

ELEMENTARY PROOFS OF SOME BASIC FACTS CONCERNING ORDER STATISTICS

By

G. HAJÓS (Budapest), member of the Academy, and A. RÉNYI (Budapest),
corresponding member of the Academy

Let ξ_1, \dots, ξ_n denote a sample of size n from a population with the distribution function $F(x)$. By other words, ξ_1, \dots, ξ_n are mutually independent random variables with the common distribution function $F(x)$. Let ξ_1^*, \dots, ξ_n^* be the same set of variables, rearranged in increasing order of magnitude, i. e.

$$\xi_k^* = R_k(\xi_1, \dots, \xi_n),$$

where $R_k(x_1, \dots, x_n)$ denotes the k 'th term of the sequence obtained by rearranging the numbers x_1, \dots, x_n in increasing order of magnitude.

The present paper deals with the *order statistic* ξ_k^* . Some basic facts will be proved by simple methods. We aim expressively to avoid the calculus and at reduction of any calculation to possibly minimal extent. As consequence, our results may be easily checked by calculation in various different ways, which we are not intended to mention.

Our results seem us mostly to be known, though we did not find some of them explicitly in the literature. We endeavoured to give an elementary and systematic treatment of our subject. Accordingly, our paper may be of methodical interest. As to the literature we refer to the bibliography compiled by S. S. WILKS [3] and by the second named author [4].

1. In order to obtain distribution-free results, i. e. results independent of the distribution function $F(x)$, we introduce $\eta_k = F(\xi_k)$ ($k=1, \dots, n$). If we suppose that $y=F(x)$ is strictly increasing and continuous, the same holds for the inverse function $x=F^{-1}(y)$ and we have¹

$$\mathbf{P}(\eta_k < x) = \mathbf{P}(\xi_k < F^{-1}(x)) = F(F^{-1}(x)) = x \quad (0 \leq x \leq 1),$$

what shows that the variables η_1, \dots, η_n are *uniformly distributed* in the interval $(0, 1)$. Putting $\eta_k^* = F(\xi_k^*)$ ($k=1, \dots, n$) we have

$$\eta_k^* = F(\xi_k^*) = F(R_k(\xi_1, \dots, \xi_n)) = R_k(F(\xi_1), \dots, F(\xi_n)) = R_k(\eta_1, \dots, \eta_n).$$

Consequently $\eta_1^*, \dots, \eta_n^*$ are order statistics of a sample of size n from a population of uniform distribution in $(0, 1)$.

¹ $\mathbf{P}(A)$ denotes the probability of the event A .

Accordingly, we confine ourselves in what succeeds to the research of the order statistics η_{jk}^* . Our results may be interpreted as facts concerning the original order statistics ξ_k^* .

2. The variables η_{jk}^* are not independent, the relation $\eta_j^* \leq \eta_k^*$ (in case $j < k$) contradicts the independency.

The joint density function of the variables $\eta_1^*, \dots, \eta_n^*$ is

$$(1) \quad f(x_1, \dots, x_n) = n! \quad (0 \leq x_1 \leq \dots \leq x_n \leq 1).$$

As a matter of fact, if E denotes any measurable subset of the n -dimensional simplex defined by the inequalities $0 \leq x_1 \leq \dots \leq x_n \leq 1$, we have

$$\mathbf{P}((\eta_1^*, \dots, \eta_n^*) \in E) = \sum \mathbf{P}((\eta_{i_1}, \dots, \eta_{i_n}) \in E),$$

where the summation is extended over all permutations i_1, \dots, i_n of the indices $1, \dots, n$ and the density function of $(\eta_{i_1}, \dots, \eta_{i_n})$ is equal to 1 at any point (x_1, \dots, x_n) of the cube $0 \leq x_k \leq 1$ ($k = 1, \dots, n$).

Considering now the case that $\eta_{jk}^* = c_k, \dots, \eta_n^* = c_n$ ($2 \leq k \leq n$) are fixed, we state that $\eta_1^*, \dots, \eta_{k-1}^*$ are order statistics of a sample of size $k-1$ from a population of uniform distribution in the interval $(0, c_k)$. In fact, this is true if $\eta_1^*, \dots, \eta_{k-1}^*$ are furnished by any given $k-1$ variables out of η_1, \dots, η_n , since these $k-1$ variables are uniformly distributed, even within the cube $0 \leq x_i \leq c_k$ ($i = 1, \dots, k-1$). Thus, by (1), the joint density function of the variables $\eta_1^*, \dots, \eta_{k-1}^*$, under condition $\eta_{jk}^* = c_k, \dots, \eta_n^* = c_n$, is

$$(2) \quad f(x_1, \dots, x_{k-1} | c_k, \dots, c_n) = \frac{(k-1)!}{c_k^{k-1}} \quad (0 \leq x_1 \leq \dots \leq x_{k-1} \leq c_k).$$

By the same argument, if $\eta_1^* = c_1, \dots, \eta_k^* = c_k$ ($1 \leq k \leq n-1$) are fixed, $\eta_{k+1}^*, \dots, \eta_n^*$ are order statistics of a sample of size $n-k$ from a population of uniform distribution in the interval $(c_k, 1)$ and the joint density function of the variables $\eta_{k+1}^*, \dots, \eta_n^*$, under condition $\eta_1^* = c_1, \dots, \eta_k^* = c_k$, is

$$(3) \quad f(x_{k+1}, \dots, x_n | c_1, \dots, c_k) = \frac{(n-k)!}{(1-c_k)^{n-k}} \quad (c_k \leq x_{k+1} \leq \dots \leq x_n \leq 1).$$

Since (2) and (3) depend only on c_k , our statements hold also under the only condition $\eta_{jk}^* = c_k$, i. e. (2) and (3) give also the values of the functions $f(x_1, \dots, x_{k+1} | c_k)$ and $f(x_{k+1}, \dots, x_n | c_k)$, and the same holds under any restriction on the non-occurring variables. By the same argument, under condition $\eta_{jk}^* = c_k$, the sets of variables $(\eta_1^*, \dots, \eta_{k-1}^*)$ and $(\eta_{k+1}^*, \dots, \eta_n^*)$ are independent. By other words, order statistics form a *Markov chain*.²

3. The joint density function of the variables $\eta_{i+1}^*, \dots, \eta_k^*$ ($1 \leq i < k \leq n$) is

$$(4) \quad f_{ik}(x_{i+1}, \dots, x_k) = \frac{n!}{i!(n-k)!} x_{i+1}^i (1-x_k)^{n-1} \quad (0 \leq x_{i+1} \leq \dots \leq x_k \leq 1).$$

² A. N. KOLMOGOROFF [1] was the first to remark this.

Indeed, the joint density function of the variables $\eta_1^*, \dots, \eta_i^*, \eta_{i+1}^*, \dots, \eta_n^*$, under the condition

$$(C_1) \quad \eta_{i+1}^* = x_{i+1}, \dots, \eta_k^* = x_k,$$

is clearly given by

$$f(x_1, \dots, x_i, x_{k+1}, \dots, x_n | x_{i+1}, \dots, x_k) = \frac{f(x_1, \dots, x_n)}{f_{ik}(x_{i+1}, \dots, x_k)}.$$

On the other hand, since the variables $\eta_1^*, \dots, \eta_i^*$ are by the Markov chain property, under condition (C₁), independent of the variables $\eta_{k+1}^*, \dots, \eta_n^*$, we have

$$\begin{aligned} & f(x_1, \dots, x_i, x_{k+1}, \dots, x_n | x_{i+1}, \dots, x_k) = \\ & = f(x_1, \dots, x_i | x_{i+1}, \dots, x_n) \cdot f(x_{k+1}, \dots, x_n | x_1, \dots, x_k). \end{aligned}$$

Taking into account (1), (2) and (3), comparison of both statements establishes (4).

By (5), taking the special case $i+1 = k$ the density function of η_{ik}^* is

$$(5) \quad f_k(x_k) = \frac{n!}{(k-1)!(n-k)!} x_k^{k-1} (1-x_k)^{n-k} = \beta_{nk}(x),$$

i. e. order statistics η_{ik}^* have Beta-distribution.

The joint density function of any $\eta_{k_1}^*, \dots, \eta_{k_r}^*$ ($1 \leq k_1 < \dots < k_r \leq n$) variables is

$$f_{k_1, \dots, k_r}(x_1, \dots, x_r) = C_{k_1, \dots, k_r} x_1^{k_1-1} (x_2-x_1)^{k_2-k_1-1} \dots (x_r-x_{r-1})^{k_r-k_{r-1}-1} (1-x_r)^{n-k_r},$$

where

$$C_{k_1, \dots, k_r} = \frac{n!}{(k_1-1)!(k_2-k_1-1)! \dots (k_r-k_{r-1}-1)!(n-k_r)!}$$

and

$$0 \leq x_1 \leq \dots \leq x_r \leq 1.$$

The proof is given immediately by the above argument, if we consider instead of (C₁) the condition

$$(C_2) \quad \eta_{k_1}^* = x_1, \dots, \eta_{k_r}^* = x_r,$$

divide (0, 1) by x_1, \dots, x_r into $r+1$ subintervals, and take into account that, by the Markov chain property, the sets of variables η_i^* lying in these subintervals are independent under condition (C₂).

4. We define $\eta_0^* = 0, \eta_{n+1}^* = 1$, and introduce the variables

$$\delta_k = \eta_{k-1}^* - \eta_{k+1}^* \quad (k = 1, \dots, n+1).$$

Since $\delta_1, \dots, \delta_n$ are obtained from $\eta_1^*, \dots, \eta_n^*$ by a measure-preserving linear transformation, their joint density functions are equal at corresponding places. By corresponding transformation $y_1 = x_1, y_k = x_k - x_{k-1}$ ($k = 2, \dots, n$) of (1) the joint density function of the random variables $\delta_1, \dots, \delta_n$ is

$$(6) \quad g(y_1, \dots, y_n) = n! \quad (y_1 \geq 0, \dots, y_n \geq 0; y_1 + \dots + y_n \leq 1).$$

We conclude from (6) that the variables $\delta_1, \dots, \delta_n$ have the same distribution. As the distribution of η_k is symmetric with respect to the point $\frac{1}{2}$, this symmetry holds also for the joint distribution of $(\eta_1^*, \dots, \eta_n^*)$. Especially, $\delta_1 = \eta_1^*$ and $\delta_{n+1} = 1 - \eta_n^*$ have equal distributions. We draw the conclusion that the variables $\delta_1, \dots, \delta_{n+1}$ are equally distributed.

Their common density function is that of η_1^* , given by (5),

$$\beta_{n1}(x) = n(1-x)^{n-1} \quad (0 \leq x \leq 1).$$

Their common mean value is, because of $\delta_1 + \dots + \delta_{n+1} = 1$,

$$(7) \quad \mathbf{M}(\delta_k) = \frac{1}{n+1} \quad (k = 1, \dots, n+1).$$

5. The variables $\delta_1, \dots, \delta_{n+1}$ are not only equally distributed, but are also *equivalent variables*, i. e. their joint distribution remains invariant under any permutation of them. This is established by (6) for permutations of $\delta_1, \dots, \delta_n$, by the above mentioned symmetry for the permutation $(\delta_1, \dots, \delta_{n+1}) \rightarrow (\delta_{n+1}, \dots, \delta_1)$, and by successive application of these for any permutation of $\delta_1, \dots, \delta_{n+1}$.

In consequence, the distribution of the *difference*

$$\eta_{i+k}^* - \eta_i^* = \delta_{i+1} + \dots + \delta_{i+k} \quad (0 \leq i < i+k \leq n+1)$$

depends only on k , is therefore equal to that of η_k^* , has the density function (5),

and, by (7), the mean value $\frac{k}{n+1}$.

Especially, the *range* $\mathcal{A} = \eta_n^* - \eta_1^*$ has the density function

$$\beta_{n, n-1}(x) = n(n-1)x^{n-2}(1-x) \quad (0 \leq x \leq 1)$$

and the mean value

$$\mathbf{M}(\mathcal{A}) = \frac{n-1}{n+1}.$$

6. As previously stated, under condition $\eta_k^* = c_k, \dots, \eta_n^* = c_n$ the variable η_i^* ($1 \leq i < k \leq n$) is the i 'th order statistic of a sample of size $k-1$ from a population of uniform distribution in the interval $(0, c_k)$. Hence, the distribution of the *quotient* $\frac{\eta_i^*}{\eta_k^*}$ does not depend even on c_k , remains therefore unaltered if $\eta_k^*, \dots, \eta_n^*$ are not fixed, and has, by (6), the density function $\beta_{i+k-1, k}(x)$. Thus, both differences and quotients of order statistics η_k^* have Beta-distribution.

The variables $\frac{\eta_k^*}{\eta_{k+1}^*}$ ($k = 1, \dots, n$) are *mutually independent*. Indeed, by above statement, the distribution of $\frac{\eta_k^*}{\eta_{k+1}^*}$ does not depend on the values of

$\eta_{k+1}^*, \dots, \eta_n^*$, is therefore independent of $\frac{\eta_{k+1}^*}{\eta_{k+2}^*}, \dots, \frac{\eta_{n-1}^*}{\eta_n^*}, \frac{\eta_n^*}{\eta_{n+1}^*} = \eta_n^*$.

Moreover the independent variables

$$\zeta_k^* = \left(\frac{\eta_k^*}{\eta_{k+1}^*} \right)^k \quad (k \leq 1, \dots, n)$$

are uniformly distributed in the interval $(0, 1)$.³ As a matter of fact, the distribution function of η_n^* is

$$\mathbf{P}(\eta_n^* < x) = \prod_{k=1}^n \mathbf{P}(\eta_k < x) = x^n \quad (0 \leq x \leq 1).$$

Correspondingly, since $\frac{\eta_k^*}{\eta_{k+1}^*}$ is the k 'th order statistic of a sample of size k from a population of uniform distribution in $(0, 1)$, we have

$$\mathbf{P}\left(\frac{\eta_k^*}{\eta_{k+1}^*} < x\right) = x^k.$$

Whence

$$\mathbf{P}(\zeta_k^* < x) = \mathbf{P}\left(\frac{\eta_k^*}{\eta_{k+1}^*} < \sqrt[k]{x}\right) = x,$$

establishing our statement.

7. We introduce the random variables

$$(8) \quad \mathcal{G}_k = -\ln \zeta_k = -k \ln \frac{\eta_k^*}{\eta_{k+1}^*} \quad (k = 1, \dots, n).$$

These are, according to our preceding result, mutually independent and equally distributed, with the distribution function

$$\mathbf{P}(\mathcal{G}_k < x) = \mathbf{P}(\zeta_k > e^{-x}) = 1 - e^{-x} \quad (x \geq 0),$$

i. e. have exponential distribution with mean value 1.

From equations (8) we get

$$(9) \quad \ln \eta_k^* = -\sum_{j=k}^n \frac{\mathcal{G}_j}{j} \quad (k = 1, \dots, n).$$

Consequently, the *logarithms* of the order statistics η_k^* form not only a Markov chain, but also an *additive chain*, i. e. they are consecutive partial sums of a sequence of mutually independent random variables.

We expressed by (9) the order statistics η_k^* as simple functions of independent and equally distributed random variables. Starting from this fact, *limit theorems* on order statistics may be obtained, by means of the central limit theorem, in a simple and straightforward way. This has been shown by the second named author [4].

(Received 4 May 1954)

³ This was proved by S. MALMQUIST [2]; his proof is rather complicated.

Bibliography

- [1] A. N. KOLMOGOROFF, Sulla determinazione empirica di una legge di distribuzione, *Gior. d. Att.*, **4** (1933), pp. 83—91.
- [2] S. MALMQUIST, On a property of order statistics from a rectangular distribution, *Skand. Aktuerietidsskrift*, **33** (1950), pp. 214—222.
- [3] S. S. WILKS, Order statistics, *Bull. Amer. Math. Soc.*, **54** (1948), pp. 6—50.
- [4] A. RÉNYI, On the theory of order statistics, *Acta Math. Acad. Sci. Hung.*, **4** (1953), pp. 191—231.

ЭЛЕМЕНТАРНОЕ ДОКАЗАТЕЛЬСТВО НЕКОТОРЫХ ОСНОВНЫХ ФАКТОВ ТЕОРИЙ ВАРИАЦИОННЫХ РЯДОВ

Г. ХАЙОШ и А. РЕНЬИ (Будапешт)

(Резюме)

В работе, которая имеет преимущественно методический характер, дается систематическое и элементарное изложение некоторых основных фактов теорий вариационных рядов. Пусть $\eta_1, \eta_2, \dots, \eta_n$ — независимые случайные величины, которые равномерно распределены в интервале $(0, 1)$, пусть $\eta_1^* \leq \eta_2^* \leq \dots \leq \eta_n^*$ — те же величины, расположенные в возрастающем порядке. Доказывается, между прочим, очень просто, что величины η_k^* ($k = 1, 2, \dots, n$) образуют цепь Маркова (теорема А. Н. Колмогорова, см. [1]), более того, что величины $\ln \eta_k^*$ ($k = 1, 2, \dots, n$) образуют аддитивную цепь Маркова, так как случайные величины $\left(\frac{\eta_k^*}{\eta_{k+1}^*}\right)^k$ ($k = 1, 2, \dots, n; \eta_{n+1}^* \equiv 1$) независимы и равномерно распределены в интервале $(0, 1)$ (теорема С. Малмквиста, см. [2]). С помощью этих фактов возможно доказательство многих предельных теорем теорий вариационных рядов на центральную предельную теорему теории вероятностей (см. [4]).