

## EGY INFORMÁCIÓELMÉLETI PROBLÉMÁRÓL

RÉNYI ALFRÉD

### Bevezetés

E dolgozatban egy olyan probléma matematikai modelljét állítjuk fel és oldjuk meg, amely a legkülönbözőbb tudományokban, valamint számos egymástól távoleső gyakorlati tevékenység során szinte mindennapos. A problémát először egy konkrét példán keresztül ismertetjük; azután rámutatunk néhány más szituációra is, ahol olyan probléma lép fel, amely ugyanarra a matematikai modellre vezet.

Tegyük fel, hogy egy sok alkatrészből álló, bonyolult, összetett berendezésben (pl. egy autó motorjában, egy elektromos hálózatban, egy elektronikus számológépben vagy más sok elemből felépített komplex berendezésben) valami hiba lép fel; a hiba kijavítására törekvő szerelő először meg kell, hogy találja a hibát; csak ez után kerülhet sor a hiba kiküszöbölésére. Tudvalevő, hogy a hiba megkeresése gyakran sokkal több időt és fáradságot igényel, mint a hiba kijavítása. Mi itt csak a hiba megkeresésének folyamatával foglalkozunk.

Jelöljük  $\mathbf{H}$ -val az összetett berendezés önálló egységét alkotó részeinek halmazát; jelöljék  $x_1, x_2, \dots, x_n$  a  $\mathbf{H}$  halmaz elemeit. A hiba tehát vagy  $x_1$ -ben, vagy  $x_2$ -ben,  $\dots$ , vagy  $x_n$ -ben van. Az egyszerűség kedvéért tegyük fel, hogy az  $x_1, x_2, \dots, x_n$  alkatrészek között csak egyetlen egy hibás van. Attól, hogy milyen jellegű a hiba (pl. kiégett egy elektroncső, valahol zárlat van, stb.) tekintsünk el. Minket csak az érdekel, hogy *melyik* a hibás alkatrész. Hogyan lát a szerelő munkához? Megpróbálja működtetni a berendezés egyes részrendszereit. Ha a berendezés egy részrendszere, amely pl. az  $x_{i_1}, x_{i_2}, \dots, x_{i_r}$  elemekből áll, nem működik, akkor a hibás alkatrész az  $x_{i_1}, x_{i_2}, \dots, x_{i_r}$  elemek között van; ha ez a részrendszer hibátlanul működik, akkor a hibás alkatrész nem az  $x_{i_1}, x_{i_2}, \dots, x_{i_r}$  elemek között, hanem a  $\mathbf{H}$  halmaznak az  $x_{i_1}, x_{i_2}, \dots, x_{i_r}$  sorozathoz nem tartozó elemei között van. Akár működik tehát a kiválasztott részrendszer, akár nem, a hibát sikerült bizonyos mértékig lokalizálni, a hiba tekintetében számbajövő alkatrészek  $\mathbf{H}$  halmazát ezen halmaz egy részhalmazára leszűkíteni. Tekintve, hogy ez a részhalmaz általában még mindig sok alkatrészt tartalmaz, a szerelő ezután egy másik részrendszerrel próbálkozik. Kellő számú részrendszer megvizsgálása alapján a hibás alkatrész egyértelműen lokalizálható. A részrendszerek kiválasztását a berendezés konstrukciója bizonyos tekintetben korlátozza; ezen korlátokon belül azonban a szerelő tetszésétől (mondhatjuk bátran: intuíciójától) függ, hogy mely részrendszereket vizsgálja meg és milyen sorrendben.

Mivel feltehetőleg egy másik szerelő másként járt volna el, úgy tekinthetjük, hogy a megvizsgált részrendszerek kiválasztása véletlenszerűen tör-

ténik. Előfordulhat az is, hogy az egyes részvizsgálatok eredménye téves és így félrevezető. (Pl. a vizsgálat céljából összekapcsolt alkatrészek között létesített kontaktus laza; a vizsgálatához felhasznált mérőműszer maga sem kifogástalanul működik stb.). Így tehát a hibakeresés folyamatára a következő matematikai modellt állíthatjuk fel:

Meg akarjuk határozni az  $x_1, x_2, \dots, x_n$  elemekből álló  $\mathbf{H}$  halmaz egy előttünk ismeretlen elemét — jelöljük ezt  $x$ -szel. E célból taláalomra kiválasztjuk a  $\mathbf{H}$  halmaz  $H_1, H_2, \dots, H_k$  részhalmozait. Tegyük fel, hogy a  $H_1, H_2, \dots, H_k$  halmazok választásai egymástól függetlenek és a  $H_1, H_2, \dots, H_k$  halmazok mindegyike ugyanakkora (tehát  $\frac{1}{2^n}$ ) valószínűséggel lehet azonos a  $\mathbf{H}$  halmaz

bármely részhalmozásával. Miután a  $H_j$  halmazt megválasztottuk, választ kapunk arra a kérdésre, hogy a keresett  $x$  elem hozzátartozik-e a  $H_j$  halmazhoz. Az így kapott válaszok azonban nem mindig felelnek meg a valóságnak; tegyük fel, hogy arra a kérdésre, hogy a  $H_j$  halmaz tartalmazza-e a keresett  $x$  elemet,  $\beta$  valószínűséggel ( $\frac{1}{2} < \beta \leq 1$ ) a valóságnak megfelelő választ, és  $1 - \beta$  valószínűséggel hamis választ kapunk. Tegyük fel továbbá, hogy az egyes kérdésekre kapott válaszok helyes vagy hibás volta egymástól független. A matematikai probléma mármost a következő: ha  $n$  és  $\beta$  értéke ismert, és elő van írva egy  $\alpha$  szám ( $0 < \alpha < 1$ ), milyen nagyra kell választani  $k$  értékét (tehát hány taláalomra választott részhalmozatot kell megvizsgálni) ahhoz, hogy a kapott (részben helyes, részben téves) válaszokból legalább  $\alpha$  valószínűséggel identifikálni tudjuk a keresett  $x$  elemet?

Mielőtt a fentiekben megfogalmazott matematikai probléma megoldását megadnánk, megemlítünk néhány más, az összetett berendezésekben fellépő hiba keresésétől látszólag nagyon is távoleső szituációt, amely ugyanerre a matematikai modellre vezethető vissza. Természetesen a felsorolt példák mindegyike rendelkezik bizonyos sajátos vonásokkal, amelyek a fentemlített leegyszerűsített modellben nem tükröződnek. Egy-egy konkrét szituáció leírására a matematikai modell természetesen megfelelően módosítható. Mi itt első lépésként csak a fentiekben ismertetett leegyszerűsített modellel foglalkozunk, mivel az a szóbanforgó, egymástól sok szempontból különböző konkrét szituációk bizonyos közös vonásait ragadja meg.<sup>1</sup> További vizsgálatok tárgyát kell, hogy képezze a fentemlített alapmodell olyan módosításainak vizsgálata, amely révén egy-egy konkrét szituáció sajátosságait pontosabban lehet visszaadni.<sup>2</sup>

Második példaként tekintsük az orvosi diagnózis példáját. Amikor egy orvos egy betegről diagnózist állít fel, lényegében hasonló módon jár el, mint a szerelő, aki a hibát keresi. E példában célszerűbb a  $\mathbf{H}$  halmaz elemeinek az összes lehetséges betegségeket tekinteni (és nem pl. a beteg egyes szerveit, hiszen számos betegség nem egy szervnek, hanem az egész szervezetnek a megbetegedése). A beteget diagnosztizáló orvos úgy jár el, hogy a betegen bizonyos vizsgálatokat végez (megméri a hőmérsékletet, vérnyomást, véréjszűlyedést,

<sup>1</sup> Az [1], [2] és [3] dolgozatokban már foglalkoztunk a jelen dolgozat tárgyát képező probléma azon speciális esetével, amikor  $\beta = 1$ , vagyis amikor a kapott válaszok kivétel nélkül helyesek, továbbá ezen speciális eset másirányú általánosításával. Jelen dolgozatban éppen az a novum, hogy megengedjük azt is, hogy a válaszok egy része hibás legyen.

<sup>2</sup> A [3] dolgozatban a  $\beta = 1$  esetre nézve már foglalkoztunk a modellel olyan módosításával, amelynél a különböző részhalmozatok nem egyenlő valószínűséggel kerülnek kiválasztásra. Ez a  $\beta < 1$  esetben is megvizsgálandó.

kopogtatja a tüdejét, megnézi a torkát, különböző laboratóriumi vizsgálatokat végeztet, stb.). Az egyes vizsgálatok eredményeképpen a számításba jövő betegségek halmaza egyre jobban leszűkül. Itt is számolni kell azzal, hogy egyes részvizsgálatok eredményei tévesek és félrevezetőek. Itt nemcsak arról van szó, hogy pl. a laboratóriumi vizsgálat során téves műszerleolvasás történhet, hanem arról is, hogy vannak olyan tünetek, amelyek bizonyos betegségnél legtöbbször fellépnek, kivételes esetekben azonban hiányozhatnak. A fent felállított matematikai modell az orvosi diagnózis-felállítás folyamatának kétségtelenül csak nagyon leegyszerűsített képe, azonban a szóbanforgó folyamat bizonyos vonásait helyesen tükrözi.

Harmadik példaként vizsgáljuk meg, hogy hogyan jár el a kémikus, amikor valamilyen ismeretlen anyag kémiai összetételét kívánja megállapítani. Az egyszerűség kedvéért gondoljunk a kvalitatív kémiai analízisre, mivel a kvantitatív analízis olyan problémákat vet fel, amelyek bonyolultabb matematikai modellre vezetnek. Ez esetben a  $H$  halmaz elemeit az összes számbajövő vegyületek képezik. Az ismeretlen anyag összetételének megállapítása céljából a kémikus az anyagot bizonyos vizsgálatoknak veti alá (pl. bizonyos savak, lúgok, stb. hatásának teszi ki, lakmusz-papír-próbát végez stb.). Minden egyes részvizsgálat eredménye szűkíti a lehetőségek halmazát, míg végül kellő számú vizsgálat után a lehetőségek számát sikerül egyre redukálni. Nyilvánvaló, hogy e példában sem elhanyagolható az a lehetőség, hogy egy-egy részvizsgálat eredménye hibás.

A kvalitatív kémiai analízis példája alkalmas arra is, hogy rámutassunk, hogy a gyakorlatban sokszor nem egyetlen ismeretlent, hanem egyidejűleg több ismeretlent kell meghatározni; hiszen a megvizsgálandó anyag legtöbbször nem homogén, hanem több különböző anyagból áll. A  $\beta = 1$  speciális esetre vonatkozólag a [3] dolgozatban a probléma ilyenirányú általánosítását is megvizsgáltuk. A szimultán meghatározás problémájára a  $\beta < 1$  esetben más alkalommal kívánunk visszatérni.

A kémiai analízis példáján keresztül rámutathatunk a modell egy másik általánosítási lehetőségére is. Egyes részvizsgálatok ugyanis a számításbajövő lehetőségek halmazát nem két, hanem kettőnél több részhalmazra bontják, vagyis az egyes adalékok nem egyszerű dichotómiák, hanem általánosabb felbontások. A  $\beta = 1$  esetet illetőleg [3]-ban ezzel az általánosabb modellel is foglalkoztunk. Itt egyelőre az egyszerű dichotómiákra szorítkozunk; az általános osztályozás modelljét a  $\beta < 1$  esetben is érdemes volna azonban megvizsgálni.

Negyedik példaként tekintsük egy vizsgálóbíró munkáját, aki egy bűntény elkövetőjét akarja megtalálni. Ez esetben a halmaz a gyanúba kerülő egyénekből áll, míg részvizsgálatoknak ez esetben az egyes tanúk kihallgatásai, a bűnjelek megvizsgálása, és más vizsgálatok (pl. ujjlenyomat-vizsgálat, stb.) tekintendők. Mindezek az adatok egyre szűkebbre szorítják a kört a valódi bűnös körül. Az, hogy ez esetben az egyes adatok (pl. tanúvallomások) hitelessége gyakran kérdéses, teljesen nyilvánvaló.

E példák számát szinte korlátlanul lehetne folytatni. Befejezésül még csak egy példát említünk meg, amely mintegy középúton van a fentebb említett gyakorlati problémák és ezek matematikai modellje között, mégpedig az ún. Bar-Kochba játék példáját. Ez a játék ugyanis már maga is mintegy leegyszerűsített modellje a fentebb említett problémáknak. A Bar-Kochba játékban tudvalevőleg két játékos vesz részt — nevezzük őket  $A$ -nak és  $B$ -nek.

Az  $A$  játékos gondol valamire, a  $B$  játékos pedig arra törekszik, hogy kitalálja, hogy  $A$  mire gondolt. E célból kérdéseket tehet fel  $A$ -nak, azonban csak olyan kérdés van megengedve, amely igen-nel vagy nem-mel megválaszolható. A kérdésekre kapott válaszokból kell  $B$ -nek kitalálnia, hogy  $A$  mire gondolt. Ha  $\mathbf{H}$ -val jelöljük azon „dolgok” (személyek, tárgyak, fogalmak stb.) halmazát, amelyekre az  $A$  játékos gondolhatott, minden kérdés, amit  $B$  feltehet, megfogalmazható olyan alakban is, hogy hozzátartozik-e a gondolt „dolog” a  $\mathbf{H}$  halmaz egy bizonyos  $H$  részhalmazához. Mint már [3]-ban rámutattunk, a kérdések kiválasztása még gyakorlott játékos esetében is bizonyos mértékig véletlenszerűnek tekinthető. Míg [3]-ban csak a  $\beta = 1$  esetet tárgyaltuk, vagyis feltettük, hogy az  $A$  játékos minden kérdésre a valóságnak megfelelően válaszol, a Bar-Kochba játék tényleges lefolyásának jobban megfelel az hipotézis, hogy a válaszok bizonyos százaléka (a kérdés félreértése, vagy bizonyos tények ismeretének hiánya folytán) téves. (Lehetséges persze a Bar-Kochba játéknak egy olyan variánsa is, ahol az  $A$  játékosnak jogában áll a kérdések egy részére szándékosan téves választ adni, (pl. azon megszorítással, hogy a valóságnak meg nem felelően megválaszolt kérdések száma a játék során soha nem haladhatja meg az összes, az adott időpontig feltett kérdések számának pl. az ötödrészét).<sup>3</sup>

Hangsúlyozni szeretnénk, hogy a tárgyalt probléma mind gyakorlati, mind pedig elméleti szempontból elsősorban akkor érdekes, ha  $n$  (a  $\mathbf{H}$  halmaz elemeinek száma) igen nagy szám. A következőkben ezt mindig fel fogjuk tenni.

Végül még csak azt szeretnénk hangsúlyozni, hogy olyan esetekben, amikor egy komplex berendezésben fellépő hiba keresését, az orvosi diagnózist, a kémiai analízist stb. képzett szakember végzi, az általa követett eljárás erősen eltér az általunk választott modelltől, amennyiben az egyes részvizsgálatok megválasztása csak részben függ a véletlentől, nagyobb részben viszont azt az illető szakember szaktudása, tapasztalatai, intuíciója szabják meg. Ha azonban az említett vagy azokhoz hasonló feladatok elvégzésére készített kibernetikai berendezésre gondolunk, annak működése már aligha képzelhető el másképpen mint az általunk választott modellnek megfelelően. Egy hibakereső gép konstruálása viszont a kibernetika mai állása mellett teljes mértékben a lehetőségek határán belül van, bár tudomásom szerint ezideig ilyen gép nem készült. Vizsgálataink eredményei felhasználásra kerülhetnek egy ilyen hibakereső kibernetikai berendezés konstrukciójánál.

Az 1. §-ban a fentebb megfogalmazott matematikai probléma megoldását adjuk meg. A 2. §-ban megmutatjuk, hogyan függ össze a vizsgált probléma a matematikai statisztika „diszkriminációs probléma” elnevezés alatt ismert problémakörével. A 3. §-ban viszont azt mutatjuk meg, hogyan függ össze az általunk vizsgált probléma az információelmélet szokásos kérdésfeltevésével (a zajos csatornán keresztül való információátvitel problémájával).

## 1. §. A probléma megoldása

E §-ban a következő problémával foglalkozunk. Legyen  $\mathbf{H}$  egy  $n$  elemű halmaz. Legyenek  $H_1, H_2, \dots, H_k$  a  $\mathbf{H}$  halmaz találmára, egymástól függetle-

<sup>3</sup> Természetesen ez a megszorítás nem teljesen felel meg a következőkben tárgyalt modell feltevéseinek, azonban nem is áll attól túl távol.

nül kiválasztott részhalmazai. Tegyük fel, hogy  $H_j$  ugyanakkora (tehát  $\frac{1}{2^n}$ ) valószínűséggel lehet azonos  $\mathbf{H}$  bármely részhalmazával.<sup>4</sup> Legyen  $x$  a  $\mathbf{H}$  halmaz egy ismeretlen eleme. Tegyük fel, hogy arra a kérdésre, hogy hozzátartozik-e  $x$  a  $H_j$  részhalmazhoz,  $\beta$  valószínűséggel helyes, és  $1 - \beta$  valószínűséggel hamis választ kapunk ( $\frac{1}{2} < \beta < 1$ ), mégpedig oly módon, hogy az egyes kérdésekre kapott válaszok helyes ill. hamis volta független a többi kérdésre kapott válasz helyes ill. hamis voltától. Kérdés: milyen nagyra kell  $k$  értékét választani ahhoz, hogy a kapott válaszokból legalább  $\alpha$  valószínűséggel meg tudjuk állapítani, hogy a  $\mathbf{H}$  halmaz melyik  $x$  eleméről van szó? ( $0 < \alpha < 1$ ).

Feltesszük, hogy a fenti problémában szereplő  $\beta$  és  $\alpha$  megadott rögzített számok, és elsősorban a keresett  $k$  értéknek  $n$ -től való függését vizsgáljuk, ha  $n \rightarrow +\infty$ . Annak valószínűségét, hogy az ismeretlen  $x$  elemet  $k$  válaszból meg lehet állapítani,  $P_{nk}$ -val fogjuk jelölni; valójában  $P_{nk}$  attól is függ, hogy a kapott válaszok kiértékelését milyen módszer alapján végezzük el. A következőkben a válaszok kiértékelésének egy plauzibilis konkrét módját fogjuk alkalmazni. Arra a kérdésre, hogy miért éppen ezt az eljárást választottuk, még a 2. §-ban visszatérünk.

Mielőtt eredményünket megfogalmazzunk, rá szeretnénk mutatni, hogy az a feltevés, hogy  $\frac{1}{2} < \beta \leq 1$  nem jelenti az általánosság megszorítását. Ugyanis a  $\beta = \frac{1}{2}$  esetben nyilvánvalóan a válaszok semmi információt nem nyújtanak, míg a  $0 \leq \beta < \frac{1}{2}$  eset visszavezethető az  $\frac{1}{2} < \beta \leq 1$  esetre oly módon, hogy minden választ az ellenkezővel helyettesítünk.

Vezessük be a következő jelölést:

$$(1.1) \quad I(\beta) = \beta \log_2 \frac{1}{\beta} + (1 - \beta) \log_2 \frac{1}{1 - \beta}.$$

Másszóval jelöljük  $I(\beta)$ -val a  $\beta$  és  $1 - \beta$  valószínűségekből álló kéttagú valószínűségeloszlás entrópiáját.  $I(\beta)$  felfogható, mint egy válasz helyes voltára vonatkozó bizonytalanság mértékszám. Ily módon egy-egy válasz legfeljebb  $1 - I(\beta)$  információt nyújt. Mivel az  $x$  elem meghatározásához  $\log_2 n$  információ szükséges, heurisztikus megfontolás szerint csak akkor remélhetjük, hogy a  $k$  válasz elegendő  $x$  meghatározására, ha  $k(1 - I(\beta)) \geq \log_2 n$ , vagyis, ha

$$(1.2) \quad k \geq \frac{\log_2 n}{1 - I(\beta)}.$$

A heurisztikus megfontolással nyert (1.2) alatti alsó korlát nincs is távol a helyes eredménytől, ugyanis érvényes a következő

**Tétel.** Ha  $\frac{1}{2} < \beta < 1$  és

$$(1.3) \quad k(n) = \frac{\log_2 n + y \sqrt{\log_2 n} + o(\sqrt{\log_2 n})}{1 - I(\beta)},$$

<sup>4</sup> Ez azt jelenti, hogy elvben  $H_j$  az üres halmaz vagy maga  $\mathbf{H}$  is lehet, továbbá azt, hogy a  $H_1, H_2, \dots, H_k$  halmazok között elvben ugyanaz a halmaz többször is előfordulhat. A továbbiakban tárgyalt nagyságrendi viszonyok mellett azonban ezen lehetőségek valószínűségei olyan elenyészően kicsinyek lesznek, hogy e lehetőségeket felesleges eleve kizárni; ha e lehetőségeket kizárnánk, ez az  $n \rightarrow +\infty$  határesetre nyert eredményeket egyáltalán nem befolyásolná.

ahol  $y$  tetszőleges rögzített valós szám, és  $I(\beta)$  az (1.1) által definiált mennyiség, akkor<sup>5</sup>

$$(1.4) \quad \lim_{n \rightarrow +\infty} P_{n,k(n)} = \Phi\left(\frac{y}{\sigma}\right),$$

ahol

$$(1.5) \quad \Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{u^2}{2}} du$$

és

$$(1.6) \quad \sigma = \sqrt{\frac{\beta(1-\beta)}{I(\beta)}} \log_2 \frac{\beta}{1-\beta}.$$

*Megjegyzés:* A tétel eredménye úgy is fogalmazható, hogy ha  $\alpha$  tetszőleges előírt valószínűség ( $0 < \alpha < 1$ ) és

$$(1.3') \quad k(n) = \frac{\log_2 n + \sigma \Phi^{-1}(\alpha) \sqrt{\log_2 n} + o(\sqrt{\log_2 n})}{1 - I(\beta)},$$

akkor

$$(1.4') \quad \lim_{n \rightarrow +\infty} P_{n,k(n)} = \alpha.$$

Ez a fogalmazás világossá teszi, hogy az  $\alpha$  valószínűségi szint megválasztásától  $k(n)$  csak viszonylag kevésé függ, hiszen a főtág (1.3')-ben független  $\alpha$ -tól.

**Bizonyítás.** Jelölje  $\tilde{H}_j$  azt a halmazt, amelyről a  $j$ -edik válasz azt állítja, hogy  $x$ -et tartalmazza. Tehát  $\tilde{H}_j = H_j$ , ha  $x \in H_j$  és a  $j$ -edik válasz helyes, vagy ha  $x \notin H_j$  és a  $j$ -edik válasz hamis, míg  $\tilde{H}_j = \bar{H}_j$ , ha  $x \notin H_j$  és a  $j$ -edik válasz helyes, vagy ha  $x \in H_j$  és a  $j$ -edik válasz hamis.<sup>6</sup> Jelöljék  $x_1, x_2, \dots, x_n$  a  $\mathbf{H}$  halmaz elemeit. Legyen  $\varepsilon_{hj} = 1$ , ha a  $j$ -edik válasz összhangban van azzal a hipotézissel, hogy  $x = x_h$  és  $\varepsilon_{hj} = 0$ , ha a  $j$ -edik válasz ellentmond annak a hipotézisnek, hogy  $x = x_h$  ( $h = 1, 2, \dots, n; j = 1, 2, \dots, k$ ).

Másszóval legyen

$$(1.7) \quad \varepsilon_{hj} = \begin{cases} 1 & \text{ha } x_h \in \tilde{H}_j \\ 0 & \text{ha } x_h \notin \tilde{H}_j. \end{cases}$$

Tegyük fel egy pillanatra, hogy  $x = x_1$ . Ez esetben, ha  $\delta_2, \delta_3, \dots, \delta_n$  tetszőleges, a 0 és 1 elemekből álló sorozat akkor

$$(1.8) \quad \mathbf{P}(\varepsilon_{1j} = 1, \varepsilon_{2j} = \delta_2, \varepsilon_{3j} = \delta_3, \dots, \varepsilon_{nj} = \delta_n) = \frac{\beta}{2^{n-1}}$$

és

$$(1.9) \quad \mathbf{P}(\varepsilon_{1j} = 0, \varepsilon_{2j} = \delta_2, \varepsilon_{3j} = \delta_3, \dots, \varepsilon_{nj} = \delta_n) = \frac{1-\beta}{2^{n-1}}.$$

<sup>5</sup> A kiértékelési módszer a bizonyítás során részletezendő alkalmas választása mellett.

<sup>6</sup> Itt  $\bar{H}_j$  a  $H_j$  halmaznak  $\mathbf{H}$ -ra vonatkozó kiegészítő halmazát jelöli.

Ugyanis  $\varepsilon_{1j} = 1$  akkor és csak akkor áll fenn, ha a  $j$ -edik válasz helyes, hiszen ha a  $j$ -edik válasz helyes, akkor  $\tilde{H}_j$  aszerint egyenlő  $H_j$ -vel vagy  $\bar{H}_j$ -sal, hogy  $x \in H_j$  vagy  $x \notin H_j$ , míg ha a  $j$ -edik válasz hamis, akkor ennek fordítottja igaz. Másrészt viszont, ha  $\delta_2, \dots, \delta_n$  tetszőleges rögzített, a 0 és 1 elemekből álló sorozat, akkor az  $\varepsilon_{2j} = \delta_2, \varepsilon_{3j} = \delta_3, \dots, \varepsilon_{nj} = \delta_n$  feltevések  $\varepsilon_{1j}$  értékével együtt egyértelműen meghatározzák a  $\tilde{H}_j$  halmazt. Ha  $\tilde{H}_j$  adva van, akkor  $H_j$  nem lehet más mint  $\tilde{H}_j$  vagy  $\bar{\tilde{H}}_j$ . Ilyen módon (1.8) és (1.9) azonnal következnek.

Jelölje  $E_{nk}$  azt az  $n \times k$  elemű mátrixot, amelynek  $h$ -adik sorának  $j$ -edik eleme  $\varepsilon_{hj}$  ( $h = 1, 2, \dots, n; j = 1, 2, \dots, k$ ) vagyis legyen

$$(1.10) \quad E_{nk} = \begin{pmatrix} \varepsilon_{11} & \varepsilon_{12} & \dots & \varepsilon_{1k} \\ \vdots & \vdots & & \vdots \\ \varepsilon_{n1} & \varepsilon_{n2} & \dots & \varepsilon_{nk} \end{pmatrix}.$$

(1.8) és (1.9) úgy interpretálhatók (figyelembevéve, hogy az  $E_{nk}$  oszlopait alkotó vektorok a  $H_j$  halmazok megválasztásának feltételezett függetlensége folytán függetlenek), hogy az  $E_{nk}$  véletlen mátrixot úgy származtathatjuk, hogy első sorának minden helyére egymástól függetlenül  $\beta$  valószínűséggel 1-est, ill.  $1-\beta$  valószínűséggel 0-t írunk, míg az összes többi helyre egymástól és az első sor kitöltésétől függetlenül  $\frac{1}{2}$  valószínűséggel írunk 1-est, ill. 0-t. Más szavakkal, az  $\varepsilon_{hj}$  ( $h = 1, 2, \dots, n; j = 1, 2, \dots, k$ ) valószínűségi változók teljesen függetlenek és

$$(1.11) \quad \mathbf{P}(\varepsilon_{1j} = 1) = \beta, \quad \mathbf{P}(\varepsilon_{1j} = 0) = 1 - \beta \quad (j = 1, 2, \dots, k),$$

míg ha  $2 \leq h \leq n$ , akkor

$$(1.12) \quad \mathbf{P}(\varepsilon_{hj} = 1) = \mathbf{P}(\varepsilon_{hj} = 0) = \frac{1}{2} \quad (j = 1, 2, \dots, k).$$

Ha  $x$  nem  $x_1$ -gyel, hanem pl.  $x_h$ -val egyenlő, akkor hasonlóképpen az  $E_{nk}$  mátrix  $h$ -adik sorának elemei lesznek  $\beta$  valószínűséggel 1-gyel és  $1-\beta$  valószínűséggel 0-val egyenlők, és az összes többi sor elemei  $\frac{1}{2}$  valószínűséggel lesznek 1-gyel, ill. 0-val egyenlők.

Legyen most

$$(1.13) \quad v_h = \sum_{j=1}^k \varepsilon_{hj}.$$

(Most nem tesszük fel, hogy  $x = x_1$ ).

Mármost a kapott válaszok kiértékelését a következő kézenfekvő módszer alapján végezzük: megnézzük, hogy az  $x = x_h$  hipotézisek közül ( $h = 1, 2, \dots, n$ ) melyiket támasztja alá a legtöbb válasz. Ha egyetlen ilyen  $h$  van, akkor a mellett döntünk, hogy  $x = x_h$ ; ha több ilyen  $h$  volna, akkor azt mondjuk, hogy a válaszok nem teszik lehetővé az egyértelmű döntést. Persze, még ha ezen eljárás egyértelmű is, akkor is lehet az így hozott döntés téves. Ha tehát  $P_{nk}$  jelöli annak a valószínűségét, hogy ezen döntési eljárással helyesen hatá-

rozzuk meg az ismeretlen  $x$  elemet, a fenti konvenció mellett, mely szerint  $x = x_1$ , a helyes döntés valószínűsége

$$(1.14) \quad P_{nk} = \mathbf{P}(v_1 > \max_{2 \leq h \leq n} v_h).$$

Figyelembevétel az  $E_{nk}$  véletlen mátrix elemeinek függetlenségét és az (1.11)—(1.12) összefüggéseket, azonnal adódik, hogy

$$(1.15) \quad P_{nk} = \sum_{l=1}^k \binom{k}{l} \beta^l (1-\beta)^{k-l} \left( \sum_{j=0}^{l-1} \binom{k}{j} \frac{1}{2^k} \right)^{n-1}.$$

(1.15)-ből a Moivre—Laplace tétel és a Stirling-formula segítségével teljesen kézenfekvő aszimptotikus megfontolások segítségével adódik a fenti tétel állítása. A számítás részleteit itt mellőzzük.

## 2. §. A probléma összefüggése a matematikai statisztika ismert problémáival

Az 1. tétel bizonyításából nyilvánvaló, hogy a tárgyalt probléma azonos a következő statisztikai problémával:  $n$  számú  $k$  elemű mintánk van; tudjuk azt, hogy az  $n$  minta független, közülük  $n-1$  egy olyan  $S\left(\frac{1}{2}\right)$  statisztikai sokaságból lett véve, amelyben az 1 és 0 számok valószínűsége  $\frac{1}{2}$ , míg az  $n$  minta közül az egyik egy olyan  $S(\beta)$  statisztikai sokaságból lett véve, amelyben az 1 ill. 0 számok valószínűsége  $\beta$  ill.  $1-\beta$  ( $\frac{1}{2} < \beta \leq 1$ ). Eldöntendő, hogy az  $n$  minta közül melyik származik  $S(\beta)$ -ből.

Másrészt azt is könnyen beláthatjuk, hogy az általunk választott döntési eljárás tulajdonképpen nem más, mint a maximum likelihood módszeren alapuló eljárás. Ugyanis azon feltevés mellett, hogy a  $h$ -adik minta származik  $S(\beta)$  ból (vagyis  $x = x_h$ ), az  $n$  minta együttes valószínűsége

$$(2.1) \quad W_h = \frac{\beta^{v_h} (1-\beta)^{k-v_h}}{2^{(n-1)k}}$$

és így, ha a maximum likelihood módszer szerint döntünk, úgy azt az  $x_h$ -t fogadjuk el az ismeretlen  $x$ -nek, amelyre  $W_h$  maximális, tehát azt, amelyre  $v_h$  maximális, és mi éppen így járunk el. Ez megvilágítja, hogy miért éppen ezt a kiértékelési eljárást választottuk. A problémát a következőképpen is fogalmazhatjuk. Jelölje  $\Pi_h$  azt a hipotézist, hogy az  $E_{nk}$  mátrix  $h$ -adik sora a  $S(\beta)$  sokaságból, a többi sorai a  $S\left(\frac{1}{2}\right)$  sokaságból vett minták. Az általunk vizsgált probléma úgy is jellemezhető, hogy eldöntendő, hogy a  $\Pi_1, \Pi_2, \dots, \Pi_n$  egymást kizáró hipotézisek közül melyik a helyes? E szerint a problémánk az ún. diszkriminációproblémák közé tartozik.

## 3. §. A probléma összefüggése az információelmélet ismert problémáival

Az 1. §-ban tárgyalt probléma a következőképpen is interpretálható. Legyen adva egy zajos csatorna, amelynél a lehetséges bemenő jelek  $x_1, x_2, \dots, x_n$ , míg a lehetséges kimenő jelek a  $\mathbf{H} = \{x_1, x_2, \dots, x_n\}$  halmaz összes részhalmazaiából állnak. Jelöljék ezeket  $H^{(1)}, H^{(2)}, \dots, H^{(2^n)}$ . Jelölje  $p_{hi}$  (ahol



$h = 1, 2, \dots, n; i = 1, 2, \dots, 2^n$ ) annak valószínűségét, hogy a  $H^{(i)}$  halmazt kapjuk kimenő jelként, feltéve, hogy az  $x_h$  jelet adtuk le, és tegyük fel, hogy

$$(3.1) \quad p_{hi} = \begin{cases} \frac{\beta}{2^{n-1}} & \text{ha } x_h \in H^{(i)} \\ \frac{1-\beta}{2^{n-1}} & \text{ha } x_h \notin H^{(i)}. \end{cases}$$

Tegyük fel, hogy  $k$ -szor egymásután leadjuk az  $x_h$  jelet; a feladat az, hogy a felvett  $H_1, H_2, \dots, H_k$  jelekből rekonstruáljuk  $x_h$ -t. (Vegyük észre, hogy ezen átfogalmazásnál még a  $\beta = 1$  esetnek is *zajos* csatornán keresztül való információátvitel felel meg!) Kérdés: hányszor kell az  $x_h$  jelet leadni (vagyis milyen nagyoknak kell választani  $k$  értékét), hogy a felvett jelekből  $x_h$  legalább  $\alpha$  valószínűséggel meghatározható legyen. A probléma ezen átfogalmazásának fő érdekessége abban áll, hogy elvezet a probléma egy igen messzeemenő általánosításához. A szóbanforgó kérdés, hogy tudniillik hányszor kell az  $x_h$  jelet leadni, hogy az előírt valószínűséggel dekódolható legyen, felvethető tetszőleges átmenet-valószínűségekkel bíró zajos-csatorna esetében is. E kérdés-feltevés, bár szorosan összefügg az információelmélet szokásos problematikájával, attól mégis eltér, és tudomásom szerint eddig nem képezte behatóbb vizsgálat tárgyát. Ezen általános információelméleti probléma tárgyalása azonban már túlnőne jelen dolgozat keretén.

(Beérkezett: 1962. február 27.)

IRODALOM

[1] RÉNYI, A.: „On random generating elements of a finite Boolean algebra.” *Acta Sci. Math. Szeged* **22** (1961) 75—81.  
 [2] RÉNYI, A.: „Statistical laws of accumulation of information.” *Bulletin of the International Statistical Institute, 33<sup>rd</sup> Session*, Paris, 1961, 1—7.  
 [3] RÉNYI, A.: „Az információ-akkumuláció statisztikus törvényszerűségeiről.” *MTA III. Osztályának Közleményei* **12** (1962) 15—33.

ОБ ОДНОЙ ПРОБЛЕМЕ ТЕОРИИ ИНФОРМАЦИИ

A. RÉNYI

Резюме

Вот типический пример положений, математический модель которого рассматривается в этой работе: Следует найти дефектную часть сложной сети. В таком положении, если число частей, которые могут быть дефектными, очень велико, возможный метод найти дефект следующий: Провер-

яются некоторые под-сети; если под-сеть не работает, то она содержит дефектную часть; если работает, то дефектная часть содержится в дополнительной под-сети. Если сеть состоит из  $n$  частей и проверяются подходящим образом выбранные  $k > \log_2 n$  подсети, то таким образом найдется дефектная часть. В работе рассматривается, сколько таких проверок нужно выполнять, если под-сети выбираются случайным образом, независимо друг от друга, так что при каждом выборе каждая возможная под-сеть выбирается с той же самой вероятностью. При этом предполагается, что проверка под-сетей не всегда, а только с вероятностью  $\beta$  ( $1/2 < \beta \leq 1$ ) дает правильный результат и с вероятностью  $1 - \beta$  получается ошибочный ответ на вопрос, содержит ли под-сеть дефект. Такой вопрос может иметь значение при конструировании автоматов для поиска дефекта.

Доказывается следующая теорема: *Если положим*

$$I(\beta) = \beta \log_2 \frac{1}{\beta} + (1 - \beta) \log_2 \frac{1}{1 - \beta}$$

*и если проверяются  $k = k(n)$  случайно и независимо выбранные под-сети сети, которая состоит из  $n$  частей, где  $k(n)$  зависит от  $n$  следующим образом:*

$$k(n) = \frac{\log_2 n + y \sqrt{\log_2 n} + o(\sqrt{\log_2 n})}{1 - I(\beta)}$$

*( $y$  фиксированное вещественное число), и  $P_{n,k}$  обозначает вероятность того, что из результатов проверок можно найти дефектную часть сети, то имеем*

$$\lim_{n \rightarrow +\infty} P_{n,k(n)} = \Phi\left(\frac{y}{\sigma}\right),$$

где

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{u^2}{2}} du$$

и

$$\sigma = \sqrt{\frac{\beta(1 - \beta)}{I(\beta)}} \cdot \log_2 \frac{\beta}{1 - \beta}.$$

Замечается, что имеются и другие процессы (кроме процесса нахождения дефекта), для которых выше описанная проблема может служить очень упрощенной моделью, так например процесс, с помощью которого врач найдет диагноз, или процесс химического качественного анализа и т. д.

Наконец указано, что вопрос может быть рассмотрен и как проблема математической статистики, или как частный случай следующей проблемы передачи информации: Пусть дан дискретный безпамятный канал с шумом. Передается тот же самый сигнал много раз; как можно из полученных сигналов разгадать, какой сигнал был передан?

ON A PROBLEM OF INFORMATION THEORY

A. RÉNYI

Abstract

A typical example of the situation, a mathematical model of which is discussed in this paper, is the following: one has to find that part of a complicated network which got out of order. In such a situation if the number of places of the defect is very large, a possible method to find the defect is to try whether certain sub-networks do work or not. If a sub-network (consisting of a subset of the parts of the total network) does not work, then it contains the defective part (it is supposed that only a single part is defective), while if it works, then the defective part is contained in the complementary subset. If the network consists of  $n$  parts and if  $k > \log_2 n$  suitably chosen sub-networks are tested in this way, the defective part can be determined. The question arises, how many sub-networks have to be tested, if the sub-networks are chosen completely at random. The answer to this question may be useful if one wants to construct defect-searching automata. The corresponding mathematical problem is solved under the supposition that the results of the tests of the sub-networks are correct with probability  $\beta$  only ( $1/2 < \beta \leq 1$ ) while with probability  $1-\beta$  they are false and thus misleading.

Let us put

$$(1) \quad I(\beta) = \beta \log_2 \frac{1}{\beta} + (1 - \beta) \log_2 \frac{1}{1 - \beta}.$$

It is shown that if the total number of parts any one of which may be defective is equal to  $n$ , and  $k = k(n)$  randomly and independently chosen sub-networks are tested (so that the probability of choosing any particular sub-network is the same) where

$$(2) \quad k(n) = \frac{\log_2 n + y \sqrt{\log_2 n} + o(\sqrt{\log_2 n})}{1 - I(\beta)}$$

where  $y$  is a fixed real number, then if  $P_{nk}$  denotes the probability that by evaluating the tests in a suitable way, the defective part can be identified, one has

$$(3) \quad \lim_{n \rightarrow +\infty} P_{n,k(n)} = \Phi\left(\frac{y}{\sigma}\right)$$

where

$$(4) \quad \Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{u^2}{2}} du$$

and

$$(5) \quad \sigma = \sqrt{\frac{\beta(1-\beta)}{I(\beta)}} \log_2 \frac{\beta}{1-\beta}$$

It is mentioned that the mathematical problem solved above may also be considered as a highly simplified model of other processes too, e. g. the process

applied by a physician trying to make a diagnose, or the work of a chemist who wants to analyse some material of an unknown composition, or even of a judge trying to find out the truth in some criminal case. The special case when all tests lead to reliable results (i. e.  $\beta = 1$ ) has been considered in previous papers ([1], [2], [3]) of the author. It is pointed out that from the point of view of statistics the problem is one of discrimination, while from the point of view of information theory the problem can be characterized as follows: a discrete noisy memoryless channel is given and the same symbol is transmitted several times; one has to determine this transmitted symbol from the received symbols. The special channel corresponding to the problem solved in the paper is such that the number of symbols which can be sent is  $n$ , that of symbols which can be received is  $2^n$  while the matrix of transition probabilities (the channel probability function) contains only two sorts of elements, so that in each row exactly one half of the elements belong to the first, and one half to the second sort and all columns are different.