

Módszer folytonos és bináris változókkal leírt minták osztályozására

B. Nagy András és Wolf Tamás

1. Bevezetés

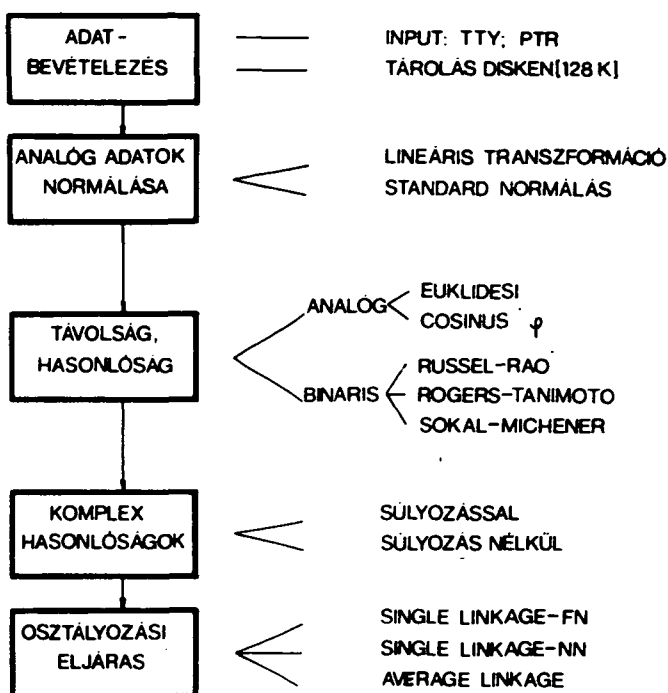
A jelfeldolgozási feladatok jelentős részénél a vizsgált objektum analóg és bináris (általában digitális) jellemzők összességével írható le. A változók - melyek a mérés eredményeként állnak rendelkezésre - többnyire folytonos értékkészlettel bírnak, míg a kérdéses objektummal kapcsolatos jegyek, tünetek megléte, illetve hiánya bináris típusu adatot eredményez. Az ilyen, ugynevezett kevert típusu változók feldolgozásának evidens példája az orvosi diagnózis alkotási folyamat. Ekkor a klinikai vizsgálatok eredményei (vérnyomás, hőmérséklet, bioelektromos jelek stb.) analóg jellegűek: az anamnézis során az orvos által feltett eldöntendő típusu kérdésekre, vagy a klinikai kérdőívekre adott IGEN/NEM válaszok mind bináris típusúak. Más területek szakemberei is gyakran találkoznak kevert változókkal leírható objektumok problematikájával (pl. a meteorológiai prognosztikánál, szociológiai vizsgálatoknál stb.).

A kevert változók feldolgozásánál is gyakran kívánatos lehet valamilyen osztályozó eljárás alkalmazása. Az irodalomból ismert cluster-algoritmusok vagy analóg (1) vagy bináris (2,3,4) változókat kezelnek. Az orvosi diagnosztikai kutatások jelenlegi szakaszában célszerűnek látszott a hierarchikus cluster-algoritmusok kiterjesztése a kevert változókkal leírható modellekre.

2. Az ABCL eljárás

A kidolgozott eljárás főbb lépései az 1. ábrán láthatók. Az ADAT-BEVÉTELEZÉS lyukszalagról vagy irógépről történhet, az utóbbi esetben az adatokról lyukszalag készül az újrafuttatás megkönnyítéséhez. Mivel a programot 16 K szavas központi tárral rendelkező TPAi gépen futtatjuk, az adatok disk-re mentése szükséges ahhoz, hogy megfelelően nagy adatbázist kezelhessünk. A 128 K szavas disk - figyelembe véve a BASIC programozási nyelv számbábrázolását és a felügyelő OS/i operációs rendszer helyigényét - mintegy 100 x 50-es analóg és 100 x 50-es bináris adatmátrix feldolgozását teszi lehetővé. A disk felhasználása természetesen a futási idő megnövekedését eredményezi, de a nem rutinszerű alkalmazás időszakában ez megengedhető.

ABCL SOFTWARE FUNKCIÓK



1. ábra

Az ANALÓG ADATOK NORMÁLÁSÁ-t kétféle eljárással végeztük. Mindkét esetben az a cél, hogy a vektorok normálása után számított távolság-értékek kompatibilisek legyenek a bináris távolságokkal, amelyek kielégítik a szokásos távolsági kritériumokat (1). Ha az objektumok száma kicsi (néhányszor 10), LINEÁRIS TRANSZFORMÁCIÓ-val (LT) az adatokat a (0,1) intervallumba transzformáljuk. Amikor a minták száma nagyobb, lehetőség van STANDARD NORMÁLÁS (SN) alkalmazására.

A TÁVOLSÁG számítást az ANALÓG változókra, vagy a vektorok EUKLIDESI távolsága vagy COSINUS φ -je alapján végezzük. A BINÁRIS hasonlóság meghatározására három módszer közül választjuk valamelyiket (2,3,4).

Az alkalmazott hasonlósági függvények értékkészlete a (0,1) tartományba esik, így az analóg és bináris hasonlóságok alapján KOMPLEX HASONLÓSÁG-ot számolhatunk. A gyakorlati feladatok többségénél a minták analóg és bináris változóinak száma eltérő. Így a belőlük számolt analóg, illetve bináris hasonlóságnak is különböző a súlya az objektumok

együttes hasonlóságában. Ezért a felhasználónak lehetősége van a dimenziók figyelembevételére, illetve bizonyos heurisztikus megfontolások eredményeképpen előálló súlyozás beiktatására.

Az OSZTÁLYOZÁSI ELJÁRÁS-ok az együttes hasonlósági elemekből adódó, ugynevezett hasonlósági mátrix elemei alapján osztályoznak. A megvalósított algoritmus ugyanakkor megengedi csak analóg vagy csak bináris típusú minták osztályozását is. Ez egyben hasznos segédeszköz lehet annak vizsgálatára, hogy egy objektumhalmaz osztályozásakor mennyivel ad több információt a kevert minták figyelembevétele. A hierarchikus osztályozási eljárások közül - melyek az objektumok eloszlására, illetve az osztályok számára vonatkozóan nem igényelnek a priori információt - három összekapcsolási módszert programoztunk be. Ezek közös tulajdonsága, hogy az osztályozás során mindig a két leghasonlóbb elemet olvasztják össze egy osztályba. Az így létrehozott új osztály tulajdonságai azonban a három módszernél eltérőek. A legtávolabbi szomszéd módszernél (Furthest Neighbour: FN) az új osztály a hozzásorolt minták hasonlóságai közül a kisebbeket tartja meg, a legközelebbi szomszéd módszer (Nearest Neighbour: NN) a nagyobbakat, az átlagos távolság szerinti összekapcsolás módszerénél (Average Linkage) pedig az új osztály és egy másik osztály hasonlósága a hozzátartozó osztályok és a másik osztály hasonlóságainak átlagával egyenlő (1,2,3). Az egyes clusterezési lépéseknél a program kiírja a kialakult osztályokat, a hozzájuk tartozó minták sorszámait és a hasonlósági szinteket. A program addig fut, míg végül egyetlen objektumot kapunk, amelybe minden minta beletartozik.

3. Az együttes hasonlóság számítása

Az együttes hasonlóság számítása az analóg és bináris hasonlóságokon alapszik. A mért analóg adatokat előzőleg normáljuk.

Jelölje x_{ki} a k-adik minta i-edik elemét,

M a minták számát,

$$x_{imin} = \min_k \{x_{ki}\}$$

$$x_{imax} = \max_k \{x_{ki}\}$$

$$\bar{x}_i = \sum_{k=1}^M x_{ki}$$

$$\sigma_i = \sqrt{\frac{\sum_k (x_{ki} - \bar{x}_i)^2}{M-1}}$$

és x'_{ki} az x_{ki} normált értékét.

a.) Lineáris transzformálásnál

$$x'_{ki} = \frac{x_{ki} - x_{i\min}}{x_{i\max} - x_{i\min}}$$

ekkor

$$0 \leq x'_{ki} \leq 1$$

b.) Standard normálásnál

$$x'_{ki} = \frac{x_{ki} - \bar{x}_i}{\sigma_i}$$

Ebben az esetben a transzformált vektorokra

$$\bar{x}'_i = 0 \quad \text{és} \quad \sigma'_i = 1 \quad \text{teljesül.}$$

A normált adatok alapján számolt analóg hasonlóságokat a 2. ábrán részletezzük.

Az alkalmazott bináris távolság fogalmak értelmezése a 3. ábrán található. A hasonlósági mértéktől megkivánjuk, hogy

$$\begin{aligned} 0 \leq S_{ij}^B \leq 1 & \quad \text{ha} \quad i \neq j & \quad /i/ \\ S_{ij}^B = 1 & \quad \text{ha} \quad i = j & \quad /ii/ \\ S_{ij}^B = S_{ji}^B & & \quad /iii/ \end{aligned}$$

teljesüljön. Az /ii/ feltétel miatt a Russell-Rao hasonlósági függvényt módosítottuk

$$S_{ij}^B = \begin{cases} \frac{n_{ij}}{L} & i \neq j = 1, 2, \dots, M \\ 1 & i = j \end{cases}$$

ANALÓG HASONLÓSÁG

$X_i^a = [x_{i1}, x_{i2}, \dots, x_{iM}]$ L normált minta
 N változók dimenziója
 d_{ij} i, j minták távolsága
 S_{ij}^a i, j -- hasonlósága

BINÁRIS HASONLÓSÁG

$X_i^b = [x_{i1}^b, \dots, x_{iL}^b]$
 L változók dimenziója
 S_{ij}^b i, j minták hasonlósága

EUKLIDESI $d_{ij} = \left[\sum_{k=1}^M (x_{ik} - x_{jk})^2 \right]^{1/2}$
 $S_{ij}^a = 1 - \frac{d_{ij}}{P \cdot D}$

COSINUS ϕ
 $h_{ij}^\phi = \frac{x_i x_j^T}{\|x_i\| \|x_j\|}$
 $S_{ij}^a = h_{ij}^\phi$

Russel-Rao $S_{ij}^b = \begin{cases} \frac{n_{ij}}{L} & i \neq j \\ 1 & i = j \end{cases}$

Roger-Tanimoto $S_{ij}^b = \frac{n_{ij} + n^j}{L + n^j}$

Sokal-Michener $S_{ij}^b = \frac{n_{ij} + n^j}{L}$

n^j "0" egyezések
 n_{ij} "1" egyezések
 n^i "0-1", "1-0" eltérések

} száma az i, j vektorok azonos bitpozícióiban nézve

LT: $D=1; P=\sqrt{N}$
 SN: $D=\text{MAX}_{i,j} \{d_{ij}\}; P=1$

2. ábra

3. ábra

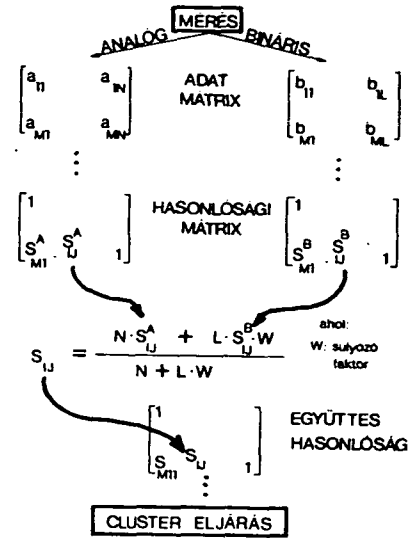
Bár a képletek formailag hasonlóak, lényeges elvi és gyakorlati alkalmazásbeli különbség van közöttük. A Russell-Rao formulát akkor célszerű választani, ha mintákban az IGEN válaszokat tekintjük fontosnak. Másszóval bizonyos tények meglétéből kívánunk következtetni az objektumok tulajdonságaira, míg a tények hiányát nem tekintjük informatívnak. A Roger-Tanimoto formula alkalmazásánál az IGEN és NEM válaszok (azaz a bináris minta "1" elemei és "0" elemei) egyforma súlyal jellemzik az objektumot. A hasonlósági képlet két minta eltérő elemeit fokozottan veszi figyelembe. A Sokal-Michener eljárás szintén egyformán jellemzőnek tekinti az IGEN és NEM válaszokat, de a két minta eltérő elemeit nem hangsúlyozza a hasonlóság számításánál.

Az együttes-hasonlósági együtthatók számítását szemlélteti a 4. ábra. A módszer alkalmazásával az /i/, /ii/ és /iii/ feltételek továbbra is teljesülnek.

4. Összefoglalás

A biológiai objektumok általános esetben analóg és bináris változók együttesével írhatók le. A szerzők által javasolt eljárás az analóg és bináris

ris hasonlóság alapján bevezeti az együttes hasonlóság fogalmát, a TPAi kiskiszámítógépen futtatott ABCL program pedig háromféle bináris hasonlóság számítás és kétfajta analóg távolságfüggvény alapján számított együttes hasonlósági mátrix elemein három típusú összekapcsolási módszer szerint képes a biológiai objektumokat osztályba sorolni. Az így előálló 18 féle clusterizációs modell az első tapasztalatok szerint új ismereteket szolgáltat a biológiai objektumok mélyebb megismeréséhez.



4. ábra

Irodalom

- (1) Duran, B.S., Odell, P.L.: Cluster Analysis. A Survey, Springer-Verlag (1974)
- (2) Sokal, R.R., Sneath, P.H.A.: Principles of Numerical Taxonomy, San Francisco. W.H. Freeman and Comp. (1963)
- (3) Winkel, P.: Numerical Taxonomic Analysis of Cirrhosis. I. Comp. and Biomed. Res. 7, 100-110 (1974)
- (4) Rogers, D.J., Tanimoto, T.T.: A Computer Program for Classifying plants, Science, Vol. 132 (oct. 21, 1960)
- (5) Young, T.Y., Calvert, T.W.: Classification, Estimation and Pattern Recognition Am. Elsevier, 1974.