

# DISZKRÉT VALÓSZÍNŰSÉG-ELOSZLÁSOK KEVERÉKÉNEK FELBONTÁSA ÖSSZETEVŐIRE

MEDGYESSY PÁL

## Bevezetés

Korábbi dolgozatainkban ([1], [2]) folytonos eloszlás-, illetve sűrűség-függvények konstans súlyokkal vett keverékének felbontásával foglalkoztunk. Ebben a cikkben diszkrét valószínűség-eloszlásokra végzünk hasonló vizsgálatokat. Tudomásunk szerint ezzel a problémával eddig nem foglalkoztak, bár fontos gyakorlati alkalmazásai vannak.

Problémánk a következőképp fogalmazható meg: legyen  $\psi_\nu(x)$  ( $-\infty < x < \infty$ ;  $\nu = a_1, a_2, \dots$ ) az  $x$  paramétertől függő diszkrét eloszlás, amelyet az  $a_1, a_2, \dots$  pontokban értelmeztünk;  $\psi_\nu(x)$  az  $x$ -nek folytonos függvénye. Legyen továbbá  $\Phi(x)$  egy eloszlásfüggvény. Bevezetjük a következő definíciót:

A  $\psi_\nu(x)$  komponensek  $\Phi(x)$  súlyfüggvénnyel való keverékének az

$$(1) \quad f_\nu = \int_{-\infty}^{\infty} \psi_\nu(x) d\Phi(x) \quad (\nu = a_1, a_2, \dots)$$

Stieltjes-integrállal értelmezett diszkrét eloszlást nevezünk.

A definíció az [1] dolgozat bevezetésében közölt, folytonos sűrűség-függvények keverékére vonatkozó definíció értelemszerű kiterjesztése. (A diszkrét eloszlás eloszlásfüggvényét tartalmazó definíció hasonlóan adható meg, semmi előnnyel sem jár azonban, ha ebből indulunk ki s ezért mellőzzük.)

Az (1)-ben szereplő  $f_\nu$  keverékeloszlás felbontásán általánosságban  $\psi_\nu(x)$  vagy  $\Phi(x)$  (esetleg mindkettő) meghatározását értjük az (1) összefüggés és az  $f_\nu$  értékek ismerete alapján. Esetenként még azt is megengedjük, hogy a keresett függvények típusát is ismerjük, csak paramétereit nem.

Ha  $\Phi(x)$  tiszta lépcsős függvény, (1) nyilván végtelen sorba megy át. Alkalmasan választott  $\psi_\nu(x)$  és tiszta lépcsős  $\Phi(x)$  függvényekkel elérhetjük, hogy  $f_\nu$  az

$$(2) \quad \left\{ \begin{array}{l} f_\nu = \sum_{k=1}^N A_k \binom{M_k}{\nu} p_k^\nu (1-p_k)^{M_k-\nu} \quad (\nu = 0, 1, \dots, \max M_k) \\ \text{illetőleg} \\ f_\nu = \sum_{k=1}^N A_k \cdot e^{-\lambda_k} \frac{\lambda_k^\nu}{\nu!} \quad (\nu = 0, 1, \dots) \end{array} \right.$$

alakot öltse  $\left( A_k > 0, \sum_{k=1}^N A_k = 1, M_k > 0, 0 < p_k < 1, \lambda_k > 0 \right.$  valós paraméterek  $\left. \right)$ . A (2) alatti kifejezések közül az elsőt különböző  $M_k, p_k$  paraméterű binomiális eloszlások konstans  $A_k$  súlyokkal vett keverékének, a másodikat pedig  $\lambda_k$  paraméterű Poisson-eloszlások  $A_k$  súlyokkal vett keverékének nevezzük.

A dolgozatunkban tárgyalt problémák erre a két speciális keverékre vonatkoznak.

Az (1) alatti alakhoz viszonyítva látjuk, hogy a keverékek komponenseinek típusát ismerjük, de paramétereit nem, a súlyfüggvényről pedig csak azt tudjuk, hogy tiszta lépcsős függvény.

A következőkben ezen keverékek felbontásával fogunk foglalkozni, ami itt az ismeretlen  $A_k, M_k, p_k, \lambda_k$  paraméterek meghatározásával ekvivalens.

Az  $A_k$  mennyiségek valószínűségszámítási interpretációjához  $\left( \sum_{k=1}^N A_k = 1 \right)$  nem szükséges ragaszkodnunk s ezért a (2) alatti alakot egyszerűen mint binomiális, illetve Poisson-eloszlások szuperpozícióját kezeljük.

Binomiális eloszlások keveréke felbontásának gyakorlati szerepével és jelentőségével egy másik dolgozatban [3] már foglalkoztunk.

## 1. §.

Az

$$(3) \quad f_\nu = \sum_{k=1}^N A_k \binom{M_k}{\nu} p_k^\nu (1-p_k)^{M_k-\nu} \quad (\nu = 0, 1, \dots, \max M_k)$$

( $A_k > 0, M_k > 0$  egész,  $0 < p_k < 1$ , azonos  $M_k, p_k$  pár nincs) keverék felbontásánál, — ami itt az  $A_k, M_k, p_k$  értékek meghatározásával ekvivalens — az [1] dolgozatban említett »szórás-csökkentési« elvet alkalmazzuk: legyenek egyelőre a  $p_k$ -k mind különbözőek, s tegyük fel, hogy a (3)-ban szereplő  $f_\nu$ -értékek ismeretében alkalmasan választott  $\mu > 1$  valós paraméterrel elő tudunk állítani egy

$$(4) \quad F_\nu(\mu) = \sum_{k=\nu}^N A_k \binom{M_k}{\nu} (\mu p_k)^\nu (1 - \mu p_k)^{M_k - \nu}$$

( $\nu = 0, 1, \dots, \max M_k$ ) eloszlást; ez (3)-tól csak abban különbözik, hogy benne  $p_k$  helyett  $\mu p_k$  szerepel ( $0 < \mu p_k \leq 1$ ).

A nem-negatív monoton növekvő  $\mu_1, \mu_2, \dots, \mu_i, \dots$  értékekkel állítsuk elő sorban a megfelelő  $F_\nu(\mu_i)$ -ket. Tegyük fel, hogy valamelyik  $\mu_i$  épp

egyenlő  $\max p_k = p_\tau$  reciprokával, vagyis  $\mu_i = \frac{1}{p_\tau}$ . Ennél a  $\mu_i$ -nél a (4)

a  $\tau$  indexű komponens egyetlen ponttá elfajult binomiális eloszlás lesz, melynek 0-adik, 1-ső,  $\dots, M_{\tau-1}$ -edik eleme zérus, az  $M_\tau$ -adik pedig  $A_\tau$ .

Ez az egyetlen pont nyilván ki fog ugrani az eloszlás többi pontja fölé és ha szomszédságában már igen kicsik az eloszlás tagjai, ordinátája jó közelítéssel megadja  $A_r$ -t, abszcisszája pedig  $M_r$ -t.  $p_r$  nyilván az alkalmazott  $\mu$  reciproka lesz.

Ezen körülmények mellett tehát egy komponens adatait meghatározhatjuk. Az  $M_r$ -ban közel  $A_r$  magasban kiugró pontot elhagyva és helyette a szomszédos adatokból interpolált új pontot vezetve be, kis hibával olyan keverékre jutunk, amelynek a kiindulónál eggyel kevesebb komponense van, és amelyben  $p_k$  helyett  $p_k/p_r$  paraméterek szerepelnek. Ezen az eljárást megismételhetjük és újabb komponenset választhatunk le, és így tovább. A felbontással nyert adatok természetesen a komponensek számát,  $N$ -et is megadják.

A gyakorlatban legtöbbször egyetlen lépés (sőt, már egyetlen  $\mu$ -értékkel való kísérlet) is elég, mert kimutatható, hogy az eljárás során az egyes komponensek maximumhelyeinek távolsága általában megnő, szórásnégyzeteik pedig általában csökkennek. Ennek a részletes tárgyalására nem térünk itt ki, mert előbbi okoskodásunk értelmében másodrendű fontosságúak<sup>1)</sup>, csak azt jegyezzük meg, hogy szomszédjaitól erősen különvált komponenset egyetlen binomiális eloszlásként kezelhetünk, ekkor pedig paramétereit közelítőleg megállapíthatjuk. A komponensek száma legtöbbször egyetlen  $\mu$ -vel való kísérletből is kiderül.

Állapodjunk meg abban, hogy a felbontást gyakorlatilag elintéztük, ha az  $f_v$  értékek ismeretében általunk megadott  $\mu$ -vel elő tudjuk állítani az  $F_v(\mu)$  mennyiségeket.

Ezt az előállítást megadja a következő

*Tétel: Ha az  $f_v$  diszkrét eloszlás binomiális eloszlások*

$$f_v = \sum_{k=1}^N A_k \binom{M_k}{v} p_k^v (1-p_k)^{M_k-v} \quad (v = 0, 1, \dots, \max M_k)$$

alakú keveréke, a komponensek szétválasztására alkalmas új

$$F_v(\mu) = \sum_{k=1}^N A_k \binom{M_k}{v} (\mu p_k)^v (1 - \mu p_k)^{M_k-v}$$

keverékeloszlást az

$$(5) \quad F_v(\mu) = \mu^v \sum_{\varrho=v}^{\max M_k} \binom{\varrho}{v} (1-\mu)^{\varrho-v} f_{\varrho}$$

összefüggés adja meg.

<sup>1)</sup> A maximumhelyek távolságának változása annak az ismert ténynek az alapján vizsgálható, hogy egy  $M_k$ ,  $\mu p_k$  paraméterű binomiális eloszlás maximuma kb.  $M_k \mu p_k$ -nál van; ilyenformán pl. két egymásra következő maximumhely távolsága  $\mu(M_{k+1} p_{k+1} - M_k p_k)$ , ami  $\mu$ -vel együtt növekedik, ha csak  $\mu$  tényezője  $\neq 0$ . (Analog módon okoskodunk a szórásnégyzetek esetében.)

*Bizonyítás:* Felhasználva az

$$(1 - \mu p_k)^{M_k - \nu} = [(1 - \mu) p_k + (1 - p_k)]^{M_k - \nu} = \\ = \sum_{\sigma=0}^{M_k - \nu} \binom{M_k - \nu}{\sigma} (1 - \mu)^\sigma p_k^\sigma (1 - p_k)^{M_k - \nu - \sigma}$$

azonosságot,

$$F_\nu(\mu) = \sum_{k=1}^N A_k \binom{M_k}{\nu} (\mu p_k)^\nu \sum_{\sigma=0}^{M_k - \nu} \binom{M_k - \nu}{\sigma} (1 - \mu)^\sigma p_k^\sigma (1 - p_k)^{M_k - \nu - \sigma}$$

Vezessük be a  $\nu + \sigma = \varrho$  új összegezési változót; ekkor

$$F_\nu(\mu) = \mu^\nu \sum_{k=1}^N A_k \sum_{\varrho=\nu}^{M_k} \binom{M_k}{\nu} \binom{M_k - \nu}{\varrho - \nu} (1 - \mu)^{\varrho - \nu} p_k^\varrho (1 - p_k)^{M_k - \varrho}.$$

Használjuk fel az

$$\binom{M_k}{\nu} \binom{M_k - \nu}{\varrho - \nu} = \binom{M_k}{\varrho} \binom{\varrho}{\nu}$$

összefüggést, és vegyük figyelembe, hogy az  $F_\nu(\mu)$  nem változik, ha  $M_k$  helyett  $\max M_k$ -ig összegezzünk. Ekkor azonban az összegezések felcserélhetők,

$$F_\nu(\mu) = \mu^\nu \sum_{\varrho=\nu}^{\max M_k} \binom{\varrho}{\nu} (1 - \mu)^{\varrho - \nu} \sum_{k=1}^N A_k \binom{M_k}{\varrho} p_k^\varrho (1 - p_k)^{M_k - \varrho},$$

vagyis, (3) alapján

$$F_\nu(\mu) = \mu^\nu \sum_{\varrho=\nu}^{\max M_k} \binom{\varrho}{\nu} (1 - \mu)^{\varrho - \nu} f_\nu$$

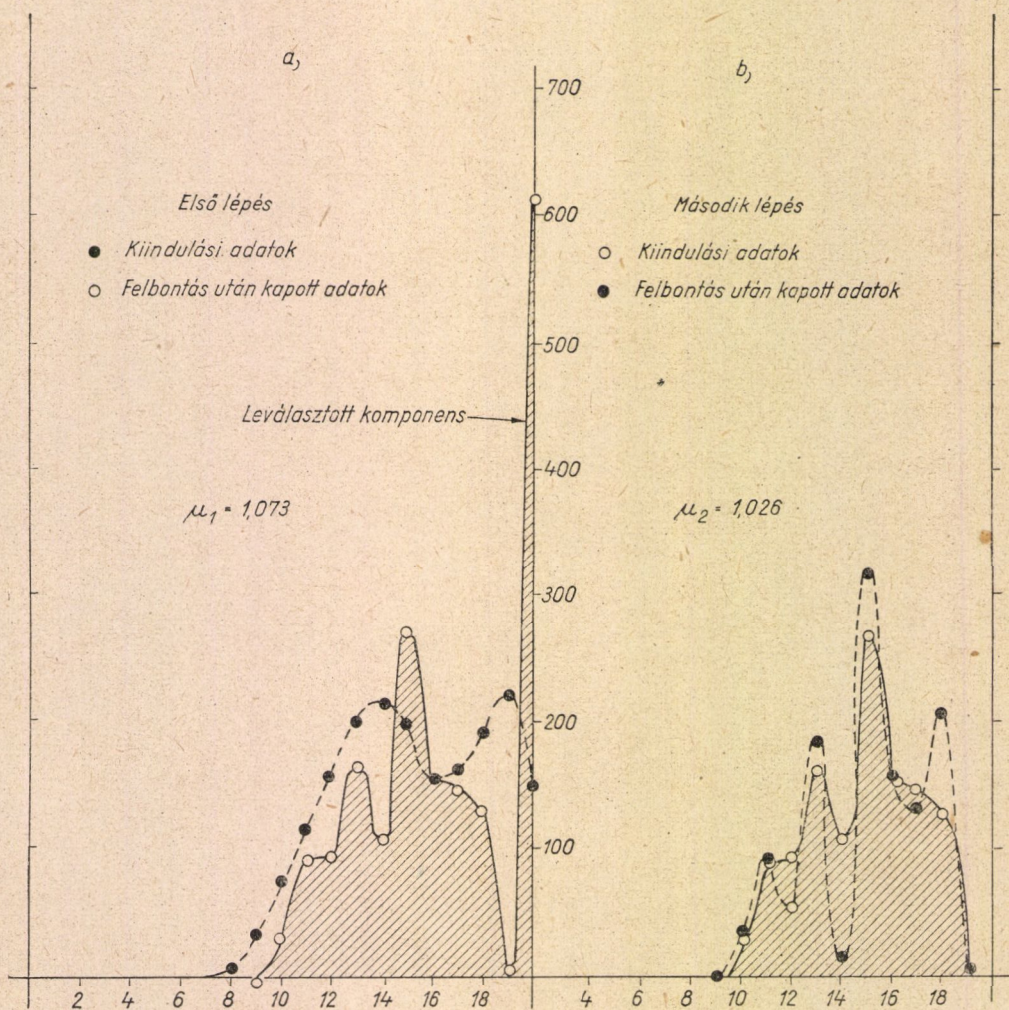
ahogy állítottuk.

Ezek után egyszerűen kínálkozik a keverék felbontásának gyakorlati módszere: Választunk valamilyen  $\mu_1, \mu_2, \dots$  sorozatot és a felbontandó binomiális keverék-eloszlás  $f_\nu$  értékei segítségével (5) felhasználásával kiszámítjuk az  $F_\nu(\mu_1), F_\nu(\mu_2), \dots$  eloszlásokat. (Ha  $\max M_k$  értéke 10–20 körül van, ez géppel gyorsan megy.)

Vegyük észre, hogy  $F_\nu(\mu)$  meghatározásához csak az  $f_\nu, \dots, f_{\max M_k}$  (a  $\nu$ -nél nagyobb indexű) értékek szükségesek. Ha  $\mu_i$  megközelíti az  $1/\max p_k$  értéket, a vonatkozó komponens egyetlen pontként kiugrik a többi közül, ahogy azt fent már részleteztük. Ennek a ténynek a megállapítása az  $F_\nu(\mu_1), F_\nu(\mu_2), \dots$  eloszlások ábrázolása segítségével könnyű. A kiugró komponens paramétereinek megállapítása ugyancsak a fentebb már tárgyalt módon történik. — A gyakorlatban az  $M_k$  értékek általában egyenlők.

(4) alakjából következik, hogy  $\mu p_k > 1$  esetén az  $F_v(\mu)$  komponensei közül egyesek negatív értékeket is felvesznek, vagyis várható, hogy maga  $F_v(\mu)$  is helyenként negatív értékeket vesz fel: Ha tehát az eljárás során valamilyen  $\mu$  értéknél már negatív tagokat is kapunk  $F_v(\mu)$ -ben, biztos, hogy túlléptük  $\mu$ -vel az  $1/\max p_k$  értéket. Ennek a fordítottja azonban már nem mondható ki, s ezért legbiztosabb eljárás több  $\mu$ -vel végzett felbontások eredményeinek összehasonlítása.

Ha a  $p_k$ -k közt egyenlők is vannak, csak az lesz a különbség, hogy a szétválasztás valamelyik lépésénél egyszerre több pont »ugrik ki«. Mivel ez fenti okoskodásunkat nem befolyásolja és az eljárás menete is rögtön látható, ezért részletesen nem tárgyaljuk.



I. ábra

Binomiális eloszlások keverékének felbontása válik szükségessé a [3] dolgozatban (jelen kötet) tárgyalt fizikai-kémiai problémánál. Ilyen eredetű keveréket és felbontását mutatjuk be az 1. ábrán. *Áttekinthetőség kedvéért* folytonos görbét fektettünk át az adatokon. A közölt felbontás  $\mu_1 = 1,073$ , majd — a kiugró komponens leválasztása után nyert eloszlásnál, második lépésként —  $\mu_2 = 1,026$ -tal történt; jól láthatók a különvált komponensek (bizonyos különválás az eredeti felvételen is látszik). Vegyük észre, hogy nem két komponens van, mint azt a kiinduló eloszlásból gondolnánk, hanem öt! A keresett paraméterek a felbontás után kapott eloszlásból:

$$p_1 \approx 0,49 \quad p_2 \approx 0,59 \quad p_3 \approx 0,68 \quad p_4 \approx 0,82 \quad p_5 \approx 0,91$$

$$A_1 \approx 140 \quad A_2 \approx 220 \quad A_3 \approx 560 \quad A_4 \approx 280 \quad A_5 \approx 610,$$

amelyek jól egyeznek az itt történetesen gyakorlatból is ismert értékekkel.

Természetesen ajánlatos ellenőrzésként az így kiszámított paraméterekkel néhány eloszlás-pontot meghatározni, és egybevetni a kísérleti adatokkal. Bemutatjuk még az

$$f_\nu = \binom{10}{\nu} 0,4^\nu (1-0,4)^{10-\nu} + 2 \binom{10}{\nu} 0,5^\nu (1-0,5)^{10-\nu} \quad (\nu = 0, 1, \dots, 10)$$

keverék felbontását is. Itt

$$N = 2, \quad A_1 = 1, \quad A_2 = 2, \quad M_1 = M_2 = 10, \quad p_1 = 0,4, \quad p_2 = 0,5.$$

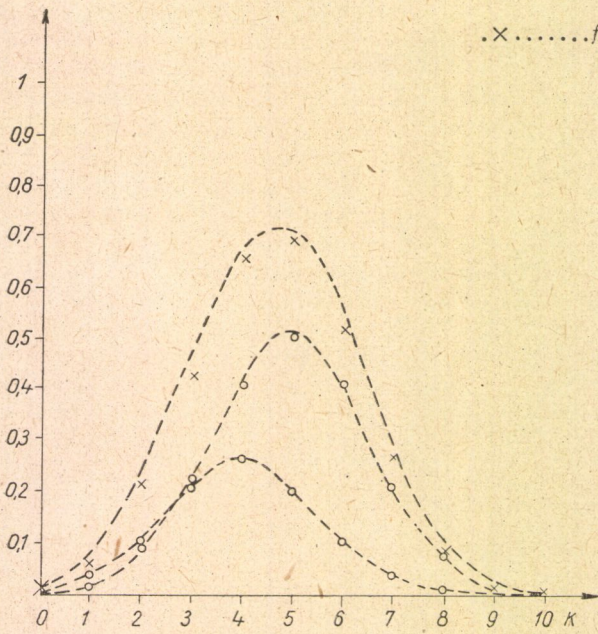
A 2. a. ábra mutatja a két komponens külön (o) és összegüket (x), a 2. b pedig összegüket (o) és a  $\mu = \frac{10}{6}$ -dal történt felbontás eredményét (x),

vagyis az  $F_\nu\left(\frac{10}{6}\right)$  értékeket (itt is Gauss-görbét fektettünk át az adatokon, a szemléltetés céljára). Az  $f_\nu$  keverék nem árul el semmit a komponensekről, az  $F_\nu\left(\frac{10}{6}\right)$  értékek viszont már élesen különvált két eloszlást mutatnak, bár még nem járunk az extrémális esethez.

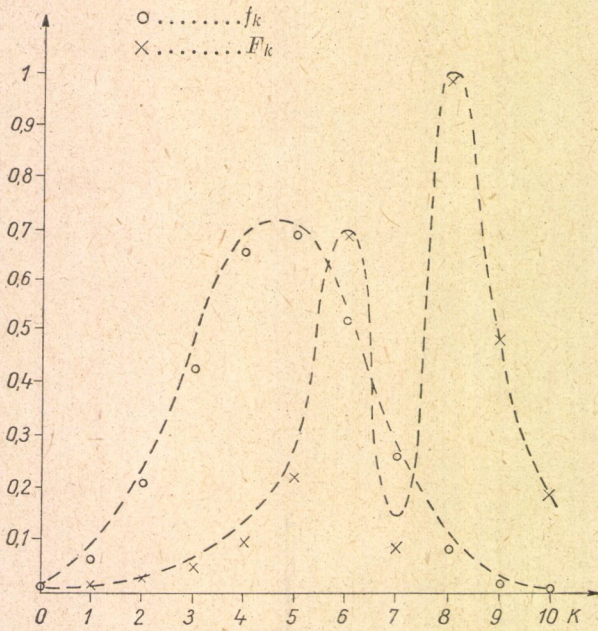
## 2. §

Az előző §-ban tárgyalt eljárás gyakorlati végrehajtásakor több nehézségre bukkanunk. Az egyik: a kísérleti adatok vagy 0, vagy  $\max M_k = M$  szomszédságában (esetleg mindkét helyen is) sokszor hibahatáron belül kicsik, és ezért a kutatók zérusnak veszik őket, hiszen nincs alap más érték hozzárendelésére.

Ha csak a bal-, illetve jobbszáron ilyenek az adatok, a nehézség még áthidalható, mert az  $F_\nu(\mu)$ -k szempontjából csak a  $\nu$ -nél *nagyobb* indexű  $f_\nu$ -értékek fontosak (vesd össze: (5)), vagyis hibás *balszéli* adatok legfeljebb az első  $F_\nu(\mu)$  értékeket rontják el, ami a felbontást alig zavarja, hibás *jobbszéli* adatoknál pedig egyszerűen tükrözzük az egész eloszlást a 0, ...,  $\max M_k$  intervallum felező



2. a. ábra



2. b. ábra

merőlegese körül, mikoris a keverék lényegileg változatlan marad, a használhatatlan adatok pedig a balszélre kerülnek, s így lényegtelené válnak. Kritikus tehát csak az az eset marad, amidőn mindkét szélen használhatatlanok az adatok. Számunkra ez azt jelenti, hogy valamilyen  $l$ -nél és  $j$ -nél az  $f_0, f_1, \dots, f_l$  és  $f_M, f_{M-1}, \dots, f_{M-j}$  értékeket is 0-nak kell vennünk, holott értékük véges, bár igen kicsiny. Ezzel azonban (5) alkalmazásakor igen nagy hibákat követhetünk el: pl. az  $F_M(\mu) = \mu^M f_M$  értéket 0-nak kell tekintenünk, holott nagy  $M$  esetén  $\mu > 1$  folytán  $\mu^M$  nagy szám és így  $F_M(\mu)$  kis  $f_M$ -nél is elég nagy marad. Egyenlő  $M_k$  értékeknél éppen ebben a pontban várhatnánk a kiugró eloszlás-értéket, az említett elhanyagolás azonban épp ezt (meg a szomszédos) adatokat teszi használhatatlanná.

Olyan eljárást kell tehát keresnünk, amely szélső adatokra *nem* támaszkodva szolgáltatja a felbontást vagy mindjárt az egyes adatokat.

Bebizonyítható, hogy a fentebb tárgyalt szétválasztási eljárás nem alkalmazható, ha nem használhatjuk fel az összes  $f_\nu$ -értékeket. Közelfekvő gondolat megpróbálni az  $f_\nu$  értékek olyan  $\Phi_s(f_1, \dots, f_M)$  ( $s = 0, 1, \dots$ ) függvényeit előállítani, amelyeknek értékét alig befolyásolják a szélső  $f_\nu$ -k értékei, amellett a  $\Phi_s(f_1, \dots, f_M)$  értékekből vagy következtetni tudunk az ismeretlen paraméterekre, vagy visszavezethetjük a problémát a 2. §-ban tárgyalt módszer alkalmazására. (Itt tulajdonképpen globálisan kezelnénk az  $f_\nu$ -ket, ami a kísérleti adatok feldolgozásánál azért előnyös, mert a hibák kevésbé befolyásolják a végeredményt.) Végeredményben az  $f_\nu$ -k sorozatát kellene áttranszformálnunk

a  $\Phi_s(f_1, \dots, f_M)$  sorozatba egy  $\Phi_s(f_1, \dots, f_M) = \sum_{k=1}^M B_k(s) \cdot f_k$  típusú transzformáció segítségével, ahol azonban a  $B_k(s)$  transzformáló sorozatot feltételeinknek megfelelően kell meghatározni. Sajnos, minden kikötésünknek eleget tevő  $B_k(s)$  sorozatot eddig nem sikerült találnunk.

*Ha tudjuk, hány komponenst várhatunk (a gyakorlatban sokszor ez az eset) és az összes  $M_k$ -k egyenlőek, sikerrel alkalmazható a következő, algebrai eljárás:*

Tekintsük ismét a keverék (3) alatti alakját ( $M_k = n$ ):

$$(6) \quad f_\nu = \sum_{k=1}^N A_k \binom{n}{\nu} p_k^\nu (1-p_k)^{n-\nu}, \quad (\nu=0, 1, \dots, n).$$

Ez felfogható magasabbfokú egyenletrendszernek az  $A_k$  és  $p_k$ -k értékek meghatározására, amely  $2N \leq n$  esetén meg is oldható. A megoldást egy már régebben kidolgozott módszer<sup>2)</sup> általánosítása adja meg: Vezessük be az új  $D_\nu, B_k, x_k$  paramétereket, amelyeket az

$$\binom{f_\nu}{\nu} = D_\nu, \quad p_k = \frac{x_k}{1+x_k}, \quad B_k = \frac{A_k}{(1+x_k)^n}$$

összefüggések definiálnak. Ekkor (6) így alakul:

$$(7) \quad D_\nu = \sum_{k=1}^N B_k \cdot x_k^\nu, \quad (\nu = 0, 1, \dots, n),$$

<sup>2)</sup> Lényegében szerepel már Gauss-nál: [4]. Lásd még: [5], p. 275 — Alkalmazását illetően lásd: [6], főleg p. 713.



ahol a  $B_k, x_k$  mennyiségek az ismeretlenek. Mivel a  $B_k, x_k$  mennyiségekből az  $A_k$  és  $p_k$  értékek egyértelműen kiszámíthatók, elég a (7) egyenlettel foglalkoznunk. Olyan megoldási módszert keresünk, amely nem használja fel a szélső  $f$ , értékeket, illetve a megfelelő indexű  $D$ , értékeket.

Tekintsünk pozitív egész számokból álló olyan

$$a_1 < a_2 < \dots < a_{N+1}$$

$$\beta_1 < \beta_2 < \dots < \beta_{N+1}$$

$$\gamma_1 < \gamma_2 < \dots < \gamma_{N+1}$$

.....

.....

$$\tau_1 < \tau_2 < \dots < \tau_{N+1}$$

.....

.....

rendszert, amelyre

$$a_2 - a_1 = \beta_2 - \beta_1 = \gamma_2 - \gamma_1 = \dots = \tau_2 - \tau_1 = \dots$$

$$a_3 - a_1 = \beta_3 - \beta_1 = \gamma_3 - \gamma_1 = \dots = \tau_3 - \tau_1 = \dots$$

$$(8) \quad \dots \dots \dots$$

$$a_{N+1} - a_1 = \beta_{N+1} - \beta_1 = \gamma_{N+1} - \gamma_1 = \dots = \tau_{N+1} - \tau_1 = \dots,$$

és tegyük fel, hogy

$$(9) \quad a_{N+1} \leq n, \quad \beta_{N+1} \leq n, \quad \gamma_{N+1} \leq n, \quad \dots, \quad \tau_{N+1} \leq n, \quad \dots$$

Szemeljük ki a (7) rendszerből először a  $\nu = a_1, a_2, \dots$  indexű egyenleteket:

$$B_1 x_1^{a_1} + \dots + B_N x_N^{a_1} = D_{a_1}$$

$$B_1 x_1^{a_2} + \dots + B_N x_N^{a_2} = D_{a_2}$$

$$(10) \quad \dots \dots \dots$$

$$B_1 x_1^{a_{N+1}} + \dots + B_N x_N^{a_{N+1}} = D_{a_{N+1}}$$

Szorozzuk ezeket sorban az egyelőre ismeretlen  $c_1, c_2, \dots, c_N$  értékkel, az utolsót 1-gyel s adjuk össze a jobb- és baloldalakat:

$$\begin{aligned}
 & B_1 x_1^{a_1} (c_1 + c_2 x_1^{a_2 - a_1} + \dots + c_N x_1^{a_N - a_1} + x_1^{a_{N+1} - a_1}) + \\
 & + B_2 x_2^{a_1} (c_1 + c_2 x_2^{a_2 - a_1} + \dots + c_N x_2^{a_N - a_1} + x_2^{a_{N+1} - a_1}) + \dots \\
 (11) \quad & \dots + B_N x_N^{a_1} (c_1 + c_2 x_N^{a_2 - a_1} + \dots + c_N x_N^{a_N - a_1} + x_N^{a_{N+1} - a_1}) = \\
 & = c_1 D_{a_1} + c_2 D_{a_2} + \dots + c_N D_{a_N} + D_{a_{N+1}}
 \end{aligned}$$

Válasszuk meg  $c_1, \dots, c_N$ -t úgy, hogy az  $x_1, x_2, \dots, x_N$  értékek épp a

$$(12) \quad c_1 + c_2 x_1^{a_2 - a_1} + \dots + c_N x_1^{a_N - a_1} + x_1^{a_{N+1} - a_1} = 0$$

egyenlet gyökei legyenek. Ekkor (11) baloldala zérus, vagyis a  $c_i$  értékek közt fenn kell állnia a

$$(13) \quad D_{a_1} c_1 + D_{a_2} c_2 + \dots + D_{a_N} c_N = -D_{a_{N+1}}$$

összefüggésnek.

Vegyük most a (7) rendszerből a  $\nu = \beta_1, \beta_2, \dots$  indexű egyenleteket, szorozzuk meg őket ismét az előbbi  $c_1, c_2, \dots$  értékekkel és adjuk össze:

$$\begin{aligned}
 & B_1 \cdot x_1^{\beta_1} (c_1 + c_2 x_1^{\beta_2 - \beta_1} + \dots + c_N x_1^{\beta_N - \beta_1} + x_1^{\beta_{N+1} - \beta_1}) + \dots \\
 (14) \quad & \dots \\
 & = c_1 D_{\beta_1} + c_2 D_{\beta_2} + \dots + c_N D_{\beta_N} + D_{\beta_{N+1}}.
 \end{aligned}$$

Mivel azonban (8) fennállásából indultunk ki,

$$c_1 + c_2 x_1^{\beta_2 - \beta_1} + \dots + x_1^{\beta_{N+1} - \beta_1} \equiv c_1 + c_2 x_1^{a_2 - a_1} + \dots + x_1^{a_{N+1} - a_1}$$

(és így tovább a többi  $x_i$  változóra).

Ezek azonban feltevés szerint zérusok, következésképp (12) baloldala is zérus és a  $c_i$  értékekre a következő újabb összefüggést kapjuk:

$$(15) \quad D_{\beta_1} c_1 + D_{\beta_2} c_2 + \dots + D_{\beta_N} c_N = -D_{\beta_{N+1}}.$$

Hasonlóan járunk el a  $\gamma_1, \gamma_2, \dots$  indexű egyenletekkel, végeredményül a

$$(16) \quad D_{\gamma_1} c_1 + D_{\gamma_2} c_2 + \dots + D_{\gamma_N} c_N = -D_{\gamma_{N+1}}$$

összefüggést kapjuk. Újabb, a (8) és (9)-nek megfelelő  $\tau_1, \tau_2, \dots$  indexrendszerhez tartozó egyenleteket választva ki (7)-ből, ismét más összefüggést nyerünk a  $c_i$ -k közt. Az eljárást mindaddig folytatjuk, míg (12), (14), (15)

és a többi egyenlet  $N$  számú egyenletet nem szolgáltat. Ezek együtt lineáris egyenletrendszert adnak a  $c_i$  értékekre.

*Tegyük fel*, hogy ez az egyenletrendszer megoldható (ennek a vizsgálatára nem térünk ki; ha a gyakorlatban valamilyen  $a_i, \beta_i, \dots$  értékrendszerrel esetleg zérus volna az egyenletrendszer determinánsa, más értékrendszerrel próbálkozunk). — A megoldással nyert  $c_i$  értékeket beírjuk a (12) egyenletbe s az így kapott  $(a_{N+1} - a_1)$ -edfokú egyismeretlenes egyenletet pontosan vagy közelítő módszerrel megoldjuk. A gyökök szolgáltatják az  $x_i$  értékeket, ezeket a (7) alatti egyenletrendszer  $N$  darab egyenletébe beírjuk és az ily módon a  $B_i$ -értékekre kapott lineáris egyenletrendszert megoldjuk. Az  $x_i$  és  $B_i$  értékekből a keresett  $p_k, A_k$  értékek könnyen meghatározhatók.

Az ismeretlenek meghatározását tehát visszavezettük két  $N$ -ismeretlenes lineáris egyenletrendszer és egy  $(a_{N+1} - a_1)$ -edfokú egyismeretlenes egyenlet megoldására, mindehhez a szükséges kísérleti adatokat elég tág határok között választhattuk ki. A kísérleti adatokból akkor kell a legkevesebb, ha  $\beta_1 = a_2, \beta_2 = a_3, \dots, \gamma_1 = a_3, \gamma_2 = a_4, \dots$ , és így tovább, vagyis ha az egyes adatok egyenlő távolságú pontokból valók. Ez azonban meg is köti a pontok választását; az eredeti (8) értékrendszer jobban biztosítja, hogy megbízható adatokat használhassunk fel. A (8) rendszert mindig az adatok megbízhatóságát figyelembevéve állítsuk össze. Természetesen csak addig használható a leírt eljárás, amíg (9) fennáll, ez azonban a gyakorlatban legtöbbször teljesül.

Általában ajánlatos többféleképpen megválasztott adatseregből többször határozni meg az ismeretlen paramétereket, s azután közepelni az eredményeket.

2–3 komponensnél ez az eljárás jól alkalmazható. Legfőbb haszna, hogy akkor is alkalmazható, ha az  $f_v$ -értékeknek csak egy része ismeretes; ez gyakran előfordul, mert a gyorsaság kedvéért nem mindig mérik ki az összes  $f_v$  értéket.<sup>3)</sup>

E paragrafus záradékaul még a következőket említjük meg:

Ha  $\max M_k$  illetve  $n$  nagy (esetleg több száz), a szóráscsökkentő eljárást megláthatja a nagy számolási munka. Az algebrai módszer ekkor is használható, feltételezve, hogy ismerjük a komponensek számát. Gondolhatnánk arra is, hogy az eloszlást lépcsősfüggvénnyé alakítjuk át (egységnyi vízszintes szakaszokat illesztünk az egyes pontokhoz), s ezt a függvényt valamilyen ortogonális függvényrendszer szerint kifejtjük, majd összefüggést keresünk a kiindulási és a csökkentett szórású keverék sorfejtésének együttműködési között. Sajnos, Fourier-sor alkalmazásakor ez az ötlet keresztülvihetetlen. Másirányú vizsgálataink még folyamatban vannak, s ezért nem időzünk tovább ennél a kérdésnél, csak azt említjük meg, hogy nagy  $\max M_k$  illetve  $n$  esetén a binomiális eloszlások már jól közelíthetők a normális sűrűségfüggvénnyel, vagyis az eloszlás pontjait összekötő folytonos görbe jó közelítéssel Gauss-függvények keverékének tekinthető, erre pedig alkalmazhatjuk az [1]-ben leírt eljárásokat.

<sup>3)</sup> A (7) magasabbfokú egyenletrendszer egyike azoknak, amelyeket »ad hoc« módszerrel meg lehet oldani, illetve egyszerűbb egyenletekre vissza lehet vezetni. Az irodalomban nem sok ilyen speciális, megoldható típus ismeretes; érdemes volna az ilyeneket és megoldási módszerüket összegyűjteni.

### 3. §.

A következőkben célkitűzésünknek megfelelően Poisson-eloszlások konstans súlyokkal vett

$$(17) \quad f_v = \sum_{k=1}^N A_k \cdot e^{-\lambda_k} \frac{\lambda_k^v}{v!} \quad (v = 0, 1, \dots)$$

keverékének  $\left( A_k > 0, \sum_{k=1}^N A_k = 1, \lambda_k > 0, \text{ a } \lambda\text{-értékek különbözök} \right)$  felbontásával foglalkozunk: az  $f_v$  értékek ismeretében meg akarjuk határozni az ismeretlen  $A_k$  és  $\lambda_k$  paramétereket.

A binomiális eloszlásnál bemutatott szóráscsökkentés, illetve algebrai módszerek lényegében itt is alkalmazhatók. Itt is elég különböző  $\lambda_k$  értékek esetével foglalkozunk.

A szóráscsökkentés módszer ebben az esetben azon az észrevételen alapszik, hogy a  $\mu < \min \lambda_k$  paraméterrel képezett

$$(18) \quad F_v(\mu) = \sum_{k=1}^N A_k \cdot e^{-(\lambda_k - \mu)} \frac{(\lambda_k - \mu)^v}{v!} \quad (v = 0, 1, \dots)$$

keverék a (17) alattitól csak abban különbözik, hogy a komponensek maximális eloszlásértéke ebben nagyobb, szórása kisebb, és  $\mu \rightarrow \min \lambda_k = \lambda_\tau$  esetén a  $\tau$  indexű komponens a  $v = 0$ -nál  $A_\tau$  magasságban jelentkező egyetlen ponttá fajul el, ami paraméterek megállapítását, komponens-leválasztást, további kezelést épp oly módon tesz lehetővé, mint ahogy azt a binomiális eloszlás esetében már részletesen leírtuk.

Tekintsük problémánkat itt is megoldottnak, ha az  $f_v$  és általunk választott, monoton növekvő  $\mu$ -értékek segítségével elő tudjuk állítani az  $F_v(\mu)$  kifejezéseket.

Erről az előállításról szól a következő

*Tétel: Poisson-eloszlások*

$$(19) \quad f_v = \sum_{k=1}^N A_k e^{-\lambda_k} \frac{\lambda_k^v}{v!}$$

keveréke és az ismeretlen  $A_k, \lambda_k$  paraméterek meghatározására szolgáló

$$(20) \quad F_v(\mu) = \sum_{k=1}^N A_k \cdot e^{-(\lambda_k - \mu)} \frac{(\lambda_k - \mu)^v}{v!}$$

( $\mu < \min \lambda_k$  általunk választott valós szám) keverék közt a következő kapcsolat áll fenn:

$$(21) \quad F_v(\mu) = e^\mu \sum_{q=0}^v \frac{(-\mu)^{v-q}}{(v-q)!} f_q$$

*Bizonyítás:* Nyilván

$$F_\nu(\mu) = e^\mu \sum_{\varrho=0}^{\nu} \binom{\nu}{\varrho} \frac{1}{\nu!} (-\mu)^{\nu-\varrho} \sum_{k=1}^N A_k \cdot e^{-\lambda_k} \lambda_k =$$

$$= e^\mu \sum_{\varrho=0}^{\nu} \frac{(-\mu)^{\nu-\varrho}}{(\nu-\varrho)!} \sum_{k=1}^N A_k \cdot e^{-\lambda_k} \frac{\lambda_k^\varrho}{\varrho!} = e^\mu \sum_{\varrho=0}^{\nu} \frac{(-\mu)^{\nu-\varrho}}{(\nu-\varrho)!} f_\varrho$$

mivel a szummációkat felcserélhetjük, q. e. d.

Vegyük észre, hogy  $F_\nu(\mu)$  meghatározásához az *első*  $f_\nu$ -kre van szükség. Ha a legelső  $f_\nu$ -k megbízhatatlanok, vagy nem mindegyiket ismerjük — viszont ismerjük a komponensek számát —, 2–3 komponensnél itt is alkalmazható algebrai eljárás, épp úgy, mint a binomiális eloszlásoknál láttuk. Az

$$(22) \quad f_\nu = \sum_{k=0}^N A_k \cdot e^{-\lambda_k} \frac{\lambda_k^\nu}{\nu!} \quad (\nu = 0, 1, \dots)$$

értékeket az  $A_k \cdot e^{-\lambda_k} = G_k$ ,  $\nu! f_\nu = H_\nu$  relációkkal bevezetett új  $G_k$ ,  $H_\nu$  paraméterek segítségével átvisszük a

$$(23) \quad \sum_{k=1}^N G_k \cdot \lambda_k^\nu = H_\nu$$

egyenletrendszerbe. Ez azonban ugyanúgy kezelhető, mint a (7) egyenletrendszer. Gyakorlati felhasználására is az ott mondottak vonatkoznak, s így itt mindezt felesleges lenne megismételni.

#### IRODALOM

[1] MEDGYESSY P.: »Valószínűség-eloszlásfüggvények keverékének felbontása összetevőikre.« *A Magyar Tudományos Akadémia Alkalmazott Matematikai Intézetének Közleményei* 2 (1953) 165—177.

[2] MEDGYESSY P.: »Újabb eredmények valószínűség-eloszlásfüggvények keverékének összetevőire bontásával kapcsolatban.« *A Magyar Tudományos Akadémia Alkalmazott Matematikai Intézetének Közleményei* 3 (1954) 155—169 (jelen kötet).

[3] MEDGYESSY P.—RÉNYI A.—TETTAMANTI K.—VINCZE I.: »A frakcionáló megosztás matematikai tárgyalása nem-teljes diffúzió esetében.« *A Magyar Tudományos Akadémia Alkalmazott Matematikai Intézetének Közleményei* 3 (1954) 81—97 (jelen kötet).

[4] C. F. GAUSS: »Methodus nova integralium valores per approximationem inveniendi.« — *Werke*, Königliche Gesellschaft der Wissenschaften, Göttingen, 1866, Band 3. 165—196.

[5] C. RUNGE—R. KÖNIG: *Vorlesungen über numerisches Rechnen*. Springer, Berlin, 1924.

[6] G. DOETSCH: »Die Elimination des Dopplereffekts bei spektroskopischen Feinstrukturen und exakte Bestimmung der Komponenten.« *Zeitschrift für Physik* 49 (1928) 705—730.

## РАЗЛОЖЕНИЕ НА КОМПОНЕНТЫ СМЕСИ ДИСКРЕТНЫХ РАСПРЕДЕЛЕНИЙ ВЕРОЯТНОСТЕЙ

П. Медеши

### Резюме

Настоящая работа посвящается следующим проблемам:

a) Известное дискретное распределение вероятностей является наложением биномиальных распределений с различными параметрами, (то есть смесью с постоянными весами), то есть оно имеет вид (3). Зная  $f_v$ , следует определить неизвестные параметры  $A_k$ ,  $M_k$ ,  $p_k$  компонентов.

b) Известное дискретное распределение является наложением (17) распределений Пуассона с различными параметрами. Следует определить неизвестные параметры  $A_k$ ,  $\lambda_k$ . Такие смеси часто встречаются при различных исследованиях. Проблема a) возникает, например, при т. н. фракционирующем разделении в химической промышленности (см. список литературы в (3)).

При решении обеих проблем используем, главным образом, т. н. »метод уменьшения дисперсии«, примененный в (1) и (2). В случае a) это приводит к следующему:

Зная значения  $f_v$  из (3), мы в состоянии представлять распределение  $F_v(\mu)$  (4). Оно отличается от (3) только тем, что все значения  $p_k$  умножены на выбранный нами параметр  $\mu$ . Пусть распределения  $\mu$  представляются при помощи возрастающих  $F_v(\mu)$ . Если параметром  $\mu$  достигается значение  $1/\max p_k$ , то компоненты  $F_v(\mu)$  с параметром вырождаются в единственную выступающую точку, из данных которой параметры соответствующих компонентов определяются хорошим приближением. В то же время уменьшается и дисперсия других компонент. После отделения вырожденных компонент, предлагаемый прием повторяется.

Величины  $F_v(\mu)$ , представляющие решение проблемы, суть линейные выражения, составленные из  $f_v$  формулой (5).

В качестве иллюстрации приводим два примера (рис. 1 и 2). На рис. 1 исходная смесь распределение изображается полными точками (это распределение характеризует результат фракционирующего разделения). Результат разложения, соответствующего значению  $\mu = 1,073$ , представляется неполными точками. Видно, что мы имеем не 2, а 5 компонент. Это подтверждается и разложением, проведенном после отщепления сильно выступающей последней компоненты (рис. 1. b). — На рис. 2. a. и 2. b. показывается разделение фиктивной смеси. Смесь имеет единственный максимум; после разделения отдельные компоненты отчетливо проявляются.

Описанный прием может найти применение лишь при известных  $f_v$ . В случае, если известна лишь часть этих данных, можно применять следующий способ: Пусть нам известно число  $N$  компонент и  $M_k = n$  ( $k = 1, 2, \dots, N$ ). Запишем для достаточного числа точек сложного распределения соотношение (6). Таким образом мы приходим к системе уравнений для  $A_k$  и  $p_k$ , которую преобразуем к форме (7).

Предположим, что известные нам  $f_v$  делают возможным записать (7) для систем  $v = \alpha_k, \beta_k, \gamma_k, \dots$  соответствующим (8) [(10), (14)]. Решение полученной таким образом системы уравнений приводится при помощи известного приема [5] к решению уравнения высшей степени и двух систем линейных уравнений. Из решений уже можно вычислить искомые  $A_k$  и  $p_k$ .

Следует заметить, что для больших  $n$ , смесь (6) с достаточной точностью представляется как смесь нормальных функций плотности. В этом случае применимы методы, изложенные в [1] и [2].

В случае разложения смеси (17) распределений Пуассона (случай b) мы также применяем изложенный метод. Основой разложения служит представление смеси (18) при помощи (17). Здесь параметр  $\mu$  следует вычитать из  $\lambda_k$ . Компоненты (18) разделяются в случае  $\mu = \min \lambda_k$  таким же образом, как и компоненты смеси a). Связь между (17) и (18) устанавливается формулой (21). При надобности можно прибегнуть и в данном случае к алгебраическим методам.

# DECOMPOSITION OF DISCRETE COMPOUND PROBABILITY DISTRIBUTIONS

P. MEDGYESSY

## Summary

This paper deals with the following problems:

a) A given discrete probability distribution is the superposition (a compound) of binomial distributions with different parameters, — i. e. it is of the form (3). The probabilities  $f_\nu$  being given, we have to determine the unknown parameters  $A_k$ ,  $M_k$ ,  $p_k$  of the components.

b) A given discrete probability distribution is a superposition (17) of Poisson-distributions with different parameters; the unknown parameters  $A_k$ ,  $\lambda_k$  are to be determined.

We often meet such compound distributions in various scientific investigations. E. g. the problem a) arises in connection with the procedure called «counter-current distribution» used in chemical industry (see the bibliography of [3]).

For the solution of both problems we primarily apply the so called «variance reduction method» used in [1] and [2]. In the case of a) the application of this method can be done as follows: From the knowledge of the probabilities  $f_\nu$  (see (3)) we can construct the distribution  $F_\nu(\mu)$  given by (4). This differs from (3) so far as all the  $p_k$ -s in it are multiplied with the parameter  $\mu$ .  $\mu$  will be chosen arbitrarily (but so that  $\mu p_k \leq 1$ ). When increasing  $\mu$  we reach  $1/\max p_k$ , those components of  $F_\nu(\mu)$  which have the parameter  $\max p_k$  degenerate into a single point (see fig. 1. a.). If these points appear separately their coordinates furnish approximate values of the parameters  $A_k$ ,  $M_k$ . — Simultaneously, the variances of all the other components will also diminish. After the separation of the degenerated components the procedure can be repeated.

The values  $F_\nu(\mu)$  yielding the solution are given by (5) as linear forms of the probabilities  $f_\nu$ .

The method is illustrated by two examples (Fig. 1. and 2.) On Fig. 1. the full points show the original compound distribution (it is the result of some counter-current distribution). The result of a transformation with  $\mu = 1,073$  is marked by circles. It is easy to see that there are not 2 but 5 components. This is also justified by a second decomposition performed after the separation of the last component. (Fig. 1. b.)

Fig. 1. a. and b. show the decomposition of an artificial compound distribution. There is only one maximum; after the decomposition the separate components present themselves very markedly.

The described method uses all the probabilities  $f_\nu$ . If they are not all available, then instead of the preceding procedure we may apply the following method [inasmuch the number  $N$  of the components is known and  $M_k = n$  ( $k = 1, 2, \dots, N$ )]. We write (6) for an adequate set of the points of the compound distribution. Then we obtain a system of equations for the  $A_k$  and  $p_k$ -s, which can be brought to the form (7). Suppose that the known  $f_\nu$ -s make possible to write (7) for the systems of values  $\nu = \alpha_k, \beta_k, \gamma_k, \dots$  corresponding to (8) (see (10), (14)). By the aid of a well-known technique [5] the solution of the system of equations thus obtained can be reduced to that of an equation of higher degree with an only unknown and of two systems of linear equations. From this solution, the  $A_k$  and  $p_k$ -s can already be determined.

We remark that in the case of larger  $n$ , the compound distribution (6) can very well be approximated by a compound of normal frequency functions; then the methods given in [1] and [2] may be applied.

Similarly, for the decomposition of a compound (17) of Poisson-distributions (case b)), the above method can be applied. The basis of a decomposition is the construction of the compound distribution (18) by the aid of (17). Here the parameter  $\mu$  will be subtracted from the  $\lambda_k$ . In the case of  $\mu = \min \lambda_k$  the components of (18) will generally separate from each other like those of the compound a). The relation between (17) and (18) is given by (21). If necessary an algebraic procedure may be applied in this case also.