

MTA SZTAKI

Azonosítási kódok statisztikai vizsgálata

Garádi János

A kórházakban ápolott személyek azonosítására általában a születési év, hó, nap, a nem, valamint az anya nevének kezdőbetűje szolgál.

Ezek az adatok jól használhatók, mivel nem változnak meg az ember élete során. A vizsgálatok azonban azt mutatták, hogy ezek az adatok önmagukban a személyek csak mintegy 30 %-át azonosítják egyértelműen. Ezért a fenti azonosító kódhoz hozzávehetjük az állandó lakóhely megyéjét és a település jellegét (pl.: falu, város, illetve Budapesten a kerület). Ezek az adatok egy éven belül a lakosság mintegy 2 %-ánál változnak meg. Ezen felül az olyan kódokat is hozzávehetjük, amelyeknek a változása a lakoságnak legfeljebb 3 %-át érinti.

A statisztikai vizsgálathoz modellként az ugynevezett cellabetöltési problémát használjuk.

Ebben a dolgozatban egy cellabetöltési problémával foglalkozunk, amely az összetartozó rekordoknak (emberek) véletlenszerűen kisorsolt azonosító számok segítségével történő azonosítása során merült fel.

A mi esetünk annyiban különbözik a klasszikus cellabetöltési problémától, hogy a különböző cellákba esés valószínűsége különböző.

A feladat leírása a következő:

Adott N számú ember, és n számú különböző azonosítószám,

melyeket az N ember egymástól függetlenül p_1, p_2, \dots, p_n valószínűséggel vesz fel (ahol $n = N$).

Kérdés:

Mennyi a hibásan azonosított emberek várható száma?
(Hibásan azonosítottak akkor tekintünk egy embert, ha található hozzá legalább még egy ember, amelynek ugyanaz az azonosítója.)

Legyen az $\xi_i^{(m)}$ valószínűségi változó jelentése a következő:

$$\xi_i^{(m)} = \begin{cases} 1, & \text{ha az } i\text{-edik azonosító pontosan} \\ & m \text{ emberhez tartozik} \\ 0, & \text{különben} \end{cases} \quad /1.1/$$

A $P(\xi_i^{(m)} = 1)$ valószínűség könnyen kiszámítható:

$$P(\xi_i^{(m)} = 1) = \binom{N}{m} p_i^m (1 - p_i)^{N-m} \quad /1.2/$$

mert az egyes emberekhez egymástól függetlenül rendeltük az azonosítókat.

Azoknak az azonosítóknak az E_m várható száma, amelyhez pontosan m ember tartozik egyenlő:

$$E_m = \sum_{i=1}^n \binom{N}{m} p_i^m (1 - p_i)^{N-m} \quad /1.3/$$

Konkrétan:

$$E_0 = \sum_{i=1}^n (1 - p_i)^N$$

$$E_1 = N \cdot \sum_{i=1}^n p_i (1 - p_i)^{N-1} \quad /1.4/$$

$$E_2 = \frac{N \cdot (N-1)}{2!} \sum_{i=1}^n p_i^2 (1 - p_i)^{N-2}$$

$$E_3 = \frac{N \cdot (N-1)(N-2)}{3!} \sum_{i=1}^n p_i^3 (1 - p_i)^{N-3}$$

Vezessük be a következő jelölést: $p_i = \frac{\alpha_i}{N}$.

Fel kell tennünk, hogy $n \gg N$ - egyébként a rosszul azonosított emberek száma tulságosan nagy lenne (a konkrét feladatban $n \approx 150 N$).

Emellett, a p_i valószínűségek, bár nem egyenlők, mégis egyenletesen kicsik, azaz ha $n \rightarrow \infty$ és $n \cdot \max_i p_i$ egy független K korlát alatt marad. Ezért /1.4/ felhasználásával:

$$E_1 \approx \sum_{i=1}^n \alpha_i \cdot e^{-\alpha_i} \quad /1.5/$$

Annak ellenére, hogy feltevésünk szerint:

$$\alpha_i < \frac{K}{n} \cdot N \approx 0, \text{ az } e^{-\alpha_i}$$

nem helyettesíthető egyszerűen 1-gyel, mert így azt kapnánk, hogy a rekordok 100 %-a helyesen azonosítható.

Az /1.4/ képletekből az előző feltevéseinket felhasználva, hogy a hibásan azonosított emberek H várható száma:

$$H = \sum_{m=2}^{\infty} m \cdot E_m$$

jól közelíthető a $2 \cdot E_2$ -vel, mert a maradék:

$$\sum_{m=3}^{\infty} m \cdot E_m = \sigma(\alpha_i^2)$$

E_2 kiszámításával már elhanyagolhatjuk az $e^{-\alpha_i}$ tényezőket, tehát:

$$H \approx 2 \cdot L_2 = 2 \cdot \frac{N \cdot (N-1)}{2!} \sum_{i=1}^2 \frac{d_i^2}{N^2} \left(1 - \frac{d_i}{N}\right)^{N-2} \approx$$

$$\approx \sum_{i=1}^{\infty} d_i^2 e^{-d_i} \approx \sum_{i=1}^{\infty} d_i^2 + 0 \left(\sum_{i=1}^{\infty} d_i^3 \right)$$

A $\sum_{i=1}^2 p_i^2$ összeg a konkrét feladatnál könnyen kiszámítható,

mert az azonosítószámot 5 - a feladat természete miatt - egymástól független azonosító egyesítésével alakítottuk ki, ezért elegendő volt külön-külön megbecsülni ezen azonosítók

$$p_{i_1}^{(1)}, p_{i_2}^{(2)}, \dots, p_{i_5}^{(5)} \text{ valószínűségeit.}$$

Ekkor a függetlenség miatt:

$$\sum_{i=1}^{\infty} p_i^2 = \prod_{j=1}^5 \sum_{i,j=1}^{\infty} (p_{ij}^{(j)})^2$$

Térjünk át a $\xi^{(m)} = \sum_{i=1}^n \xi_i^{(m)}$ valószínűségi változók szórás négyzetének kiszámítására.

$$D^2(\xi^{(m)}) = M(\xi^{(m)2}) - (M(\xi^{(m)}))^2 \quad /1.6/$$

$$M(\xi^{(m)2}) = M(\xi^{(m)}) + \sum_{\substack{i \neq j \\ i, j=1}}^n M(\xi_i \xi_j) =$$

$$= \sum_{i=1}^n \binom{N}{m} p_i^m (1 - p_i)^{N-m} + \sum_{\substack{i \neq j \\ i, j=1}}^n \binom{N}{m} \binom{N-m}{m} \cdot p_i^m \cdot (1-p_i)^m \cdot p_j^m \cdot (1-p_j)^{N-2m} \quad /1.7/$$

itt természetesen feltételeztük, hogy $N > 2 \cdot m$.

$$M \left(f^{(m)} \right)^2 = \sum_{i,j=1}^n \binom{N}{m}^2 p_i^m p_j^m (1-p_i)^{N-m} (1-p_j)^{N-m}$$

Könnyen belátható, a következő állítás:

Ha $m \geq 2$, és a $\max_i d_i \rightarrow 0$, hogy eközben

$$\frac{\max_i d_i}{\min_i d_i}$$

korlát alatt marad, akkor a

$$\frac{D^2 \left(f^{(m)} \right)}{M \left(f^{(m)} \right)} \rightarrow 1. \quad /1.8/$$

Bizonyítás:

Azt kell belátnunk, hogy $\frac{D^2 \left(f^{(m)} \right) - M \left(f^{(m)} \right)}{M \left(f^{(m)} \right)} \rightarrow 0$

Az /1.6/ és az /1.7/ alapján:

$$D^2 \left(f^{(m)} \right) - M \left(f^{(m)} \right) = \sum_{\substack{i,j=1 \\ i \neq j}}^n \binom{N}{m} \binom{N-m}{m} p_i^m (1-p_i)^m \cdot$$

$$\cdot p_j^m (1-p_i - p_j)^{N-2m} - \sum_{i,j=1}^n \binom{N}{m}^2 p_i^m (1-p_i)^m \cdot \quad /1.9/$$

$$\cdot p_j^m (1 - p_i - p_j + p_i \cdot p_j)^{N-m}$$

Először vizsgáljuk az $i = j$ esetet: Ilyen tagból n darab van a kivonandóban. Könnyen belátható, hogy összegük rendje:

$$O \left(n \cdot \max_i \alpha_i^{2m} \right)$$

Ha az $i \neq j$, akkor az /1.9/-ben a kivonást tagonként elvégezve, és alkalmazva a binomiális tételt nyerjük, hogy egy tag rendje:

$$O \left(\frac{1}{n} \cdot \max_i \alpha_i^{2m-1} \right)$$

Tehát az összeg rendje:

$$O \left(n \cdot \max_i \alpha_i^{2m-1} \right)$$

Igy adódik, hogy az $\frac{1}{M(\xi^{(m)})}$ rendje pedig: $O \left(\frac{1}{n \cdot \min_i \alpha_i^m} \right)$

Amiből következik állításunk:

$$\frac{D^2 \left(f^{(m)} \right) - M \left(f^{(m)} \right)}{M \left(f^{(m)} \right)} = O \left(\max_i \alpha_i^{m-1} \right)$$

Q. E. D.

Számításainkat két módon is alátámasztottuk számítógépes programok segítségével, reprezentatív mintavétel alapján.

Egyrészt közvetlen úton: sikerült találni a mintának egy olyan - viszonylag nagyméretű - részhalmazát, amelyben az összehasonlítható rekordok előfordulása logikailag lehetetlen volt. Ezen halmazon a logikailag megengedett azonosítók száma 300-szorosa az azonosítandó emberek számának. Ennek ellenére a rosszul azonosított emberek az összes emberek 5,5 %-ét teszik ki. Ez 1100

ilyen azonosítónak felel meg, amely pontosan 2 embert azonosít.

$$A \text{ szűrés tehát: } E \left(\xi^{(m)} \right) \approx 31.$$

A három vagy ennél több embert azonosító kódok száma már elenyészően kicsiny, azaz: 60.

A másik ellenőrzési mód közvetett: előző halmazunkhoz találtunk egy másik, vele kb. azonos méretű (N -elemű), tőle logikailag diszjunkt ember-halmazt, amelyen az azonosító kódok eloszlása (a p_i valószínűségek) - a probléma természetéből adódóan - ugyanaz.

Kérdés: mekkora a két halmaznak a hibés azonosításból adódo közös része.

A közös elemek várható száma könnyen meghatározható lenne, ha ismernénk az egyik halmaz által lefoglalt azonosítókhoz rendelt p_i valószínűségek összegét $P_{\text{össz}}$. A keresett várható érték ebben az esetben: $14 \cdot P_{\text{össz}}$.

Ha csak az azonosítók eloszlását (a p_i valószínűségeket) ismerjük, de nem ismerjük, hogy az első halmaz konkrétan mely azonosítókat foglalja le, akkor a $P_{\text{össz}}$ várható értékét tudjuk kiszámítani.

$$P_{\text{össz}} = \sum_{i=1}^n p_i \xi_i^{(1)} + \sum_{i=1}^n p_i \xi_i^{(2)} + \dots$$

Az $\xi_i^{(m)}$ definíciójét lásd /1.1/-ben.

$$E(P_{\text{össz}}) = \binom{N}{1} \sum_{i=1}^n p_i^2 (1-p_i)^{N-1} + \binom{N}{2} \sum_{i=1}^n p_i^3 (1-p_i)^{N-2} + \dots$$

Tekintettel arra, hogy az $a_i = p_i \cdot N$ számok nullához közeliek,

ezért elegendő csak az első tagot figyelembe venni, sőt az

$$(1 - p_i)^{i-1} = e^{-d_i}$$

szorzót is elhanyagolhatjuk.

Igy azt nyerjük, hogy:

$$E(p_{\text{össz}}) \approx \frac{2 \cdot E(f^{(2)})}{n-1},$$

azaz a két halmaz hibás kódolásból adódó közös részének várható elemszáma kétszerese az egyik halmaz által kétszeresen lefoglalt kódok várható számának.

Az /1.6/ összefüggéshez hasonlóan belátható, hogy

$$D^2(p_{\text{össz}}) = N \left(\sum_{i=1}^n p_i^3 - \left(\sum_{i=1}^n p_i^2 \right)^2 \right) + O\left(\frac{N}{n^3}\right)$$

ha $n \rightarrow \infty$ és $\max_i d_i \rightarrow 0$.

Ha a p_i -k egyenletességére nem kötünk ki semmilyen feltételt, akkor:

$$L^2(p_{\text{össz}}) = O\left(\frac{N}{n^2}\right), \text{ ezért a}$$

$$D^2(N \cdot p_{\text{össz}}) = O(n).$$

Ugyanígy nagyságrendű a közös rész elemszámának - a binomiális eloszlás alapján számított - szórásnégyzete, ha $p_{\text{össz}}$ ismert.

Fenti megfontolásaink érdekessége az, hogy a két halmazban közös azonosítók elemeinek számára következtetni tudunk az egyik halmaz hibásan azonosított elemeinek számából - anélkül -, hogy ismernénk az azonosító kódok eloszlását.

A közös rész elemeinek tényleges (gépi úton nyert) száma : 1300.

A számított értékeket úgy kaptuk, hogy különbözőnek tekintettük azokat a személyeket, akiknek állandó lakóhelye és foglalkozása különböző, majd megvizsgáltuk, hogy hány olyan 9 jegyű azonosító kód van, amely az ily módon különbözőnek tekintett személyek közül 1, 2, 3, stb. személyhez tartozik.

Megjegyezzük, hogy a belső vándorlás és az egy éven belüli foglalkozás változás - pl. nyugdíjazás - viszonylag nagy száma miatt ez az eljárás nem teljesen korrekt, mégis hü képet nyújt az azonosító kódok statisztikai viselkedéséről.

A fentiek alapján a személyek azonosítását egy 13 számjegyű azonosítóval végeztük.