

MTA SZTAKI

Az 1974 óta végzett kórházi morbiditásvizsgálat által fel-  
vetett statisztikai kérdések

Ratkó István - Soltész János

Bevezetés

Az MTA SZTAKI Valószínűségszámítási és Matematikai Statisztikai Osztályának munkatársai 1973 óta vesznek részt a kórházi morbiditási vizsgálatokban. Ennek kapcsán különböző matematikai statisztikai kérdések merültek fel, melyek megválaszolása - bármilyen orvosstatisztikai adatfeldolgozás alkalmával - a kapott eredmények értékelése, megbízhatósága szempontjából alapvető fontosságú. Erről részletebben [1]-ben beszámoltunk. Most csak négy problémára szeretnénk rávilágítani. A felhasznált adatok az 1972-73 évi kórházi morbiditási vizsgálat adatai.

1. Adott százaléku minta biztosítása

Mintavételnél két fontos szempontot kell figyelembe venni: a minta lehetőleg a teljes populáció egy meghatározott hányada legyen, hogy a mintából könnyen következtethessünk a teljes populációra és ugyanezért fontos az is,

hogy a minta reprezentatív legyen.

Az évenkénti kórházi morbiditási adatokból osztályonként 10 %-os mintát kell kiválasztani. A minta pontossága az elsődleges cél, még ha ez néhány többszörösen ápoltságos személy egyes adatainak elvesztését vonja is maga után.

Korábbi statisztikai vizsgálatok igazolták azt a természetes feltevést, hogy a morbiditási adatok függetlenek a születésnaptól. A betegek születésnap szerinti eloszlása, ha a hónapot nem vesszük figyelembe, egyenletesnek tekinthető fel [2]. Hány születésnap /havonként/ szükséges a 10 %-os minta biztosításához feltéve, hogy a kórházak a kijelölt napokon született valamennyi beteg fejlapját beküldik. A Moivre-Laplace tétel és a Bernstein egyenlőtlenség felhasználásával bizonyítható, hogy három nap 0,5 két nap 0,00008 négy nap 0,9974 valószínűséggel elegendő a 10 %-os mintához.

A kórházak négy születésnap betegeinek lapjait /4.-i, 14.-i, 24.-i és 6.-i/ küldik el. Ezekből egy már megírt program segítségével választja ki a számítógép a 10 %-os mintát egy egyszerű itt nem részletezett algoritmus alapján.

Meggondolásaink, számolásaink nem 10 %-os mintavételnél is alkalmazhatók.

## 2. A többszörös ápolással kapcsolatos mintavételi problémák

Azt vizsgáljuk meg először, milyen hibák adódnak, amikor a mintánk alapján a többszörösen ápolott betegek számát akarjuk megbecsülni. Most csak a speciális kérdésfeltevésből és a mintavétel sajátosságaiból adódó hibákkal fogunk foglalkozni.

A probléma a következő: tegyük fel, hogy egy A osztályon ápolott olyan esetek száma, amelyeknél a beteg 4.-én, 14.-én vagy 24.-én született kisebb, mint  $0,1 M_A$ . Ekkor a mintába beveszünk még néhány 6.-án született, A osztályon kezelt beteget. Tegyük fel továbbá, hogy a 4.-én, 14.-én vagy 24.-én született B osztályon kezelt betegek eseteinek száma nagyobb, mint  $0,1 M_B$ . Ekkor el kell hagynunk néhány 4.-én született beteg esetét. Számoljuk most össze, hogy a mintában hány olyan beteg van, akit előbb az A osztályon, majd a B osztályon kezeltek. Az ilyen betegeket az egyszerűség kedvéért AB betegeknak fogjuk hívni. A fenti feltevések esetén a mintából az AB betegekre adódó becslés valószínűleg kisebb lesz a pontos értéknél, mivel elvesznek azok a betegek, akik 4.-én születtek, de a B osztály mintájából kihagytuk őket. Akik 6.-án születtek és bekerültek az A osztály mintájába, azok is elvesznek, ugyanis a minta alapján nem lehet megállapítani, hogy őket a későbbiek so-

rán a B osztályon kezelték.

A következő modellel fogunk dolgozni:

Az ápolási esetek rendezve vannak születésnap /először vannak a 14.-én, majd a 24.-én, 4.-én, 6.-án, stb. született emberek, majd azokon belül egyéb azonosítók /születési év, hó, név, anyja neve, stb./ szerint. Így minden kétszer ápolat beteg két esete egymás mellé kerül. Ez a feltevés nem jelent megszorítást a kórházi morbiditási adatok statisztikai viselkedésére vonatkozóan. Egy beteg  $p_1, p_2, \dots, p_{20}$  valószínűséggel kerül az 1., 2., ..., 20. osztályra. Ha kétszer kezelik, akkor a második alkalommal az elsőtől függetlenül kerül  $p_1, \dots, p_{20}$  valószínűséggel a megfelelő osztályra. Egy beteg  $i$ -én  $1/30$  valószínűséggel születik / $i = 1, 2, \dots, 30$ /, függetlenül attól, hányszor és melyik osztályon kezelték.

A mintavételezést úgy végezzük, hogy elindulunk a rendezett populáció elejéről és minden az A osztályon kezelt esetet beveszünk a mintába, egészen addig, amíg  $M_A/10$  esetünk nem lesz. Jelöljük  $\xi_A$ -val az A osztály mintájába bekerülő utolsó eset sorszámát. Jelölje  $\eta$  a teljes mintába kerülő azon esetek számát, akiket mind az A, mind a B osztályon ápolat. A

$$h = \frac{E(\eta | \xi_A=a, \xi_B=b) - E\eta}{E\eta}$$



### 3. A mintavétellel kapcsolatos megbízhatósági kérdések

[1]-ben részletes leírás található a felvethető kérdésekről, a felhasznált összefüggésekről. Most csak konkrét példákon keresztül szeretnénk rávilágítani a lényegre.

a/ Leggyakrabban az a feladatunk, hogy becslést adjunk annak  $p$  valószínűségére, hogy egy beteg valamely előre adott tulajdonsággal rendelkezik, pl. a beteg Pest megyei v. adott korformájú betegséggel ápolták, stb. Mászóval ez pl. a következőt jelenti 95 %-os biztonsággal állíthatjuk, hogy a Pest megyei betegek száma 14200 és 14800 közé esik, stb. Másik példa 0,95 megbízhatósági szintű konfidencia intervallumot akarunk szerkeszteni annak  $p$  valószínűségére, hogy egy adott beteg Szabolcs megyei. 7600 elemű a mintánk, így azt kapjuk, hogy  $0,433 \leq p \leq 0,0461$  0,95 megbízhatósági szintű konfidencia intervallum. Ez azt jelenti, hogy 5 %-os biztonsággal állíthatjuk a Szabolcs megyei betegek száma 7311 és 7837 = 17000.0,0461 közé esik.

b/ Feltételezve, hogy kórházainkban évente kb. 1700000 beteget ápolnak, felmerül a kérdés, hogy adott megbízhatósági szint /adott hibavalószínűség és hibakorlát/ esetén hány százalékos mintára van szükségünk.

Ha pl. azt az eseményt vizsgáljuk, hogy a beteget 333-as korformával ápolták, akkor  $0,0000644 \leq p \leq 0,0002298$

adódik, ami "rossznak" mondható. Élesebb konfidenciaintervallumhoz juthatunk  $M$  növelésével. Ha pl. az intervallum két végpontja közötti távolságra  $0,00005$  értéket kívánjuk meg - ez olyankor fordulhat elő, amikor az  $A$  esemény valószínűsége igen kicsi, mint pl. ebben a példában is. Most azt kapjuk, hogy  $M \geq 1124000000$ , ami természetesen semmilyen mintavétellel sem érhető el, figyelembevételre Magyarország lakosainak számát. Vegyünk egy másik példát. Az esemény legyen most az, hogy a beteget a 10. osztályon ápolják. Ekkor  $0,05$  intervallum esetén az következik, hogy 1382 elemű minta is elég lenne a  $0,95$  megbízhatóságu szintű  $0,05$  hosszúságu konfidencia intervallum megadásához. Látjuk tehát, hogy adott megbízhatóságu szintű adott nagyságu konfidencia intervallumhoz más-más mintanagyság kellene. Van amikor ez problémába ütközik.

c/ Szükség lehet arra, hogy eldöntsük adott  $H_0$  feltevés /pl. a szellemi dolgozók 30 %-a infarktuszban hal meg, vagy a születésnapok eloszlása egyenletes, stb./, u.n. nullhipotézis adott szinten elfogadható-e. Megadandó továbbá a  $H_0$ -t elutasító u.n. kritikus tartomány. /Ilyen kérdésekről ld. [5]./

d/ Ha az a/ kérdést szeretnénk megválaszolni azokban az esetekben, amikor a "tulajdonság" rendre az, hogy a beteg négyjegyű BNO kódja 0001, 0002, ..., 9998, 9999

és az ott követendő eljárást alkalmaznánk most is, sok és felesleges számolást végeznénk. Ehelyett a Kolmogorov eloszlás alapján konfidenciasávot adunk az eloszlásfüggvényre. Hangsúlyozni szeretnénk, hogy tulajdonképpen az egyes valószínűségekre adunk becslést, csak más módon, mint az a/ pontban. Ugyancsak ezt az eloszlást használhatjuk annak eldöntésére, hogy kapott eredményeink mennyire egyeznek régebbi eredményeinkkel vagy külföldi eredményekkel.

e/ Homogenitás vizsgálat alkalmazása is felmerülhet. Állandó lakóhely megyéje, születési hely megyéje azonos eloszlásúnak tekinthető-e.

f/ Két tényező, amelyek egymásrahatása feltételezhető, függetlennek vehető-e, pl. beteg és beteg édesanyja keresztnévének kezdőbetűje, nem v. kor és bizonyos betegségek, keresztnév kezdőbetűje és nem, stb.

### 3. Azonosító kódok vizsgálata

Ha egy populáció egyedeinek azonosítása nem lehetséges sorszámozással, akkor az egyedeket valamilyen természetes adatuk alapján lehet megkülönböztetni egymástól. Ezek az adatok személyeknél lehet pl. néhány születési adat, stb. Ilyen adatok azonban több különböző egyednél is lehetnek azonosak /pl. egyazon napon született azonos nemű emberek/. Pl. belátható, hogy már 23 ember megkülön-



böztetésére sem elég jó azonosító az év 365 napja: ha véletlenszerűen kiválasztunk 23 embert, akkor az esetek 50 %-ában eközött a 23 ember között legalább kettőnek az év ugyanazon napján van a születésnapja /az év minden napját egyenlő valószínűnek tekintve/.

A kórházban ápoltszemélyek azonosítására bizonyos adatokat használunk fel. Kérdés: a/ ezek az adatok a személyek hány %-át azonosítják egyértelműen? b/ Hány újabb adatot kell hozzávennünk az azonosítóhoz, hogy az előbbi százalékszámot növeljük?

Az azonosítás hatásfokának növelése érdekében nyilván az azonosításra csak olyan adatokat célszerű használni, melyek nem változnak meg az ember élete során. Ilyen adat pl. a születési év, hó, nap, stb, de nem ilyen adat pl. az állandó lakóhely megyéje, annak település jellege, stb. Ennek megfelelően vizsgálatunk az alábbi adatokra terjed ki:

születési dátum	6 karakter
nem	1 karakter
beteg /leánykor/ nevének kezdőbetűi	4 karakter
anyja nevének kezdőbetűi	4 karakter
születési hely megyéje	2 karakter

Bizonyított összefüggéseinkből következik, hogy a

duplán azonosított személyek várható száma 36.

Felvetődik a kérdés, mi történik, ha valamelyik adatot kihagyjuk az azonosítóból, mennyire változik meg a rosszul azonosított emberek várható száma. Erre vonatkozik a következő táblázat.

Kihagyott adat	Duplán azonosítottak várható száma
Beteg vezetékneve	292
Születési megye	216
Beteg keresztnévének kezdőbetűje	294

#### 4. Megjegyzések

- a/ A kettőnél többször ápoltak száma elhanyagolható, az ebből adódó hiba egy nagyságrenddel kisebb, mint az általunk adott becslés hibája.
- b/ Az AB betegek számára vonatkozó becslés elég pontos lesz, ha A is és B is "nagy" osztály.
- c/ Az elmondott példák alapján a következő megállapításokat tehetjük. Bizonyos értékek - a 10 %-os mintát alapul véve - nem szolgáltatnak megbízható

eredményeket, ugyancsak vannak esetek, amikor kisebb mintából is megbízhatóan következtethetünk. Felmerülhet annak igénye, hogy a kapott táblázatokban valamilyen formában jelöljük, mely eredmények nem megbízhatóak - adott szinten. Ez azonban két problémát vet fel: megnöveli a számolási időt, csökkenti a rendszer hatékonyságát, általánosságát. Mindezek ellenére nyilvánvaló, hogy bizonyos esetekben feltétlenül szükség van erre.

Ennek és az említett egyéb kérdések alkalmazási lehetőségeinek pontos behatárolására - hol, milyen számítások elvégzésénél kell bizonyos próbákat, stb. kivitelezni - további vizsgálatokra van szükség.

- d/ Az azonosítóba túl sok adatot nem célszerű belevenni, mert ez egyrészt megnövelné a különböző helyigényeket /az adathordozókon/, másrészt meglassítaná az adatmozgatást.
- e/ Az azonosítás egy másik problémája közvetlenül a kódolással kapcsolatos. Ha természetes adatokkal azonosítunk, akkor sok esetben igen rossz hatásfoku kódokat kell használnunk. Például a "beteg neve" kétféle érték lehet, holott a felhasznált egyjegyű decimális kód tíz érték megkülönböztetését teszi le-

hetővé. Ugyanez a helyzet a születés hónapjánál és napjánál is, de még az olyan látszólag teljesen kihasznált kódnál is, mint a születés éve, hiszen pl. a kórházi ápoltak között bizonyos viszonylag szűk korosztályba tartozó betegek nagy számban fordulhatnak elő /pl. szülő nők/. Így pl. a születési dátumból és nemből álló hétjegyű azonosító közel sem ad 10 millió féle értéket, hanem csak néhány tizezernyit.

#### I r o d a l o m

- [1] Králmi A., Ratkó I., Ruda M., Soltész J.: A statisztikai adatfeldolgozás matematikai és számítástechnikai problémái, MTA SZTAKI, Tanulmányok, 70/1977.
- [2] Garádi J., Krámlí A., Ratkó I., Ruda M.: Statisztikai és számítástechnikai módszerek alkalmazása kórházi morbiditás vizsgálatokban, MTA SZTAKI, Tanulmányok, 35/1975.
- [3] Rényi A.: Valószínűségszámítás, Tankönyvkiadó, Bp., 1966.
- [4] Tomkó J.: A Markov folyamatok elemei és néhány operációkutatási vonatkozása, Bolyai J. Matematikai Társulat kiadványa, Bp., 1968.

- [5] Arató M.: Fejezetek a matematikai statisztikából  
számítógépes alkalmazásokkal I., MTA SZTAKI,  
Tanulmányok, 42/1975.