

# AZ INFORMÁCIÓ-AKKUMULÁCIÓ STATISZTIKUS TÖRVÉNYSZERŰSÉGEIRŐL

Írta: RÉNYI ALFRÉD

## Bevezetés

A tudományos kutatás során és a gyakorlati életben egyaránt gyakran előfordul, hogy egy minket érdeklő tényállásra vonatkozó információ nem egyszerre jut teljes egészében birtokunkba, hanem azt részletekben szerezzük meg; az egyes, külön-külön keveset mondó adatokból — mintegy mozaikkövekből — lehet azután egy összképet kialakítani és a tényállást egészében felderíteni. *A teljes információt tehát részinformációk akkumulációja útján szerezzük meg.* Az információgyűjtés ezen folyamata során gyakran előfordul, hogy az egyes adatok részben átfedik egymást, és egy-egy új adat olyan információt is nyújt, amely az előzőleg szerzett adatokban már benne foglaltatik és csak részben ad valójában új felvilágosítást. Más szóval, az információ részletekben való gyűjtésénél általában bizonyos „redundancia” lép fel. Ezt a redundanciát csak úgy lehet kiküszöbölni, ha az információ gyűjtése igen jól átgondolt terv szerint történik; legtöbbször azonban ez olyan mértékben bonyolítja az információ összegyűjtését, hogy nem is célszerű erre törekedni.

Az információgyűjtés szóban forgó folyamata általában bizonyos véletlen elemeket is tartalmaz. Többnyire a véletlentől függ, hogy milyen részletadatok és milyen sorrendben jutnak az információt gyűjtő birtokába\*.

Hogy az elmondottakat konkrétabbá tegyük, vizsgáljunk meg néhány példát. Első példaként tekintünk a kvalitatív kémiai analízis problémáját. Valamely anyagról meg akarjuk állapítani annak kémiai összetételét. E célból több vizsgálatnak, próbának vetjük alá az illető (homogénnek feltételezett) anyag egy-egy adagját. Minden egyes próba szűkíti a lehetőségek halmazát, míg végül elegendő számú próba után a lehetőségek halmaza egyetlenegy lehetőségre szűkül össze. Az egyes vizsgálatok részben egymást átfedő információt is nyújtanak. A vizsgálatok sorrendje befolyásolja azt, hogy milyen gyorsan sikerül a kérdést eldönteni. Persze, ha már eleve van bizonyos hipotézisünk, ez megkönnyíti az analízist, azonban minket az az eset érdekel, amikor a vizsgálat megkezdésekor még semmilyen támpontunk nincsen, és így a véletlennek is van bizonyos szerepe abban, hogy milyen sorrendben végezzük el az egyes vizsgálatokat.

Hasonló situációval áll szemben az orvos is, amikor egy beteg diagnózisát kívánja megállapítani, e célból a beteget különböző vizsgálatoknak veti alá és ezek

\* A valóságos helyzetet rendkívül bonyolítja, hogy az ember gyakran jut téves információkhoz is, vagyis a véletlen sokszor véletlen hibák (elírások, félreértések stb.) formájában is közrejátsszik; ez a helyzet akkor, ha — az információelmélet szokásos terminológiájával élve — az egyes adatok „zajos” csatornán át jutnak el az információ gyűjtőjéhez. Bár az ilyen hamis adatok kiküszöbölése (pl. az adatok közötti ellentmondások felfedése útján) gyakran az információ feldolgozásának legnehezebb részét képezi, a kérdésnek ezzel az oldalával itt nem kívánunk foglalkozni; e kérdésre más alkalommal kívánunk visszatérni.

eredményeinek összevetése útján állapítja meg, hogy milyen betegségben szenved az illető beteg. Itt nagy szerepet játszik persze az orvos intuíciója és gyakorlata, amelyre támaszkodva már eleve csak kisszámú lehetőségre kell korlátoznia a figyelmét. Az orvosi tapasztalat és intuíció szerepét azonban úgy foghatjuk fel, hogy az orvos a beteg panaszai, a betegről szerzett benyomásai alapján és a beteg szervezetének esetleges előzetes kezelésekből származó ismeretéből, még mielőtt a laboratóriumi vizsgálatokhoz folyamodna, máris nagy mennyiségű információ felett rendelkezik; vagyis tulajdonképpen az információgyűjtés nem a laboratóriumi vizsgálattal kezdődik, hanem abban a pillanatban, amikor az orvos az illető beteggel foglalkozni kezd (és ennek következtében gyakran laboratóriumi vizsgálatokra nincs is szükség). Minket e példából elsősorban az érdekel, hogy az orvos a végleges diagnózishoz nagyszámú részletinformáció összevetése révén jut el. Az, hogy milyen információk jutnak az orvos birtokába és milyen sorrendben, bizonyos mértékig a véletlentől is függ. Pl. az, hogy a beteg nem felejt-e ki valamit betegsége tüneteinek felsorolásából (esetleg azért, mert ő maga orvosilag tájékozatlan lévén, az illető, őt viszonylag kevésbé zavaró tünetnek nem tulajdonít jelentőséget), erősen a véletlentől függ.

Harmadik példaként tekintsük egy vizsgálóbíró munkáját, aki egy peres ügyben kívánja a tényállást felderíteni; e célból számos tanút hallgat ki, és azok vallomásaiból igyekszik a történeteket rekonstruálni. Legtöbbször az a helyzet, hogy az egyes tanúk a szóban forgó eseményeknek csak egy-egy részletéről bírnak tudomással; továbbá az egyes tanúk által szolgáltatott felvilágosítások legtöbbször részben átfedik egymást\*. A bíró feladata pl. állhat abban, hogy egyének egy bizonyos *H* halmazából kiválassza azt az egyént, aki egy bizonyos cselekményt elkövetett. Azt, hogy kiket és különösen, hogy milyen sorrendben hallgat ki a bíró, legtöbbször tartalmaz véletlen elemeket. (Pl. előfordulhat, hogy a megidézett tanúk közül egyesek akadályoztatás miatt csak későbbi időpontban vagy egyáltalán nem hallgathatók ki.)

További példaként vizsgáljuk meg az ún. „Bar—Kochba” játékot, amely a szóban forgó helyzetnek olyan messzemenően leegyszerűsített esete, hogy ennek alapján azonnal eljuthatunk a probléma matematikai modelljének felállításához. Az említett játékban ugyanis a minket érdeklő probléma szinte „laboratóriumi” tisztaságban áll előttünk. A Bar—Kochba játékban két játékos vesz részt; nevezzük őket *A*-nak és *B*-nek. Az *A* játékos gondol valamire, a *B* játékos pedig kérdéseket tesz fel *A*-nak és a kérdéseire kapott válaszok alapján igyekszik kitalálni, hogy *A* mire gondolt\*\*. Azt, hogy mire gondolhat *A*, a játék szabályai nem korlátozzák lényegesen: *A* gondolhat egy meghatározott személyre, egy tárgyra, egy elvont fogalomra stb. *B* csak olyan kérdéseket tehet fel, amelyekre igennel vagy nemmel lehet válaszolni, *A* pedig köteles minden feltett és megválaszolható kérdésre legjobb tudása szerint helyesen válaszolni\*\*\*. Csak ha a kérdés nem válaszolható meg, vagy *A* a helyes választ nem ismeri, áll *A*-nak jogában nem válaszolni egy feltett kérdésre.

\* Mint már fentebb említettük, azt a (gyakorlatilag egyáltalán nem elhanyagolható) lehetőséget, hogy egyes tanúk (szándékosan vagy tévedésből) a valóságnak nem megfelelő vallomást tesznek, itt nem vesszük tekintetbe.

\*\* A játékot szokták úgy is játszani, hogy többen együtt gondolnak valamire és felváltva válaszolnak a kérdezőnek; a mi szempontunkból ez nem jelent lényeges különbséget.

\*\*\* A téves információkat, amelyekkel a probléma matematikai tárgyalásának egyszerűsítése céljából itt amúgy sem kívánunk foglalkozni, a Bar—Kochba játék esetében tehát a játék szabályai eleve kizárják.

A Bar-Kochba játéknál a gyakorlott kérdező általában bizonyos rendszer szerint szokott kérdezni, a kérdések megválasztását és sorrendjét azonban még a gyakorlott játékosnál is kétségtelenül többé-kevésbé a véletlentől függő tényezők befolyásolják, míg olyan kérdező esetében, aki a játékban kevés tapasztalattal bír, a kérdések megválasztása teljesen véletlenszerűnek tekinthető.

A Bar—Kochba játék egy lehetséges módosításával is érdemes foglalkozni, amelyet „többértékű Bar—Kochba játék”-nak nevezhetünk. Ez abban különbözik a játék szokásos formájától, hogy a kérdezőnek jogában áll olyan kérdéseket is feltenni, amelyekre nem lehet igen-nel vagy nem-mel válaszolni, hanem a lehetséges válaszok száma kettőnél több. Ez esetben azonban a kérdezőnek fel kell sorolni a kérdésre adható összes lehetséges válaszokat, és az  $A$  játékos meg kell, hogy mondja, hogy az általa gondolt dolog esetében a kérdező által felsorolt válaszok közül melyik a helyes (másszóval, hogy az általa gondolt dolog a kérdező által felsorolt, egymást kizáró kategóriák közül melyikbe tartozik).

A „többértékű Bar—Kochba” esetében célszerű az alternatívák számát előre korlátozni (pl. úgy, hogy legfeljebb 3 vagy legfeljebb 4 alternatíva adható meg).

Ahhoz, hogy a Bar—Kochba játéktól eljussunk egy matematikailag kezelhető, leegyszerűsített modellhez, csak egyetlen egyszerűsítést kell bevezetni, mégpedig a gondolható dolgok halmazát kell egy — a játékosok által a játék megkezdése előtt elfogadott megállapítással előre lerögzített — véges  $H$  halmazra korlátozni. Az, hogy a gondolható dolgok halmaza véges, amennyiben az elemszámot nem korlátozzuk, nem jelent gyakorlatilag megszorítást. (Lehetséges megszorítás pl., hogy csak olyasmire lehet gondolni, ami, mint külön címszó, egy bizonyos lexikonban szerepel, vagy pl. ami egy meghatározott szótárban szereplő 5 szóval definiálható stb.).

Ebben az esetben minden egyes kérdésnek egyértelműen megfelel a  $H$  halmaz egy  $H'$  részhalmaza oly módon, hogy a kérdés azzal ekvivalens, hogy a gondolt dolog hozzátartozik-e a  $H'$  részhalmazhoz?

Könnyen megadhatunk egy elméletileg optimális, bár gyakorlatilag nemigen számításba jövő „stratégiát” a kérdező részére. Ha ugyanis a gondolható dolgok  $H$  halmaza  $n$  elemű,  $2^{s-1} < n \leq 2^s$  ( $s = 1, 2, \dots$ ), és  $H$  elemeit a  $0, 1, \dots, n-1$  számokkal számozzuk meg; ez esetben (feltéve, hogy a számozásban a játékosok megállapodtak, vagy egy egyértelmű természetes számozás adva van) elegendő, ha a  $B$  játékos a gondolt dolog sorszámára kérdez. Más szóval azt is feltehetjük, hogy  $A$  valójában a  $0, 1, \dots, n-1$  számok valamelyikére gondolt, és ezt a *gondolt számot* akarja  $B$  kitalálni. Mármint a  $0, 1, \dots, n-1$  számokat a kettes számrendszerben felírva, azok mindegyike legfeljebb  $s$  diadikus jegyet tartalmaz, és így  $B$  minden esetben célhoz ér a következő  $s$  kérdéssel; „A gondolt dolog sorszámát a kettes számrendszerben felírva, abban a  $j$ -edik jegy zérus-e? ( $j = 0, 1, \dots, s-1$ )”.

Az  $n$  elemű  $H$  halmaz gondolt eleme tehát  $s = \{\log_2 n\}$  kérdésre adott válaszból egyértelműen kitalálható; itt  $\{x\}$  a legkisebb olyan egész számot jelöli, amely nem kisebb  $x$ -nél. Azt, hogy ennél kevesebb kérdéssel nem juthatunk célhoz, a következő egyszerű megfontolással láthatjuk be. Ha csak olyan kérdést tehet  $B$  fel, amelyre a válasz igen vagy nem, akkor minden elképzelhető kérdéshez hozzárendelhető a  $H$  halmaznak egy  $H'$  részhalmaza, amely  $H$  azon elemeiből áll, amelyek esetében a válasz a szóban forgó kérdésre igenlő. Ha mármint  $f(n)$ -nel jelöljük azon kérdések minimális számát, amelyek minden körülmények között elégségesek egy  $n$  elemű halmaz egy elemének kitalálásához, akkor először is nyilvánvaló, hogy  $f(n)$  monoton

nem-csökkenő függvény, másrészt, ha az első kérdésnek megfelelő halmaz  $m$ -elemű, akkor még további  $f(m)$  ill.  $f(n-m)$  kérdésre van szükség aszerint, hogy az első kérdésre igenlő vagy tagadó választ kaptunk; ennél fogva

$$f(n) \cong 1 + \text{Min}_{0 \leq m \leq n} [\text{Max}(f(m), f(n-m))].$$

Tehát

$$f(n) \cong 1 + f\left(\left\{\frac{n}{2}\right\}\right).$$

Ha azonban  $2^{s-1} < n \leq 2^s$ , akkor  $\left\{\frac{n}{2}\right\} > 2^{s-2}$ . Ezt a megfontolást  $(s-1)$ -szer megismételve, azt kapjuk, hogy

$$f(n) \cong s - 1 + f(2);$$

mármint nyilvánvaló, hogy  $f(2) = 1$ , tehát

$$f(n) \cong s;$$

de  $s$  alkalmas kérdés bizonyosan célhoz vezet, tehát  $f(n) = s$ , ha  $2^{s-1} < n \leq 2^s$ .

Hasonlóképpen látható be, hogy a többértékű Bar—Kochba játéknál, ha az egy kérdésre adandó válaszok megengedett száma  $r$ , a szükséges kérdések száma  $\left\{\frac{\log_2 n}{\log_2 r}\right\}$ , és egy optimális kérdésrendszer pl. az, hogy sorra megkérdezzük, hogy  $H$  gondolt elemének sorszámát az  $r$ -alapú számrendszerben felírva a 0-adik, első, második, ... jegy a 0, 1, ...,  $r-1$  számok közül melyik?

A mondottakból látható, hogy ha egy-egy kérdésre  $r$  alternatív válasz adható, mindig feltehetjük, hogy a gondolható dolgok  $H$  halmazának elemszáma  $r$ -nek hatványa; ha ugyanis  $r^{s-1} < n \leq r^s$ , semmit nem veszítünk, ha a  $H$  halmazt  $r^s - n$  „fiktív” elemmel  $r^s$  eleművé egészítjük ki. Tegyük tehát fel, hogy  $n = r^s$ . Mármint könnyen belátható, hogy ez esetben a fentebb ismertetett optimális kérdésrendszerek a következő két okból optimálisak: a) először is minden válasz teljes egészében új információt nyújt, azaz redundancia nem lép fel; ugyanis egy számot az  $r$ -alapú számrendszerben felírva az egyes jegyek egymástól függetlenül vehetik fel a 0, 1, ...,  $r-1$  értékeket. b) Egy olyan kérdésre adott válasz, amely kérdésre elvileg  $r$ -féle válasz lehetséges, az információelmélet elemei szerint maximálisan  $\log_2 r$  „bit” információt tartalmaz és az említett optimális kérdésrendszer esetében ezt a maximális információt minden egyes kérdésre ténylegesen meg is kapjuk. Ugyanis a 0, 1, ...,  $r^s - 1$  számokat az  $r$ -alapú számrendszerben felírva, ezek közül pontosan  $r^{s-1}$  számnak lesz a  $j$ -edik jegye éppen  $i$  ( $j = 0, 1, \dots, s-1$ ;  $i = 0, 1, \dots, r-1$ ). Ha tehát a 0, 1, ...,  $r^s - 1$  számok mindegyikére ugyanolyan valószínűséggel gondolt az  $A$  játékos (amit feltehetünk, hiszen ez éppen azt jelenti, hogy semmilyen előzetes tudomásunk sincs arról, hogy  $A$  melyik számra gondolt), akkor arra a kérdésre, hogy a gondolt dolog sorszámát az  $r$ -alapú számrendszerben felírva, mi a  $j$ -edik jegy,  $\frac{1}{r}$  valószínűséggel kapjuk azt a választ, hogy ez a jegy  $i$  ( $i = 0, 1, \dots, r-1$ ). Ez

esetben tehát  $\log_2 r$  információt kapunk minden egyes válaszból és így  $s$  válaszból összesen  $s \log_2 r = \log_2 n$  információt kapunk, és éppen ennyi információra van szükség ahhoz, hogy egy  $n$  elemű halmaz egy elemét egyértelműen jellemezni tudjuk.

Mielőtt továbbmennénk, vegyük még észre az említett optimális kérdésrendszer következő sajátosságát: teljesen mindegy, hogy a kérdéseket egymás után egyenként tesszük-e fel vagy egyszerre. Míg a tényleges Bar—Kochba játéknál előbb megvárjuk az első kérdésre adott választ és csak azután tesszük fel a második kérdést, s. i. t., és első pillanatra úgy tűnhet, hogy ez előnyt jelent (vagyis, hogy érdemes a második kérdés megválasztását függővé tenni az első kérdésre kapott választól, s. i. t.), valójában ez az előny illuzórikus, hiszen az említett optimális kérdésrendszer esetében semmilyen hátrányt nem jelent, ha az összes kérdéseket egyszerre tesszük fel, hiszen akármilyen is az első kérdésre a válasz, úgyis ugyanazt a kérdést tesszük fel másodiknak. Pl.  $r=2$  esetében az első kérdés úgy szól, hogy a gondolt szám páros-e, a második úgy, hogy a gondolt számból kivonva e szám mod 2 vett maradékát és a különbséget 2-vel osztva a hányados páros-e? Ha a második kérdés feltevése előtt már rendelkezésre áll az első kérdésre adott válasz, ez csak annyi különbséget jelent, hogy ha a válasz az volt, hogy a szám páros, a második kérdést úgy fogalmazhatjuk, hogy a szám fele páros-e, míg ha azt hallottuk, hogy a szám páratlan, másodszorra azt kérdezzük, hogy a gondolt számból egyet levonva a kapott páros szám fele páros-e? Ez azonban csak a második kérdés *fogalmazását* érinti, a szükséges kérdések *számát* azonban nem befolyásolja.

Más szavakkal kifejezve: teljesen mellékes, hogy a kapott válaszokból nyert információkat minden egyes válasz után „összegezzük”-e, vagy az „összegezést” csak akkor kezdjük meg, ha már mind az  $s$  válasz rendelkezésünkre áll.

Egy optimális kérdésrendszer az  $r=2$ ,  $n=2^s$  esetben a következőképpen jellemezhető. Ha  $H_1, H_2, \dots, H_s$  jelölik az egyes kérdéseknek megfelelő részhalmazokat (vagyis a  $j$ -edik kérdés azt kérdezi, hogy a gondolt dolog hozzátartozik-e a  $H$  halmaz  $H_j$  részhalmazához) és  $\bar{H}_j$  jelöli a  $H_j$  halmaz kiegészítő halmazát  $H$ -ra nézve, akkor a  $H_j$  és  $\bar{H}_j$  halmazok  $2^{s-1}$  elemből állnak, a  $H_i H_j$ ,  $H_i \bar{H}_j$ ,  $\bar{H}_i H_j$  és  $\bar{H}_i \bar{H}_j$  halmazok ( $i \neq j$ ) mindegyike  $2^{s-2}$  elemből áll és általában ha  $\tilde{H}_j$  a  $H_j$  és  $\bar{H}_j$  halmazok közül az egyiket jelöli, akkor az összes  $\tilde{H}_{j_1} \tilde{H}_{j_2} \dots \tilde{H}_{j_i}$  alakba írható halmazok ( $j_1 < j_2 < \dots < j_i$ ) mindegyike  $2^{s-i}$  elemből áll. Speciálisan tehát a  $\tilde{H}_1 \tilde{H}_2 \dots \tilde{H}_s$  halmazok mindegyike egyetlen elemet tartalmaz. Mivel ilyen alakú halmaz pontosan  $2^s$  darab van, ezek tehát a  $H$  halmaz elemeihez kölcsönösen egyértelműen vannak hozzárendelve; könnyen belátható, hogy az összes lehetséges  $2^s!$  ilyen hozzárendelésnek megfelel egy optimális kérdésrendszer, vagyis az optimális kérdésrendszer száma éppen  $2^s!$ .

A tényleges Bar—Kochba játéknál az említett optimális kérdésrendszerek azonban nem alkalmazhatók, kivéve ha a gondolható dolgok halmazának elemszáma igen kicsiny, már pusztán azért sem, mert ha  $n$  nagy, a gondolható dolgok halmaza elemeinek megszámozása igen fáradságos volna és e számozás fejbentartása gyakorlatilag nem vihető keresztül. (Ha ez a rendszer ténylegesen alkalmazható volna, a játék elveszítene minden érdekességét.) Így tehát felmerül a kérdés, hogy mennyivel több kérdést kell feltenni akkor, ha nem optimális kérdésrendszert használunk. Ezt a problémát azon extrém feltevés mellett fogjuk megválaszolni, amikor nemcsak hogy nem optimális kérdésrendszer szerint kérdezzünk, hanem kérdéseinkben semmiféle rendszert nem alkalmazunk, vagyis a kérdéseket teljesen taláalomra választjuk meg. Arra a meglepő eredményre jutottunk, hogy a szükséges kérdések száma általában nem sokkal lesz több, mint az optimális kérdésrendszer esetében, pl.  $n$  értékétől függetlenül az esetek kb. 99%-ában a kérdéseket taláalomra választva legfeljebb 7-tel több kérdésre van szükség, mint optimális kérdésrendszer esetében.

Az 1. §-ban a közönséges Bar—Kochba játék modelljével, a 2. §-ban a többértékű Bar—Kochba játék modelljével foglalkozunk.

A tárgyalás során azzal az általánosabb esettel is foglalkozunk, amikor a lehetséges válaszok egy-egy kérdésre eleve (a priori) nem egyformán valószínűek. Ha egy kérdésre a többértékű Bar—Kochbánál  $r$  válasz lehetséges és ezek valószínűségei a  $\mathfrak{S} = \{p_1, p_2, \dots, p_r\}$  valószínűségeloszlás tagjai, akkor Shannon képlete szerint

egy válasz  $I_1(\mathfrak{S}) = \sum_{i=1}^r p_i \log_2 \frac{1}{p_i}$  információt nyújt és így várható, hogy legalább

$\left\{ \frac{\log_2 n}{I_1(\mathfrak{S})} \right\}$  kérdést kell feltenni. Ez esetben is kiszámítjuk, hogy ha a kérdéseket minden rendszer nélkül, találmásra tesszük fel, mennyivel több kérdés szükséges.

A 3. és 4. §-okban egy, az előbbivel rokon másik problémát vizsgálunk. Ennek a problémának a modellje a Bar—Kochba játék egy másik általánosítása, amelyet „szimultán Bar—Kochba játéknak” nevezhetünk. E játékban nem egy, hanem több játékos — pl. az  $A_1, A_2, \dots, A_l$  játékosok mindegyike — gondol a  $H$  halmaz egy-egy általa választott elemére (feltesszük, hogy nincs két játékos, aki  $H$  ugyanazon elemére gondolna) és a kérdező  $B$  játékosnak ki kell találnia, hogy az  $A_1, A_2, \dots, A_l$  játékosok rendre mire gondoltak. E célból  $B$  kérdéseket tesz fel szimultán az  $A_1, A_2, \dots, A_l$  játékosoknak, amelyekre ezek mindegyike — az általa gondolt dologra vonatkozólag — a valóságnak megfelelően válaszol. Aszerint, hogy csak igennel vagy nemmel megválaszolható kérdések vannak-e megengedve, vagy olyanok, amelyekre  $r > 2$  válasz lehetséges, beszélünk közönséges ill. többértékű szimultán Bar—Kochba játékról. A matematikai tárgyalás szempontjából lényeges különbséget jelent, hogy  $l$  kis rögzített szám-e (pl.  $l=2$  vagy  $l=3$ ) vagy pedig  $l$  nagy szám, pl.  $l=n$ . Az  $l=n$  esetben a  $H$  halmaz minden egyes elemére gondolt az  $A_1, A_2, \dots, A_l$  játékosok közül egy, tehát lényegében a  $H$  halmaz elemeinek egy ismeretlen permutációját kell  $B$ -nek kitalálnia. Az  $l$  szám nagyságától függetlenül  $B$ -nek a szimultán Bar—Kochba esetében valójában azt a függvényt kell meghatároznia, amely a  $A_1, A_2, \dots, A_l$  játékosok mindegyikéhez hozzárendeli a  $H$  halmaz egy elemét, vagyis egy véges halmaznak egy másik véges halmazra való leképezését kell „kitalálnia”. Abban az esetben, amikor  $l=n$ , tehát a  $H$  halmaz minden egyes elemére gondolt az  $A_1, A_2, \dots, A_l$  játékosok közül legalább egy,  $B$ -nek olyan kérdéseket kell feltennie, amelyekre adott válaszok  $H$  minden egyes elemét egyidejűleg jellemzik.\* Mivel egy olyan kérdésnek, amelyre  $r$  válasz lehetséges, egyértelműen megfeleltethető a  $H$  halmaz  $r$  részre való felosztása, tehát  $B$  akkor tudja kitalálni, hogy az  $A_1, A_2, \dots, A_l$  játékosok mire gondoltak, ha az általa feltett kérdéseknek megfelelő felosztások  $H$  elemeit teljesen szeparálják, azaz nincs két olyan eleme  $H$ -nak, amely minden egyes felosztásnál ugyanabba az osztályba tartozik.\*\*

Ilyen módon a szimultán Bar—Kochba játék problémája lényegében azonos a

\* A közönséges Bar—Kochbával kapcsolatban említett optimális kérdésrendszerek egyben a szimultán Bar—Kochbára nézve is optimálisak; ha azonban a kérdéseket véletlenszerűen választjuk, a közönséges és a szimultán Bar—Kochba között lényeges eltérés van: a szimultán Bar—Kochbánál az  $l=n$  esetben kb. kétszerannyi kérdést kell feltenni, mint az  $l=1$  esetben.

\*\* A rövidség kedvéért a következő terminológiát használjuk: amennyiben az egyes részinformációk a  $H$  halmaz két részre vágásának felelnek meg („igen vagy nem” válasz), azt mondjuk, hogy az információt *dichotomiákból* nyerjük (dichotomia = kettévágás), míg ha az egyes osztályozásoknál kettőnél több osztály is felléphet, azt mondjuk, hogy *általános osztályozásokból* nyerjük az információt.

következő problémával: Egy könyvtár könyveit különböző szempontok szerint osztályozzuk kategóriákba; milyen feltételek mellett fogják ezek az osztályozások együttvéve egyértelműen jellemezni a könyvtár minden egyes könyvét. Vagy tekintsük a kvalitatív kémiai analízis példáját. E példát, amelyet már fentebb említettünk, szintén lehet úgy módosítani, hogy az a „szimultán” Bar—Kochba esetének feleljen meg. Tegyük ugyanis fel, hogy nem egyetlen ismeretlen anyagot, hanem többféle anyagot kell egyidejűleg analizálnunk, mégpedig úgy, hogy ugyanazokkal a vizsgálatokkal állapítsuk meg ezek mindegyikének kémiai minőségét. Tehát a szóban forgó anyagokat többféle próbának vetjük alá és feljegyezzük, hogy melyik hogyan reagált azokra. Minden egyes vizsgálatnál a lehetséges reakciók az összes kémiai vegyületek egy-egy osztályozását adják meg. Kérdés, hogyan kell az elvégzendő vizsgálatok sorozatát összeállítani ahhoz, hogy a reakciókból egyértelműen meghatározható legyen, milyen anyagokkal állunk szemben. Hasonló probléma merül fel a képességvizsgálatoknál alkalmazott „teszt”-eknél is, stb.

A szóban forgó probléma érdekessége, hogy megoldásában a  $\mathfrak{S} = \{p_1, p_2, \dots, p_r\}$  valószínűségeloszláshoz tartozó információnak nem a Shannon-féle  $I_1(\mathfrak{S}) = \sum_{i=1}^r p_i \log_2 \frac{1}{p_i}$  mértékszám, hanem egy másik információ-mérték, az  $I_2(\mathfrak{S}) = \log_2 \frac{1}{\sum_{i=1}^r p_i^2}$  ún. *másodrendű információ-mértékszám* játszik szerepet. Az  $I_2(\mathfrak{S})$

információ-mértékszám a szerző [1] és [2] dolgozataiban bevezetett információ-mértékszámok közé tartozik. Az említett dolgozatokban a  $\mathfrak{S} = \{p_1, p_2, \dots, p_r\}$  valószínűségeloszláshoz tartozó információmennyiség  $\alpha$ -adrendű  $I_\alpha(\mathfrak{S})$  mértékszámát ( $\alpha > 0$ ) az  $\alpha \neq 1$  esetben a következőképpen definiáltuk:

$$I_\alpha(\mathfrak{S}) = \frac{1}{1-\alpha} \log_2 \left( \sum_{i=1}^r p_i^\alpha \right),$$

megjegyezve, hogy ha  $\alpha \rightarrow 1$ , akkor  $I_\alpha(\mathfrak{S})$  az  $I_1(\mathfrak{S}) = \sum_{i=1}^r p_i \log_2 \frac{1}{p_i}$  Shannon-féle mértékszámhoz konvergál, amely tehát úgy tekinthető, hogy szintén az  $\alpha$ -adrendű információ-mértékszámok közé tartozik, mégpedig az  $\alpha = 1$  esetnek felel meg.

Megjegyezzük még, hogy az  $r=2, p_1=p_2=\frac{1}{2}$  speciális esettel foglalkoznak a szerző [3] és [4] dolgozatai. [3]-ban a kérdést a Boole-algebra terminológiájával fogalmaztuk meg. A többértékű szimultán Bar—Kochba játék problémája a Boole-algebra nyelvén a következőképpen fogalmazható meg: a  $H$  halmaznak a  $B$  játékos által feltett kérdéseknek megfelelő részhalmazai milyen feltételek mellett generálják (a halmazelméleti műveletekre nézve) a  $H$  összes részhalmazából álló Boole-algebrát?

Az e dolgozatban tárgyalt eredményeket 1961 nyarán előadtam a Michigan State University-n tartott szemináriumomban. A szeminárium résztvevői közül elsősorban H. RUBIN professzor kapcsolódott be a szóban forgó témakör vizsgálatába. Értékes hozzászólásaiért, amelyekből sok ösztönzést merítettem, ezúton is köszönetet mondok. A szeminárium egy másik résztvevője, J. FOX, a következő módosítást vetette fel a közönséges szimultán Bar—Kochba játéknak: a kérdések feltevése során  $B$  mindig csak olyan kérdéseket tesz fel, amelyek biztosan nyújtanak új infor-

mációt az összes, az előző kérdések által még ki nem talált, gondolt dolgokra vonatkozólag. Más szóval, a  $j$ -edik kérdésnek olyannak kell lennie, hogy a  $H$  halmaznak az előző  $j-1$  kérdésnek megfelelő két részre osztásai egyesítésével létrejött felbontásának minden egynél több elemet tartalmazó osztályát a  $j$ -edik kérdésnek megfelelő felosztás két valódi részre bontsa fel. A probléma ezen módosítását illetőleg H. RUBIN ért el érdekes eredményeket.

### 1. §. Információ-akkumuláció egy ismeretlenre vonatkozó dichotomiák esetében

E §-ban először a következő problémát oldjuk meg.

Legyen  $H$  egy  $n$  elemű halmaz. Válasszuk ki találoomra  $H$ -nak  $k$  számú rész-halmazát — ezeket jelöljük  $H_1, H_2, \dots, H_k$  —, oly módon, hogy az egyes rész-halmazok kiválasztásai függetlenek és minden választásnál  $H$  minden egyes rész-halmaza (beleértve az üres halmazt és magát  $H$ -t is) ugyanakkora, tehát  $\frac{1}{2^n}$  valószínűséggel

kerülhet kiválasztásra.\* Kiszámítandó annak a valószínűsége — ezt  $P_{nk}$ -val jelöljük —, hogy kijelölve  $H$  egy  $x$  elemét, ez az elem egyértelműen meg legyen határozva annak megadása által, hogy a  $H_1, H_2, \dots, H_k$  halmazok közül melyeknek eleme és melyeknek nem eleme. Meg kívánjuk vizsgálni továbbá, hogy mekkorára kell  $k$  értékét választani, hogy a szóban forgó  $P_{nk}$  valószínűség egy megadott  $\alpha$  ( $0 < \alpha < 1$ ) számot meghaladjon.

Jelöljük  $H$  elemeit  $x_1, x_2, \dots, x_n$  és legyen

$$(1.1) \quad \varepsilon_{jh} = \begin{cases} 1 & \text{ha } x_h \in H_j \\ 0 & \text{ha } x_h \notin H_j \end{cases} \quad (j=1, 2, \dots, k, h=1, 2, \dots, n).$$

Feltevésünk szerint a  $H_j$  halmaz  $\frac{1}{2^n}$  valószínűséggel lesz azonos  $H$  bármely rész-halmazával; mivel továbbá  $H$ -nak nyilván  $2^{n-1}$  olyan rész-halmaza van, amely  $x_h$ -t tartalmazza és ugyancsak  $2^{n-1}$  olyan rész-halmaza, amely  $x_h$ -t nem tartalmazza, tehát

$$(1.2) \quad \mathbf{P}(\varepsilon_{jh} = 1) = \mathbf{P}(\varepsilon_{jh} = 0) = \frac{2^{n-1}}{2^n} = \frac{1}{2}.$$

Hasonlóképpen látható be az is, hogy az  $\varepsilon_{jh}$  ( $j=1, 2, \dots, k; h=1, 2, \dots, n$ ) valószínűségi változók függetlenek; ugyanis először is az  $(\varepsilon_{j1}, \varepsilon_{j2}, \dots, \varepsilon_{jn})$  vektorok függetlenek, mert feltevésünk szerint a  $H_1, H_2, \dots, H_k$  halmazokat egymástól függetlenül választjuk; másrészt viszont, ha  $\delta_1, \delta_2, \dots, \delta_n$  tetszőleges, a 0 és 1 számokból álló számsorozat, akkor

$$(1.3) \quad \mathbf{P}(\varepsilon_{j1} = \delta_1, \varepsilon_{j2} = \delta_2, \dots, \varepsilon_{jn} = \delta_n) = \frac{1}{2^n} = \prod_{h=1}^n \mathbf{P}(\varepsilon_{jh} = \delta_h),$$

\* Ily módon előfordulhat, hogy ugyanazt a rész-halmazt többször is kiválasztjuk, ennek valószínűsége azonban rendkívül csekély lesz, ha  $n$  nagy szám, mivel vizsgálatainkban  $k$  nagyságrendje  $\log n$ , és ahhoz, hogy ugyanazon halmaz többszöri kiválasztásának a valószínűsége számottevő legyen,  $k$ -nak legalább  $\sqrt[2]{2^n}$  nagyságrendűnek kellene lennie.



ugyanis  $H$ -nak egyetlenegy olyan részhalmaza van, amely tartalmazza azon  $x_h$ -kat, amelyekre  $\varepsilon_{jh} = 1$  és nem tartalmazza azon  $x_h$ -kat, amelyekre  $\varepsilon_{jh} = 0$ . Tehát az  $\varepsilon_{jh}$  ( $j=1, 2, \dots, k; h=1, 2, \dots, n$ ) valószínűségi változók függetlenek és mindegyik a 0 és 1 értékeket  $\frac{1}{2}$  valószínűséggel veszi fel. E változókat rendezzük el egy  $k \times n$  elemű mátrix alakjában, oly módon, hogy  $\varepsilon_{jh}$  a mátrix  $j$ -edik sorában a  $h$ -adik helyre kerüljön. Jelöljük e mátrixot  $E_{kn}$ -nel. Könnyen belátható, hogy az a követelmény, hogy az  $x$  elemet egyértelműen meghatározza az, hogy a  $H_1, H_2, \dots, H_k$  halmazok közül melyeknek eleme, az  $x = x_i$  esetben azt jelenti, hogy az  $E_{kn}$  mátrixnak az  $i$ -edik oszlopa különbözzék e mátrix összes többi oszlopától. Mivel az  $E_{kn}$  mátrix oszlopvektorai függetlenek és mindegyikük a  $k$ -dimenziós és a 0 és 1 komponensekből álló

összes lehetséges  $2^k$  vektorral  $\frac{1}{2^k}$  valószínűséggel egyenlő, a keresett  $P_{nk}$  valószínűség nyilván

$$(1.4) \quad P_{nk} = \left(1 - \frac{1}{2^k}\right)^{n-1},$$

függetlenül attól, hogy mi  $i$  értéke és hogy milyen elemek állnak  $E_{kn}$   $i$ -edik oszlopában. Másrészt az  $e^{-\frac{x}{1-x}} < 1 - x < e^{-x}$  ( $0 < x < 1$ ) egyenlőtlenség szerint

$$e^{-\frac{n-1}{2^k-1}} < P_{nk} < e^{-\frac{n-1}{2^k}}.$$

Ha  $2^k \cong n$ , akkor  $\frac{n-1}{2^k-1} \cong \frac{n}{2^k}$ , ezért

$$(1.5) \quad e^{-\frac{n}{2^k}} < P_{nk} < e^{-\frac{n-1}{2^k}}, \text{ hacsak } 2^k \cong n.$$

Más szóval, ha  $\frac{1}{e} < \alpha < 1$ , és  $k \cong \log_2 n + \log_2 \frac{1}{\log \frac{1}{\alpha}}$ , akkor  $P_{nk} > \alpha$ , míg ha

$k \cong \log_2(n-1) + \log_2 \frac{1}{\log \frac{1}{\alpha}}$ , akkor  $P_{nk} < \alpha$ . Ilyen módon bebizonyítottuk a következő tételt:

1a. TÉTEL. Találomra válasszuk ki egy  $n$  elemű  $H$  halmaz  $k$  részalmazát oly módon, hogy az egyes részalmazokat egymástól függetlenül választjuk és minden részalmaz kiválasztásakor  $H$  összes részalmazai ugyanakkora valószínűséggel jönnek tekintetbe. Jelölje  $P_{nk}$  annak valószínűségét, hogy  $H$  egy tetszőleges rögzített  $x$  eleme egyértelműen meg lesz határozva azáltal, hogy megadjuk, hogy a kiválasztott  $k$  részalmaz közül melyek tartalmazzák  $x$ -et. Ha  $\frac{1}{e} < \alpha < 1$  és  $k \cong \log_2 n + \log_2 \frac{1}{\log \frac{1}{\alpha}}$ ,

akkor  $P_{nk} \cong \alpha$ , míg ha  $k \cong \log_2(n-1) + \log_2 \frac{1}{\log \frac{1}{\alpha}}$ , akkor  $P_{nk} \cong \alpha$ .

*Megjegyzés.* Ha  $\alpha = 0,99$ , akkor  $\log_2 \frac{1}{\log \frac{1}{\alpha}} < \log_2 100 < 7$ , tehát ha a Bar—

Kochba játéknál kérdéseinket találomra választjuk meg, és  $n$  dolog valamelyikére lehet gondolni, a gondolt dolgot  $\{\log_2 n\} + 7$  találomra megválasztott kérdés alapján 99%-ot meghaladó valószínűséggel ki fogjuk találni. Meglepő, hogy a szükséges kérdések száma mindig 7-tel lesz nagyobb, mint ha optimális módon tennénk fel a kérdéseket, függetlenül attól, hogy mekkora  $n$ , vagyis függetlenül attól, hogy mi a minimális kérdészám. Ha például a játékosok abban állapotnak meg, hogy egy 6-jegyű telefonszámot kell kitalálni, akkor ez 27 találomra feltett kérdés alapján 99% valószínűséggel sikerül ( $2^{20} > 10^6 > 2^{19}$ ).

Könnnyen belátható, hogy a részhalmazok választására tett feltevésünkéből következik, hogy általában a kiválasztott részhalmazok mindegyike a  $H$  halmaz elemeinek körülbelül a felét fogja tartalmazni.\*

Felmerül a kérdés, hogy hogyan módosul a helyzet, ha a részhalmazok kiválasztási szabályát úgy választjuk, hogy a kiválasztott részhalmazok túlnyomórészt kb.  $pn$  elemből álljanak, ahol  $0 < p < 1$  és  $p \neq \frac{1}{2}$ .

Egy ilyen kiválasztási szabályt könnyen nyerhetünk, ha továbbra is feltesszük, hogy az  $\varepsilon_{jh}$  valószínűségi változók függetlenek, azonban ezek eloszlását úgy választjuk, hogy

$$(1.6) \quad \mathbf{P}(\varepsilon_{jh} = 1) = p \quad \text{és így} \quad \mathbf{P}(\varepsilon_{jh} = 0) = q = 1 - p$$

legyen. Ez annak felel meg, hogy a  $H_j$  részhalmazt úgy választjuk meg, hogy  $H$  minden egyes  $x_h$  eleméről sorsolás útján döntjük el, hogy belekerüljön-e  $H_j$ -ba, és feltesszük, hogy a sorsolás úgy történik, hogy ennek valószínűsége  $j$  és  $h$  minden értékére  $p$ -vel egyenlő. Ez esetben is elkészítjük az  $\varepsilon_{jh}$  számokból az  $E_{kn}$  mátrixot, és az, hogy az ismeretlen  $x$  elem meg van határozva azáltal, hogy megadjuk, hogy a  $H_1, H_2, \dots, H_k$  halmazok közül melyek tartalmazzák  $x$ -et, az  $x = x_i$  esetben továbbra is azzal az állítással ekvivalens, hogy az  $E_{kn}$  mátrix  $i$ -edik oszlopa különbözik-e a mátrix összes többi oszlopától. Most azonban ezen esemény  $P_{nk}$  valószínűsége már függ attól, hogy az  $i$ -edik oszlopban hány 1-es áll. A teljes valószínűség tételét alkalmazva a  $P_{nk}$  valószínűsége, a

$$(1.7) \quad P_{nk} = \sum_{v=0}^k \binom{k}{v} p^v q^{k-v} (1 - p^v q^{k-v})^{n-1},$$

kifejezést nyerjük. Alkalmazva a Moivre—Laplace tételt, némi számolással adódik, hogy érvényes a következő tétel:

2a. TÉTEL. *Egy  $n$  elemű  $H$  halmaznak válasszuk ki találomra, egymástól függetlenül  $k$  részhalmazát, oly módon, hogy  $H$  minden egyes elemére vonatkozólag, egymástól függetlenül, sorsolással döntjük el, hogy az illető elem beletartozzék-e vagy nem a szóban forgó részhalmazba és ezen lehetőségek valószínűségei  $p$  és  $q = 1 - p$  legyenek ( $0 < p < 1$ ;  $p \neq \frac{1}{2}$ ). Jelölje  $P_{nk}$  annak a valószínűségét, hogy  $H$  egy tetszőlegesen meg-*

\* Pontosabban: mindegyik halmaz elemszáma nagy valószínűséggel az  $\frac{n \pm \sqrt{(2+\delta)n \log \log n}}{2}$  határok közé fog esni, ahol  $\delta > 0$  tetszőleges, és  $n \geq n_0(\delta)$ , ahol  $n_0(\delta)$  csak  $\delta$ -tól függ.

választott rögzített  $x$  elemét egyértelműen jellemezze az, hogy a kiválasztott részhalmazok közül melyek tartalmazzák  $x$ -et. Ez esetben, bevezetve az  $I_1(\{p, q\}) = -p \log_2 \frac{1}{p} + q \log_2 \frac{1}{q}$  jelölést, ha  $n \rightarrow +\infty$  és

$$(1.8) \quad k = k(n) = \frac{\log_2 n + y \sqrt{\log_2 n} + o(\log_2 n)}{I_1(\{p, q\})},$$

akkor

$$(1.9) \quad \lim_{n \rightarrow +\infty} P_{n, k(n)} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{y/\sigma} e^{-t^2/2} dt,$$

ahol

$$(1.10) \quad \sigma = \sqrt{\frac{pq}{I_1(\{p, q\})}} \cdot |\log p/q|.$$

*Megjegyzés.* Ha (1.10)-ben  $p = \frac{1}{2}$ , akkor  $\sigma = 0$ ; ez esetben a tétel állítása akkor marad érvényes, ha  $y \neq 0$  és  $y/\sigma$  alatt  $+\infty$  ill.  $-\infty$  értendő, aszerint, hogy  $y > 0$ , ill.  $y < 0$ .

## 2. §. Információ-akkumuláció egy ismeretlenre vonatkozó általános osztályozások esetében

E §-ban az 1. §-ban tárgyalt probléma következő általánosításával foglalkozunk. Az  $n$  elemű  $H$  halmazt  $k$ -szor egymásután találomra  $r$  darab (számozott) részre bontjuk, ( $r > 2$ ) oly módon, hogy az egyes felbontások egymástól függetlenül történnek és minden felbontás ugyanakkora (tehát  $\frac{1}{r^n}$ ) valószínűséggel lesz egyenlő az összes lehetséges  $r$ -részre való felbontások mindegyikével. Kiszámítandó annak a valószínűsége — ezt  $P_{nk}$ -val jelöljük —, hogy kijelölve  $H$  egy tetszőleges rögzített  $x$  elemét, ez az elem egyértelműen meg legyen határozva annak megadása által, hogy a  $k$  felbontás mindegyikénél  $x$  hányadik részhalmazba tartozik.

Nyilván a szóban forgó felbontások mindegyike jellemezhető azon függvény által, amely  $H$  minden eleméhez hozzárendeli a felbontás azon részhalmazának sorszámát, amelyhez az illető elem tartozik. Ez a függvény (és ennek következtében a számbajövő felbontás) nyilván  $r^n$ -féleképpen adható meg. Jelölje  $\varepsilon_{jh}$  a  $j$ -edik felbontáshoz tartozó ilyen függvényt, tehát legyen

$$(2.1) \quad \varepsilon_{jh} = l,$$

ha  $x_h$  a  $j$ -edik felbontásnál az  $l$ -edik osztályba tartozik ( $l = 1, 2, \dots, r$ ;  $j = 1, 2, \dots, k$ ;  $h = 1, 2, \dots, n$ ). Ugyanúgy, mint az 1. §-ban, beláthatjuk, hogy az  $\varepsilon_{jh}$  változók függetlenek és

$$(2.2) \quad \mathbf{P}(\varepsilon_{jh} = l) = \frac{1}{r} \quad (j = 1, 2, \dots, k; h = 1, 2, \dots, n; l = 1, 2, \dots, r).$$

Elkészítve újból az  $\varepsilon_{jh}$  számokból az  $E_{kn}$  mátrixot, az  $x = x_i$  esetben  $P_{nk}$  annak valószínűségével lesz egyenlő, hogy  $E_{kn}$   $i$ -edik oszlopa az összes többi oszloptól különbözik. Így nyerjük, hogy

$$(2.3) \quad P_{nk} = \left(1 - \frac{1}{r^k}\right)^{n-1},$$

tehát, ha  $r^k > n$ , akkor

$$(2.4) \quad e^{-\frac{n}{r^k}} < P_{nk} < e^{-\frac{n-1}{r^k}}.$$

Ily módon az 1. tétel következő általánosítását kaptuk:

1b. TÉTEL. Ha az  $n$  elemű  $H$  halmaznak taláломra kiválasztjuk  $k$  számú  $r$  (számozott) részre való felbontását ( $r \geq 2$ ) oly módon, hogy az egyes felbontásokat egymástól függetlenül választjuk ki úgy, hogy az  $r^n$  lehetséges felbontás mindegyike ugyanolyan valószínűséggel kerülhet kiválasztásra és  $P_{nk}$ -val jelöljük annak a valószínűségét, hogy  $H$  egy tetszőleges rögzített  $x$  eleme egyértelműen meg legyen határozva általa, hogy megadjuk, hogy a  $k$  felbontás mindegyikénél melyik osztályba tartozik,

továbbá, ha  $\frac{1}{e} < \alpha < 1$  és  $k \geq \frac{\log_2 n + \log_2 \left(\log \frac{1}{\alpha}\right)^{-1}}{\log_2 r}$ , akkor  $P_{nk} \geq \alpha$ , míg ha  $k \leq \frac{\log_2(n-1) + \log_2 \left(\log \frac{1}{\alpha}\right)^{-1}}{\log_2 r}$ , akkor  $P_{nk} \leq \alpha$ .

*Megjegyzés.* Ha  $r=2$ , akkor az 1b. tételből speciális esetként az 1a. tételt kapjuk.

A mondott feltételek mellett minden egyes felbontás minden egyes részhalmaza 1-hez közeli valószínűséggel  $\frac{n \pm O(\sqrt{n \log \log n})}{r}$  elemet fog tartalmazni. Ha azt akarjuk, hogy a felbontások részhalmazai ne álljanak körülbelül ugyanannyi elemből, hanem elemszámaik a  $p_1 n, p_2 n, \dots, p_r n$  számokhoz essenek közel, ahol  $\mathfrak{S} = \{p_1, p_2, \dots, p_r\}$  egy tetszőleges  $r$  elemű valószínűségeloszlás, azaz  $p_l > 0$  ( $l=1, 2, \dots, r$ ) és  $\sum_{l=1}^r p_l = 1$ , akkor a felbontások véletlenszerű megválasztását úgy kell módosítani, hogy

$$(2.5) \quad \mathbf{P}(\varepsilon_{jh} = l) = p_l \quad (l=1, 2, \dots, r; j=1, 2, \dots, k; h=1, 2, \dots, n)$$

legyen.

Ez esetben az 1. §-ban alkalmazott megközelítéshez hasonló módon a

$$(2.6) \quad P_{nk} = \sum_{\sum_{l=1}^r k_l = k} \frac{k!}{k_1! k_2! \dots k_r!} p_1^{k_1} p_2^{k_2} \dots p_r^{k_r} (1 - p_1^{k_1} p_2^{k_2} \dots p_r^{k_r})^{n-1}$$

képletet nyerjük. Ebből némi számolással adódik a következő

2b. TÉTEL. Ha egy  $n$  elemű  $H$  halmaz elemeit  $k$ -szor taláalomra  $r \equiv 2$  számozott osztályra bontjuk fel, oly módon, hogy az egyes felbontások egymástól függetlenek és  $H$  minden eleme minden egyes felbontás  $l$ -edik osztályába  $p_l$  valószínűséggel esik  $\left( p_l > 0, l = 1, 2, \dots, r; \sum_{l=1}^r p_l = 1 \right)$  és  $P_{nk}$  jelöli annak a valószínűségét, hogy  $H$  egy tetszőleges rögzített  $x$  eleme egyértelműen meg legyen határozva azáltal, hogy megadjuk, hogy a  $k$  felbontás mindegyikénél  $x$  hányadik osztályba tartozik, továbbá, ha

$$(2.7) \quad k = k(n) = \frac{\log_2 n + y \sqrt{\log_2 n} + o(\sqrt{\log_2 n})}{I_1(\mathfrak{S})},$$

ahol  $\mathfrak{S}$  jelöli a  $(p_1, p_2, \dots, p_r)$  eloszlást, amelyről feltesszük, hogy elemei nem mind egyenlők, és

$$(2.8) \quad I_1(\mathfrak{S}) = \sum_{l=1}^r p_l \log_2 \frac{1}{p_l},$$

akkor

$$(2.9) \quad \lim_{n \rightarrow +\infty} P_{n, k(n)} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{y/\sigma} e^{-t^2/2} dt,$$

ahol

$$(2.10) \quad \sigma = \sqrt{\frac{\sum_{l=1}^r p_l \log_2^2 p_l - \left( \sum_{l=1}^r p_l \log_2 \frac{1}{p_l} \right)^2}{I_1(\mathfrak{S})}}.$$

Megjegyzés. Ha  $r=2, p_1=p, p_2=q$  ( $p+q=1, p \neq \frac{1}{2}$ ), akkor a (2.10) által definiált  $\sigma$  értéke  $\sqrt{\frac{pq}{I_1(\{p, q\})} |\log p/q|}$ , és így a 2b. tétel az  $r=2$  speciális esetben a

2a. tételre redukálódik. Ha viszont  $r$  tetszőleges, viszont  $p_1=p_2=\dots=p_r=\frac{1}{r}$ , akkor a (2.10) által definiált  $\sigma$  értéke 0; a 2b. tétel állítása ez esetben is érvényes marad  $y \neq 0$ -ra, ha  $y/\sigma$ -t  $+\infty$ -nek ill.  $-\infty$ -nek értelmezzük, aszerint, hogy  $y > 0$  vagy  $y < 0$ .

### 3. §. Szimultán információ-akkumuláció véletlen dichotomiák esetében

E §-ban először a következő problémával foglalkozunk: Egy  $n$  elemű  $H$  halmaznak egyidejűleg  $l$  különböző ismeretlen elemét kívánjuk meghatározni oly módon, hogy egymástól függetlenül taláalomra kiválasztjuk  $H$   $k$  részalmazát (úgy, hogy minden választásnál minden részalmaz ugyanakkora valószínűséggel kerül kiválasztásra) és feltesszük, hogy az  $l$  ismeretlen elem mindegyikét illetőleg megállapítható, hogy azok a kiválasztott részalmazok közül melyekhez tartoznak hozzá. Jelöljük  $P_{nkl}$ -vel annak a valószínűségét, hogy ezek az adatok egyértelműen meghatározzák a  $H$  halmaz szóban forgó  $l$  elemét. Jelöljék újból  $x_1, x_2, \dots, x_n$  a  $H$  halmaz elemeit,  $H_1, H_2, \dots, H_k$  a kiválasztott részalmazokat és legyen  $\varepsilon_{jh} = 1$  vagy 0 aszerint, hogy  $x_h$  eleme-e  $H_j$ -nek vagy nem ( $j = 1, 2, \dots, k; h = 1, 2, \dots, n$ ). Akkor  $P_{nkl}$  nyilván

egyenlő annak valószínűségével, hogy az  $\varepsilon_{jh}$  elemekből alkotott  $E_{kn}$  mátrixnak az  $l$  ismeretlen elemnek megfelelő oszlopa egymástól, valamint az összes többi oszlop-tól különbözzék. Ilyen módon azt kapjuk, hogy

$$(3.1) \quad P_{nkl} = \prod_{v=0}^{l-1} \left(1 - \frac{v}{2^k}\right) \left(1 - \frac{l}{2^k}\right)^{n-l}.$$

Egyszerűen adódik (3.1)-ből, hogy ha  $n, k$  és  $l$  úgy tartanak végtelenhez, hogy  $l = o(n)$  és  $k = \log_2 nl + \log_2 \frac{1}{\log \frac{1}{\alpha}} + o(1)$ , ahol  $0 < \alpha < 1$ , akkor

$$(3.3) \quad \lim P_{nkl} = \alpha.$$

Ezzel szemben, ha  $n$  és  $l$  úgy tartanak végtelenhez, hogy

$$(3.4) \quad l \sim \vartheta n \quad (0 < \vartheta \leq 1),$$

akkor (3.3) abban az esetben teljesül, ha

$$(3.5a) \quad k = \log_2 nl + \log_2 \frac{1 - \frac{\vartheta}{2}}{\log \frac{1}{\alpha}} + o(1),$$

vagyis, ha

$$(3.5b) \quad k = 2 \log_2 n + \log_2 \frac{\vartheta(2 - \vartheta)}{2 \log \frac{1}{\alpha}} + o(1).$$

Speciálisan, ha

$$(3.6) \quad l = n,$$

akkor (3.3) fennállásához kell, hogy

$$(3.7) \quad k = 2 \log_2 n + \log_2 \frac{1}{2 \log \frac{1}{\alpha}} + o(1)$$

legyen.

Érvényes tehát a következő tétel:

3a. TÉTEL. Ha az  $n$  elemű  $H$  halmaznak egymástól függetlenül taláalomra kiválasztjuk  $k$  részhalmazát, oly módon, hogy minden választásnál  $H$   $2^n$  részhalmaza közül mindegyik ugyanolyan valószínűséggel kerülhet kiválasztásra, és  $P_{nk}$  jelöli annak valószínűségét, hogy a kiválasztott  $k$  részhalmaz  $H$  bármely két elemét elválasztja egymástól, azaz, hogy nincs  $H$ -nak két olyan eleme, amelyek az összes kiválasztott részhalmazzal ugyanolyan relációban állnak (ugyanazoknak a részhalmazoknak

elemei, ill. nem elemei) és  $n$  és  $k$  oly módon tartanak végtelenhez, hogy

$$k = 2 \log_2 n + \log_2 \frac{1}{2 \log \frac{1}{\alpha}} + o(1),$$

ahol  $0 < \alpha < 1$ , akkor

$$\lim_{n \rightarrow +\infty} P_{nk} = \alpha.$$

Összehasonlítva a 3a. tételt az 1a. tétellel, azt látjuk, hogy egy  $n$  elemű halmaz összes elemeinek szeparálásához kerekén kétszer annyi dichotomia szükséges, mint  $e$  halmaz egyetlen elemének szeparálásához. E problémánál valójában egy ismeretlen  $n$ -edrendű permutációt kell meghatároznunk. Ehhez legalább  $\log_2 n! \sim n \log_2 n$  „bit” információra van szükség; egy dichotomia a permutáció minden egyes elemére nézve maximálisan 1 bit, az egész permutációra nézve  $n$  bit információt adhat. Más szóval, ha véletlen dichotomiákból határozzuk meg az ismeretlen permutációt, kb. 100%-os redundancia lép fel.\*

#### 4. §. Szimultán információ-akkumuláció, általános osztályozások esetében

Ha az  $n$  elemű  $H$  halmaz  $l$  különböző elemét akarjuk meghatározni és minden egyes rész-információ abban áll, hogy találomra kiválasztjuk a  $H$  halmaz egy  $r$  részre való felbontását (úgy, hogy minden választásánál az összes ilyen felbontások egyformán valószínűek, és az egyes választások egymástól függetlenek), és ezután értesülünk arról, hogy az  $l$  ismeretlen elem mindegyike minden egyes felbontásnál melyik részhalmazhoz tartozik, akkor —  $P_{nkl}$ -vel jelölve annak a valószínűségét, hogy  $k$  ilyen felbontás elegendő az  $l$  ismeretlen elem meghatározására — az előző §-ban alkalmazott megfontoláshoz hasonlóan beláthatjuk, hogy

$$(4.1) \quad P_{nkl} = \prod_{v=0}^{l-1} \left(1 - \frac{v}{r^k}\right) \left(1 - \frac{l}{r^k}\right)^{n-l}.$$

Ha tehát

$$(4.2) \quad l = o(n),$$

akkor a

$$(4.3) \quad k = \frac{\log_2 nl + \log_2 \frac{1}{\log \frac{1}{\alpha}} + o(1)}{\log_2 r}$$

\* Ha  $l > n$ , viszont megengedett az, hogy — a szimultán Bar—Kochba terminológiáját használva — több játékos gondoljon ugyanarra a dologra, akkor sem szükséges több dichotomia mint az  $l = n$  esetben, hiszen  $2 \log_2 n + \log_2 \frac{1}{2 \log \frac{1}{\alpha}}$  dichotomia körülbelül  $\alpha$  valószínűséggel elégséges

a  $H$  halmaz minden elemének szeparálásához és teljesen mellékes, hogy hány játékos gondolt ugyanarra a dologra, hiszen akik ugyanarra gondoltak, mindig ugyanazokat a válaszokat adják és így, ha  $B$  kitalálta, hogy egyikük mire gondolt, akkor a többiről is tudja ugyanezt.

esetben lesz

$$(4.4) \quad \lim P_{nkl} = \alpha, \quad (0 < \alpha < 1),$$

míg ha

$$(4.5) \quad l \sim \vartheta n, \quad (0 < \vartheta \leq 1),$$

akkor a

$$(4.6) \quad k = \frac{2 \log_2 n + \log_2 \frac{\vartheta(2-\vartheta)}{2 \log \frac{1}{\alpha}} + o(1)}{\log_2 r}$$

esetben lesz (4.4) érvényes.

Ezek után azt váránk, hogy ha az  $r$ -része osztásnál az egyes részek elemszámai körülbelül úgy viszonylanak egymáshoz, mint a  $p_1, p_2, \dots, p_r$  számok

( $\sum_{i=1}^r p_i = 1$ ), akkor a fenti eredmények csak annyiban módosulnak, hogy  $\log_2 r$  helyébe  $I_1(\mathfrak{S})$  lép. A tüzetesebb vizsgálat azonban azt mutatja, hogy ez nem igaz, ugyanis  $\log_2 r$  helyébe nem a Shannon-féle  $I_1(\mathfrak{S}) = \sum_{i=1}^r p_i \log_2 \frac{1}{p_i}$  elsőrendű információ-mennyiség, hanem az  $I_2(\mathfrak{S}) = \log_2 \left( \frac{1}{\sum_{i=1}^r p_i^2} \right)$  másodrendű információ-mennyiség

lép. (Persze a  $p_1 = p_2 = \dots = p_r = 1/r$  esetben  $I_1(\mathfrak{S}) = I_2(\mathfrak{S}) = \log_2 r$ .) Ezt a következőképpen láthatjuk be: Jelöljék  $\pi_1, \pi_2, \dots, \pi_k$  valamilyen sorrendben az olyan  $k$ -tényezős szorzatokat, amelyek minden tényezője a  $p_1, p_2, \dots, p_r$  számok egyike

(sorrend is számít!). Nyilván a  $\pi_i$  számok mind  $p_1^{k_1} p_2^{k_2} \dots p_r^{k_r}$  alakúak, ahol  $\sum_{i=1}^r k_i = k$

és ez a szám a  $\pi_i$  számsorozatban  $\frac{k!}{k_1! k_2! \dots k_r!}$ -szor fordul elő.

Az egyszerűség kedvéért csak a legérdekesebb esetet tekintjük, vagyis azt, amikor  $l = n$ .

Ez esetben  $P_{nk}$ -val jelölve annak valószínűségét, hogy a  $k$  felbontás  $H$  összes elemeit szeparálja,

$$(4.7) \quad P_{nk} = \sum' \pi_{i_1} \pi_{i_2} \dots \pi_{i_n},$$

ahol az összegezés az összes különböző számokból álló  $(i_1, i_2, \dots, i_n)$ -szám  $n$ -esekre terjesztendő ki.

A (4.7) jobb oldalán álló összeg aszimptotikus viselkedésének megvizsgálásához éppen arra a szita-formulára van szükség, amelyet egy előző dolgozatunkban [5] bizonyítottunk be. Jelölje  $B_{ij}$  azt az eseményt, hogy az  $E_{kn}$  mátrix  $i$ -edik és  $j$ -edik oszlopa azonos. E jelölés mellett  $P_{nk}$  annak a valószínűségével egyenlő, hogy a  $B_{ij}$  ( $1 \leq i < j \leq n$ ) események egyike sem következik be. Jelölje  $G$  azt a gráfot, amelynek szögpontjai az  $(i, j)$  számpárok ( $1 \leq i < j \leq n$ ) és amelyben az  $(i, j)$  és  $(i', j')$  számpároknak megfelelő szögpontok akkor és csak akkor vannak összekötve egy éllel,



ha az  $i, j, i', j'$  számok nem mind különbözők (tehát, ha  $i=i'$  vagy  $i=j'$  vagy  $j=i'$  vagy  $j=j'$ ; nyilvánvaló, hogy e négy lehetőség közül  $i < j$  és  $i' < j'$  miatt legfeljebb egy állhat fenn). Jelölje  $\Gamma^*$  a  $G$  szögpontjai halmazának azon részhalmazai összességét, amelyek nem tartalmaznak egyetlen  $G$ -ben éllel összekötött szögpontpárt sem, és  $\Gamma^{**}$  a  $G$  szögpontjai halmazának azon részhalmazai összességét, amelyek legfeljebb egy  $G$ -ben éllel összekötött szögpontpárt tartalmaznak. Legyen

$$(4.8) \quad S_d^* = \Sigma^* P(B_{i_1 j_1} B_{i_2 j_2} \dots B_{i_d j_d}),$$

ill.

$$(4.9) \quad S_d^{**} = \Sigma^{**} P(B_{i_1 j_1} B_{i_2 j_2} \dots B_{i_d j_d}),$$

ahol  $\Sigma^*$  ill.  $\Sigma^{**}$  azt jelöli, hogy az összegezés csak az olyan  $(i_1, j_1), \dots, (i_d, j_d)$  szám-pár  $d$ -esekre terjesztendő ki, amelyek  $\Gamma^*$ -hoz, ill.  $\Gamma^{**}$ -hoz tartoznak. Ez esetben az [5]-ben bebizonyított tételünk szerint  $s$  bármely nemnegatív egész értékre

$$(4.10) \quad 1 - S_1^{**} + S_2^* - S_3^{**} + \dots - S_{2s+1}^{**} \leq P_{nk} \leq 1 - S_1^* + S_2^{**} - \dots + S_{2s}^{**}.$$

Azonban nyilvánvalóan

$$(4.11) \quad S_d^* = \left( \sum_{i=1}^{rk} \pi_i^2 \right)^d \frac{n(n-1) \dots (n-2d+1)}{d! 2^k},$$

továbbá

$$(4.12) \quad S_d^{**} = S_d^* \left( 1 + \frac{8 \left( \sum_{i=1}^{rk} \pi_i^3 \right) d(d-1)}{3 \left( \sum_{i=1}^{rk} \pi_i^2 \right)^2 (n-2d+1)} \right).$$

Ha mármost

$$(4.13) \quad k = \frac{2 \log_2 n + \log \frac{1}{2 \log \frac{1}{\alpha}} + o(1)}{\log_2 \frac{1}{\sum_{v=1}^r p_v^2}},$$

akkor figyelembe véve, hogy

$$(4.14) \quad \sum_{i=1}^{rk} \pi_i^2 = \left( \sum_{v=1}^r p_v^2 \right)^k,$$

azt kapjuk, hogy  $d$  minden rögzített értékére

$$(4.15) \quad S_d^* \sim \frac{\left( \log \frac{1}{\alpha} \right)^d}{d!},$$

továbbá figyelembe véve, hogy

$$(4.16) \quad \sum_{i=1}^{r^k} \pi_i^3 = \left( \sum_{v=1}^r p_v^3 \right)^k,$$

és ezért, bevezetve a  $\bar{p} = \max_{1 \leq v \leq r} p_v$  jelölést,

$$(4.17) \quad \sum_{i=1}^{r^k} \pi_i^3 < \bar{p}^k \left( \sum_{v=1}^r p_v^2 \right)^k,$$

azt kapjuk, hogy

$$(4.18) \quad S_d^{**} = S_d^* (1 + o(n^{-\epsilon})),$$

ahol

$$(4.19) \quad \varrho = \frac{\log_2 \frac{\sum_{v=1}^r p_v^2}{\bar{p}^2}}{\log_2 \frac{1}{\sum_{v=1}^r p_v^2}} > 0.$$

Tehát  $s$  minden nemnegatív egész értékre

$$(4.20) \quad \sum_{d=0}^{2s+1} \frac{(-1)^d \left( \log \frac{1}{\alpha} \right)^d}{d!} \leq \underline{\lim} P_{nk} \leq \overline{\lim} P_{nk} \leq \sum_{d=0}^{2s} \frac{(-1)^d \left( \log \frac{1}{\alpha} \right)^d}{d!}.$$

Ilyen módon, elvégezve az  $s \rightarrow +\infty$  határátmenetet nyerjük, hogy

$$(4.21) \quad \lim P_{nk} = e^{-\log \frac{1}{\alpha}} = \alpha.$$

Ezzel tehát bebizonyítottuk a következő tételt:

**3b. TÉTEL.** *Ha egy  $n$  elemű  $H$  halmaz elemeit egymástól függetlenül  $k$ -szor  $r$  (számozott) osztályba soroljuk, oly módon, hogy  $H$  elemei minden egyes osztályba-sorolásnál egymástól függetlenül  $p_v$  valószínűséggel kerülnek a  $v$ -edik osztályba ( $v=1, 2, \dots, r$ ), és  $P_{nk}$  jelöli annak a valószínűségét, hogy a szóban forgó  $k$  felosztás  $H$  bármely két elemét elválasztja, vagyis, hogy ne legyen  $H$ -nak két olyan különböző eleme, amelyek mindegyik választott felosztásnál ugyanabba az osztályba esnek, továbbá, ha  $n$  és  $k$  úgy tartanak végtelenhez, hogy*

$$(4.22) \quad k = \frac{2 \log_2 n + \log_2 \frac{1}{\log \frac{1}{\alpha}} + o(1)}{\log_2 \frac{1}{\sum_{v=1}^r p_v^2}}$$

ahol  $0 < \alpha < 1$ , akkor

$$(4.23) \quad \lim P_{nk} = \alpha.$$

Nyilvánvaló, hogy a 3b. tétel speciális esetként tartalmazza (az  $r=2$ ,  $p_1=p_2=\frac{1}{2}$  esetben) a 3a. tételt.

#### IRODALOMJEGYZÉK

- [1] RÉNYI ALFRÉD: Az információelmélet néhány alapvető kérdése, *MTA III. Oszt. Közl.* **10** (1960) 251–282.
- [2] „ On measures of entropy and information, *Proceedings of the 4th Berkeley Symposium on Mathematical Statistics and Probability*, University of California Press, 1961, Vol. I. 547–561.
- [3] „ On random generating elements of a finite Boolean algebra, *Acta Sci. Math. Szeged*, **22** (1961) 75–81.
- [4] „ Statistical laws of accumulation of information, *Bulletin of the International Statistical Institute, 33rd. Session of the ISI in Paris*, 1961. pp. 1–7.
- [5] „ Egy általános módszer valószínűségszámítási tételek bizonyítására és annak néhány alkalmazása, *MTA III. Oszt. Közl.* **11** (1961) 79–105.

(Beérkezett: 1961. XI. 9.)