \\ &= T_{ni} - \frac{1}{4} T_{ni}^2 \int_0^1 \frac{2(1 - \lambda)}{(1 + \frac{1}{2} \lambda T_{ni})^2} d\lambda \end{aligned}$$

and

$$(b) \quad \Lambda_n = W_n - \frac{1}{4} \sum_{i=1}^n T_{ni}^2 \int_0^1 \frac{2(1 - \lambda)}{(1 + \frac{1}{2} \lambda T_{ni})^2} d\lambda.$$

Set  $T_{ni}^\delta \equiv T_{ni} 1_{\{|T_{ni}| \leq \delta\}}$  for  $\delta > 0$ . From the normal convergence criterion (see, e.g., Loève (1977, section 23.5, page 328)),  $L(W_n | P_n) \rightarrow N(-\sigma^2/4, \sigma^2)$  and the UAN condition (2) hold if and only if

$$(c) \quad \sum_{i=1}^n P_n(|T_{ni}| > \delta) \rightarrow 0,$$

$$(d) \quad \sum_{i=1}^n E(T_{ni}^\delta) \rightarrow -\frac{1}{4} \sigma^2,$$

and

$$(e) \quad \sum_{i=1}^n \text{Var}(T_{ni}^\delta) \rightarrow \sigma^2.$$

where all expectations and variances are under  $P_n$ . Note that



$$\int_0^1 2(1-\lambda) d\lambda = 1 \text{ and}$$

$$P_n \{ \max_{1 \leq i \leq n} |T_{ni}| > \delta \} \leq \sum_{i=1}^n P_n(|T_{ni}| > \delta) \rightarrow 0 \quad \text{by (c).}$$

Thus for any  $0 < \eta < 1$  there is an  $N = N(\eta)$  such that, for  $n \geq N$ ,  $P_n(S_n) > 1 - \eta$  with  $S_n \equiv [ \max_{1 \leq i \leq n} |T_{ni}| \leq \eta ]$ . It follows that, on  $S_n$

$$\sup_{\lambda} \max_{1 \leq i \leq n} | (1 + \frac{1}{2} \lambda T_{ni})^{-2} - 1 | \leq 3\eta$$

and hence

$$\max_{1 \leq i \leq n} | \int_0^1 \frac{2(1-\lambda)}{(1 + \frac{1}{2} \lambda T_{ni})^2} d\lambda - 1 | \leq 3\eta.$$

Also, since  $T_{ni} = T_{ni}^\eta$  for  $i = 1, \dots, n$ , on  $S_n$

$$| \sum_{i=1}^n T_{ni}^2 \int_0^1 \frac{2(1-\lambda)}{(1 + \frac{1}{2} \lambda T_{ni})^2} d\lambda - \sum_{i=1}^n T_{ni}^2 | \leq 3\eta \sum_{i=1}^n T_{ni}^2 = 3\eta \sum_{i=1}^n (T_{ni}^\eta)^2$$

so that

$$\begin{aligned} & | \sum_{i=1}^n T_{ni}^2 \int_0^1 \frac{2(1-\lambda)}{(1 + \frac{1}{2} \lambda T_{ni})^2} d\lambda / \sum_{i=1}^n (T_{ni}^\eta)^2 - 1 | \\ & \leq 3\eta \quad \text{on } S_n. \end{aligned}$$

Thus in order to prove the lemma it suffices to show that, under  $P_n$  as  $n \rightarrow \infty$  and subsequently  $\eta \downarrow 0$ ,

$$(f) \quad \sum_{i=1}^n (T_{ni}^\eta)^2 \rightarrow_P \sigma^2.$$

To prove (f) it suffices, by Chebychev's inequality, to show that

$$(g) \quad \sum_{i=1}^n E[(T_{ni}^\eta)^2] \rightarrow \sigma^2$$

and

$$(h) \quad \limsup_{\eta \rightarrow 0} \limsup_{n \rightarrow \infty} \sum_{i=1}^n \text{Var}[(T_{ni}^\eta)^2] = 0.$$

But by virtue of (e), (g) is equivalent to

$$(i) \quad \sum_{i=1}^n [ET_{ni}^\eta]^2 \rightarrow 0.$$

We first prove (i) and hence (g): If  $\eta > 2$ , then  $T_{ni}^\eta \leq T_{ni}$  since  $T_{ni} \geq -2$  a.s. by definition of  $T_{ni}$ . Therefore

$$ET_{ni}^\eta \leq ET_{ni} = 2E\left\{\frac{g_{ni}^{1/2}}{f_{ni}^{1/2}}(X_{ni})\right\} - 2 = -d_H^2(P_{ni}, Q_{ni}) \leq 0.$$

Thus for  $\eta > 2$

$$\sum_{i=1}^n (-ET_{ni}^\eta)^2 \leq \max_{1 \leq i \leq n} (-ET_{ni}^\eta) \sum_{i=1}^n (-ET_{ni}^\eta) \rightarrow 0$$

since

$$\sum_{i=1}^n (-ET_{ni}^\eta) \rightarrow \frac{1}{4} \sigma^2 \quad \text{by (d)}$$

and

$$\max_{1 \leq i \leq n} (-ET_{ni}^\eta) \rightarrow 0 \quad \text{by the UAN condition (2).}$$

Now note that if (i) holds for any  $\eta > 2$ , it holds for all  $\eta > 0$ : Since

$$\sum_{i=1}^n E[(T_{ni}^\eta)^2] \leq \sum_{i=1}^n E[(T_{ni}^\gamma)^2] \quad \text{for } \eta < \gamma$$

and, by (e) both  $\sum_{i=1}^n \text{Var}(T_{ni}^\eta) \rightarrow \sigma^2$  and  $\sum_{i=1}^n \text{Var}(T_{ni}^\gamma) \rightarrow \sigma^2$ , it follows that

$$\begin{aligned} \sum_{i=1}^n [ET_{ni}^\eta]^2 &= \sum_{i=1}^n \{E[(T_{ni}^\eta)^2] - \text{Var}(T_{ni}^\eta)\} \\ &\leq \sum_{i=1}^n \{E[(T_{ni}^\gamma)^2] - \text{Var}(T_{ni}^\gamma) + \text{Var}(T_{ni}^\gamma) - \text{Var}(T_{ni}^\eta)\} \\ &= \sum_{i=1}^n [ET_{ni}^\gamma]^2 + \sum_{i=1}^n \text{Var}(T_{ni}^\gamma) - \sum_{i=1}^n \text{Var}(T_{ni}^\eta) \\ &\rightarrow 0 + \sigma^2 - \sigma^2 = 0, \end{aligned}$$

completing the proof of (i) and hence (g).

To prove (h), note that

$$\sum_{i=1}^n \text{Var}[(T_{ni}^\eta)^2] \leq \sum_{i=1}^n E[(T_{ni}^\eta)^4] \leq \eta^2 \sum_{i=1}^n E[(T_{ni}^\eta)^2].$$

Then, by (g)

$$\limsup_{n \rightarrow \infty} \sum_{i=1}^n \text{Var}[(T_{ni}^\eta)^2] \leq \eta^2 \sigma^2,$$

and hence (h) holds.  $\square$

We conclude this section with a proof of the uniform LAN for regular models stated in proposition 2.1.2. Before proving the proposition we need one more lemma.

**Lemma 5.** Suppose that  $\mathbf{P} = \{P_\theta : \theta \in \Theta\}$  is regular and set

$$T_n \equiv 2 \left\{ \frac{s(\theta + n^{-1/2}t)}{s(\theta)}(X) - 1 \right\}.$$

Then:

$$(24) \quad E_\theta T_n = -\frac{1}{n} \frac{1}{4} t^T I(\theta) t + o(n^{-1}),$$

$$(25) \quad E_\theta T_n^2 = \frac{1}{n} t^T I(\theta) t + o(n^{-1}),$$

$$(26) \quad E_\theta (T_n - \frac{t^T}{\sqrt{n}} \dot{I}(X, \theta))^2 = o(n^{-1}),$$

$$(27) \quad P_\theta (|T_n| \geq \varepsilon) = o(n^{-1})$$

uniformly in  $|t| \leq M$  and  $\theta \in K \subset \Theta$ ,  $K$  compact.

**Proof.** By B of proposition A.5.3, for  $K$  compact

$$\begin{aligned} & \sup_{\theta \in K} \frac{1}{|h|^2} \|s(\theta + h) - s(\theta) - h^T \dot{s}(\theta)\|^2 \\ &= \sup_{\theta \in K} \frac{1}{|h|^2} \|h^T (\int_0^1 \dot{s}(\theta + uh) du - \dot{s}(\theta))\|^2 \\ &\leq \int_0^1 \sup_{\theta \in K} \|\dot{s}(\theta + uh) - \dot{s}(\theta)\|^2 du \end{aligned}$$

$$(a) \quad \rightarrow 0 \quad \text{as } h \rightarrow 0$$

by continuity of the map  $\theta \rightarrow \dot{s}(\theta)$ , which implies that  $\theta \rightarrow \dot{s}(\theta)$  is uniformly continuous on compact subsets  $K$  of  $\Theta$ . Therefore

$$\begin{aligned} & E_\theta (T_n - \frac{t^T}{\sqrt{n}} \dot{I}(X, \theta))^2 \\ &= \int \left( 2 \left( \frac{s(\theta + n^{-1/2}t)}{s(\theta)} - 1 \right) - \frac{t^T}{\sqrt{n}} 2 \frac{\dot{s}(\theta)}{s(\theta)} \right)^2 s^2(\theta) d\mu \\ &\leq 4 \|s(\theta + n^{-1/2}t) - s(\theta) - \frac{t^T}{\sqrt{n}} \dot{s}(\theta)\|^2 \end{aligned}$$

$$(b) \quad = o\left(\frac{|t|^2}{n}\right) = o(n^{-1})$$

uniformly in  $\theta \in K$  and  $|t| \leq M$  by (a), so (26) holds. Since

$$E_\theta \left( \frac{t^T}{\sqrt{n}} \dot{I}(X, \theta) \right)^2 = \frac{1}{n} t^T I(\theta) t,$$

(25) follows from (26). Furthermore, with  $A = [s(\theta) > 0]$ ,

$$E_\theta T_n^2 = 4 \int [s(\theta + \frac{t}{\sqrt{n}}) - s(\theta)]^2 1_A d\mu$$

$$= -4E_{\theta}T_n - 4 \int 1_A s^2(\theta + \frac{t}{\sqrt{n}}) d\mu$$

$$(c) \quad = -4E_{\theta}T_n - o(n^{-1})$$

holds uniformly in  $\theta \in K$  and  $|t| \leq M$  by the proof of proposition A.5.3.E and (a). Hence (24) follows from (25).

To prove (27), we first show that

$$(d) \quad \lim_{\lambda \rightarrow \infty} \sup_{\theta \in K} E_{\theta} |\dot{I}(X, \theta)|^2 1_{[|\dot{I}(X, \theta)| > \lambda]} = 0.$$

Suppose that (d) fails. Then there exist  $\theta_n \in K$  and  $\lambda_n \rightarrow \infty$  such that

$$(e) \quad E_{\theta_n} |\dot{I}(X, \theta_n)|^2 1_{[|\dot{I}(X, \theta_n)| > \lambda_n]} > \delta > 0.$$

Since  $K$  is compact, we may assume without loss of generality that  $\theta_n \rightarrow \theta \in K$ . Then, using the continuity of  $I(\theta)$  from proposition A.5.3.A, it follows that

$$E_{\theta_n} (t^T \dot{I}(X, \theta_n))^2 = t^T I(\theta_n) t \rightarrow t^T I(\theta) t = E_{\theta} (t^T \dot{I}(X, \theta))^2.$$

Therefore

$$E_{\theta_n} \dot{I}_i(X, \theta_n)^2 \rightarrow E_{\theta} \dot{I}_i(X, \theta)^2 \quad \text{for } i = 1, \dots, k,$$

and adding these up yields

$$E_{\theta_n} |\dot{I}(X, \theta_n)|^2 \rightarrow E_{\theta} |\dot{I}(X, \theta)|^2.$$

Thus by lemma A.7.2.B the sequence of random variables  $\{|\dot{I}(X, \theta_n)|\}$  is uniformly integrable; hence

$$\lim_{n \rightarrow \infty} E_{\theta_n} |\dot{I}(X, \theta_n)|^2 1_{[|\dot{I}(X, \theta_n)| > \lambda_n]} = 0$$

which contradicts (e). Therefore (d) holds.

Now we use (d) to establish (27):

$$\begin{aligned} P_{\theta}(|T_n| \geq \varepsilon) &\leq P_{\theta}(|T_n - \frac{t^T}{\sqrt{n}} \dot{I}(X, \theta)| \geq \frac{1}{2}\varepsilon) \\ &\quad + P_{\theta}(|\frac{t^T}{\sqrt{n}} \dot{I}(X, \theta)| \geq \frac{1}{2}\varepsilon) \\ &\leq \frac{4}{\varepsilon^2} E_{\theta} |T_n - \frac{t^T}{\sqrt{n}} \dot{I}(X, \theta)|^2 \\ &\quad + \frac{4}{\varepsilon^2} E_{\theta} |\frac{t^T}{\sqrt{n}} \dot{I}(X, \theta)|^2 1_{[|(t^T/\sqrt{n})\dot{I}(X, \theta)| \geq \varepsilon/2]} \\ &\leq o(n^{-1}) + \frac{4|t|^2}{n\varepsilon^2} E_{\theta} |\dot{I}(X, \theta)|^2 1_{[|\dot{I}(X, \theta)| > \varepsilon\sqrt{n}/2|t|]} \\ &= o(n^{-1}) \end{aligned}$$

uniformly in  $\theta \in K$  and  $|t| \leq M$  by (26) and (d). □

**Proof of proposition 2.1.2.** Let  $T_{ni}$  denote the  $n$  i.i.d. copies of the random variable  $T_n$  defined in lemma 5 corresponding to  $X_1, \dots, X_n$ , and for  $\varepsilon < 1$  define

$$A_n \equiv \{ \max_{1 \leq i \leq n} |T_{ni}| < \varepsilon \}.$$

Then, on the event  $A_n$  it follows by expansion of  $\log(1+x)$  that

$$\begin{aligned} L_n\left(\theta + \frac{t}{\sqrt{n}}\right) - L_n(\theta) &= 2 \sum_{i=1}^n \log\left(1 + \frac{1}{2} T_{ni}\right) \\ &= 2 \sum_{i=1}^n \left\{ \frac{1}{2} T_{ni} - \frac{1}{8} T_{ni}^2 + \alpha_{ni} \frac{1}{12} |T_{ni}|^3 \right\} \\ &= \sum_{i=1}^n T_{ni} - \frac{1}{4} \sum_{i=1}^n T_{ni}^2 + \frac{1}{6} \sum_{i=1}^n \alpha_{ni} |T_{ni}|^3 \\ &= t^T S_n(\theta) - \frac{1}{2} t^T I(\theta) t \\ &\quad + \sum_{i=1}^n T_{ni} - \left( t^T S_n(\theta) - \frac{1}{4} t^T I(\theta) t \right) \\ &\quad - \frac{1}{4} \left( \sum_{i=1}^n T_{ni}^2 - t^T I(\theta) t \right) \\ &\quad + \frac{1}{6} \sum_{i=1}^n \alpha_{ni} |T_{ni}|^3, \end{aligned}$$

where  $|\alpha_{ni}| \leq 1$ . To prove (2.1.13) it therefore suffices to show that

(a)  $R_{n1}(\theta, t) \equiv \sum_{i=1}^n T_{ni} - \{S_n(\theta) - \frac{1}{4} t^T I(\theta) t\} \rightarrow 0,$

(b)  $R_{n2}(\theta, t) \equiv \sum_{i=1}^n T_{ni}^2 - t^T I(\theta) t \rightarrow 0,$

(c)  $R_{n3}(\theta, t) \equiv \sum_{i=1}^n \alpha_{ni} |T_{ni}|^3 \rightarrow 0,$

in  $P_\theta$ -probability and

(d)  $P_\theta(A_n^c) \rightarrow 0$

uniformly in  $\theta \in K \subset \Theta$  for  $K$  compact and  $|t| \leq M$ .

We first prove (d):

$$\begin{aligned} P_\theta(A_n^c) &\leq \sum_{i=1}^n P_\theta(|T_{ni}| \geq \varepsilon) \\ &= n P_\theta(|T_n| \geq \varepsilon) = o(1) \end{aligned}$$

uniformly in  $\theta \in K$ ,  $|t| \leq M$  by (27) of lemma 5, so (d) holds. To prove (b), write

$$R_{n2}(\theta, t) = \sum_{i=1}^n T_{ni}^2 - \frac{1}{n} \sum_{i=1}^n (t^T \dot{\mathbf{i}}(X_i, \theta))^2 \\ + \frac{1}{n} \sum_{i=1}^n (t^T \dot{\mathbf{i}}(X_i, \theta))^2 - t^T I(\theta) t,$$

where

$$E_{\theta} \left| \sum_{i=1}^n T_{ni}^2 - \frac{1}{n} \sum_{i=1}^n (t^T \dot{\mathbf{i}}(X_i, \theta))^2 \right| \\ \leq n E_{\theta} \left| T_n^2 - \left( \frac{t^T}{\sqrt{n}} \dot{\mathbf{i}}(X, \theta) \right)^2 \right| = o(1)$$

uniformly in  $\theta \in K$ ,  $|t| \leq M$  as a consequence of (26) of lemma 5, and where

$$\frac{1}{n} \sum_{i=1}^n (t^T \dot{\mathbf{i}}(X_i, \theta))^2 = \frac{1}{n} \sum_{i=1}^n t^T \dot{\mathbf{i}}(X_i, \theta) \dot{\mathbf{i}}^T(X_i, \theta) t \rightarrow t^T I(\theta) t$$

in  $P_{\theta}$ -probability uniformly in  $\theta \in K$ ,  $|t| \leq M$  by (d) of the proof of lemma 5, and theorem A.7.3. Thus (b) holds. Now (c) follows since

$$P_{\theta} \left( \sum_{i=1}^n |T_{ni}|^3 > \varepsilon \right) \leq P_{\theta} \left( \max_{1 \leq i \leq n} |T_{ni}| \sum_{i=1}^n T_{ni}^2 > \varepsilon \right) \\ = P_{\theta} \left( \max_{1 \leq i \leq n} |T_{ni}| \sum_{i=1}^n T_{ni}^2 > \varepsilon, \sum_{i=1}^n T_{ni}^2 \leq 1 + t^T I(\theta) t \right) \\ + P_{\theta} \left( \max_{1 \leq i \leq n} |T_{ni}| \sum_{i=1}^n T_{ni}^2 > \varepsilon, \sum_{i=1}^n T_{ni}^2 > 1 + t^T I(\theta) t \right) \\ \leq P_{\theta} \left( \max_{1 \leq i \leq n} |T_{ni}| > \frac{\varepsilon}{1 + t^T I(\theta) t} \right) \\ + P_{\theta} \left( \sum_{i=1}^n T_{ni}^2 > 1 + t^T I(\theta) t \right) \\ \rightarrow 0$$

uniformly in  $\theta \in K$ ,  $|t| \leq M$ , by (d) and (b). Finally to prove (a), note that

$$E_{\theta} [R_{n1}(\theta, t)]^2 = \text{Var}_{\theta}(R_{n1}(\theta, t)) + [E_{\theta} R_{n1}(\theta, t)]^2$$

$$\begin{aligned}
 &= n \operatorname{Var}_{\theta}[T_n - \frac{t^T}{\sqrt{n}} \dot{\mathbf{i}}(X, \theta)] + [nE_{\theta} T_n + \frac{1}{4} t^T I(\theta) t]^2 \\
 &= n \left\{ E_{\theta} [T_n - \frac{t^T}{\sqrt{n}} \dot{\mathbf{i}}(X, \theta)]^2 - (E_{\theta} T_n)^2 \right\} + o(1) \\
 &= o(1)
 \end{aligned}$$

uniformly in  $\theta \in K, |t| \leq M$ , by (24) and (26) of lemma 5.

To prove (2.1.14) of proposition 2.1.2, let  $f: R^k \rightarrow R$  be bounded and continuous. We want to show that for  $K$  compact  $\subset \Theta$ ,

(e) 
$$\sup_{\theta \in K} |E_{\theta} f(S_n(\theta)) - E_{\theta} f(Z)| \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

where  $Z \sim N_k(0, I(\theta))$ . Suppose that (e) fails. Then there exist  $\theta_n \in K$  such that

(f) 
$$\limsup_{n \rightarrow \infty} |E_{\theta_n} f(S_n(\theta_n)) - E_{\theta_n} f(Z)| > 0.$$

But since  $K$  is compact, we may assume without loss of generality that  $\theta_n \rightarrow \theta \in K$ .

Now let  $t \in R^k$  and apply the bound of theorem A.7.4 to the mean 0 variance 1 random variables  $Y_i \equiv t^T \dot{\mathbf{i}}(X_i, \theta_n) / \sigma_n$  where  $\sigma_n^2 \equiv t^T I(\theta_n) t$  to get, by (d) of the proof of lemma 5,

$$L(t^T S_n(\theta_n) / \sigma_n | P_{\theta_n}) \rightarrow N(0, 1)$$

and hence, by continuity of  $I(\theta)$  from proposition A.5.3.A,

$$L(t^T S_n(\theta_n) | P_{\theta_n}) \rightarrow N(0, t^T I(\theta) t),$$

which, in view of the Cramér-Wold device, contradicts (f).

We conclude the proof of proposition 2.1.2 by noting that (2.1.15) follows from (2.1.13), contiguity (proposition 2.1.3), and continuity and symmetry of  $I(\theta)$ , since under  $P_{\theta}$

$$\begin{aligned}
 0 &= L_n\left(\theta + \frac{t}{\sqrt{n}}\right) - L_n(\theta) + L_n\left(\theta + \frac{t+h}{\sqrt{n}}\right) - L_n\left(\theta + \frac{t}{\sqrt{n}}\right) \\
 &\quad - \left\{ L_n\left(\theta + \frac{t+h}{\sqrt{n}}\right) - L_n(\theta) \right\}
 \end{aligned}$$

(g) 
$$\begin{aligned}
 &= t^T S_n(\theta) - \frac{1}{2} t^T I(\theta) t + h^T S_n\left(\theta + \frac{t}{\sqrt{n}}\right) - \frac{1}{2} h^T I\left(\theta + \frac{t}{\sqrt{n}}\right) h \\
 &\quad - (t+h)^T S_n(\theta) + \frac{1}{2} (t+h)^T I(\theta) (t+h) + o_P(1)
 \end{aligned}$$

$$= h^T \left\{ S_n\left(\theta + \frac{t}{\sqrt{n}}\right) - S_n(\theta) + I(\theta) t \right\} + o_P(1).$$

□

### A.10 THE MASTER THEOREM FOR ASYMPTOTIC GENERALIZED $M$ -ESTIMATES

In this section we collect the proofs of several results presented in sections 7.2 through 7.4. For ease of reference we repeat some of the definitions, conditions, and statements of results.

Suppose that  $\mathbf{M}_0 \supset \mathbf{P}$  and all distributions with finite support as in section 7.1 and 7.2. Suppose that  $W_n, W$  map  $R^m \times \mathbf{M}_0 \rightarrow R^m$  and that

(i)  $W_n(v, P) = W(v, P) + o(1)$  for all  $P \in \mathbf{M}_0$ , all  $v$ ,

(ii)  $W(v(P), P) = 0$  for all  $P \in \mathbf{P}$ .

Let  $\mathbf{W}_n(v) = W_n(v, P_n)$ . If

$$\mathbf{W}_n(\hat{v}_n) = 0,$$

then  $\hat{v}_n$  is a *generalized  $M$ -estimate* (or *GM-estimate* for short) of  $v(P)$ . If

$$(1) \quad \mathbf{W}_n(\hat{v}_n) = o_p(n^{-1/2}) \quad \text{for all } P \in \mathbf{P},$$

then  $\hat{v}_n$  is an *asymptotic generalized  $M$ -estimate* (or *AGM-estimate* for short) of  $v(P)$ . Define

$$\mathbf{V}_n(v) \equiv \sqrt{n} \{ \mathbf{W}_n(v) - W(v, P) \}$$

for  $v \in \mathbf{N}$  open  $\subset R^m$ . Let  $\mathbf{Q}$  be a model with  $\mathbf{Q} \supset \mathbf{P}$ . Here are the key assumptions:

(GM0) There exists  $v : \mathbf{Q} \rightarrow R^m$  such that  $v(P)$  satisfies  $W(v(P), P) = 0$  for all  $P \in \mathbf{Q}$ .

(GM1) For any  $\varepsilon_n \downarrow 0$  we have

$$\sup \left\{ \frac{|\mathbf{V}_n(v) - \mathbf{V}_n(v(P))|}{1 + \sqrt{n}|v - v(P)|} : |v - v(P)| \leq \varepsilon_n \right\} = o_p(1).$$

(GM2) There is a function  $\psi : \mathbf{X} \times \mathbf{Q} \rightarrow R^m$  with  $\int \psi(x, P) dP(x) = 0$  and  $|\psi(\cdot, P)| \in L_2(P)$  such that

$$\mathbf{W}_n(v(P)) = n^{-1} \sum_{i=1}^n \psi(X_i, P) + o_p(n^{-1/2}).$$

(GM3)  $W(\cdot, P) = (W_1(\cdot, P), \dots, W_m(\cdot, P))^T$  is differentiable with derivative  $\dot{W}(v, P) \equiv [(\partial/\partial v_j) W_i(v, P)]_{m \times m}$  and  $\dot{W}(P) \equiv \dot{W}(v(P), P)$  is nonsingular.

We introduce the model  $\mathbf{Q}$  to emphasize that even though  $W_n$  may be motivated by features of  $\mathbf{P}$ , the *AGM-estimates* corresponding to  $W_n$  can be thought of as estimates of parameters on a larger model.

**Theorem 1. (AGM-estimates).** Suppose  $P \in \mathbf{Q}$ . Let  $\hat{v}_n$  be an *AGM-estimate* of  $v(P)$  on  $\mathbf{Q}$ . If  $\hat{v}_n$  is consistent and if (GM0)–(GM3) hold, then

$$\sqrt{n}(\hat{v}_n - v(P)) = -\frac{1}{\sqrt{n}} \sum_{i=1}^n \dot{W}^{-1}(P) \psi(X_i, P) + o_p(1)$$



$$\rightarrow_d N(0, \Sigma(\hat{v}, P)) \quad \text{as } n \rightarrow \infty$$

where

$$\Sigma(\hat{v}, P) = \dot{W}^{-1}(P) E[\psi(X, P) \psi^T(X, P)] [\dot{W}^{-1}(P)]^T.$$

If (GM1) holds for  $\varepsilon_n = O(n^{-1/2})$  only and if  $\hat{v}_n$  is  $\sqrt{n}$ -consistent, this asymptotic linearity remains valid.

**Proof.** By (GM0) and (1),

$$\begin{aligned} \text{(a)} \quad \mathbf{V}_n(v(P)) + \sqrt{n} \{W(\hat{v}_n, P) - W(v(P), P)\} \\ = \mathbf{V}_n(v(P)) - \mathbf{V}_n(\hat{v}_n) + o_p(1). \end{aligned}$$

Dividing both sides of (a) by  $1 + \sqrt{n} |\hat{v}_n - v(P)|$ , and applying (GM1) yields, since  $\hat{v}_n$  is consistent,

$$\frac{\mathbf{V}_n(v(P)) + \sqrt{n} \{W(\hat{v}_n, P) - W(v(P), P)\}}{1 + \sqrt{n} |\hat{v}_n - v(P)|} = o_p(1)$$

and hence, by (GM2),

$$\text{(b)} \quad \frac{n^{-1/2} \sum_{i=1}^n \psi(X_i, P) + \sqrt{n} \{W(\hat{v}_n, P) - W(v(P), P)\}}{1 + \sqrt{n} |\hat{v}_n - v(P)|} = o_p(1).$$

We now complete the proof by arguing as in Huber (1967). Fix  $\varepsilon > 0$  small, and let  $M^2 \equiv 2 \int |\psi(x, P)|^2 dP(x) / \varepsilon < \infty$  by (GM2). Then it follows from (GM2) and (b) that, with probability at least  $1 - \varepsilon$  for  $n$  sufficiently large, the following two inequalities hold:

$$\text{(c)} \quad |n^{-1/2} \sum_{i=1}^n \psi(X_i, P)| < M$$

and

$$\begin{aligned} \text{(d)} \quad |n^{-1/2} \sum_{i=1}^n \psi(X_i, P) + \sqrt{n} \{W(\hat{v}_n, P) - W(v(P), P)\}| \\ \leq \varepsilon (1 + \sqrt{n} |\hat{v}_n - v(P)|). \end{aligned}$$

Since  $\hat{v}_n$  is consistent, (GM3) implies that with probability converging to 1 we have

$$\text{(e)} \quad |W(\hat{v}_n, P) - W(v(P), P)| \geq \alpha |\hat{v}_n - v(P)|$$

for some  $\alpha = \alpha(P) > 0$ . Combining (c)–(e) yields

$$\begin{aligned} \varepsilon (1 + \sqrt{n} |\hat{v}_n - v(P)|) &\geq |\sqrt{n} \{W(\hat{v}_n, P) - W(v(P), P)\}| - M \\ &\geq \alpha \sqrt{n} |\hat{v}_n - v(P)| - M \end{aligned}$$

and hence, for  $\epsilon < \alpha$  and with probability of at least  $1 - \epsilon$ ,

$$(f) \quad \sqrt{n} |\hat{v}_n - v(P)| \leq \frac{M + \epsilon}{\alpha - \epsilon},$$

so that  $\hat{v}_n$  is  $\sqrt{n}$ -consistent. Finally, since  $\epsilon M = O(\epsilon^{1/2})$ , (d) and (f) imply that

$$(g) \quad \sqrt{n} \{W(\hat{v}_n, P) - W(v(P), P)\} = -n^{-1/2} \sum_{i=1}^n \psi(X_i, P) + o_p(1),$$

and the first part of the theorem follows from (g) and (GM3). The modifications of this proof needed for the second part are simple. □

We now specialize theorem 1 to the case of M-estimates. Suppose that  $\psi : X \times R^m \rightarrow R^m$  is such that  $\int |\psi(x, v)|^2 dP(x) < \infty$  and  $\int \psi(x, v(P)) dP(x) = 0$ , and define

$$W_n(v) \equiv \frac{1}{n} \sum_{i=1}^n \psi(X_i, v) = \int \psi(x, v) dP_n(x),$$

and

$$W(v, P) \equiv \int \psi(x, v) dP(x) = E_P \psi(X, v).$$

Note that (GM0) is satisfied by the assumptions on  $\psi$ :

$$(M0) \quad W(v(P), P) = 0.$$

The key hypothesis (GM1) of theorem 1 becomes:

(M1) For any  $\epsilon_n \downarrow 0$  we have

$$\sup_{|v - v(P)| \leq \epsilon_n} \frac{|n^{-1/2} \sum_{i=1}^n \{\psi(X_i, v) - \psi(X_i, v(P)) - E_P[\psi(X, v) - \psi(X, v(P))]\}|}{1 + \sqrt{n} |v - v(P)|} = o_p(1).$$

Note that

$$(M2) \quad \sqrt{n} W_n(v(P)) = n^{-1/2} \sum_{i=1}^n \psi(X_i, v(P)),$$

so that (GM2) holds trivially. Finally, (GM3) becomes

(M3)  $W(v, P) = \int \psi(x, v) dP(x)$  is differentiable and the derivative at  $v(P)$ ,  $\dot{W}(P)$ , is nonsingular.

These identifications prove the following corollary of theorem 1.

**Corollary 1. (M-estimates).** Suppose that (M1) and (M3) hold,  $\hat{v}_n$  is consistent, and is an asymptotic M-estimate: i.e.  $\sqrt{n} W_n(\hat{v}_n) = o_p(1)$ . Then  $\hat{v}_n$  is asymptotically linear with influence function  $-\dot{W}^{-1}(P) \psi(\cdot, P)$ .

Suppose that  $\dot{W}_n \equiv [(\partial/\partial v_j) W_{ni}(v)]$  exists where  $W_n \equiv (W_{n1}, \dots, W_{nm})^T$  and that:

(U) For some sequence  $\{\varepsilon_n\}$  with  $\varepsilon_n \downarrow 0$ ,  $\varepsilon_n n^{1/2} \rightarrow \infty$ ,

$$\sup\{|\dot{W}_n(v) - \dot{W}(P)| : |v - v(P)| < \varepsilon_n\} = o_p(1).$$

Let  $T_n^{(0)}$  be a preliminary estimate and define

$$(2) \quad T_n^{(j+1)} = T_n^{(j)} - \dot{W}_n^{-1}(T_n^{(j)}) W_n(T_n^{(j)}), \quad j = 0, 1, \dots$$

**Theorem 2. (Iteration).** Suppose (GM0), (GM2), (GM3), and (U) hold. Then:

- A. With probability converging to 1,  $W_n(v)$  has a unique root  $v_n^{(\infty)}$  in  $\{v : |v - v(P)| \leq \varepsilon_n\}$ .  $v_n^{(\infty)}$  is  $\sqrt{n}$ -consistent and asymptotically linear with influence function  $-\dot{W}^{-1}(P)\psi(\cdot, v(P))$ .
- B. If there exists an estimator  $T_n^{(0)}$  satisfying  $P(|T_n^{(0)} - v(P)| < \varepsilon_n) \rightarrow 1$  for  $\{\varepsilon_n\}$  as in (U) and if the Newton-Raphson iteration (2) starts at  $T_n^{(0)}$ , then

$$P(T_n^{(\infty)} \text{ exists and equals } v_n^{(\infty)}) \rightarrow 1, \quad \text{as } n \rightarrow \infty.$$

**Proof.** We begin by establishing two preliminary results. First

$$(a) \quad \begin{aligned} a_n &\equiv \sup\{\|I - \dot{W}^{-1}(P)\dot{W}_n(v)\| : |v - v(P)| \leq \varepsilon_n\} \\ &\leq \sup\{\|\dot{W}^{-1}(P)\| \|\dot{W}(P) - \dot{W}_n(v)\| : |v - v(P)| \leq \varepsilon_n\} \\ &= o_p(1) \quad \text{by (U) and (GM3)}. \end{aligned}$$

Note that (U) and (GM3) imply also that

$$P\{\dot{W}_n^{-1}(v) \text{ exists and } \|\dot{W}_n^{-1}(v)\| \leq 2\|\dot{W}^{-1}(P)\|\}$$

$$\text{for all } v \text{ with } |v - v(P)| \leq \varepsilon_n \} \rightarrow 1.$$

Hence, just as in (a),

$$(b) \quad \begin{aligned} b_n &\equiv \sup\{\|I - \dot{W}_n^{-1}(v)\dot{W}_n(v')\| : |v - v(P)| \leq \varepsilon_n, \\ &\quad |v' - v(P)| \leq \varepsilon_n\} \\ &= o_p(1). \end{aligned}$$

**Proof of A:** Let  $h_n(v) = v - \dot{W}^{-1}(P)W_n(v)$ , and let  $v_n^{(0)} = v(P)$ ,  $v_n^{(j+1)} = h_n(v_n^{(j)})$ ,  $j \geq 1$ .

Suppose that

$$|v_n^{(j)} - v(P)| \leq (1 - a_n)^{-1} |\dot{W}^{-1}(P)W_n(v(P))| \equiv (1 - a_n)^{-1} A_n,$$

where  $A_n = O_p(n^{-1/2})$  by (GM2). Since

$$(c) \quad \begin{aligned} v_n^{(j+1)} - v(P) &= v_n^{(j)} - v(P) - \dot{W}^{-1}(P)W_n(v_n^{(j)}) \\ &= v_n^{(j)} - v(P) - \dot{W}^{-1}(P)(W_n(v_n^{(j)}) - W_n(v(P))) \end{aligned}$$

$$\begin{aligned}
& - \dot{W}^{-1}(P) \mathbf{W}_n(v(P)) \\
& = (I - \dot{W}^{-1}(P) \dot{\mathbf{W}}_n(\tilde{v}_n^{(j)}))(v_n^{(j)} - v(P)) \\
& \quad - \dot{W}^{-1}(P) \mathbf{W}_n(v(P)),
\end{aligned}$$

for some intermediate point  $\tilde{v}_n^{(j)}$ , it follows from (a) and (U) that

$$|v_n^{(j+1)} - v(P)| \leq a_n \frac{A_n}{1 - a_n} + A_n = \frac{A_n}{1 - a_n}.$$

Since  $v_n^{(0)} - v(P) = 0$ , induction yields

$$(d) \quad \sup_{j \geq 1} |v_n^{(j)} - v(P)| \leq (1 - a_n)^{-1} A_n = O_p(n^{-1/2}).$$

Now consider

$$\begin{aligned}
(e) \quad & |h_n(v_n^{(j+1)}) - h_n(v_n^{(j)})| \\
& = |v_n^{(j+1)} - v_n^{(j)} - \dot{W}^{-1}(P)(\mathbf{W}_n(v_n^{(j+1)}) - \mathbf{W}_n(v_n^{(j)}))| \\
& = |(I - \dot{W}^{-1}(P) \dot{\mathbf{W}}_n(\tilde{v}_n^{(j)}))(v_n^{(j+1)} - v_n^{(j)})| \\
& \leq a_n |v_n^{(j+1)} - v_n^{(j)}|,
\end{aligned}$$

where  $\tilde{v}_n^{(j)}$  is another intermediate point, which, by (d), satisfies  $|\tilde{v}_n^{(j)} - v(P)| = O_p(n^{-1/2})$ . It follows from (a) and (e) that  $h_n$  is a contraction on the set  $\{v : |v - v(P)| \leq (1 - a_n)^{-1} A_n\}$ , and hence that  $v_n^{(j)} \rightarrow v_n^{(\infty)}$  as  $j \rightarrow \infty$  on a set with probability arbitrarily close to 1 for large  $n$ . Further  $v_n^{(\infty)}$  satisfies

$$(f) \quad v_n^{(\infty)} - v(P) = O_p(n^{-1/2}) = o_p(\varepsilon_n).$$

Moreover,  $v_n^{(\infty)}$  is a fixed point of  $h_n$  and satisfies  $\mathbf{W}_n(v_n^{(\infty)}) = 0$ . The linearity of  $v_n^{(\infty)}$  follows from theorem 1.

**Proof of B:** Consider the sequence (2). Assume that  $|T_n^{(j)} - v(P)| \leq \varepsilon_n$ . Then

$$\begin{aligned}
(g) \quad & |T_n^{(j+1)} - v_n^{(\infty)}| = |T_n^{(j)} - v_n^{(\infty)} - \dot{\mathbf{W}}_n^-(T_n^{(j)})(\mathbf{W}_n(T_n^{(j)}) - \mathbf{W}_n(v_n^{(\infty)}))| \\
& = |(I - \dot{\mathbf{W}}_n^-(T_n^{(j)}) \dot{\mathbf{W}}_n(\tilde{T}_n^{(j)}))(T_n^{(j)} - v_n^{(\infty)})| \\
& \leq b_n |T_n^{(j)} - v_n^{(\infty)}|,
\end{aligned}$$

where  $\tilde{T}_n^{(j)}$  is another intermediate point. Hence

$$\begin{aligned}
|T_n^{(j+1)} - v(P)| & \leq b_n |T_n^{(j)} - v(P)| + |v_n^{(\infty)} - v(P)| \\
& \leq b_n \varepsilon_n + o_p(\varepsilon_n) \\
& \leq \varepsilon_n
\end{aligned}$$

on the event  $\{|v_n^{(\infty)} - v(P)| \leq \varepsilon_n/2\} \cap \{b_n < \frac{1}{2}\}$ . This event has probability converging to 1 by (b) and (f). We obtain from (g) that, for  $n$  sufficiently large,

$$|T_n^{(j)} - v_n^{(\infty)}| \leq b_n^j |\tilde{v}_n - v_n^{(\infty)}| \rightarrow 0 \quad \text{as } j \rightarrow \infty$$

with probability converging to 1. □

By a theorem of Brown (1985) to be discussed below, convexity assumptions can be coupled with much weaker conditions than (GM0)–(GM3) to yield the conclusion of theorem 1.

Let  $N$  be an open convex subset of  $R^m$ . Let  $W$  be the class of all functions  $W : N \rightarrow R^m$  such that for all  $u \in R^m, t \in N$ , the maps  $\lambda \rightarrow u^T W(t + \lambda u)$  from  $\{\lambda \in R : t + \lambda u \in N\}$  to  $R$  are monotone nondecreasing. Let  $W_0 \subset W$  be the subclass of functions which have a unique root in  $N$ . Assume that (GM0) and (GM3) hold, and let  $v_0 \equiv v(P)$ . Suppose that  $M_0 \supset Q \cup \{\text{all realizations of the empirical measures } P_n, n \geq 1\}$ . Here  $Q \supset P$ . Fix  $P \in P$  and  $W : N \times Q \rightarrow R^m$ , and consider the assumptions:

- (C1)  $P(W(\cdot, P_n) \in W_0) \rightarrow 1$  as  $n \rightarrow \infty$ .
- (C2) For each fixed  $\tau \in R^m$ ,  
 $\sqrt{n}(W(v_0 + n^{-1/2}\tau, P_n) - W(v_0, P_n)) = \dot{W}(P)\tau + o_p(1)$ .
- (C3)  $W(v_0, P_n) = \int \psi(x, P) dP_n + o_p(n^{-1/2})$  where  $|\psi| \in L_2(P)$ ,  
 $\int \psi(x, P) dP(x) = 0$ .

We will also use a strengthening of (C1):

$$(C1') \quad W(\cdot, Q) \in W_0 \text{ for all } Q \in M_0.$$

The following useful result is due to Brown (1985) and Ritov (1987).

**Theorem 3. (Convexity).** Suppose that  $X_1, \dots, X_n$  are i.i.d.  $P \in Q$ . Suppose that (GM0), (GM3), (C1'), (C2), and (C3) hold. Then  $\hat{v}_n$  corresponding to  $W_n(v) \equiv W(v, P_n) = 0$  is uniquely defined and asymptotically linear with influence function  $-\dot{W}^{-1}(P)\psi(\cdot, P)$ . If the hypothesis (C1') is replaced by (C1), then the AGM-estimate  $\hat{v}_n$  exists and asymptotic linearity continues to hold.

Theorem 3 is a result similar to theorem 7.3.2, but under very weak hypotheses. It can also be viewed as an alternative to theorem 1 with weak smoothness hypotheses counterbalanced by the strong assumptions (C1) or (C1'). The proof hinges on theorem A.7.8.

**Proof.** Let  $U_n(\tau) = \sqrt{n}W(v_0 + n^{-1/2}\tau, P_n)$ . Since  $U_n(\cdot) - U_n(0) \in W$  and  $\tau \rightarrow \dot{W}(P)\tau$  is continuous, (C2) and theorem A.7.8 imply that

$$(a) \quad \sup_{|\tau| \leq M} |U_n(\tau) - U_n(0) - \dot{W}(P)\tau| = o_p(1).$$

We proceed to check the conditions of the second part of theorem 1 for  $W_n(v, P_n) = W(v, P_n)$ . Note that

$$(b) \quad V_n(v) - V_n(v_0) = U_n(n^{1/2}(v - v_0)) - U_n(0) \\ - n^{1/2}[W(v, P) - W(v_0, P)].$$

So, by (a) and differentiability of  $W$ ,

$$\begin{aligned} & \sup\{|V_n(v) - V_n(v_0)| : n^{1/2}|v - v_0| \leq M\} \\ &= o_p(1) \\ & \quad + \sup\{|n^{1/2}(W(v_0 + n^{-1/2}\tau, P) - W(v_0, P)) - \dot{W}(P)\tau| : |\tau| \leq M\} \\ &= o_p(1) \quad \text{by (GM3)}. \end{aligned}$$

Hence condition (GM1) of theorem 1 follows, at least for  $\epsilon_n = O(n^{-1/2})$ . Suppose (C1') holds. We now need only to verify the  $\sqrt{n}$ -consistency of  $\hat{v}_n$  which is uniquely defined by  $W(\hat{v}_n, P_n) = 0$ . It is enough to show that for all  $\epsilon > 0$  there exists an  $M = M(\epsilon)$ , so that, for  $n$  sufficiently large

$$(c) \quad P(A_n) > 1 - \epsilon,$$

where

$$A_n = [\inf\{\tau^T U_n(\tau) : |\tau| = M\} > 0].$$

This is enough since  $U_n(\tau) \in \mathbf{W}$  implies  $\tau^T U_n(\lambda\tau)$  is increasing in  $\lambda$  and hence that, on  $A_n$ ,  $\tau^T U_n(\tau) > 0$  for all  $|\tau| \geq M$ , and hence  $\sqrt{n}|\hat{v}_n - v_0| < M$ . (The argument is spelled out in the proof of theorem 7.5.3.)

We obtain from (a) that

$$(d) \quad \sup\{|\tau^T U_n(\tau) - \tau^T U_n(0) - \tau^T \dot{W}(P)\tau| : |\tau| = M\} = o_p(1).$$

Let  $A(P) = \frac{1}{2}(\dot{W}(P) + \dot{W}^T(P))$ . Then

$$(e) \quad \tau^T A(P)\tau = \tau^T \dot{W}(P)\tau, \quad \text{for all } \tau.$$

$W(\cdot, P) \in \mathbf{W}_0$  implies that  $A(P)$  is symmetric, positive definite. In particular its minimum eigenvalue  $\lambda_1(P)$  is strictly positive.

Fix  $\epsilon > 0$ . By (C3), for  $M$  sufficiently large and all  $n$ :

$$(f) \quad P(|U_n(0)| < \frac{1}{2}\lambda_1(P)M) \geq 1 - \frac{1}{2}\epsilon.$$

We obtain from (d)–(f) that with probability at least  $1 - \epsilon$  for all  $n$  and  $M$  sufficiently large

$$\inf\{\tau^T U_n(\tau) : |\tau| = M\} > \frac{1}{3}\lambda_1(P)M^2.$$

Hence (c) holds, and  $\hat{v}_n$  is  $\sqrt{n}$ -consistent. If (C1) holds then (c) still holds and implies that the minimum of  $|W(\cdot, P_n)|$  is assumed on  $\{v : \sqrt{n}|v - v_0| < M\}$ . To complete the proof we need only to show therefore, that the minimizer is an AGM-estimate. That is, it suffices to show that  $\inf_{\tau} |U_n(\tau)| = o_p(1)$ . Let  $\tau_n = -\dot{W}^{-1}(P)U_n(0)$ . We obtain from (a) that

$$\inf_{\tau} |U_n(\tau)| \leq |U_n(\tau_n)| = o_p(1),$$

and the theorem follows.  $\square$

# List of Symbols

## LATIN ALPHABET

Symbol	Page	Meaning or explanation
$a(\cdot, \theta)$	88	transformations parametrized by $\theta$
$\mathbf{A}$	88	group of measurable transformations
$\mathcal{A}$	46, 475	a $\sigma$ -field of subsets of $\mathbf{X}$ or $\Omega$
$A$	81, 434	compensator of the counting process $N$
$B$	197	Brownian motion
$B_0$	192	Brownian bridge process on $[0, 1]$
$\mathcal{B}$	1	Borel $\sigma$ -field
$\mathcal{B}_0$	104	a sub- $\sigma$ -field of the Borel $\sigma$ -field $\mathcal{B}$
$\mathcal{B}_T$	431	invariant $\sigma$ -field
$\mathbf{B}$	176, 414	Banach space, $b \in \mathbf{B}$
$\mathbf{B}^*$	177, 415	dual of Banach space $\mathbf{B}$ , $b^* \in \mathbf{B}^*$
$\mathbf{B}(\mathbf{X}, \mathbf{Y})$	416	set of bounded linear operators from $\mathbf{X}$ to $\mathbf{Y}$
$c_\theta$	155	density of $C_\theta$
$C_\theta$	155	copula function; i.e., d.f. on $[0, 1]^2$
$Cov$		covariance
$C[0, 1]$	414	continuous functions on $[0, 1]$
df		distribution function
$d_H(\cdot, \cdot)$	464	Hellinger distance
$d_K(\cdot, \cdot)$	465	Kolmogorov distance
$d_{pr}(\cdot, \cdot)$	465	Prohorov distance
$d_r(\cdot, \cdot)$	466	Mallows distance
$d_v(\cdot, \cdot)$	464	variational distance
$D(\cdot, \cdot, \cdot), D_n(\cdot, \cdot)$	299, 336	minimization or contrast function
$D[-\infty, \infty]$	191	space of right-continuous functions with left limits
$D_n(\cdot)$	325	$D(\cdot, \mathcal{P}_n)$
$\mathbf{D}$	325	class of minimization functions
$E, E_P, E_\theta$	18	expectation (under $P$ or $P_\theta$ )
$F, F_\theta$	22	distribution function of $P, P_\theta$
$\mathbf{F}_n$	192	empirical df
$\mathcal{F}_i$	81	a sub- $\sigma$ -field
$g$	1, 55	infinite-dimensional parameter
$G$	1, 55	infinite-dimensional parameter
$\mathbf{G}$	1, 55	infinite-dimensional parameter space
$\mathbf{H}$	2, 415	usually a Hilbert space
i.i.d.		independent identically distributed

Symbol	Page	Meaning or explanation
$I(\theta)$	11, 13	Fisher information matrix
$I_*(\theta)$	214	efficient information matrix
$I(P   v, P)$	63	information matrix for the parameter $v$ in the model $P$ at $P \in \mathbf{P}$
$I_v^{-1}$	178, 184	inverse information covariance function or functional
$I^{-1}(P   v, P)$	23, 46, 63, 178	inverse of $I(P   v, P)$ or inverse information covariance functional for $v$
$j(\cdot, \theta)$	89	Jacobian for a parametric group model
$J$	61	identity matrix
$K(\cdot, \cdot)$	129, 162	kernel function
$l, l$	11	$\log p(\theta)$ , log-likelihood of $P_\theta$
$\dot{l}, \dot{l}(\theta)$	13	score function
$\dot{l} = \dot{l}_g$	160, 202	score operator
$I^*, I_i^*$	28, 70, 213	efficient score function
$\tilde{I}, \tilde{I}(\cdot, P   v, P)$	23, 62, 179	efficient influence function for the parameter $v$ in the model $P$ at $P \in \mathbf{P}$
$\tilde{I}_g$	210	efficient influence operator for $g$
$\tilde{I}_g^T \tilde{I}_g$	211	inverse information operator for $g$
$\tilde{I}_v, \tilde{I}(P   v, P)$	178	efficient influence operator
$\tilde{I}_v^T \tilde{I}_v$	179	information bound operator
$l^\infty(T)$	184	the space of bounded functions on $T$ with the supremum norm
$L$	78, 421	the $L$ -operator
$L_1(\mu)$	12, 414	Banach space of $\mu$ -integrable functions
$L_2(\mu)$	12, 414	Hilbert space of $\mu$ -square integrable functions
$L_2^0(P)$	52	Hilbert space of random variables with mean 0 and finite variance under $P$ and covariance as inner product
$L_n(\theta)$	16	log-likelihood of a sample of size $n$ under $\theta$
$L, L_\theta$	12	law or distribution (under $P_\theta$ )
$(M, d)$	471	metric space
$M$	81, 434	a martingale
$\mathbf{M}$	1	the collection of all probability measures on $(X, \mathcal{B})$
$M_\mu$	11	the collection of all probability measures in $\mathbf{M}$ dominated by $\mu$
$\mathcal{M}_B$	475	the Borel $\sigma$ -field of $M$
$N(\mu, \sigma^2)$	5	normal distribution with mean $\mu$ , variance $\sigma^2$
$N(\mu, \Sigma)$	7, 16	normal distribution with mean vector $\mu$ , covariance matrix $\Sigma$



Symbol	Page	Meaning or explanation
$N(\cdot)$	417	null space of an operator
$p(\theta)$	11	$dP_\theta/d\mu$ , density of $P_\theta$
$P, P_0$	1	distribution
$P(f)$	199	$\int f dP$
$P_n$	5, 41, 465	empirical distribution
$\mathbf{P}, \mathbf{P}_0$	1	model, collection of distributions
$\mathbf{P}_1(\eta), \mathbf{P}_1(G)$	27, 55	submodel of $\mathbf{P}$ with $\eta$ (or $G$ ) fixed and the first component of the parameter varying
$\mathbf{P}_2(v)$	32, 55	submodel of $\mathbf{P}$ with $v$ fixed and the second component of the parameter varying
$\dot{\mathbf{P}}$	50	tangent space of $\mathbf{P}$
$\mathbf{P}^0$	50	tangent set of $\mathbf{P}$
$\dot{\mathbf{P}}_i$	55	tangent space of $\mathbf{P}_i$
$\mathbf{Q}$	46, 143	regular parametric submodel of $\mathbf{P}$ , or latent model
$r$	103	regression function, $r \in \mathbf{R}$
$r(\theta)$	15	proxy of $l(\theta)$
$R$	78, 420	real numbers, or the $R$ -operator
$R^m$	2	$m$ -dimensional real Euclidean space
$R^+$	7, 269	the positive reals
$\mathbf{R}$	103	collection of regression functions
$\mathbf{R}(\cdot)$	76, 417	range of an operator
$s(\theta)$	12	$\sqrt{p(\theta)}$
$S_n(\theta)$	16	normed score function of a sample of size $n$
$\mathbf{S}$	48	embedding of $\mathbf{P}$ into $L_2(\mu)$
$\dot{\mathbf{S}}, \dot{\mathbf{S}}$	49, 50, 51	tangent set, tangent space of $\mathbf{S}$
$T = \{T_n\}$	18	sequence of estimates/estimators of $v(P)$
$\mathbf{T}$	153, 292, 430	a group of transformations $T$
$Var$		variance
$\mathbf{V}_n(\cdot)$	311	
$w_i$	113	stratum weight function in biased sampling
$W_i$	113	stratum weight in biased sampling
$W(\cdot, \cdot), W_n(\cdot, \cdot)$	299, 309	estimating function or function used to define an estimating equation
$\dot{W}(P)$	312	matrix of derivatives of $W$
$W_n(\cdot)$	309	$W_n(\cdot, P_n)$
$\mathbf{W}$	325	class of estimating functions
$X, X_i$	1	observations/data/rv's on $\mathbf{X}$
$\mathbf{X}$	1	sample space
$\mathbf{X}$	471	random function with values in a metric space $M$
$\mathbf{Z}, \mathbf{Z}_0$	181	limit processes

## GREEK ALPHABET

Symbol	Page	Meaning or explanation
$\gamma(P)$	9, 53, 222	constraint
$\delta$	273	small positive real or realization of a 0 -1 valued random variable
$\Delta$	9, 207	indicator random variable
$\Delta_\theta$	24	limit random variable in convolution theorem
$\Delta_0$	181	limit process in convolution theorem
$\Delta(\cdot, \cdot)$	163	Green's function
$\varepsilon$	5	error, or small positive real
$\eta$	1	Euclidean nuisance parameter
$H$	27	parameter space of $\eta$
$\theta$	1, 27	parameter, $\theta = (v, \eta) \in R^k$ , $v \in R^m$ , $\eta \in R^{k-m}$
$\Theta$	1, 27	Euclidean parameter space
$\hat{\theta}, \hat{\theta}_n, \theta_n^*$	11, 43, 44	estimators of $\theta$
$\lambda$	9	Lebesgue measure, or hazard function/rate
$\lambda_i$	113	selection probability in biased sampling
$\Lambda$	20, 77	set of all probability measures on $R^m$ , or cumulative hazard function
$\Lambda_n$	499	log-likelihood ratio
$\mu$	11	dominating measure
$v, v(P)$	1, 17, 176	parameter of interest
$\hat{v}_n$	302	estimator of $v$
$\dot{v}$	57, 177	pathwise derivative of $v$
$N$	27	parameter space for $v$
$\pi_r$	184	projection or evaluation map
$\Pi(\cdot   \cdot)$	425	projection
$\Pi_0(\cdot   \cdot)$	50	projection within $L_2(P_0)$
$\Sigma(P, T)$	18	covariance matrix of limit distribution
$\tau$	153	strictly increasing, continuous transformation from $R$ to $R$
$\Psi$	19, 95, 180	influence function
$\Psi_y$	180	influence operator
$\dot{\Psi}$	202	pathwise derivative of $\Psi$
$\dot{\omega}(\theta)$	90	see (i)(c) of section 4.2

## SPECIAL SYMBOL

Symbol	Page	Meaning or explanation
$\sim$	5	distributed according to
$=_d$	22	the same in distribution as

Symbol	Page	Meaning or explanation
$ \cdot $	12	Euclidean norm in $R$ or $R^k$
$\ \cdot\ $	12, 414	norm in Hilbert or Banach space
$\ \cdot\ _0$	50, 70	norm in $L_2(P_0)$
$\ \cdot\ _\infty$	414	supremum norm
$\ \cdot\ _*$	356	norm
$[\cdot]$	49	closed linear span
$\langle \cdot, \cdot \rangle$	414, 415	inner product in Hilbert space, or "duality relation" in Banach space
$\langle \cdot, \cdot \rangle_0$	50	inner product in $L_2(P_0)$
$\bar{A}$	49	closure of the set $A$
$\perp$	417, 425	is orthogonal to; as superscript, orthogonal complement
$A^T$	416	transpose of an operator, matrix or vector $A$
$A^-$	205	generalized inverse of matrix $A$
$\rightarrow$	468	convergence (of several types)
$\Rightarrow$	477	weak convergence for nonmeasurable "r.v.'s"
$f^*$	475	measurable covering function of $f$
$f_*$	475	measurable "undercovering" function of $f$
$\triangleleft, \triangleleft \triangleright$	17, 499	contiguity
$\square$		Halmos: end of proof or example
$\wedge$	147	minimum
$\vee$	165	maximum
$\langle_m$	269	majorization symbol



# Bibliography

- Adams, R. A. (1975). *Sobolev Spaces*, Academic Press, New York.
- Alexander, K. S. (1984). "Probability inequalities for empirical processes and a law of the iterated logarithm," *Ann. Probability* **12**, 1041–1067. Correction: *Ann. Probability* **15** (1987), 428–430.
- Andersen, E. B. (1973). *Conditional Inference and Models for Measuring*, Mentalhygienisk Forlag, Copenhagen.
- Andersen, N. T. and Dobrić, V. (1987). "The central limit theorem for stochastic processes," *Ann. Probability* **15**, 164–177.
- Andersen, P. K., Borgan, O., Gill, R. D., and Keiding, N. (1988). "Censoring, truncation, and filtering in statistical models based on counting processes," *Contemporary Mathematics* **80**, 19–60. American Mathematical Society, Providence.
- Andersen, P. K. and Gill, R. D. (1982). "Cox's regression model for counting processes: a large sample study," *Ann. Statist.* **10**, 1100–1120.
- Anderson, T. W. (1955). "The integral of a symmetric unimodal function over a symmetric convex set and some probability inequalities," *Proc. Amer. Math. Soc.* **6**, 170–176.
- Anderson, T. W. (1984). "Estimating linear statistical relationships," *Ann. Statist.* **12**, 1–45.
- Aranda-Ordaz, F. J. (1983). "An extension of the proportional-hazards model for grouped data," *Biometrics* **39**, 100–107.
- Araujo, A. and Giné, E. (1980). *The Central Limit Theorem for Real and Banach Valued Random Variables*, Wiley, New York.
- Aronszajn, N. (1950). "Theory of reproducing kernels," *Trans. Amer. Math. Soc.* **68**, 337–404.
- Averbukh, V. I. and Smolyanov, O. G. (1967). "The theory of differentiation in linear topological spaces," *Russian Math. Surveys* **22**, (6) 201–258.
- Averbukh, V. I. and Smolyanov, O. G. (1968). "The various definitions of the derivative in linear topological spaces," *Russian Math. Surveys* **23**, (4) 67–113.
- Ayer, M., Brunk, H. D., Ewing, G. M., Reid, W. T., and Silverman, E. (1955). "An empirical distribution function for sampling with incomplete information," *Ann. Math. Statist.* **26**, 641–647.
- Bahadur, R. R. (1960). "On the asymptotic efficiency of tests and estimates," *Sankhya* **22**, 229.
- Bahadur, R. R. (1967). "Rates of convergence of estimates and test statistics," *Ann. Math. Statist.* **38**, 303–325.
- Bailey, K. R. (1984). "Asymptotic equivalence between the Cox estimator and the

- general ML estimators of regression and survival parameters in the Cox model," *Ann. Statist.* **12**, 730-736.
- Bajamonde, A. (1991). "On efficient and robust estimation in semiparametric linear regression with missing observations," Ph.D. dissertation, Univ. California, Berkeley.
- Barlow, R. E., Bartholomew, D. J., Bremner, J. M., and Brunk, H. D. (1972). *Statistical Inference under Order Restrictions*, Wiley, New York.
- Barndorff-Nielsen, O. (1978). *Information and Exponential Families in Statistical Theory*, Wiley, New York.
- Basawa, I. V. and Koul, H. L. (1988). "Large-sample statistics based on quadratic dispersion," *Int. Statist. Rev.* **56**, 199-219.
- Bauer, H. (1972). *Probability Theory and Elements of Measure Theory*, Holt, Rinehart, and Winston, New York.
- Begun, J. M. (1987). "Estimates of relative risk," *Metrika* **34**, 65-82.
- Begun, J. M., Hall, W. J., Huang, W. M., and Wellner, J. A. (1983). "Information and asymptotic efficiency in parametric-nonparametric models," *Ann. Statist.* **11**, 432-452.
- Bennett, S. (1983). "Analysis of survival data by the proportional odds model," *Statist. Med.* **2**, 273-277.
- Beran, R. (1974). "Asymptotically efficient adaptive rank estimates in location models," *Ann. Statist.* **2**, 63-74.
- Beran, R. (1977a) "Estimating a distribution function," *Ann. Statist.* **5**, 400-404.
- Beran, R. (1977b) "Minimum Hellinger distance estimates for parametric models," *Ann. Statist.* **5**, 445-463.
- Beran, R. (1978). "An efficient and robust adaptive estimator of location," *Ann. Statist.* **6**, 292-313.
- Berk, R. H. (1972). "A note on sufficiency and invariance," *Ann. Math. Statist.* **43**, 647-650.
- Berk, R. and Bickel, P. (1968). "On invariance and almost invariance," *Ann. Math. Statist.* **39**, 1573-1576.
- Berman, A. and Plemmons, R. J. (1979). *Nonnegative Matrices in the Mathematical Sciences*, Academic Press, New York.
- Bhanja, J. and Ghosh, J. K. (1991). "Efficient estimation with many nuisance parameters," preprint, Indian Statistical Institute, Calcutta.
- Bhattacharya, P. K. (1983). "Justification for a K-S type test for the slope of a truncated regression," *Ann. Statist.* **11**, 697-701.
- Bhattacharya, P. K., Chernoff, H., and Yang, S. S. (1983). "Nonparametric estimation of the slope of a truncated regression," *Ann. Statist.* **11**, 505-514.
- Bickel, P. J. (1981). "Quelques aspects de la statistique robuste," *Lecture Notes in Math.* Vol. **876**, 1-72. Springer-Verlag, Berlin.
- Bickel, P. J. (1982). "On adaptive estimation," *Ann. Statist.* **10**, 647-671.
- Bickel, P. J. (1985). "Discussion of papers on semiparametric models," *Proc. of 1985 ISI Meeting, Amsterdam*.
- Bickel, P. J. (1986). "Efficient testing in a class of transformation models," *Papers on Semiparametric Models at the ISI Centenary Session, Amsterdam*, 63-81, Report MS-R8614, Centrum voor Wiskunde en Informatica, Amsterdam.
- Bickel, P. J. and Freedman, D. A. (1981). "Some asymptotic theory for the bootstrap," *Ann. Statist.* **9**, 1196-1217.
- Bickel, P. J. and Hodges, J. L. (1967). "The asymptotic theory of Galton's test and a related simple estimate of location," *Ann. Math. Statist.* **38**, 73-89.

- Bickel, P. J. and Klaassen, C. A. J. (1986). "Empirical Bayes estimation in functional and structural models, and uniformly adaptive estimation of location," *Adv. Appl. Math.* 7, 55–69.
- Bickel, P. J. and Lehmann, E. L. (1975). "Descriptive statistics for nonparametric models," I, II, *Ann. Statist.* 3, 1038–1069.
- Bickel, P. J. and Ritov, Y. (1987). "Efficient estimation in the errors in variables model," *Ann. Statist.* 15, 513–540.
- Bickel, P. J. and Ritov, Y. (1988). "Estimating integrated squared density derivatives: sharp best order of convergence estimates," *Sankhya* 50, 381–393.
- Bickel, P. J. and Ritov, Y. (1991). "Large sample theory of estimation in biased sampling regression models, I," *Ann. Statist.* 19, 797–816.
- Bickel, P. J., Ritov, Y., and Wellner, J. A. (1991). "Efficient estimation of linear functionals of a probability measure  $P$  with known marginal distributions," *Ann. Statist.* 19, 1316–1346.
- Billingsley, P. (1968). *Convergence of Probability Measures*, Wiley, New York.
- Billingsley, P. (1971). *Weak Convergence of Measures: Applications in Probability*, Society for Industrial and Applied Mathematics, Philadelphia.
- Billingsley, P. (1986). *Probability and Measure*, 2nd ed., Wiley, New York.
- Birgé, L. (1983). "Approximation dans les espaces métriques et théorie de l'estimation," *Z. Wahrsch. Gebiete* 65, 181–237.
- Birgé, L. (1987a) "Estimating a density under order restrictions: nonasymptotic minimax risk," *Ann. Statist.* 15, 995–1012.
- Birgé, L. (1987b) "On the risk of histograms for estimating decreasing densities," *Ann. Statist.* 15, 1013–1022.
- Birgé, L. (1989). "The Grenander Estimator: a nonasymptotic approach," *Ann. Statist.* 17, 1532–1549.
- Birnbaum, Z. W., Esary, J. D., and Marshall, A. W. (1966). "A stochastic characterization of wear-out for components and systems," *Ann. Math. Statist.* 37, 816–825.
- Blumenthal, S. (1967). "Proportional sampling in life length studies," *Technometrics* 9, 205–218.
- Box, G. E. P. and Cox, D. R. (1964). "An analysis of transformations," *J. Roy. Statist. Soc. Ser. B* 36, 211–252.
- Boyles, R. A., Marshall, A. W., and Proschan, F. (1985). "Inconsistency of the maximum likelihood estimator of distributions having increasing failure rate average," *Ann. Statist.* 13, 413–417.
- Breiman, L. and Friedman, J. H. (1985). "Estimating optimal transformations for multiple regression and correlation," *J. Amer. Statist. Assoc.* 80, 580–598.
- Breslow, N. (1974). "Covariance analysis of censored survival data," *Biometrika* 30, 89–99.
- Breslow, N. and Crowley, J. (1974). "A large sample study of the life table and product limit estimates under random censorship," *Ann. Statist.* 2, 437–453.
- Breslow, N. E. and Day, N. E. (1980). *The Analysis of Case-Control Studies*, International Agency for Research on Cancer, Lyon.
- Bretagnolle, J. and Huber, C. (1979). "Estimation des densités: risque minimax," *Z. Wahrsch. Gebiete* 47, 119–137.
- Brillinger, D. R. (1977). "The identification of a particular nonlinear time series system," *Biometrika* 64, 509–515.
- Brillinger, D. R. (1982). "A generalized linear model with "Gaussian" regressor variables," *A Festschrift for Erich L. Lehmann, In Honor of His Sixty-Fifth Birthday*, 97–114, Wadsworth, Belmont.

- Brown, B. M. (1985). "Multiparameter linearization theorems," *J. Roy. Statist. Soc. Ser. B* **47**, 323–331.
- Brown, L. D. (1986). *Fundamentals of Statistical Exponential Families: with Applications in Statistical Decision Theory*, IMS Lecture Notes-Monograph Series, Vol. 9 Institute of Mathematical Statistics, Hayward.
- Bruni, C. and Koch, G. (1985). "Identifiability of continuous mixtures of unknown Gaussian distributions," *Ann. Probability* **13**, 1341–1357.
- Buckley, J. and James, I. (1979). "Linear regression with censored data," *Biometrika* **66**, 429–436.
- Buja, A. (1990). "Remarks on functional canonical variates, alternating least squares methods and ACE," *Ann. Statist.* **18**, 1032–1069.
- Burkholder, D. L. (1962). "Successive conditional expectations of an integrable function," *Ann. Math. Statist.* **33**, 887–893.
- Burkholder, D. L. and Chow, Y. S. (1961). "Iterates of conditional expectation operators," *Proc. Amer. Math. Soc.* **12**, 490–495.
- Cambanis, S., Huang, S., and Simons, G. (1981). "On the theory of elliptically contoured distributions," *J. Mult. Anal.* **11**, 368–385.
- Campbell, G. (1981). "Nonparametric bivariate estimation with randomly censored data," *Biometrika* **68**, 417–422.
- Campbell, G. (1982). "Asymptotic properties of several nonparametric multivariate distribution function estimators under random censorship," in *Survival Analysis* (J. Crowley and R. A. Johnson, ed.), IMS Lecture Notes-Monograph Series, Vol. 2, pp. 243–256, Institute of Mathematical Statistics, Hayward.
- Carroll, R. J. (1982). "Adapting for heteroscedasticity in linear models," *Ann. Statist.* **10**, 1224–1233.
- Causey, B. D. (1972). "Sensitivity of raked contingency table totals to changes in problem conditions," *Ann. Math. Statist.* **43**, 656–658.
- Chamberlain, G. (1986). "Asymptotic efficiency in semi-parametric models with censoring," *Jour. of Econometrics* **32**, 189–218.
- Chang, M. N. (1990). "Weak convergence of a self-consistent estimator of the survival function with doubly censored data," *Ann. Statist.* **18**, 391–404.
- Chang, M. N. and Yang, G. L. (1987). "Strong consistency of a nonparametric estimator of the survival function with doubly censored data," *Ann. Statist.* **15**, 1536–1547.
- Chen, H. (1985). "Data smoothing in analysis of covariance," preprint, State University of New York at Stony Brook.
- Chen, H. (1988). "Convergence rates for parametric components in a partly linear model," *Ann. Statist.* **16**, 136–146.
- Choi, S. (1989). "On Asymptotically Optimal Tests," Ph. D. dissertation, University of Rochester, Rochester.
- Choi, S. and Hall, W. J. (1988). "On asymptotically optimal tests," Technical Report 8805, Department of Statistics, University of Rochester, Rochester.
- Chou, C. S. and Meyer, P. A. (1974). "La representation des martingales relatives a un processus ponctuel discret," *C. R. Acad. Sci. Paris* **A278**, 1561–1563.
- Chou, C. S. and Meyer, P. A. (1975). "Sur la representation des martingales comme integrales stochastiques dans les processus ponctuels," *Lecture Notes in Math.* Vol. **465** 226–236, Springer-Verlag, Berlin.
- Chow, Y. S., Robbins, H., and Siegmund, D. (1971). *Great Expectations: The Theory of Optimal Stopping*, Houghton Mifflin, Boston.
- Chow, Y. S. and Teicher, H. (1978). *Probability Theory: Independence, Interchangeability, Martingales*, Springer-Verlag, Heidelberg.



- Chung, K. L. (1951). "The strong law of large numbers," *Proc. Second Berkeley Symp. Math. Statist. Prob.* 341–352. Univ. California Press, Berkeley.
- Chung, K. L. (1974). *A Course in Probability Theory*, 2nd ed., Academic Press, New York.
- Clayton, D. (1978). "A model for association in bivariate life tables and its application in epidemiological studies of familial tendency in chronic disease incidence," *Biometrika* 65, 141–151.
- Clayton, D. and Cuzick, J. (1985a) "Multivariate generalizations of the proportional hazards model," *J. Roy. Statist. Soc. Ser. A* 148, 82–117.
- Clayton, D. and Cuzick, J. (1985b) "An approach to inference for rank-regression models with right-censored data," preprint.
- Clayton, D. and Cuzick, J. (1986). "The semi-parametric Pareto model for regression analysis of survival times," *Papers on Semiparametric Models at the ISI Centenary Session, Amsterdam*, Report MS-R8614, Centrum voor Wiskunde en Informatica, Amsterdam.
- Cochran, W. G. (1977). *Sampling Techniques*, 3rd ed., Wiley, New York.
- Cohen, A. C. and Whitten, B. J. (1980). "Estimation in the three-parameter lognormal distribution," *J. Amer. Statist. Assoc.* 75, 399–404.
- Cohn, D. L. (1980). *Measure Theory*, Birkhäuser, Boston.
- Coleman, R. (1979). "An Introduction to Mathematical Stereology," Memoirs No. 3, Department of Theoretical Statistics, Institute of Mathematics Univ. Aarhus.
- Cosslett, S. R. (1981). "Maximum likelihood estimator for choice-based samples," *Econometrica* 49, 1289–1316.
- Cosslett, S. R. (1987). "Efficiency bounds for distribution free estimators of the binary choice and the censored regression models," *Econometrica* 55, 559–585.
- Cox, D. D. (1988). "Approximation of method of regularization estimators," *Ann. Statist.* 16, 694–712.
- Cox, D. D. and O'Sullivan, F. (1990). "Asymptotic analysis of penalized likelihood and related estimators," *Ann. Statist.* 18, 1676–1695.
- Cox, D. R. (1969). "Some sampling problems in technology," (N. L. Johnson and H. Smith, Jr., eds.), *New Developments in Survey Sampling*, 506–527. Wiley-Interscience, New York.
- Cox, D. R. (1972). "Regression models and life tables (with discussion)," *J. Roy. Statist. Soc. Ser. B* 34, 187–220.
- Cox, D. R. (1975). "Partial likelihood," *Biometrika* 62, 269–276.
- Cox, D. R. and Reid, N. (1987). "Parameter orthogonality and approximate conditional inference, (with discussion)," *J. Roy. Statist. Soc. Ser. B* 49, 1–39.
- Cramér, H. (1946). *Mathematical Methods of Statistics*, Princeton University Press, Princeton.
- Csiszar, I. and Tusnády, G. (1984). "Information geometry and alternating minimization procedures," *Statist. Decisions* 1, 205–237.
- Csörgő, M., Csörgő, S., Horváth, L. (1986). *An Asymptotic Theory for Empirical Reliability and Concentration Processes*, Lecture Notes in Statist. Vol. 33, Springer-Verlag, New York.
- Cuzick, J. (1985). "Asymptotic properties of censored linear rank tests," *Ann. Statist.* 13, 133–141.
- Cuzick, J. (1988). "Rank regression," *Ann. Statist.* 16, 1369–1389.
- Dabrowska, D. M. (1988). "Kaplan-Meier estimate on the plane," *Ann. Statist.* 16, 1475–1489.

- Dabrowska, D. M. (1989). "Kaplan-Meier estimate on the plane: weak convergence, LIL, and the bootstrap," *J. Mult. Anal.* **29**, 308-325.
- Dabrowska, D. M. and Doksum, K. (1988). "Partial likelihood in transformation models with censored data," *Scand. J. Statist.* **15**, 1-23.
- Das Gupta, S. (1980). "Brunn-Minkowski inequality and its aftermath," *J. Mult. Anal.* **10**, 296-318.
- De Boor, C. (1978). *A Practical Guide to Splines*, Springer-Verlag, New York.
- Deming, W. E. and Stephan, F. F. (1940). "On a least squares adjustment of a sampled frequency table when the expected marginal totals are known," *Ann. Math. Statist.* **11**, 423-444.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). "Maximum likelihood from incomplete data via the EM algorithm," *J. R. Stat. Soc. Ser. B* **39**, 1-38.
- Deutsch, F. (1983). "Von Neumann's alternating method: the rate of convergence," in *Approximation Theory IV* 427-434. Academic Press, New York.
- Devroye, L. and Györfi, L. (1985). *Nonparametric Density Estimation: The  $L_1$  View*, Wiley, New York.
- Diaconis, P. and Freedman, D. A. (1986). "On the consistency of Bayes estimates (with discussion)," *Ann. Statist.* **14**, 1-67.
- Dieudonné, J. (1960). *Foundations of Modern Analysis*, Academic Press, New York.
- Dinse, G. E. and Lagakos, S. W. (1982). "Nonparametric estimation of lifetime and disease onset distributions from incomplete observations," *Biometrics* **38**, 921-932.
- Dionne, L. (1981). "Efficient nonparametric estimators of parameters in the general linear hypothesis," *Ann. Statist.* **9**, 457-460.
- Doksum, K. A. (1987). "An extension of partial likelihood methods for proportional hazard models to general transformation models," *Ann. Statist.* **15**, 325-345.
- Donoho, D. L. and Liu, R. C. (1987). "Geometrizing rates of convergence, I," Technical Report 137, Department of Statistics, Univ. of California, Berkeley.
- Donoho, D. L. and Liu, R. C. (1991a) "Geometrizing rates of convergence, II," *Ann. Statist.* **19**, 633-667.
- Donoho, D. L. and Liu, R. C. (1991b) "Geometrizing rates of convergence, III," *Ann. Statist.* **19**, 668-701.
- Doss, H., Freitag, S., and Proschan, F. (1989). "Estimating jointly system and component reliability using a mutual censorship approach," *Ann. Statist.* **17**, 764-782.
- Duan, N. and Li, K.-C. (1989). "Regression analysis under link violation," *Ann. Statist.* **17**, 1009-1052.
- Duan, N. and Li, K.-C. (1991). "Slicing regression: a link-free regression method," *Ann. Statist.* **19**, 505-530.
- Dudley, R. M. (1966). "Weak convergence of probabilities on nonseparable metric spaces and empirical measures on Euclidean spaces," *Ill. J. Math.* **10**, 109-126.
- Dudley, R. M. (1968). "Distance of probability measures and random variables," *Ann. Math. Statist.* **39**, 1563-1572.
- Dudley, R. M. (1978). "Central limit theorems for empirical measures," *Ann. Probability* **6**, 899-929. Correction: *Ann. Probability* **7** (1979), 909-911.
- Dudley, R. M. (1984). *A Course on Empirical Processes*, Lecture Notes in Math., Vol. 1097, 1-142. Springer-Verlag, New York.
- Dudley, R. M. (1985). "An Extended Wichura Theorem, Definitions of Donsker Class, and Weighted Empirical Distributions," *Lecture Notes in Math.* Vol. 1153, 141-178. Springer-Verlag, New York.
- Dudley, R. M. (1987). "Universal Donsker classes and metric entropy," *Ann. Probability* **15**, 1306-1326.

- Dudley, R. M. (1989). *Real Analysis and Probability*, Wadsworth and Brooks/Cole, Pacific Grove.
- Dudley, R. M. and Philipp, W. (1983). "Invariance principles for sums of Banach space valued random elements and empirical processes," *Z. Wahrsch. Gebiete* **62**, 509–552.
- Dunford, N. and Schwartz, J. T. (1958). *Linear Operators, Part I*, Interscience, New York.
- Dvoretzky, A., Kiefer, J., and Wolfowitz, J. (1956). "Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator," *Ann. Math. Statist.* **27**, 642–669.
- Efron, B. (1977). "The efficiency of Cox's likelihood function for censored data," *J. Amer. Statist. Assoc.* **72**, 557–565.
- Efron, B. and Johnstone, I. M. (1990). "Fisher's information in terms of the hazard rate," *Ann. Statist.* **18**, 38–62.
- Ekeland, I. and Temam, R. (1976). *Convex Analysis and Variational Problems*, North-Holland, Amsterdam.
- Ekeland, I. and Turnbull, T. (1983). *Infinite-Dimensional Optimization and Convexity*, University of Chicago, Chicago.
- Elbers, C. and Ridder, G. (1982). "True and spurious duration dependence," *Rev. Econ. Stud.* **99**, 403–409.
- Elliott, R. J. (1982). *Stochastic Calculus and Applications*, Springer-Verlag, New York.
- Engle, R. F., Granger, C. W. J., Rice, J., and Weiss, A. (1986). "Semiparametric estimates of the relation between weather and electricity sales," *J. Amer. Statist. Assoc.* **81**, 310–320.
- Fan, J. (1991). "On the optimal rates of convergence for nonparametric deconvolution problem," *Ann. Statist.* **19**, 1257–1272.
- Faraway, J. J. (1992). "Smoothing in Adaptive Estimation," *Ann. Statist.* **20**, 414–427.
- Farrell, R. H. (1972). "On the best obtainable asymptotic rates of convergence in estimation of a density function at a point," *Ann. Math. Statist.* **43**, 170–180.
- Feller, W. (1971). *An Introduction to Probability Theory and Its Applications: Vol. II*, Wiley, New York.
- Ferguson, T. S. (1982). "An inconsistent maximum likelihood estimate," *J. Amer. Statist. Assoc.* **77**, 831–834.
- Fernholz, L. (1983). *Von Mises Calculus for Statistical Functionals*, Lecture Notes in Statistics. Vol. **19**, Springer-Verlag, New York.
- Filippova, A. A. (1962). "Mises' theorem on the asymptotic behavior of functionals of empirical distribution functions and its statistical applications," *Th. Prob. Appl.* **7**, 24–57.
- Fisher, R. A. (1922). "On the mathematical foundations of theoretical statistics," *Phil. Trans. Roy. Soc. London A*, **222**, 309–368, reprinted in R. A. Fisher (1950), *Contributions to Mathematical Statistics* 10.309–10.368, Wiley, New York.
- Fisher, R. A. (1925). "Theory of statistical estimation," *Proc. Cambridge Phil. Soc.* **22**, 700–725. reprinted in R. A. Fisher (1950), *Contributions to Mathematical Statistics* 11.700–11.725, Wiley, New York.
- Fisher, R. A. (1956). *Statistical Methods and Scientific Inference*, Oliver and Boyd, Edinburgh.
- Flett, T. M. (1980). *Differential Analysis*, Cambridge University Press, Cambridge.
- Franchetti, C. and Light, W. A. (1986). "On the von Neumann alternating algorithm in Hilbert space," *J. Math. Anal. Appl.* **114**, 305–314.
- Frank, M. J. (1979). "On the simultaneous associativity of  $F(x,y)$  and  $x + y - F(x,y)$ ," *Aequationes Math.* **19**, 194–226.

- Friedman, J. H. and Stuetzle, W. (1981). "Projection pursuit regression," *J. Amer. Statist. Assoc.* **76**, 817-823.
- Friedrichs, K. (1937). "On certain inequalities and characteristic value problems for analytic functions and for functions of two variables," *Trans. Amer. Math. Soc.* **41**, 321-364.
- Fuller, W. A. and Rao, J. N. K. (1978). "Estimation for a linear regression model with unknown diagonal covariance matrix," *Ann. Statist.* **6**, 1149-1158.
- Gaenssler, P. (1983). *Empirical Processes*, IMS Lecture Notes-Monograph Series, Vol. 3, Institute of Mathematical Statistics, Hayward.
- Geman, S. and Hwang, C.-R. (1982). "Nonparametric maximum likelihood estimation by the method of sieves," *Ann. Statist.* **10**, 401-414.
- Genest, C. (1987). "Frank's family of bivariate distributions," *Biometrika* **74**, 549-555.
- Genest, C. and MacKay, R. J. (1986a) "The joy of copulas: bivariate distributions with uniform marginals," *Amer. Statistician* **40**, 280-283.
- Genest, C. and MacKay, R. J. (1986b) "Copules archimediennes et familles de lois bidimensionnelles dont les marges sont donnees," *Canad. J. Statist.* **14**, 145-159.
- Gill, R. D. (1980). *Censoring and Stochastic Integrals*, Math. Centre Tract **124**, Mathematisch Centrum, Amsterdam.
- Gill, R. D. (1983). "Large sample behaviour of the product-limit estimator on the whole line," *Ann. Statist.* **11**, 49-58.
- Gill, R. D. (1985a) "Discussion of the paper by D. Clayton and J. Cuzick," *J. Roy. Statist. Soc. Ser. A* **148**, 82-117.
- Gill, R. D. (1985b) "Marginal partial likelihood," preprint, Math Center, Amsterdam.
- Gill, R. D. (1988). "Non- and semiparametric maximum likelihood estimators and the von-Mises method (part II)," preprint.
- Gill, R. D. (1989). "Non- and semiparametric maximum likelihood estimators and the von-Mises method (part I)," *Scand. J. Statist.* **16**, 97-128.
- Gill, R. D., Vardi, Y., and Wellner, J. A. (1988). "Large sample theory of empirical distributions in biased sampling models," *Ann. Statist.* **16**, 1069-1112.
- Giné, E. and Zinn, J. (1984). "Some limit theorems for empirical processes," *Ann. Probability* **12**, 929-998.
- Giné, E. and Zinn, J. (1986). *Lectures on the Central Limit Theorem for Empirical Processes*, Lect. Notes in Math., Vol. **1221**, 50-113. Springer-Verlag, Berlin.
- Gleser, L. J. (1981). "Estimation in a multivariate "errors in variables" regression model: large sample results," *Ann. Statist.* **9**, 24-44.
- Godambe, V. P. (1976). "Conditional likelihood and unconditional optimum estimating equations," *Biometrika* **63**, 277-284.
- Gokhale, D. V. and Kullback, S. (1978). *The Information in Contingency Tables*, Dekker, New York.
- Golub, G. H. and Van Loan, C. F. (1983). *Matrix Computations*, Johns Hopkins Univ. Press, Baltimore.
- Gong, G. and Samaniego, F. J. (1981). "Pseudo maximum likelihood estimation: theory and applications," *Ann. Statist.* **9**, 861-869.
- Goodman, L. A. (1985). "The analysis of cross-classified data having ordered and/or unordered categories: Association models, correlation models, and asymmetry models for contingency tables with or without missing entries," *Ann. Statist.* **13**, 10-69.
- Graves, L. M. (1946). *The Theory of Functions of Real Variables*, McGraw-Hill, New York.
- Graybill, F. A. (1983). *Matrices with Applications to Statistics*, Wadsworth, Belmont.

- Green, C. D. (1969). *Integral Equation Methods*, Nelson, London.
- Greenwood, P. and Wefelmeyer, W. (1991). "Cox's estimator is regular," in *Statistical Inference in Stochastic Processes* (Prabhu and Basawa, ed.), Marcel Dekker, New York.
- Grenander, U. (1956). "On the theory of mortality measurement, Part II," *Skand. Akt.* **39**, 125–153.
- Grenander, U. (1981). *Abstract Inference*, Wiley, New York.
- Griffiths, D. (1980). "Interval estimation for the three-parameter lognormal distribution via the likelihood function," *Appl. Statist.* **29**, 58–68.
- Groeneboom, P. (1985). "Estimating a monotone density," *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer II*, (L. M. Le Cam and R. A. Olshen, ed.), 539–555. Wadsworth, Belmont.
- Groeneboom, P. (1988). "Asymptotics for incomplete censored observations," Technical Report, 87–18 Department of Mathematics, University of Amsterdam.
- Groeneboom, P. (1991). "Nonparametric maximum likelihood estimators for interval censoring and deconvolution," Technical Report, 91–53, Faculty of Technical Mathematics and Informatics, Delft University of Technology.
- Groeneboom, P. and Wellner, J. A. (1992). *Information Bounds and Nonparametric Maximum Likelihood Estimation*, Birkhäuser, New York.
- Gumbel, E. J. (1960). "Bivariate exponential distributions," *J. Amer. Statist. Assoc.* **55**, 698–707.
- Haberman, S. J. (1979). *Analysis of Qualitative Data*, **2**, Academic Press, New York.
- Haberman, S. J. (1984). "Adjustment by minimum discriminant information," *Ann. Statist.* **12**, 971–988. Correction: *Ann. Statist.* **14** (1986), 358.
- Hájek, J. (1970). "A characterization of limiting distributions of regular estimates," *Z. Wahrsch. Gebiete* **14**, 323–330.
- Hájek, J. (1972). "Local asymptotic minimax and admissibility in estimation," *Proc. Sixth Berkeley Symp. Math. Statist. Prob.* **1**, 175–194. Univ. California Press, Berkeley.
- Hájek, J. and Šidák, Z. (1967). *Theory of Rank Tests*, Academia, Prague.
- Hall, P. and Titterton, D. M. (1984). "Efficient nonparametric estimation of mixture proportions," *J. R. Statist. Soc. Ser. B* **46**, 465–473.
- Hall, W. J. and Loynes, R. M. (1977). "On the concept of contiguity," *Ann. Probability* **5**, 278–282.
- Hall, W. J. and Mathiason, D. J. (1990). "On large-sample estimation and testing in parametric models," *Int. Statist. Rev.* **58**, 77–97.
- Hall, W. J. and Wellner, J. A. (1980a) "Estimation of mean residual life," Unpublished technical report, University of Rochester,
- Hall, W. J. and Wellner, J. A. (1980b) "Confidence bands for a survival curve from censored data," *Biometrika* **67**, 133–143.
- Hall, W. J. and Wellner, J. A. (1981). "Mean residual life," *Statistics and Related Topics*, (M. Csörgő, D. A. Dawson, J.N.K. Rao, and A.K.Md.E. Saleh, ed.), 169–184. North Holland, Amsterdam.
- Halmos, P. R. (1950). *Measure Theory*, Van Nostrand, Princeton, NJ.
- Halmos, P. R. (1982). *A Hilbert Space Problem Book*, 2nd ed., Springer-Verlag, New York.
- Halperin, I. (1962). "The product of projection operators," *Acta. Sci. Math.* **23**, 96–99.
- Hammerstrom, T. (1978). "On Asymptotic Optimality Properties of Tests and Estimates in the Presence of Increasing Numbers of Nuisance Parameters," Ph. D. dissertation, Univ. California, Berkeley.

- Hampel, F. R. (1974). "The influence curve and its role in robust estimation," *J. Amer. Statist. Assoc.* **62**, 1179–1186.
- Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. J., and Stahel, W. A. (1986). *Robust Statistics; The Approach Based on Influence Functions*, Wiley, New York.
- Han, A. K. (1987). "Nonparametric analysis of a generalized regression model," *J. Econometrics* **35**, 303–316.
- Hanley, J. A. and Parnes, M.N. (1983). "Nonparametric estimation of a multivariate distribution in the presence of censoring," *Biometrics* **39**, 129–139.
- Härdle, W. and Stoker, T. (1989). "Investigating smooth multiple regression by the method of average derivatives," *J. Amer. Statist. Assoc.* **84**, 986–995.
- Hardy, G. H., Littlewood, J. E., and Pólya, G. (1952). *Inequalities*, 2nd ed., Cambridge University Press, Cambridge.
- Has'minskii, R. Z. and Ibragimov, I. A. (1983). "On asymptotic efficiency in the presence of an infinite-dimensional nuisance parameter," *Lecture Notes in Math.*, Vol. **1021**, 195–229. Springer-Verlag, New York.
- Has'minskii, R. Z. and Ibragimov, I. A. (1991). "Asymptotically normal families of distributions and efficient estimation," *Ann. Statist.* **19**, 1681–1724.
- Hastie, T. J. and Tibshirani, R. J. (1990). *Generalized Additive Models*, Chapman and Hall, London.
- Hathaway, R. J. (1985). "A constrained formulation of maximum-likelihood estimation for normal mixture distributions," *Ann. Statist.* **13**, 795–800.
- Hausman, J. A. and Wise, D. A. (1982). "Stratification on endogenous variables and estimation," in *Structural Analysis of Discrete Data: With Econometric Applications* (C. Manski and D. McFadden, ed.), 365–391, M.I.T. Press, Cambridge.
- Heckman, J. J. (1976). "The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models," *Ann. Econ. Social Meas.* **5**, 475–492.
- Heckman, J. J. (1979). "Sample selection bias as a specification error," *Econometrica* **47**, 153–161.
- Heckman, J. J. and Singer, B. (1984). "A method for minimizing the impact of distributional assumptions in economic studies for duration data," *Econometrica* **52**, 271–320.
- Herstein, I. N. (1964). *Topics in Algebra*, Blaisdell, New York.
- Hewitt, E. and Stromberg, K. (1965). *Real and Abstract Analysis*, Springer-Verlag, New York.
- Hill, B. M. (1963). "The three-parameter lognormal distribution and Bayesian analysis of a point-source epidemic," *J. Amer. Statist. Assoc.* **58**, 72–84.
- Hille, E. (1969). *Lectures on Ordinary Differential Equations*, Addison-Wesley, Reading.
- Hjort, N. L. (1990). "Nonparametric Bayes estimators based on beta processes in models for life history data," *Ann. Statist.* **18**, 1259–1294.
- Hochstadt, H. (1973). *Integral Equations*, Wiley, New York.
- Hoffmann-Jørgensen, J. (1984). "Stochastic Processes on Polish Spaces," Unpublished manuscript.
- Horváth, L. (1983). "The rate of strong uniform consistency for the multivariate product-limit estimator," *J. Mult. Anal.* **13**, 202–209.
- Hougaard, P. (1984). "Life table methods for heterogeneous populations: distributions describing the heterogeneity," *Biometrika* **71**, 75–83.
- Hougaard, P. (1986a) "Survival models for heterogeneous populations derived from stable distributions," *Biometrika* **73**, 387–396.

- Hougaard, P. (1986b) "A class of multivariate failure time distributions," *Biometrika* **73**, 671–678. Correction: *Biometrika* **75** (1988), 395.
- Hsieh, D. A. and Manski, C. F. (1987). "Monte Carlo evidence on adaptive maximum likelihood estimation of a regression," *Ann. Statist.* **15**, 541–551.
- Huang, W. M. (1982). "Parameter Estimation when there are Nuisance Functions," Ph.D. dissertation, University of Rochester, Rochester.
- Huang, W. M. (1984). "On effective score estimation in semiparametric models," preprint, Department of Mathematics, Lehigh University, Bethlehem, Pennsylvania.
- Huber, P. J. (1964). "Robust estimation of a location parameter," *Ann. Math. Statist.* **35**, 73–101.
- Huber, P. J. (1967). "The behavior of maximum likelihood estimates under nonstandard conditions," Proc. Fifth Berk. Symp. Math. Statist. and Prob. **1**, 221–233. Univ. California Press, Berkeley.
- Huber, P. J. (1981). *Robust Statistics*, Wiley, New York.
- Huber, P. J. (1985). "Projection pursuit," *Ann. Statist.* **13**, 435–525.
- Ibragimov, I. A. and Has'minskii, R. Z. (1981). *Statistical Estimation: Asymptotic Theory*, Springer, New York.
- Ibragimov, I. A. and Has'minskii, R. Z. (1982). "Estimation of distribution density belonging to a class of entire functions," *Theor. Prob. Applic.* **27**, 551–562.
- Ibragimov, I. A., and Has'minskii, R. Z. (1991). "Asymptotically normal families of distributions and efficient estimation," *Ann. Statist.* **19**, 1681–1724.
- Ireland, C. T. and Kullback, S. (1968). "Contingency tables with given marginals," *Biometrika* **55**, 179–188.
- Jagers, P., Odén, A., and Trulsson, L. (1985). "Post-stratification and ratio estimation: usages of auxiliary information in survey sampling and opinion polls," *Int. Statist. Rev.* **53**, 221–238.
- Jain, N. C. (1977). "Central Limit Theorem and Related Questions In Banach Space," *Proc. Symp. Pure Math.* **31**, 55–65. American Mathematical Society, Providence.
- James, I. R. and Smith, P. J. (1984). "Consistency results for linear regression with censored data," *Ann. Statist.* **12**, 590–600.
- Jameson, G. J. O. (1974). *Topology and Normed Spaces*, Chapman and Hall, London.
- Jewell, N. P. (1982). "Mixtures of exponential distributions," *Ann. Statist.* **10**, 479–484.
- Jewell, N. P. (1985). "Least squares regression with data arising from stratified samples of the dependent variable," *Biometrika* **72**, 11–21.
- Jewell, N. P. and Quesenberry, C. P. (1986). "Regression analysis based on stratified samples," *Biometrika* **73**, 605–614.
- Jin, K. (1992). "Empirical smoothing parameter selection in adaptive estimation," *Ann. Statist.* **20**, 1844–1874.
- Johnson, M. E. (1987). *Multivariate Statistical Simulation*, Wiley, New York.
- Jörgens, K. (1982). *Linear Integral Operators*, Pitman, London. Transl. by G. F. Roach.
- Kagan, A. M., Linnik, Y. V., and Rao, C. R. (1973). *Characterization Problems of Mathematical Statistics*, Wiley, New York.
- Kalbfleisch, J. D. and Prentice, R. L. (1980). *The Statistical Analysis of Failure Time Data*, Wiley, New York.
- Kantorovich, L. V. and Akilov, G. P. (1982). *Functional Analysis*, 2nd ed., Pergamon Press, Oxford.
- Kanwal, R. P. (1971). *Linear Integral Equations*, Academic Press, New York.

- Kaplan, E. L. and Meier, P. (1958). "Nonparametric estimation from incomplete observations," *J. Amer. Statist. Assoc.* **53**, 457-481.
- Karlin, S., and Studden, W. J. (1966). *Tchebycheff Systems: With Applications in Analysis and Statistics*, Wiley Interscience, New York.
- Karlin, S., and Taylor, H. M. (1975). *A First Course in Stochastic Processes*, 2nd ed., Academic Press, New York.
- Kass, R. E. (1989). "The geometry of asymptotic inference," *Statist. Sci.* **4**, 188-234.
- Kato, T. (1976). *Perturbation Theory of Linear Operators*, 2nd. ed., Springer-Verlag, Berlin.
- Kayalar, S. and Weinert, H. L. (1988). "Error bounds for the method of alternating projections," *Math. Control Signals Syst.* **1**, 43-59.
- Keiding, N. and Gill, R. D. (1990). "Random truncation models and Markov processes," *Ann. Statist.* **18**, 582-602.
- Kendall, M. G. and Stuart A. (1979). *The Advanced Theory of Statistics, 2, Inference and Relationship*, 4th ed., Griffin, London.
- Kersting, G. D. (1978). "Die Geschwindigkeit der Glivenko-Cantelli-Konvergenz gemessen in der Prohorov-Metrik," *Math. Z.* **163**, 65-102.
- Kiefer, J. (1961). "On large deviations of the empiric d.f. of vector chance variables and a law of iterated logarithm," *Pacific J. Math.* **11**, 649-660.
- Kiefer, J. and Wolfowitz, J. (1956). "Consistency of the maximum likelihood estimator in the presence of infinitely many nuisance parameters," *Ann. Math. Statist.* **27**, 887-906.
- Kiefer, J. and Wolfowitz, J. (1959). "Asymptotic minimax character of the sample distribution function for vector chance variables," *Ann. Math. Statist.* **30**, 463-489.
- Kiefer, J. and Wolfowitz, J. (1976). "Asymptotically minimax estimation of concave and convex distribution functions," *Z. Wahrsch. Gebiete* **34**, 73-85.
- Kiefer, J. and Wolfowitz, J. (1985). "Note on asymptotic efficiency of m.l. estimators in nonparametric problems," *Collected Papers of Jack Carl Kiefer, Vol. II*, (L. D. Brown, I. Olkin, J. Sacks, H. P. Wynn, eds.), 567-575. Springer-Verlag, New York.
- Kimeldorf, G. and Sampson, A. R. (1975a) "One parameter families of bivariate distributions with fixed marginals," *Comm. Statist.* **4**, 293-301.
- Kimeldorf, G. and Sampson, A. R. (1975b) "Uniform representation of bivariate distributions," *Comm. Statist.* **4**, 617-627.
- Klaassen, C. A. J. (1987). "Consistent estimation of the influence function of locally asymptotically linear estimates," *Ann. Statist.* **15**, 1548-1562.
- Klaassen, C. A. J. (1989). "Efficient estimation in the Cox model for survival data," *Proceedings of the Fourth Prague Symposium on Asymptotic Statistics* (P. Mandl and M. Hušková, eds.), 313-319. Charles University, Prague.
- Klaassen, C. A. J. (1993). "Efficient estimation in the Clayton-Cuzick model for survival data," preprint, University of Amsterdam, Amsterdam.
- Klaassen, C. A. J., Van der Vaart, A. W., and Van Zwet, W. R. (1988). "On estimating a parameter and its score function, II," *Statistical Decision Theory and Related Topics IV*, **2**, (S. S. Gupta and J. O. Berger, eds.), 281-288. Springer-Verlag, New York.
- Klaassen, C. A. J. and Van Zwet, W. R. (1985). "On estimating a parameter and its score function," *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer*, Vol. II, (L. M. Le Cam and R. A. Olshen, eds.), 827-839. Wadsworth.
- Klefsjö, B. (1983). "Some tests against aging based on the total time on test transform," *Comm. Statist.* **12**, 907-927.



- Kober, H. (1939). "A theorem on Banach spaces," *Compositio Math.* 7, 135-140.
- Kodell, R. L., Shaw, G. W., and Johnson, A. M. (1982). "Nonparametric joint estimators for disease resistance and survival functions in survival sacrifice experiments," *Biometrics* 38, 43-58.
- Koshevnik, Yu. A. and Levit, B. Ya. (1976). "On a non-parametric analogue of the information matrix," *Theor. Prob. Applic.* 21, 738-753.
- Koul, H. and Susarla, V. (1983). "Adaptive estimation in linear regression," *Statist. Decisions* 11, 379-400.
- Koul, H., Susarla, V., and Van Ryzin, J. (1981). "Regression analysis with randomly right-censored data," *Ann. Statist.* 9, 1276-1288.
- Kraft, C. H. and Le Cam, L. M. (1956). "A remark on the roots of the maximum likelihood equation," *Ann. Math. Statist.* 27, 1174-1177.
- Kumon, M. and Amari, S. (1984). "Estimation of a structural parameter in the presence of a large number of nuisance parameters," *Biometrika* 71, 445-459.
- Lai, T. L. and Ying, Z. (1988). "Stochastic integrals of empirical-type processes with applications to censored regression," *J. Multivariate Anal.* 27, 334-358.
- Lai, T. L. and Ying, Z. (1991a) "Estimating a distribution function with truncated and censored data," *Ann. Statist.* 19, 417-442.
- Lai, T. L., and Ying, Z. (1991b) "Rank regression methods for left-truncated and right-censored data," *Ann. Statist.* 19, 531-556.
- Lai, T. L., Ying, Z., and Zheng, Z. (1987). "Asymptotic properties of a class of adaptive statistics with applications to regression analysis of censored data," Technical Report, Dept. Statistics, Columbia Univ., New York.
- Laird, N. (1978). "Nonparametric maximum likelihood estimation of a mixing distribution," *J. Amer. Statist. Assoc.* 73, 805-811.
- Lambert, D. and Tierney, L. (1984). "Asymptotic properties of maximum likelihood estimates in the mixed Poisson model," *Ann. Statist.* 12, 1388-1399.
- Lancaster, H. O. (1969). *The Chi-squared Distribution*, Wiley, New York.
- Langholz, B. and Thomas, D. (1986). "Residual relative risk functions as a diagnostic tool in modelling cohort data over time," preprint, USC Department of Preventive Medicine.
- Lawless, J. F. (1982). *Statistical Models and Methods for Lifetime Data*, Wiley, New York.
- Le Cam, L. (1956). "On the asymptotic theory of estimation and testing hypotheses," *Proc. Third Berkeley Symp. Math. Statist. Prob.* 1, 129-156. Univ. California Press, Berkeley.
- Le Cam, L. (1960). "Locally asymptotically normal families of distributions," *Univ. California Publ. Statist.* 3, 37-98.
- Le Cam, L. (1964). "Sufficiency and approximate sufficiency," *Ann. Math. Statist.* 35, 1419-1455.
- Le Cam, L. (1969). *Theorie Asymptotique de la Decision Statistique*, Les Presses de l'Universite de Montreal, Montreal.
- Le Cam, L. (1970). "On the assumptions used to prove asymptotic normality of maximum likelihood estimates," *Ann. Math. Statist.* 41, 802-828.
- Le Cam, L. (1972). "Limits of experiments," *Proc. Sixth Berkeley Symp. Math. Statist. Prob.* 1, 245-261. Univ. California Press, Berkeley.
- Le Cam, L. (1973). "Convergence of estimates under dimensionality restrictions," *Ann. Statist.* 1, 38-53.
- Le Cam, L. (1979). "On a theorem of J. Hájek," *Contributions to Statistics: Jaroslav Hájek Memorial Volume*, (J. Jurecková, ed.) 119-135. Reidel, Dordrecht.

- Le Cam, L. (1985). "Sur l'approximation de familles de mesures par des familles gaussiennes," *Ann. Inst. H. Poincaré* **21**, 225–287.
- Le Cam, L. (1986). *Asymptotic Methods in Statistical Decision Theory*, Springer-Verlag, New York.
- Le Cam, L. (1990). "Maximum likelihood: an introduction," *Int. Statist. Rev.* **58**, 153–171.
- Le Cam, L. and Schwartz, L. (1960). "A necessary and sufficient condition for the existence of consistent estimates," *Ann. Math. Statist.* **31**, 140–150.
- Le Cam, L. and Yang, G. L. (1988). "On the preservation of local asymptotic normality under information loss," *Ann. Statist.* **16**, 483–520.
- Le Cam, L. and Yang, G. L. (1990). *Asymptotics in Statistics; Some Basic Concepts*, Springer-Verlag, New York.
- Lehmann, E. L. (1983). *Theory of Point Estimation*, Wiley, New York.
- Lehmann, E. L. (1986). *Testing Statistical Hypotheses*, 2nd ed., Wiley, New York.
- Leurgans, S., Tsai, W.Y., and Crowley, J. (1982). "Freund's bivariate exponential distribution and censoring," *Survival Analysis*, (J. Crowley and R. A. Johnson, eds.), 230–242. IMS.
- Levit, B. Ya. (1975). "Conditional estimation of linear functionals," *Probl. Inform. Transmission* **11**, 39–54.
- Levit, B. Ya. (1978). "Infinite-dimensional informational lower bounds," *Theory Prob. Applic.* **23**, 388–394.
- Lin, D. Y. and Wei, L. J. (1989). "The robust inference for the Cox proportional hazards model," *J. Amer. Statist. Assoc.* **84**, 1074–1078.
- Lindley, D.V. and Singpurwalla, N. D. (1986). "Multivariate distributions for the life lengths of components of a system sharing a common environment," *J. Appl. Prob.* **23**, 418–431.
- Lindsay, B. G. (1980). "Nuisance parameters, mixture models, and the efficiency of partial likelihood estimators," *Phil. Trans. Royal Soc.* **296**, 639–665.
- Lindsay, B. G. (1982). "Conditional score functions: some optimality results," *Biometrika* **69**, 503–512.
- Lindsay, B. G. (1983a) "Efficiency of the conditional score in a mixture setting," *Ann. Statist.* **11**, 486–497.
- Lindsay, B. G. (1983b) "The geometry of mixture likelihoods: A general theory," *Ann. Statist.* **11**, 86–94.
- Lindsay, B. G. (1983c) "The geometry of mixture likelihoods, II: The exponential family," *Ann. Statist.* **11**, 783–792.
- Lindsay, B. G. (1985). "Using empirical partially Bayes inference for increased efficiency," *Ann. Statist.* **13**, 914–931.
- Liptser, R. S. and Shirayayev, A. N. (1977). *Statistics of Random Processes I*, Springer-Verlag, New York.
- Little, R. J. A. and Rubin, D. B. (1987). *Statistical Analysis with Missing Data*. Wiley, New York.
- Loève, M. (1963). *Probability Theory*, 3rd ed., Van Nostrand Reinhold, New York.
- Loève, M. (1977). *Probability Theory, I*, 4th ed., Springer-Verlag, New York.
- Loève, M. (1978). *Probability Theory, II*, 4th ed., Springer-Verlag, New York.
- Luenberger, D. G. (1969). *Optimization by Vector Space Methods*, Wiley, New York.
- Maguluri, G. (1986). "Inference for a Bivariate Survival Function," Ph.D. dissertation, Univ. of Rochester, Rochester.
- Maguluri, G. (1993). "Semiparametric inference for association in a bivariate survival function," *Ann. Statist.* **21**, to appear.

- Manski, C. (1984). "Adaptive estimation of non-linear regression models," *Econometric Rev.* 3, 145-194.
- Manski, C. F. and Lerman, S. R. (1977). "The estimation of choice probabilities from choice-based samples," *Econometrica* 45, 1977-1988.
- Manski, C. F. and McFadden, D. (eds.) (1981). *Structural Analysis of Discrete Data with Applications*, MIT Press, Cambridge.
- Maronna, R. A. (1976). "Robust M-estimators of multivariate location and scatter," *Ann. Statist.* 4, 51-67.
- Marshall, A. W. and Olkin, I. (1979). *Inequalities: Theory of Majorization and Its Applications*, Academic Press, New York.
- Marshall, A. W. and Olkin, I. (1988). "Families of multivariate distributions," *J. Amer. Statist. Assoc.* 83, 834-841.
- McCullagh, P. (1984). "On the elimination of nuisance parameters in the proportional odds model," *J. Roy. Statist. Soc. Ser. B* 46, 250-256.
- McCullagh, P. and Nelder, J. A. (1983). *Generalized Linear Models*, Chapman and Hall, London.
- McDonald, J. (1983). "Periodic smoothing of time series," Project Orion Technical Report 17, Dept. Statist., Stanford Univ., Stanford.
- Meilijson, I. (1981). "Estimation of the lifetime distribution of the parts from the autopsy statistics of the machine," *J. Appl. Prob.* 18, 829-838.
- Millar, P. W. (1979). "Asymptotic minimax theorems for the sample distribution function," *Z. Wahrsch. Gebiete* 48, 233-252.
- Millar, P. W. (1983). "The Minimax Principle in Asymptotic Statistical Theory," *Lecture Notes in Math.*, Vol. 976, 75-265, Springer-Verlag, New York.
- Millar, P. W. (1985). "Non-parametric applications of an infinite dimensional convolution theorem," *Z. Wahrsch. Gebiete* 68, 545-556.
- Millar, R. (1989). "Estimation of Mixing and Mixed Distributions," Ph.D. dissertation Univ. of Washington, Seattle.
- Miller, R. G. (1976). "Least squares regression with censored data," *Biometrika* 63, 449-464.
- Miller, R. and Halpern, J. (1982). "Regression with censored data," *Biometrika* 69, 521-531.
- Morgenstern, D. (1956). "Einfache Beispiele zweidimensionaler Verteilungen," *Mitt. Math. Statist.* 8, 234-235.
- Mosteller, F. (1968). "Association and estimation in contingency tables," *J. Amer. Statist. Assoc.* 63, 1-28.
- Nakano, H. (1953). *Spectral Theory in the Hilbert Space*, Japan Soc. for Promotion of Sci., Tokyo.
- Newey, W. K. (1990). "Semiparametric efficiency bounds," *J. Appl. Econometrics* 5, 99-135.
- Neyman, J. and Scott, E. (1948). "Consistent estimates based on partially consistent observations," *Econometrica* 16, 1-32.
- Neyman, J. and Scott, E. L. (1951). "On certain methods of estimating the linear structural relation between two variables," *Ann. Math. Statist.* 22, 352-361.
- Novinger, W. P. (1972). "Mean convergence in  $L_p$ -spaces," *Proc. Amer. Math. Soc.* 34, 627-628.
- Oakes, D. (1982). "A model for association in bivariate survival data," *J. Roy. Statist. Soc. Ser. B* 44 414-422.
- Oakes, D. (1986). "Semiparametric inference in a model for association in bivariate survival data," *Biometrika* 73, 353-361.

- Oakes, D. (1989). "Bivariate survival models induced by frailties," *J. Amer. Statist. Assoc.* **84**, 487-493.
- Oosterhoff, J. and Van Zwet, W. R. (1979). "A note on contiguity and Hellinger distance," *Contributions to Statistics: Jaroslav Hájek Memorial Volume* (J. Jurecková, ed.), 157-166. Reidel, Dordrecht.
- Ossiander, M. (1987). "A central limit theorem under metric entropy with  $L_2$  bracketing," *Ann. Probability* **15**, 897-919.
- Park, B. U. (1990). "Efficient estimation in the two-sample semiparametric location-scale model," *Probab. Th. Rel. Fields* **86**, 21-39.
- Parke, W. R. (1986). "Pseudo maximum likelihood estimation: the asymptotic distribution," *Ann. Statist.* **14**, 355-357.
- Parthasarathy, K. R. (1967). *Probability Measures on Metric Spaces*, Academic Press, New York.
- Patil, G. P. and Rao, C. R. (1977). "The weighted distributions: a survey of their applications," *Applications of Statistics* (P. R. Krishnaiah, ed.), North-Holland, Amsterdam.
- Perlman, M. (1972). "On the strong consistency of approximate maximum likelihood estimators," *Proc. Sixth Berkeley Symp. Math. Statist. Probab.* **1**, 263-281.
- Petrov, V. V. (1975). *Sums of Independent Random Variables*, Springer-Verlag, New York.
- Pfanzagl, J. (1979). "Nonparametric minimum contrast estimators," *Selecta Statist. Canad.* **5**, 105-140.
- Pfanzagl, J. (1988). "Consistency of maximum likelihood estimators for certain nonparametric families, in particular: mixtures," *J. Statist. Planning Inference*, **19**, 137-158.
- Pfanzagl, J. (1989). "A supplement to the convolution theorem," Technical Report, **125**, University of Cologne, Cologne.
- Pfanzagl, J. (1990). *Estimation in Semiparametric Models: Some Recent Developments*, Lecture Notes in Statist., Vol. **63**, Springer-Verlag, New York.
- Pfanzagl, J. (1991). "Incidental versus random nuisance parameters," Technical Report, University of Cologne, Cologne.
- Pfanzagl, J. (with the assistance of W. Wefelmeyer) (1982). *Contributions to a General Asymptotic Statistical Theory*, Lecture Notes in Statist. Vol. **13**, Springer-Verlag, New York.
- Pfanzagl, J. (with the assistance of W. Wefelmeyer) (1985). *Asymptotic Expansions for General Statistical Models*, Lecture Notes in Statist. Vol. **31**, Springer-Verlag, New York.
- Pitman, E. J. G. (1979). *Some Basic Theory for Statistical Inference*, Chapman and Hall, London.
- Plackett, R. L. (1965). "A class of bivariate distributions," *J. Amer. Statist. Assoc.* **60**, 516-522.
- Pollard, D. (1982). "A central limit theorem for empirical processes," *J. Austral. Math. Soc. A* **33**, 235-248.
- Pollard, D. (1984). *Convergence of Stochastic Processes*, Springer-Verlag, New York.
- Pollard, D. (1985). "New ways to prove central limit theorems," *Econometric Theory* **1**, 295-314.
- Pollard, D. (1990). *Empirical Processes: Theory and Applications*, NSF-CBMS Regional Conference Series in Probability and Statistics **2**, Inst. Math. Statist. and Ameri. Statist. Assoc., Hayward.
- Pons, O. (1986). "A test of independence for two censored survival times," *Scand. J. Statist.* **13**, 173-185.

- Pontryagin, L. S. (1966). *Topological Groups*, 2nd ed., Gordon and Breach, London.
- Prakasa Rao, B. L. S. (1983). "Estimation of a unimodal density," *Sankhya Ser. A* **31**, 23-36.
- Pratt, J. W. (1960). "On interchanging limits and integrals," *Ann. Math. Statist.* **31**, 74-77. Acknowledgement: *Ann. Math. Statist.* **37** (1966), 1407.
- Prentice, R. L. and Pyke, R. (1979). "Logistic disease incidence models and case-control studies," *Biometrika* **66**, 403-411.
- Prentice, R. L. and Self, S. G. (1983). "Asymptotic distribution theory for Cox-type regression models with general risk form," *Ann. Statist.* **11**, 804-812.
- Prohorov, Yu. V. (1956). "Convergence of random processes and limit theorems in probability theory," *Theor. Prob. Appl.* **1**, 157-214.
- Pyke, R. (1965). "Spacings," *J. Roy. Statist. Soc. Ser. B* **27**, 395-449.
- Rao, P. V., Schuster, E. F., and Littell, R. C. (1975). "Estimation of shift and center of symmetry based on Kolmogorov-Smirnov statistics," *Ann. Statist.* **3**, 862-873.
- Reed, M. and Simon, B. (1972). *Methods of Modern Mathematical Physics I: Functional Analysis*, Academic Press, New York.
- Reeds, J. A. (1976). "On the Definition of von Mises Functionals," Ph.D. dissertation, Harvard University, Cambridge.
- Reeds, J. A. (1985). "Asymptotic number of roots of Cauchy location likelihood equations," *Ann. Statist.* **13**, 775-784.
- Reiersøl, O. (1950). "Identifiability of a linear relation between variables which are subject to error," *Econometrica* **18**, 375-389.
- Rényi, A. (1959). "On measures of dependence," *Acta Math. Acad. Sci. Hungar.* **10**, 441-451.
- Ridder, G. and Verbakel, W. (1983). "On the estimation of the proportional hazards model in the presence of unobserved heterogeneity," preprint.
- Ritov, Y. (1984). "Efficient and unbiased estimation in nonparametric linear regression with censored data," unpublished manuscript, Dept. Statist., Univ. California, Berkeley.
- Ritov, Y. (1986). "On the deconvolution of a mixture of normal distributions," preprint, Univ. California, Berkeley.
- Ritov, Y. (1987). "Tightness of monotone random fields," *J. Roy. Statist. Soc. Ser. B* **49**, 331-333.
- Ritov, Y. (1990). "Estimation in a linear regression model with censored data," *Ann. Statist.* **18**, 303-328.
- Ritov, Y. and Bickel, P. J. (1990). "Achieving information bounds in non and semi-parametric models," *Ann. Statist.* **18**, 925-938.
- Ritov, Y. and Fygenson, M. (1990). "A monotone estimating equation for censored regression," preprint.
- Ritov, Y. and Wellner, J. A. (1988). "Censoring, martingales, and the Cox model," *Contemp. Math.*, Vol. **80**, 191-219. American Mathematical Society, Providence.
- Robbins, H. (1956.) "An empirical Bayes approach to statistics," *Proc. Third Berkeley Symp. Math. Statist. Prob.* **1**, 157-163. Univ. California Press, Berkeley.
- Robbins, H. and Zhang, C.-H. (1989). "Estimating the superiority of a drug to a placebo when all and only those patients at risk are treated with the drug," *Proc. Nat. Acad. Sci. USA* **86**, 3003-3005.
- Robinson, P. M. (1988a) "Root  $n$  consistent semiparametric regression," *Econometrica* **56**, 931-954.
- Robinson, P. M. (1988b) "Semiparametric economics: a survey," *J. Appl. Econometrics* **3**, 35-51.

- Rockafellar, R. T. (1970). *Convex Analysis*, Princeton University Press, Princeton.
- Rota, G.-C. (1962). "An 'alternierende Verfahren' for general positive operators," *Bull. Amer. Math. Soc.* **68**, 95–102.
- Roussas, G. G. (1972). *Contiguity of Probability Measures: Some Applications in Statistics*, Cambridge University Press, Cambridge.
- Rudin, W. (1966). *Real and Complex Analysis*, McGraw-Hill, New York.
- Rudin, W. (1973). *Functional Analysis*, McGraw-Hill, New York.
- Rüschendorf, L. (1985). "Projections and iterative procedures," *Multivariate Analysis VI* (P. R. Krishnaiah, ed.), 485–493, North-Holland, Amsterdam.
- Ruud, P. (1986). "Consistent estimation of limited dependent variable models despite misspecification of distribution," *J. Econometrics* **32**, 157–187.
- Sacks, J. (1975). "An asymptotically efficient sequence of estimators of a location parameter," *Ann. Statist.* **3**, 285–298.
- Sasieni, P. (1989). "Beyond the Cox Model: Extensions of the Model and Alternative Estimators," Ph.D. dissertation Univ. Washington, Dept. Biostatistics, Seattle.
- Sasieni, P. (1992). "Information bounds for the conditional hazard ratio in a nested family of regression models," *J. Roy. Statist. Soc. Ser. B.* **54**, 617–635.
- Savage, L. J. (1976). "On rereading R. A. Fisher," *Ann. Statist.* **4**, 441–500.
- Schay, G. (1974). "Nearest random variables with given distributions," *Ann. Probability* **2**, 163–166.
- Schick, A. (1986). "On asymptotically efficient estimation in semiparametric models," *Ann. Statist.* **14**, 1139–1151.
- Schick, A. (1987). "A note on the construction of asymptotically linear estimates," *J. Statist. Planning Inference* **16**, 89–105.
- Scholz, F. W. (1974). "A comparison of efficient location estimators," *Ann. Statist.* **2**, 1323–1326.
- Scholz, F. W. (1980). "Towards a unified definition of maximum likelihood," *Canad. J. Statist.* **8**, 193–203.
- Schuster, E. F. (1973). "On the goodness-of-fit problem for continuous symmetric distributions," *J. Amer. Statist. Assoc.* **68**, 713–715. Corrigenda (1974). *J. Amer. Statist. Assoc.* **69**, 288.
- Schuster, E. F. (1975). "Estimating the distribution function of a symmetric distribution," *Biometrika* **62**, 631–635.
- Schweizer, B. and Sklar, A. (1983). *Probabilistic Metric Spaces*, North-Holland, Amsterdam.
- Scott, A. J. and Wild, C. J. (1985). "Fitting logistic models in case-control studies," *Proc., ISI Centenary Meeting* **2**, 12.3.1–12.3.15. ISI, Amsterdam.
- Self, S. G. and Prentice, R. L. (1982). "Commentary on Andersen and Gill's Cox regression model for counting processes: a large sample study," *Ann. Statist.* **10**, 1121–1124.
- Serfling, R. J. (1980). *Approximation Theorems of Mathematical Statistics*, Wiley, New York.
- Serre, J.-P. (1977). *Linear Representations of Finite Groups*, Springer-Verlag, New York.
- Sheehy, A. (1987). "Kullback-Leibler Estimation of Probability Measures with an Application to Clustering," Ph.D. dissertation, Dept. of Statist., Univ. of Washington, Seattle.
- Sheehy, A. (1988). "Kullback-Leibler constrained estimation of probability measures," Technical Report 137, Dept. of Statist., Univ. of Washington, Seattle.

- Sheehy, A. and Wellner, J. A. (1988). "Uniformity in  $P$  of some limit theorems for empirical measures and processes," Technical Report 134, Dept. Statistics. Univ. of Washington, Seattle.
- Sheehy, A. and Wellner, J. A. (1992). "Uniform Donsker classes of functions," *Ann. Probability* 20, 1983–2030.
- Shiryayev, A. N. (1984). *Probability*, Springer-Verlag, New York.
- Shorack, G. R., and Wellner, J. A. (1986). *Empirical Processes with Applications to Statistics*, Wiley, New York.
- Silverman, B. (1985). *Density Estimation*, Chapman and Hall, London.
- Skorokhod, A. V. (1956). "Limit theorems for stochastic processes," *Theor. Probability Appl.* 1, 261–290.
- Smith, K. T., Solomon, D. C., and Wagner, S. L. (1977). "Practical and mathematical aspects of the problem of reconstructing objects from radiographs," *Bull. Amer. Math. Soc.* 83, 1227–1270.
- Spiegelman, C. (1979). "On estimating the slope of a straight line when both variables are subject to error," *Ann. Statist.* 7, 201–206.
- Staniswalis, J. G. (1989). "The kernel estimate of a regression function in likelihood-based models," *J. Amer. Statist. Assoc.* 84, 276–283.
- Stein, C. (1956). "Efficient nonparametric testing and estimation," *Proc. Third Berkeley Symp. Math. Statist. Prob.* 1, 187–195, Univ. California Press, Berkeley.
- Stoker, T. M. (1986). "Consistent estimation of scaled coefficients," *Econometrica* 54, 1461–1481.
- Stoker, T. M. (1991). "Lectures on Semiparametric Econometrics," CORE Lecture Series, CORE Foundation, Louvain la Neuve, Belgium.
- Stone, C. J. (1975). "Adaptive maximum likelihood estimators of a location parameter," *Ann. Statist.* 3, 267–284.
- Stone, C. J. (1977). "Consistent nonparametric regression (with discussion)," *Ann. Statist.* 5, 595–645.
- Stone, C. J. (1980). "Optimal rates of convergence for nonparametric estimators," *Ann. Statist.* 8, 1348–1360.
- Stone, C. J. (1985). "Additive regression and other nonparametric models," *Ann. Statist.* 13, 689–705.
- Stone, C. J. (1990). "Large-sample inference for log-spline models," *Ann. Statist.* 18, 717–741.
- Straf, M. L. (1972). "Weak convergence of stochastic processes with several parameters," *Proc. Sixth Berk. Symp. Math. Statist. and Prob.*, 2, 187–221. Univ. Calif. Press, Berkeley.
- Strassen, V. (1965). "The existence of probability measures with given marginals," *Ann. Math. Statist.* 36, 423–439.
- Strasser, H. (1985). *Mathematical Theory of Statistics: Statistical Experiments and Asymptotic Decision Theory*, De Gruyter, New York.
- Sweeting, T. J. (1980). "Uniform asymptotic normality of the maximum likelihood estimator," *Ann. Statist.* 8, 1375–1381. Correction: *Ann. Statist.* 10 (1982), 320.
- Thomas, D. (1986). "Use of auxiliary information in fitting nonproportional hazards models," *Modern Statistical Methods in Chronic Disease Epidemiology* (S. H. Moolgavkar and R. L. Prentice, eds.), 197–210. Wiley, New York.
- Tierney, L. and Lambert, D. (1984). "Asymptotic efficiency of estimators of functionals of mixed distributions," *Ann. Statist.* 12, 1380–1387.
- Tikhonov, A. and Arsenin, V. (1977). *Solutions of Ill-Posed Problems*, Wiley, New York.

- Topsoe, F. (1967). "Preservation of weak convergence under mappings," *Ann. Math. Statist.* **38**, 1661-1665.
- Tricomi, F. G. (1957). *Integral Equations*, Interscience, New York.
- Tsai, W.-Y. and Crowley, J. (1985). "A large sample study of generalized maximum likelihood estimators from incomplete data via self-consistency," *Ann. Statist.* **13**, 1317-1334.
- Tsai, W.-Y., Jewell, N. P., and Wang, M. C. (1987). "A note on the product limit estimator under right censoring and left truncation," *Biometrika* **74**, 883-886.
- Tsiatis, A. A. (1981). "A large sample study of Cox's regression model," *Ann. Statist.* **9**, 93-108.
- Tsiatis, A. A. (1990). "Estimating regression parameters using linear rank tests for censored data," *Ann. Statist.* **18**, 354-372.
- Tsui, K., Jewell, N. P., and Wu, C.F.J. (1988). "A nonparametric approach to the truncated regression problem," *J. Amer. Statist. Assoc.* **83**, 785-792.
- Turnbull, B. W. (1974). "Nonparametric estimation of a survivorship function with doubly censored data," *J. Amer. Statist. Assoc.* **69**, 169-173.
- Turnbull, B. W. (1976). "The empirical distribution function with arbitrarily grouped, censored, and truncated data," *J. Roy. Statist. Soc. Ser. B* **38**, 290-295.
- Turnbull, B. W. and Mitchell, T. J. (1984). "Nonparametric estimation of the distribution of time to onset for specific diseases in survival/sacrifice experiments," *Biometrics* **40**, 41-50.
- Van der Vaart, A. W. (1988a) *Statistical Estimation in Large Parameter Spaces*, CWI Tract 44, Centrum voor Wiskunde en Informatica, Amsterdam.
- Van der Vaart, A. W. (1988b) "Estimating a real parameter in a class of semiparametric models," *Ann. Statist.* **16**, 1450-1474.
- Van der Vaart, A. W. (1988c) "Efficiency and Hadamard Differentiability," Technical Report 143, Dept. of Statist., Univ. of Washington, Seattle.
- Van der Vaart, A. W. (1989). "On the asymptotic information bound," *Ann. Statist.* **17**, 1487-1500.
- Van der Vaart, A. W. (1991). "On differentiable functionals," *Ann. Statist.* **19**, 178-204.
- Van der Vaart, A. W. and Wellner, J. A. (1990). "Prohorov and continuous mapping theorems in the Hoffmann-Jørgensen weak convergence theory with applications to convolution and asymptotic minimax theorems," Technical Report 157, Dept. of Statist., Univ. of Washington, Seattle.
- Van Eeden, C. (1970). "Efficiency-robust estimation of location," *Ann. Math. Statist.* **41**, 172-181.
- Van Zuijlen, M. C. A. (1978). "Properties of the empirical distribution function for independent non-identically distributed random variables," *Ann. Probability* **6**, 250-266.
- Varadarajan, V. S. (1958). "Convergence in Distribution of Stochastic Processes," Thesis, Calcutta University,
- Varadarajan, V. S. (1961). "Measures on topological spaces," *Mat. Sb.* **55**, (97) 35-100. English translation 1965. Providence: American Mathematical Society Translations, Series 2, **48**, 161-228.
- Vardi, Y. (1982). "Nonparametric estimation in the presence of length bias," *Ann. Statist.* **10**, 616-620.
- Vardi, Y. (1985). "Empirical distributions in selection bias models," *Ann. Statist.* **13**, 178-205.



- Vardi, Y., Shepp, L. A., and Kaufman, L. (1985). "A statistical model for positron emission tomography," *J. Amer. Statist. Assoc.* **80**, 8-37.
- Vardi, Y. and Zhang, C.-H. (1992). "Large sample study of empirical distributions in a random-multiplicative censoring model," *Ann. Statist.* **20**, 1022-1039.
- Vaupel, J. W., Manton, K. G., and Stallard, E. (1979). "The impact of heterogeneity in individual frailty on the dynamics of mortality," *Demography* **16**, 439-454.
- Vitale, R. A. (1979). "Regression with given marginals," *Ann. Statist.* **7**, 653-658.
- Von Mises, R. (1947). "On the asymptotic distribution of differentiable statistical functionals," *Ann. Math. Statist.* **18**, 309-348.
- Von Neumann, J. (1950). *Functional Operators, Vol. II: The Geometry of Orthogonal Spaces*, Ann. of Math. Stud. **22**, Princeton University Press, Princeton.
- Wachter, K. W. and Trussell, J. (1982). "Estimating historical heights," *J. Amer. Statist. Assoc.* **77**, 279-293.
- Wald, A. (1943). "Tests of statistical hypotheses concerning several parameters when the number of observations is large," *Trans. Amer. Math. Soc.* **54**, 426-482.
- Wald, A. (1949). "Note on the consistency of maximum likelihood estimate," *Ann. Math. Statist.* **20**, 595-601.
- Walter, G. G. and Blum, J. R. (1984). "A simple solution to a nonparametric maximum likelihood estimation problem," *Ann. Statist.* **12**, 372-379.
- Wang, J.-L. (1985). "Strong consistency of approximate maximum likelihood estimators with applications in nonparametrics," *Ann. Statist.* **13**, 932-946.
- Wang, J.-L. (1986). "Asymptotically minimax estimators for distributions with increasing failure rate," *Ann. Statist.* **14**, 1113-1131.
- Wang, J.-L. (1987). "Estimators of a distribution function with increasing failure rate average," *J. Statist. Planning Inference* **16**, 415-427.
- Wang, M.-C., Jewell, N. P., Tsai, W.-Y. (1986). "Asymptotic properties of the product limit estimate under random truncation," *Ann. Statist.* **14**, 1597-1605.
- Watson, G. S. (1983). *Statistics on Spheres*, Wiley-Interscience, New York.
- Weiss, L. and Wolfowitz, J. (1970). "Asymptotically efficient non-parametric estimators of location and scale parameters," *Z. Wahrsch. Gebiete* **16**, 134-150.
- Wellner, J. A. (1982). "Asymptotic optimality of the product limit estimator," *Ann. Statist.* **10**, 595-602.
- Wellner, J. A. (1989). "Discussion of: 'Non- and semiparametric maximum likelihood estimators and the von Mises method (Part I)' by Richard Gill," *Scand. J. Statist.* **16**, 124-127.
- Wellner, J. A. (1992). "Empirical processes in action; a review," *Int. Statist. Rev.* **60**, 247-269.
- Wellner, J. A. (1993). "Estimation in semiparametric copula models via efficient scores," Unpublished manuscript.
- Wichura, M. (1969). "Inequalities with applications to the weak convergence of random processes with multi-dimensional time parameters," *Ann. Math. Statist.* **40**, 681-687.
- Wichura, M. (1970). "On the construction of almost uniformly convergent random variables with given weakly convergent image laws," *Ann. Math. Statist.* **41**, 284-291.
- Wiener, N. (1955). "On the factorization of matrices," *Comment. Math. Helv.* **29**, 97-110.
- Wilks, S. S. (1962). *Mathematical Statistics*, Wiley, New York.
- Williams, W. H. (1978). *A Sampler on Sampling*, Wiley, New York.
- Witting, H. and Nölle, G. (1970). *Angewandte Mathematische Statistik*, Teubner, Stuttgart.
- Wolfowitz, J. (1957). "The minimum distance method," *Ann. Math. Statist.* **28**, 75-88.

- Wolfowitz, J. (1974). "Asymptotically efficient non-parametric estimators of location and scale parameters, II," *Z. Wahrsch. Gebiete* **30**, 117–128.
- Wong, W. H. (1986). "Theory of partial likelihood," *Ann. Statist.* **14**, 88–123.
- Wong, W. H. and Severini, T. A. (1991). "On maximum likelihood estimation in infinite dimensional parameter spaces," *Ann. Statist.* **19**, 603–632.
- Woodroffe, M. (1985). "Estimating a distribution function with truncated data," *Ann. Statist.* **13**, 163–177. Correction: *Ann. Statist.* **15** (1987), 883.
- Yang, G. L. (1978). "Estimation of a biometric function," *Ann. Statist.* **6**, 112–116.
- Yosida, K. (1974). *Functional Analysis*, 4th ed., Springer-Verlag, New York.

# Author Index

- Adams, R. A., 354, 527  
Akilov, G. P., 371, 372, 448, 456, 459, 537  
Alexander, K. S., 467, 468, 527  
Amari, S., 141, 539  
Andersen, E. B., 132, 527  
Andersen, N. T., 478, 482, 527  
Andersen, P. K., 294, 331, 335, 527  
Anderson, T. W., 27, 127, 466, 527  
Araujo, A., 186, 187, 527  
Aronszajn, N., 438, 448, 527  
Arsenin, V., 344, 545  
Averbukh, V. I., 456, 527  
Ayer, M., 207, 210, 527
- Bahadur, R. R., 325, 349, 527  
Bajamonde, A., 147, 528  
Barlow, R. E., 269, 528  
Barndorff-Nielsen, O., 327, 528  
Bartholomew, D. J., 269, 528  
Bauer, H., 469, 492, 528  
Begun, J. M., xviii, 3, 32, 46, 76, 80, 81, 205, 210, 216, 294, 528  
Beran, R., 41, 76, 192, 403, 457, 528  
Berk, R. H., 327, 431, 528  
Bhanja, J., 126, 528  
Bhattacharya, P. K., 8, 115, 240, 254, 528  
Bickel, P. J., xviii, 6, 17, 23, 47, 96, 98, 111, 115, 126, 138, 139, 154, 164, 226, 227, 235, 321, 347, 362, 381, 382, 383, 384, 395, 396, 402, 431, 466, 527, 528, 529, 543  
Billingsley, P., 179, 184, 228, 373, 469, 472, 478, 482, 484, 485, 489, 496, 503, 529  
Birgé, L., 266, 529  
Birnbaum, Z. W., 269, 529  
Borgan, O., 527  
Box, G. E. P., 153, 529  
Boyles, R. A., 269, 529  
Breiman, L., 437, 440, 529  
Bremner, J. M., 269, 528  
Breslow, N. E., 276, 294, 360, 529  
Bretagnolle, J., 336, 529  
Brillinger, D. R., 109, 529  
Brown, B. M., 328, 519, 530  
Brown, L. D., 346, 530  
Bruni, C., 126, 530
- Brunk, H. D., 207, 210, 269, 527, 528  
Buckley, J., 147, 152, 388, 530  
Burkholder, D. L., 439, 530
- Cambanis, S., 89, 530  
Campbell, G., 290, 530  
Carroll, R. J., 112, 530  
Causey, B. D., 227, 530  
Chang, M. N., 286, 287, 530  
Chen, H., 111, 530  
Chernoff, H., 8, 115, 240, 254, 528  
Choi, S., 76, 530  
Chou, C. S., 434, 530  
Chow, Y. S., 439, 476, 503, 530  
Chung, K. L., 470, 531  
Clayton, D., 8, 153, 156, 158, 531  
Cochran, W. G., 113, 531  
Cohen, A. C., 15, 531  
Cohn, D. L., 431, 484, 531  
Cosslett, S. R., 8, 114, 342, 531  
Cox, D. R., 10, 29, 79, 80, 153, 529  
Cramér, H., 43, 300, 303, 531  
Crowley, J., 276, 286, 360, 529, 540, 546  
Csörgő, M., 198, 531  
Csörgő, S., 198, 531  
Cuzick, J., 153, 158, 172, 320, 321, 323, 531
- Dabrowska, D. M., 290, 320, 531, 532  
Das Gupta, S., 466, 532  
De Boor, C., 346, 532  
Deming, W. E., 227, 532  
Dempster, A. P., 144, 147, 532  
Deutsch, F., 438, 532  
Devroye, L., 177, 336, 532  
Diaconis, P., 345, 532  
Dieudonné, J., 49, 142, 532  
Dionne, L., 96, 532  
Dobrič, V., 478, 482, 527  
Doksum, K. A., 153, 320, 532  
Donoho, D. L., 48, 177, 532  
Duan, N., 110, 532  
Dudley, R. M., 176, 179, 180, 200, 386, 387, 466, 472, 475, 476, 477, 478, 482, 483, 484, 485, 532, 533  
Dunford, N., 423, 533  
Dvoretzky, A., 192, 533

- Efron, B., 424, 533  
 Ekeland, I., 349, 533  
 Elliott, R. J., 434, 533  
 Engle, R., 87, 107, 533  
 Esary, J. D., 269, 529  
 Ewing, G. M., 207, 210, 527
- Fan, J., 267, 533  
 Faraway, J., 403, 533  
 Farrell, R. H., 336, 533  
 Feller, W., 262, 533  
 Ferguson, T. S., 43, 533  
 Femholz, L., 358, 533  
 Filippova, A. A., 4, 310, 313, 363, 533  
 Fisher, R. A., 11, 12, 533  
 Flett, T. M., 456, 533  
 Franchetti, C., 438, 533  
 Frank, M. J., 157, 533  
 Freedman, D. A., 345, 466, 523, 532  
 Friedman, J. H., 107, 437, 440, 529, 534  
 Friedrichs, K., 438, 448, 534  
 Fuller, W. A., 112, 534  
 Fygenson, M., 329, 543
- Gaenssler, P., 482, 484, 534  
 Geman, S., 325, 349, 534  
 Genest, C., 156, 534  
 Ghosh, J. K., 126, 528  
 Gill, R. D., 247, 252, 276, 294, 313, 331, 335, 336, 357, 362, 363, 391, 455, 527, 534, 538  
 Giné, E., 186, 187, 200, 478, 527, 534  
 Gleser, L. J., 138, 534  
 Godambe, V. P., 129, 534  
 Gokhale, D. V., 227, 534  
 Gong, G., 311, 534  
 Goodman, L. A., 227, 534  
 Granger, C. W. J., 87, 107, 533  
 Graves, L. M., 460, 534  
 Graybill, F. A., 84, 534  
 Greenwood, P., 172, 534  
 Grenander, U., 266, 344, 535  
 Griffiths, D., 15, 535  
 Groeneboom, P., 207, 210, 266, 535  
 Gumbel, E. J., 157, 535  
 Györfi, L., 177, 336, 532
- Haberman, S. J., 69, 223, 227, 535  
 Hájek, J., 11, 15, 21, 93, 95, 176, 233, 407, 460, 461, 470, 499, 506, 535  
 Hall, W. J., xviii, 3, 32, 40, 46, 76, 80, 81, 198, 205, 210, 216, 276, 294, 424, 503, 528, 530, 535  
 Halmos, P. R., 91, 419, 436, 535  
 Halperin, I., 443, 447, 535  
 Halpern, J., 541  
 Hammerstrom, T., 134, 535
- Hampel, F. R., 4, 19, 536  
 Han, A. K., 158, 536  
 Härdle, J., 110, 536  
 Hardy, G. H., 423, 536  
 Has'minskii, R. Z., xviii, 3, 5, 17, 27, 41, 43, 48, 107, 139, 177, 309, 460, 461, 466, 467, 536, 537  
 Hausman, J. A., 9, 536  
 Heckman, J. J., 113, 536  
 Herstein, I. N., 230, 536  
 Hewitt, E., 361, 416, 536  
 Hill, B. M., 15, 536  
 Hille, E., 160, 162, 163, 165, 536  
 Hjort, N. L., 345, 536  
 Hodges, J. L., 21, 23, 528  
 Hoffmann-Jørgensen, J., 176, 179, 475, 478, 536  
 Horváth, L., 198, 531, 536  
 Hougaard, P., 157, 536, 537  
 Hsieh, D. A., 96, 403, 537  
 Huang, S., 89, 530  
 Huang, W. M., xviii, 3, 32, 46, 76, 80, 81, 205, 210, 216, 294, 528, 537  
 Huber, P. J., 4, 6, 19, 41, 87, 300, 302, 303, 310, 311, 312, 316, 325, 515, 537  
 Huber, C., 336, 529  
 Hwang, C.-R., 325, 349, 534
- Ibragimov, I. A., xviii, 3, 5, 17, 27, 41, 43, 48, 107, 139, 176, 309, 460, 461, 466, 467, 536, 537  
 Ireland, C. T., 227, 537
- Jagers, P., 122, 537  
 James, I. R., 147, 152, 388, 530, 537  
 Jameson, G. J. O., 490, 537  
 Jewell, N. P., 9, 115, 240, 247, 266, 337, 338, 537  
 Jin, K., 403, 537  
 Johnstone, I. M., 424, 533  
 Jörgens, K., 419, 537
- Kagan, A. M., 98, 537  
 Kalbfleisch, J. D., 147, 537  
 Kantorovich, L. V., 371, 372, 448, 456, 459, 537  
 Kaplan, E. L., 276, 538  
 Karlin, S., 434, 538  
 Kass, R. E., 29, 538  
 Kato, T., 438, 451, 452, 538  
 Kayalar, S., 438, 443, 538  
 Keiding, N., 247, 527, 538  
 Kendall, M. G., 128, 538  
 Kersting, G. D., 465, 538  
 Kiefer, J., 7, 41, 43, 126, 192, 266, 267, 325, 339, 465, 533, 538  
 Kimeldorf, G., 156, 538

- Klaassen, C. A. J., 126, 127, 171, 382, 396, 528, 537  
 Klefsjö, B., 269, 538  
 Kober, H., 436, 452, 539  
 Koch, G., 126, 530  
 Koshevnik, Yu. A., 3, 9, 46, 58, 223, 539  
 Koul, H. L., 96, 147, 527, 538  
 Kraft, C. H., 43, 539  
 Kullback, S., 227, 534, 536  
 Kumon, M., 141, 539  
  
 Lai, T. L., 285, 539, 542  
 Laird, N. M., 144, 147, 338, 532, 539  
 Lambert, D., 264, 267, 539, 545  
 Lawless, J. F., 147, 539  
 Le Cam, L., xviii, 3, 4, 11, 17, 21, 42, 43, 44, 176, 457, 461, 499, 539, 540  
 Lehmann, E. L., 17, 21, 43, 83, 84, 88, 131, 132, 136, 235, 264, 320, 326, 341, 529, 540  
 Lerman, S. R., 8, 114, 541  
 Levit, B. Ya., xviii, 3, 5, 9, 46, 58, 69, 223, 539, 540  
 Li, K.-C., 110, 532  
 Light, W. A., 438, 533  
 Lindsay, B. G., 7, 126, 128, 132, 134, 135, 141, 142, 266, 338, 540  
 Linnik, Y. V., 98, 537  
 Liptser, R. S., 434, 435, 540  
 Littell, R. C., 23, 543  
 Little, R. J. A., 144, 540  
 Littlewood, J. E., 423, 536  
 Liu, R. C., 48, 177, 532  
 Loève, M., 51, 433, 469, 470, 506, 540  
 Loynes, R. M., 503, 535  
 Luenberger, D. G., 456, 540  
  
 MacKay, R. J., 156, 534  
 Maguluri, G. M., 156, 540  
 Manski, C. F., 8, 96, 105, 114, 403, 539, 537, 541  
 Maronna, R. A., 307, 402, 541  
 Marshall, A. W., 8, 157, 269, 529, 530, 541  
 Mathiason, D. J., 40, 535  
 McCullagh, P., 85, 541  
 McDonald, J., 107, 541  
 McFadden, D., 541  
 Meier, P., 276, 538  
 Meyer, P. A., 434, 530  
 Millar, P. W., 176, 182, 183, 195, 263, 264, 271, 369, 541  
 Millar, R., 267, 541  
 Müller, R. G., 9, 147, 541  
 Morgenstern, D., 157, 541  
 Mosteller, F., 227, 541  
  
 Nelder, J. A., 85, 541  
 Newey, W. K., 88, 541  
 Neyman, J., 126, 127, 541  
 Nölle, G., 499, 547  
 Novinger, W. P., 470, 541  
  
 Oakes, D., 8, 156, 157, 541, 542  
 Odén, A., 122, 537  
 Olkin, I., 8, 157, 269, 541  
 Oosterhoff, J., 504, 505, 506, 542  
 Osslander, M., 200, 542  
  
 Park, B. U., 101, 404, 542  
 Parke, W. R., 311, 542  
 Parthasarathy, K. R., 478, 542  
 Perlman, M., 325, 542  
 Petrov, V. V., 470, 471, 542  
 Pfanzagl, J., xviii, 3, 4, 46, 58, 64, 91, 99, 126, 129, 266, 299, 338, 542  
 Philipp, W., 476, 484, 533  
 Plackett, R. L., 157, 542  
 Pollard, D., 192, 200, 312, 362, 485, 542  
 Pölya, G., 423, 536  
 Pontryagin, L. S., 91, 100, 543  
 Prakasa Rao, B. L. S., 266, 336, 543  
 Pratt, J. W., 489, 543  
 Prentice, R. L., 113, 147, 537, 543, 544  
 Prohorov, Yu. V., 472, 543  
 Proschan, F., 269, 529, 532  
 Pyke, R., 113, 424, 543  
  
 Rao, C. R., 98, 537, 542  
 Rao, J. N. K., 112, 534  
 Rao, P. V., 23, 543  
 Reed, M., 280, 418, 440, 453, 543  
 Reeds, J. A., 313, 316, 455, 543  
 Reid, N., 29, 531  
 Reid, W. T., 207, 210, 527  
 Reiersøl, O., 128, 543  
 Rényi, A., 442, 543  
 Rice, J., 87, 107, 533  
 Ritov, Y., 47, 111, 115, 138, 139, 148, 153, 226, 227, 267, 285, 328, 329, 347, 362, 383, 384, 388, 389, 391, 394, 519, 529, 543  
 Robbins, H., 126, 309, 476, 530, 543  
 Robinson, P. M., 88, 111, 543  
 Rockafellar, R. T., 338, 544  
 Ronchetti, E. M., 4, 536  
 Rota, G.-C., 439, 544  
 Rousseeuw, P. J., 4, 536  
 Rubin, D. B., 144, 147, 542, 540  
 Rudin, W., 206, 208, 415, 416, 417, 418, 419, 420, 423, 425, 439, 453, 544  
 Ruud, P., 110, 544  
  
 Sacks, J., 76, 544

- Samaniego, F. J., 311, 534  
 Sampson, A. R., 156, 538  
 Sasieni, P., 82, 335, 544  
 Savage, L. J., 12, 544  
 Schay, G., 466, 544  
 Schick, A., 111, 396, 403, 412, 544  
 Scholz, F. W., 339, 401, 544  
 Schuster, E. F., 23, 195, 543, 544  
 Schwartz, J. T., 423, 533  
 Schwartz, L., 540  
 Schweizer, B., 156, 544  
 Scott, E. L., 126, 127, 541  
 Serfling, R. J., 330, 544  
 Serre, J.-P., 234, 544  
 Severini, T. A., 313, 336, 357, 363, 365, 548  
 Sheehy, A., 69, 223, 335, 484, 544, 545  
 Shiriyayev, A. N., 434, 435, 489, 540, 545  
 Shorack, G. R., 23, 81, 192, 197, 200, 276,  
 315, 334, 359, 422, 424, 435, 436, 472,  
 502, 545  
 Sidák, Z., 15, 93, 95, 233, 407, 461, 470,  
 499, 506, 535  
 Siegmund, D., 476, 530  
 Silverman, B., 348, 545  
 Silverman, E., 207, 210, 527  
 Simon, B., 280, 418, 440, 453, 543  
 Simons, G., 89, 530  
 Sklar, A., 156, 544  
 Skorokhod, A. V., 472, 485, 545  
 Smith, K. T., 443, 447, 545  
 Smith, P. J., 152, 537  
 Smolyanov, O. G., 456, 527  
 Solomon, D. C., 443, 447, 545  
 Spiegelman, C., 139, 545  
 Stahel, W. A., 4, 536  
 Staniswalis, J. G., 345, 545  
 Stephan, F. F., 227, 532  
 Stein, C., xviii, 46, 99, 100, 545  
 Stoker, T. M., 88, 110, 536, 545  
 Stone, C. J., 76, 112, 177, 345, 350, 353,  
 403, 545  
 Strassen, V., 20, 465, 466, 545  
 Stromberg, K., 361, 416, 536  
 Stuart A., 128, 538  
 Studden, W. J., 538  
 Stuetzle, W., 107, 534  
 Susarla, V., 96, 147, 539  
  
 Taylor, H. M., 434, 538  
 Teicher, H., 503, 530  
 Temam, R., 349, 533  
 Tierney, L., 264, 267, 539, 545  
 Tikhonov, A., 344, 545  
 Topsoe, F., 484, 546  
 Tricomi, F. G., 160, 162, 163, 165, 289, 546  
 Trulsson, L., 122, 537  
 Tsai, W.-Y., 247, 286, 540, 546, 547  
  
 Tsiatis, A. A., 80, 152, 285, 294, 320, 329,  
 331, 333, 389, 546  
 Turnbull, B. W., 286, 546  
  
 Van der Laan, M., 292  
 Van der Vaart, A. W., 3, 5, 102, 126, 127,  
 131, 176, 180, 182, 183, 187, 201, 202,  
 204, 210, 247, 248, 266, 270, 418, 420,  
 457, 461, 475, 538, 546  
 Van Eeden, C., 76, 546  
 Van Ryzin, J., 147, 539  
 Van Zuijlen, M. C. A., 323, 546  
 Van Zwet, W. R., 127, 504, 505, 506, 538,  
 541  
 Varadarajan, V. S., 472, 546  
 Vardi, Y., 9, 115, 252, 341, 342, 362, 363,  
 380, 534, 546, 547  
 Von Mises, R., 313, 547  
 Von Neumann, J., 436, 438, 547  
  
 Wagner, S. L., 443, 447, 545  
 Wald, A., 17, 325, 349, 547  
 Wang, J.-L., 271, 547  
 Wang, M.-C., 247, 546, 547  
 Watson, G. S., 230, 547  
 Wefelmeyer, W., xviii, 3, 4, 46, 91, 99, 129,  
 172, 535, 542  
 Weinert, H. L., 438, 443, 538  
 Weiss, A., 87, 107, 533  
 Weiss, L., 101, 547  
 Wellner, J. A., xviii, 3, 5, 23, 33, 46, 76, 80,  
 81, 148, 176, 180, 183, 192, 197, 198,  
 200, 205, 207, 210, 216, 226, 227, 252,  
 276, 294, 315, 334, 335, 347, 359, 362,  
 363, 382, 422, 424, 435, 436, 472, 475,  
 484, 502, 528, 529, 530, 534, 535, 543,  
 545, 546, 547  
 Whitten, B. J., 15, 531  
 Wichura, M., 472, 547  
 Wiener, N., 438, 448, 547  
 Wilks, S. S., 97, 547  
 Wise, D. A., 9, 536  
 Witting, H., 499, 547  
 Wolfowitz, J., 7, 41, 43, 101, 126, 192, 266,  
 267, 325, 339, 533, 538, 547, 548  
 Wong, W. H., 313, 336, 357, 363, 365, 548  
 Woodroffe, M., 247, 548  
  
 Yang, G. L., 3, 4, 198, 286, 287, 461, 530,  
 531, 540, 548  
 Yang, S. S., 8, 115, 240, 254, 528  
 Ying, Z., 285, 541, 539  
 Yosida, K., 419, 548  
  
 Zhang, C-H., 126, 309, 380, 543, 547  
 Zinn, J., 200, 478, 534

# Subject Index

- accelerated time model, 147
- adaptation, 2, 27, 29, 94, 236, 238, 239, 278, 287
- adaptive estimate, 29, 75, 76, 96, 98, 99, 101, 111, 112, 133, 238
- adaptive estimation, 237–240, 306–309
- adjoint operator, 416
- almost invariant, 430
- alternating conditional expectations (ACE), 118, 140, 226, 255, 437, 440
- alternating projections, 167, 436
- analytic continuation, 25, 65
- Anderson's lemma, 27, 187, 466
- angle between subspaces, 168, 438
- annihilator, 417
- ANOVA, 85–86, 127, 134
- asymptotically
  - efficient, 43, 63, 77, 182
  - linear, 19, 20, 180
  - locally regular, 21, 180, 356
  - measurable, 481
  - minimax, 183
  - optimal, 182
  - orthogonal, 499
  - regular, 180
  - tight, 478
  - uniformly regular, 18
- asymptotic generalized  $M$  (AGM) estimate, 310–312, 325, 328, 329, 336, 337, 356, 363, 383, 514
- asymptotic generalized minimum contrast (AGMC) estimate, 356
- asymptotic  $M$  (AM) estimate, 301–303, 310, 311, 316
- asymptotic minimax theorem, 27, 183, 192, 195
- asymptotic optimality theorem, 26, 27, 182
- asymptotic relative efficiency, 225, 297
  
- ball  $\sigma$ -field, 482
- Banach space, 176, 335, 414
- Banach-Steinhaus theorem, 206, 260, 263, 419
- Banach's theorem, 206, 260, 263, 418
- Basu's theorem, 136
- Bernoulli pairs, 128
- Bernstein's inequality, 361
  
- biased sampling, 8, 86, 113, 240, 340, 346, 362
  - regression, 115, 383
- Bickel-Hodges estimate, 23
- bivariate distribution, 223, 225, 289, 346, 360
- bivariate normal distribution, 32, 36, 157, 174
- bootstrap, 4
- bound
  - information, 23, 33, 61, 76, 141, 176
  - (local asymptotic) minimax, 27, 183, 192, 195
- boundary conditions, 162
- bounded Lipschitz function, 386
- Box-Cox model, 86, 153, 172, 292
- bowl-shaped, 26, 187, 466
- Brownian bridge, 23, 192, 194, 199
  - covariance kernel, 162
  - P-Brownian bridge, 200, 223, 227
- Brownian motion, 197, 243, 275
  
- calculus
  - martingale, 81, 246, 276, 434
- canonical gradient, 58, 178
- Carathéodory theorem, 338
- case control, 113
- Cauchy distribution, 316, 317, 373
- censoring, 10, 271, 335
  - bivariate, 289
  - contrast with truncation, 252
  - double, 285
  - indicator, 204, 207
  - interval, 144
  - random, 272, 276, 283, 342
  - right, 86, 144, 147, 272, 329, 342, 388
- chain rule, 456
- choice-based sampling, 8, 114
- closure of sum space, 438, 440
- coefficient of variation, 9, 68
- compatible metric, 41, 298, 309, 337, 465
- compensator, 245, 256
- completely monotone, 262
- concave,
  - Schur, 269, 271
- conditional,
  - expectation, 430, 434
  - hazard function, 10, 77

- conditional (*continued*)
  - likelihood, 35, 308
  - log-likelihood, 132, 321
- confidence region, 4
- consistent, 18, 335
  - $\sqrt{n}$ -, 18, 41, 42
  - uniformly, 18, 20
  - uniformly- $\sqrt{n}$ , 18
- constrained models, 9, 53, 68, 193, 198, 221
- constraints, 53, 221
  - finitely many, 53, 222
  - infinitely many, 223
  - linear, 222
  - monotonicity, 261, 269
  - nonlinear, 222
- construction of estimators, 41, 298, 325, 370
- contiguity, 17, 395, 479, 498, 499, 504
- contiguous, 17, 122, 499
- continuous mapping theorem, 472, 479, 484
- contraction, 372, 378, 379
- convergence
  - almost uniformly, 483
  - in distribution, 471
  - in law, 468, 471
  - in  $\mu$ -measure, 468
  - in (outer) probability, 468, 483
  - in  $r$ -th mean, 468
  - weak, 468, 471, 475, 477
- convex, 300, 325, 326, 329, 348, 349, 353
- convex parametrization, 305, 308
- convexity, 473, 519
  - existence and uniqueness, 325, 519
- convolution theorem
  - finite dimensional, 24, 27, 63, 270
  - infinite dimensional, 180, 182, 270
- copula, 155, 204
  - Archimedean, 156
  - bivariate normal, 157, 174
  - Clayton-Oakes, 8, 156, 173, 295
  - Frank, 157, 173, 295
  - Hougaard, 157
  - Morgenstern, 157, 174
  - Plackett's constant odds, 157
  - regression, 158
- core model, 135, 144, 154, 171
- counting process, 81, 245
  - martingale, 81, 148, 245, 421, 422, 435
  - retro, 246
- covariates, 85, 146
- Cox
  - estimate, 75, 77, 79, 153, 172, 320, 330
  - model, 1, 9, 77, 82, 86, 88, 153, 171, 217, 292-293, 305, 320, 330, 335, 382, 420
  - partial likelihood, 10
- Cramér-Rao information bound, 12, 15, 39
- Cramér-Wold device, 513
- cumulative hazard, 78, 217, 241, 272, 293
  - backward, retro, 244, 254, 255
  - curve, 49
    - one-sided, 263, 265, 270
  - cylinder measure, 186, 368
- Darmois-Skitovich theorem, 98
- data-splitting, 396
- delta-method, 357, 384
- density estimation, 48, 345, 349, 366
- density functions, 11, 48, 345
  - completely monotone, 261
- derivatives, 453
- differentiability
  - $B, \Gamma$ -, 336
  - compactly, 454
  - Fréchet, 12, 454-457
  - Gâteaux, 374, 453-457
  - Hadamard, 336, 364, 383, 454-457
  - Hellinger, 12, 202
  - non, 265
  - ordinary, 13
  - one-sided, 263
  - pathwise, 57, 58, 177, 201, 314, 456
  - pathwise weak, 177
  - in quadratic mean, 457
- differentiable
  - function, 177, 201
- differential equation, 162
- dihedral group, 230
- dihedrally symmetric, 230, 234
- discretized estimator, 44, 316, 385-388, 401, 402
- discretization, 44, 316
- distance. *See* metric
- distribution
  - absolutely continuous, 338
  - Bernoulli, 85, 111, 128, 441
  - binomial, 327
  - bivariate, 223, 225, 289, 346, 360
  - bivariate normal, 32, 36, 157, 174
  - Cauchy, 316, 317, 373
  - completely monotone, 261
  - copula, 295
  - cyclically symmetric, 230, 234
  - dihedrally symmetric, 230, 234
  - Dirichlet, 97, 345
  - double exponential, 373
  - elliptical, 89, 96, 237
  - exponential, 84, 86, 152, 154, 261, 264, 292, 318, 320
  - exponential family, 14, 52, 84, 126, 326, 341, 345
  - extreme value, 152
  - gamma, 85, 238
  - Gaussian, 262, 264, 307
  - Gumbel, 157
  - Kolmogorov-Smirnov, 194



- least favorable, 61
- logistic, 408
- log-normal, 15
- multinomial, 12, 33, 84
- negative binomial, 85
- normal, Gaussian, 32, 36, 50, 84, 262, 264, 307
- Pareto, 153, 154, 292
- Poisson, 85, 262, 264
- Rényi's, 442
- rotationally symmetric, 230
- Schur-convex/concave, 269, 271
- spherically symmetric, 229, 232
- symmetric, 53, 55, 58, 75, 229, 233, 235, 304, 340, 343
- uniform, 50, 134, 155, 261, 263, 266, 441, 481
- Weibull, 15
- Donsker, 200, 360
- dual space, 415
- efficiency, 2, 182
  - weak, 182
- efficient
  - estimation, 43, 391, 405
  - estimators, 24, 41, 44, 63, 139, 147, 182, 394
  - influence function, 3, 23, 30, 32–33, 39, 62, 72, 91, 120, 148, 179, 194, 394
  - influence operator, 178, 210
  - locally, 24
  - score function, 3, 28, 30, 32–33, 70, 76, 90, 149, 158, 165
  - score operator, 210
  - uniformly, 24
- elliptic, 89, 96, 305, 307, 399, 401
- EM algorithm, 86
- empirical distribution (function), 2, 19, 22, 122, 194, 301, 309, 315, 320, 338, 360
- empirical measure, 200, 465, 467
- empirical process, 200, 467
- errors in variables, 7, 86, 127, 135, 305
- estimate/estimator, 2, 17
  - AGM, 310–312, 325, 328, 329, 336, 337, 356, 363, 383, 514
  - AGMC, 356
  - AM, 301–303, 310, 311, 316
  - GM, 299, 309, 311, 317, 319–320, 325, 326, 336, 337, 370, 391–394, 514
  - GMC, 299, 325, 329, 332, 336, 337, 348–349
  - M, 4, 299–311, 317, 325, 326, 328, 341, 384, 516
  - MC, 299, 325, 326, 328, 331, 332
- adaptive, 29, 75, 76, 96, 98, 99, 101, 111, 112, 133, 238
- asymptotically efficient, 43, 77
  - asymptotically linear, 19, 39
  - $B_0^*$  linear, 356
  - Bickel-Hodges, 23
  - conditional maximum likelihood, 132, 134
  - Cox, 75, 77, 79, 153, 172, 320, 330
  - $\sqrt{n}$ -consistent, 43–44, 298
  - efficient, 24, 40–41, 44, 63, 139, 147, 182, 394
  - existence, 41, 325, 370
  - Gaussian regular, 39, 356
  - Grenander, 266
  - Hodges-Lehmann, 310, 315
  - inconsistent, 43, 269
  - Kaplan-Meier, 276, 282, 285, 343, 358, 388, 389
  - least squares, 35, 109, 112, 137, 305
  - linear, 356
  - locally regular, 21, 46, 356
  - maximum likelihood, 11, 15, 35, 37, 43, 138, 147, 299, 304, 328, 337, 338, 343, 392, 401
  - minimum distance, 22, 41, 299, 318
  - Nelson-Aalen, 276, 358
  - one-step, 43, 316, 371, 381
  - partial likelihood, 79, 80, 331
  - preliminary, 41, 146, 395, 396, 402
  - product limit, 247, 276, 358, 389
  - reduced sample, 252, 282
  - regular, 17, 21, 39, 179, 180, 356
  - superefficient, 21
  - Stein's, 22
  - uniformly asymptotically linear, 20, 46
  - uniformly  $\sqrt{n}$ -consistent, 18, 42, 44
  - uniformly Gaussian regular, 18, 20, 24, 46
  - uniformly efficient, 24, 44
  - weakly  $B_0^*$  efficient, 391
  - weakly  $B_0^*$  Gaussian regular, 356
  - weakly  $B_0^*$  linear, 356
  - weakly regular, 181
- estimating equations, 312, 391
- estimation (of)
  - density, 48, 345, 349, 366
  - distribution function, 191, 192, 195, 216, 217, 223, 225, 274, 276, 284
  - cumulative hazard function, 196, 274, 276
  - influence function, 301, 394–398, 402
  - information, 409
  - joint distribution, 346, 360
  - mean, 67, 68
  - mean residual life, 197, 198
  - median, 303
  - probability measure, 199, 222, 226, 227, 235
  - score function, 395, 396, 398, 400, 402
  - symmetric distribution function, 193
- Euler's equation, 98
- evaluation map, 184
- exchangeable, 101, 135, 229, 233, 270, 432

- existence of estimates, 41, 325, 370  
 exponential distribution, 318  
 exponential family, 14, 52, 84, 126, 326, 341, 346  
     curved, 14  
 exponential inequality, 467  
 exponential mixture model, 134, 261, 264, 337, 399, 401  
 exponential proportional hazards model, 86  
 extension, 41, 298
- factorization theorem, 131  
 failure rate  
     increasing, 269, 270  
     increasing, average, 269, 271  
 filtration, 245  
 finite sample, 127  
 Fisher information, 13, 139, 330, 389, 424, 460, 461  
     for location, 15, 55, 75, 147, 399, 407  
     matrix, 13, 91  
     for scale, 134  
 frailty, 153  
 Fréchet differentiable, 12, 454–457  
 Fredholm, 140, 288  
 function  
     differentiability of, 177, 178, 201  
     efficient influence, 3, 23, 30, 33, 39, 62, 72, 91, 120, 148, 179, 194, 394  
     efficient score, 28, 90  
     Green's, 163, 165, 169, 173, 294, 382  
     influence, 3, 19, 38, 39, 40, 180, 395  
     score, 13, 28, 70, 116  
 functional, 57, 415  
     linear, 415  
     log-likelihood, 300
- Gaussian  
     linear regression model, 35  
     regularity, 46, 357  
     shift experiment, 4, 17  
 general linear model, 95  
 generalized inverse, 316  
 generalized  $M$  ( $GM$ )-estimate, 299, 309, 311, 317–319, 325, 326, 336, 337, 370, 391–394, 514  
 generalized minimum contrast ( $GMC$ )-estimate, 299, 325, 329, 332, 336, 337, 348  
 geometry, 27, 30  
 gradient, 58  
     canonical, 58, 178  
 gram matrix, 428  
 Green's function, 163, 165, 169, 173, 294, 382  
 Grenander estimator, 266  
 group,  
     Lie, 234  
     location, 235  
     models, 83, 88  
     parametric, 83, 88  
     semiparametric, 88, 229, 234  
     transformation, 86, 153, 430  
 group models, 83, 88  
 nonparametric, 229  
 parametric, 83, 88  
 semiparametric, 88, 229, 234
- Haar measure, 231, 431  
 Hardy's inequality, 161, 422, 423  
 Has'minskii-Ibragimov model, 139, 147, 263, 267, 308, 348, 372, 392  
 hazard rate / function, 9, 77, 78, 156  
 Hellinger metric, 301, 464  
 Hessian matrix, 43  
 heteroscedastic, 104, 105  
 Hilbert-Schmidt, 208, 440  
 Hilbert space, 4, 5, 48, 415  
 Hodges-Lehmann estimator, 310, 315, 400  
 Hoeffding's formula, 320  
 Hoffmann-Jørgensen-Dudley theory, 179, 475
- identifiability, 6, 126, 238  
 identifiable, 6, 98, 108, 112, 114, 127, 128, 139, 273, 309  
 i.i.d., 4  
 inconsistent, 43, 269  
 independence, 227  
 indicator censoring, 204, 207  
 inequality  
     Anderson's, 27, 187, 466  
     exponential, 467  
     Hardy's, 161, 422, 423  
     information, 26, 27  
 influence function, 3, 19, 38, 39, 40, 180, 395  
     efficient, 3, 23, 30, 33, 39, 62, 72, 91, 120, 148, 179, 194, 394  
     estimation of, 394, 395, 396  
     inefficient, 19, 38, 203  
     operator, 177, 180  
 information, 28, 32–33, 63, 158, 165  
     bound, 23, 32–33, 61, 76, 141, 176  
     bound operator, 179  
     efficient, 62, 214  
     Fisher, 13, 139, 330, 389, 424, 460, 461  
     inequality, 26, 27  
     inverse, 178  
     matrix, 13, 91  
     operator, 206, 210  
 inseparable space, 184  
 integrability  
     Bochner, 393  
     uniform, 466, 468, 469, 503  
 invariant, 320, 431  
     almost, 430  
 inverse information covariance function, 184,

- 192, 218
- inverse information covariance functional, 178, 181, 184
- inverse information operator, 211
- inverse operator, 418
- boundedness of, 418
- invertibility, 418
- isometry, 418, 420, 421
- iteration, 372, 377
- Newton-Raphson, 43, 316, 318, 321, 517
- iterative proportional fitting, 227
- jackknife, 4
- joint distribution-transformation model, 154, 204, 294, 380
- Kaplan-Meier estimator, 276, 282, 285, 343, 358, 388, 389
- Kolmogorov existence theorem, 186
- Kullback-Leibler divergence, 227
- law of large numbers, 45, 324, 333, 470
- least favorable submodel, 61, 62, 76, 134, 135, 139, 395
- Le Cam's
- first lemma, 499
- second lemma, 500
- third lemma, 40, 66, 189, 190, 479, 503
- fourth lemma (lemma A.9.4), 503
- likelihood, 13, 499
- conditional, 35, 308
- log, 11, 43, 499
- maximum, 11, 43, 299
- nonparametric maximum, 126, 252, 266, 267, 338–344, 380, 383, 384, 388, 392
- linearity
- asymptotic, 20, 180
- $B_0^*$ , 356
- locally, 21, 46
- weakly  $B_0^*$ , 356
- linear span, 49
- closed, 49
- linking, 85, 87
- Lipschitz,  $\rho$ -, 42, 309
- local asymptotic
- linearity, 21, 46
- minimax theorem, 27, 183
- normality (LAN), 4, 16, 508
- locally
- efficient, 24
- Gaussian regular, 21, 46, 47
- regular, 21, 46, 356
- location
- models, 1, 5, 109
- parameter, 1
- symmetric, 53, 55, 58, 75, 235, 304, 305, 343, 398, 400, 403
- logistic partial spline, 111
- logistic regression, 85, 111, 113, 118, 329
- log-likelihood, *See* likelihood
- $L_p$ -continuity theorem, 56, 407, 416
- loss
- quadratic, 26
- zero-one, 27
- lower semicontinuous, 187
- marginal distributions
- known, 223, 225
- unknown, 155, 156
- martingale, 359, 422, 434
- calculus, 81, 246, 276, 434
- counting process, 81, 148, 245, 421, 422, 435
- Doob's, 434
- operator, 421
- theory, 81, 245, 276, 422
- transform, 81, 246, 276, 435
- master theorem, 311, 312, 514
- maximum likelihood, 11, 43, 299 *See also* likelihood
- conditional, 308
- inconsistent, 269
- nonparametric, 126, 252, 266, 267, 339–344, 380, 383, 384, 388, 392
- parametric, 11, 43
- penalized, 301, 339, 344, 347, 348, 353, 369, 392
- regularized, 301, 339, 343, 344, 348, 363, 372, 380, 392
- mean
- constrained, 67
- residual life, 197
- residual life operator, 78, 420
- unconstrained, 68
- measurability, 179
- ball  $\sigma$ -field, 482
- Borel, 471, 475
- measurable covering function, 475
- measure
- Haar, 231, 431
- outer, 476
- median, 303
- $M$ -estimator/estimate, 4, 299–311, 317, 325, 326, 328, 341, 384, 516
- asymptotic (AM), 301–303, 310, 311, 316
- asymptotic generalized (AGM), 310–312, 325, 328, 329, 336, 337, 356, 363, 383, 514
- generalized (GM), 299, 309, 311, 317–319, 325, 326, 336, 337, 370, 391–394, 514
- metrics, 464
- compatible, 41, 298, 309, 337, 465
- Hellinger, 301, 464
- $L_2$ , 12, 199, 415

- metrics (*continued*)  
 Kolmogorov, 22, 300, 465  
 Mallows, 466  
 Prohorov, 465  
 uniform/supremum, 23, 475  
 variational, 464  
 minimum contrast (*MC*) estimates, 299, 325, 326, 328, 331, 332  
 minimum Cramér-von Mises distance estimation, 318  
 minimum distance estimates, 22, 41, 299, 318  
 missing  
   covariates, 146  
   data, 143, 271  
   at random, 144  
 mixture models, 261  
   exponential, 134, 261, 264, 337, 399, 401  
   Gaussian location, 262, 264  
   Gaussian scale, 262, 264  
   Poisson, 262, 264  
   two-sample scale, 134  
   uniform, 261, 263, 266  
 model, 1 *See also* list of examples  
   accelerated time, 147  
   Bernoulli pairs, 128  
   biased sampling, 8, 87, 113, 240, 340, 346, 362, 383  
   Box-Cox, 86, 153, 172, 292  
   censored linear regression, 9, 284, 388, 394, 420  
   choice-based sampling, 8, 114  
   Clayton-Cuzick, 170, 173, 292, 296  
   Clayton-Oakes, 8, 156, 173, 295  
   constraint defined, 9, 53, 68, 192, 198, 221  
   copula, 155–158, 173, 174, 204, 229, 295  
   Cox proportional hazards, 1, 9, 77, 82, 86, 88, 153, 171, 217, 292, 305, 320, 330, 335, 382, 420  
   elliptic, 89, 96, 305, 307, 399, 401  
   errors in variables, 7, 86, 127, 135, 305  
   exponential lifetime, 84  
   exponential mixture, 134, 261, 264, 337, 399, 401  
   group, 83, 88, 234, 391  
   Has'minskii-Ibragimov, 139, 147, 263, 267, 308, 348, 372, 392  
   linear, 89, 95, 328  
   linear regression, 6, 35, 147, 153, 158, 329  
   location, 1, 5, 109  
   location and scale, 29, 100, 112, 404  
   lognormal, 15  
   logistic regression, 85, 111, 113, 118, 329  
   measurement, 5, 84  
   missing data, 143, 271  
   mixture, 7, 85, 101, 125, 141, 204, 261, 305, 308, 391  
   Neyman-Scott, 127, 133  
   nonlinear regression, 104  
   nonparametric, 1, 2, 52, 144, 229, 338, 339  
   normal convolution, 309  
   paired comparisons, 5  
   parametric, 11, 13, 42, 83, 88  
   regression, 85, 103, 106, 125, 239, 305, 347, 353, 369, 383, 399, 401, 403  
   regression-transformation, 154, 319, 380  
   Reiersøl, 128, 138, 308  
   selection bias, 9  
   semiparametric, 1, 2, 87, 88, 229, 234  
   simultaneous equation, 125, 127  
   symmetric location, 53, 55, 58, 75, 235, 304, 305, 343, 398, 400, 403  
   transformation, 153, 204, 292, 294, 380, 382, 403  
   translation, 15, 292  
   truncated regression, 8, 115, 119, 253  
   two-sample, 99, 101, 134, 404  
 monotone density, 261  
 monotonicity constraints, 261, 269  
 multinomial distribution, 12, 33, 84  
 Nelson-Aalen estimator, 276, 358  
 Newton-Leibnitz theorem, 459  
 Newton-Raphson iteration, 43, 316, 318, 321, 517  
 Neyman-Scott models, 127, 133  
 nonparametric, 339  
   core model, 144  
   group model, 229, 319  
   maximum likelihood estimation (NPMLE), 126, 252, 266, 267, 339–344, 380, 383, 384, 388, 392  
   model, 1, 2  
 nonsingular matrix, 12  
 norm, 336, 349, 356, 414  
   Banach, 336  
   Sobolev, 347, 350, 366  
   supremum, 184, 191, 199  
 nuisance parameter, 1, 2, 27, 74, 301  
 null space, 417  
 odds ratio, 114, 128  
 one-step estimates, 43, 316, 371, 381  
 operator  
   adjoint, 416  
   bounded, 415  
   compact, 440  
   efficient influence, 178, 210  
   Hilbert-Schmidt, 208, 440  
   idempotent, 426  
   influence, 177, 180  
   information, 206, 210  
   inverse, 418  
   inverse information, 179, 211  
   linear, 415, 426

- $L$ -, 78, 249, 276, 420
- null-space, 202, 417
- $R$ -, 78, 249, 276, 420
- range, 202, 417, 419
- score, 76, 77, 202, 210
- self-adjoint, 418, 426
- square root of, 418
- transpose, 416
- optimality theorem, 26, 27, 182
- orthocomplement, 30
- orthogonal
  - complement, 417, 418, 425
  - projection, 425
- orthogonality, 238, 425
- paired observations, 101
- parameter, 2
  - abstract, 298
  - $B$ -valued, 177, 179
  - Euclidean, 1, 2
  - general, 177
  - incidental, 126
  - infinite-dimensional, 177, 221, 356
  - location, 1
  - nuisance, 1, 2, 27, 74, 301
  - of interest, 27
  - pathwise differentiable, 57, 60, 177, 456
  - scale, 1
  - structural, 126
- parametric
  - model, 11, 13, 42, 83, 88
  - submodel, 46, 70
- partial likelihood, 10, 79
- partial likelihood estimator, 79, 80, 331
- partial spline regression, 107, 110, 240
- pathwise differentiable, 57, 60, 177, 456
- penalized maximum likelihood, 301, 339, 344, 347, 348, 353, 369, 392
- perfect functions, 485
- periodic function regression, 107, 110, 240
- Poisson mixture model, 262, 264
- polar decomposition, 418
- portmanteau theorem, 478
- Pratt's theorem, 479, 489
- predictable variation, 246, 422, 435
- preliminary estimator, 41, 146, 395, 396, 402
- product limit estimator, 247, 276, 358, 389
- product measures
  - contiguity, 498, 504
- Prohorov's theorem, 24, 472, 480, 481
- projection, 425
  - of influence functions, 31, 66
  - of score functions, 31, 71
  - operator, 425
  - pursuit, 87, 107, 109, 240, 305
  - theorem, 425
- proportional hazards model
  - Cox, 1, 9, 77, 82, 86, 88, 153, 171, 217, 293, 305, 320, 330, 335, 382, 420
  - exponential, 86
  - Pythagorean theorem, 425
  - quadratic loss, 26
  - range (of an operator), 201, 417, 419
  - ranks, 320
  - rate of convergence, 176, 335, 336
    - $\sqrt{n}$ -, 176, 336
    - not  $\sqrt{n}$ -, 47, 48, 176, 207, 266, 335
  - reduced sample estimators, 252, 282
  - regression
    - biased sampling, 115, 383
    - censored, 147, 284, 388
    - linear, 6, 35, 147
    - logistic, 85, 111
    - missing observations, 146
    - models, 85, 103, 106, 125, 239, 305, 347, 353, 369, 383, 399, 401, 403
    - nonlinear, 104
    - nonparametric, 347, 353, 369
    - partial spline, 107, 110, 240
    - periodic function, 107
    - projection pursuit, 107, 109
    - transformation, 154, 319, 380
    - truncated, 8, 115, 119, 240, 253
  - regular
    - estimate/estimator, 17, 21, 39, 179, 180, 356
    - Gaussian, 46, 357
    - locally, 21, 46, 180, 356
    - model, 11, 12
    - parametric (sub)model, 2, 11, 13, 46, 70, 395
    - parametrization, 12, 458
    - point, 12
    - uniformly, 18, 20
    - uniformly Gaussian, 18, 20, 24, 46
    - weakly, 181
    - weakly  $B_0^*$ , 356
  - regularity
    - preservation, 461
  - regularization, 343, 344, 348
  - regularized (NP)MLE, 301, 339, 343, 345, 348, 363, 372, 380, 392
  - Reiersøl model, 128, 138, 308
  - relatively compact, 471
  - reparametrization, 29
  - residual life function, 197
  - residual life operator, 78, 420
  - reverse cumulative hazard, 255
  - Riesz representation theorem, 58, 178, 416
  - risk set, 331
  - robustness, 4, 311
  - sampling,
    - biased, 86, 113, 240, 340

- sampling (*continued*)  
  stratified, 115, 122  
saturated tangent space, 252, 272, 282, 287  
scale parameter, 1  
Schur concave, 269, 271  
score function, 13, 70, 116  
  efficient, 3, 28, 30, 33, 70, 76, 90, 149, 158, 165  
  estimator of, 395, 396, 398, 400, 402  
score operator, 76, 77, 202, 210  
selection probability, 113, 240  
semiparametric, 1, 2, 87, 88, 229, 234  
  group model, 89, 234  
separable, 180, 184, 477  
sieves, 301, 337, 339, 343–345, 349, 362, 392  
Skorokhod theorem, 386, 472, 485  
Slutsky's theorem, 472, 485  
space  
  Banach, 414  
  dual, 415  
  Hilbert, 4, 5, 48, 415  
  normed linear, 414  
  null, 417  
  pre-Hilbert, 414  
  range, 417  
  spline, 350  
  Sobolev, 347, 350, 354, 366, 369  
  sum, 436, 441  
  tangent, 2, 16, 49, 50, 93, 16, 158, 165, 201, 272  
spectral decomposition, 203, 280  
spectral theorem, 448  
spline, 345, 346, 350  
state space, 86, 143  
Stein's estimator, 22  
Strassen's theorem, 465  
stratified sampling, 115, 122  
stratum, 8, 121, 240  
stratum weight (function), 113, 121, 240  
strongly complete,  $G$ -, 129  
strongly unimodal, 328  
Sturm-Liouville, 160, 162, 165, 167, 169, 173, 382  
subconvex, 182, 466  
submodel, 46, 70  
  least favorable, 61, 62, 76, 134, 135, 139, 395  
  regular parametric, 46, 70  
sufficient, 127, 130, 308  
sumspace, 436, 441  
superefficient, 21  
supremum norm, 184, 191, 199  
surface, 49  
symmetric  
  constructions, 402  
  coordinatewise, 229, 233  
  cyclically, 230, 234  
  dihedrally, 230, 234  
  location model, 53, 55, 58, 75, 235, 304, 305, 343, 398, 400, 403  
  permutation, 229, 233  
  rotationally, 230  
  spherically, 229, 232  
  symmetry, 308, 431  
tangent  
  set, 50  
  space, 2, 49, 50, 93, 116, 158, 165, 201, 272  
testing, 4  
Tietze extension theorem, 490  
tight, 47, 471, 479  
  asymptotically, 478  
  uniformly, 471  
transformations, 86, 153, 430  
translation equivariant, 23  
transpose operator, 416  
truncated regression, 8, 115, 119, 240, 253  
truncation, 87, 115, 252  
  random, 240, 247, 420  
two-sample, 99, 101, 134, 404  
  location and scale, 100  
  scale mixture, 101  
uniform asymptotic negligibility, 470, 500  
uniformly  
  asymptotic negligibility, 470, 500  
  asymptotically efficient, 24  
  asymptotically linear, 19, 20, 46  
  consistent, 18, 20  
  consistent,  $\sqrt{n}$ -, 18, 42, 44  
  efficient, 24, 44  
  Gaussian regular, 18, 20, 24, 46  
  integrable, 466, 468, 469, 503  
  regular, 18, 20  
U-statistic, 314, 330  
Vapnik-Chervonenkis, 200  
Vardi's model, 115, 122  
vectors, 12  
Vitali's theorem, 416, 469  
V-statistic, 334  
weak convergence, 468, 471, 475, 477  
weakly  
  asymptotically linear, 180  
  differentiable, 177  
  efficient, 182  
  regular, 181  
Weibull, 15  
weight function, 113  
Wronskian, 165, 173  
zero-one loss, 27