

KOMPARACE PODMÍNĚNÝCH DISTRIBUCÍ KONTINGENČNÍ TABULKY
KLASICKÝM MNOHOROZMĚRNÝM ŠKÁLOVÁNÍM

Jan Řehák
ÚFS ČSAV Praha
Ivana Loučková
KAS FFUP Olomouc

Podmíněně řádkové distribuce mají v dvourozměrných /i vícerozměrných/ kontingenčních tabulkách samostatný význam. Při odmítnutí hypotézy jejich homogenity pomocí některého ze známých testů je další analytickou úlohou určit vztahy mezi distribucemi a to nejen jejich seskupení do homogenních podskupin, ale také podle stupně jejich vzájemné heterogenity. Jedním z možných postupů je grafické znázornění distribucí jako bodů v eukleidovském prostoru /co nejmenší dimenze/, přičemž bližší body ukazují na podobnější a vzdálenější body na rozdílnější distribuce.

Úlohu můžeme formulovat takto:

Je dáno R distribucí $f_{(r)} = (f_{k/r})_{k=1}^K$, $r=1, \dots, R$ a míra vzdálenosti $D(f_{(r)}, f_{(r')}) = D_{rr'}$; máme nalézt R bodů $x_r = (x_{r1}, x_{r2}, \dots, x_{rM}) \in E_M$ a dimenzi M tak, aby

$$(1) \quad D(f_{(r)}, f_{(r')}) = \left[\sum_m (x_{rm} - x_{r'm})^2 \right]^{\frac{1}{2}} = d(x_r, x_{r'})$$

respektive pro dané L nalézt takové body $x_r = (x_{r1}, \dots, x_{rL})$, aby

$$(2) \quad \frac{\sum_r \sum_{r'} \left[\sum_m^L (x_{rm} - x_{r'm})^2 \right]^{\frac{1}{2}}}{\sum_r \sum_{r'} D^2(f_{(r)}, f_{(r')})} = \min .$$

Tuto úlohu lze řešit pomocí metody klasického mnohorozměrného škálování, mají-li $D^2(f, g)$ tvar kvadratické formy.

Zobrazení $f_{(r)} \longrightarrow x_r$, $x_r \in E_M$ hledáme pomocí algoritmu kanonického rozkladu symetrických čtvercových matic /viz metoda DISTAN; J. Řehák, I. Loučková [1983]/. Pro obecnější míry nepodobnosti resp. podobnosti lze využít algoritmu MINISSA mnohorozměrného škálování /viz J.C. Lingoes [1973] /.

1. Vzdálenosti distribucí

Vzdálenosti mezi řádkovými distribucemi můžeme určit podle definice typu proměnné pomocí matice skóru $D = \|d_{ij}\|$ /viz J. Řehák, B. Řeháková [1979], [1982] /.

Pro **n o m i n á l n í** proměnné / p r o s t o u klasifikací/

$$(3) \quad D_N^2(f_{(r)}, f_{(r')}) = \sum_{k=1}^K (f_{k/r} - f_{k/r'})^2 .$$

Pro **o r d i n á l n í** proměnné / u s p o ř á d a n o u klas./

$$(4) \quad D_0^2(f_{(r)}, f_{(r')}) = \sum_{k=1}^K (f_{k/r} - f_{k/r'})^2, \quad f_{k/r} = \text{distribuční funkce.}$$

Pro zobecněný typ proměnné

$$(5) \quad D^2(f_{(r)}, f_{(r')}) = (f_{(r)} - f_{(r')})' D (f_{(r')} - f_{(r)}).$$

Pro kardinální typ ztrácí úloha smysl, neboť distribuce jsou uspořádány pomocí průměrů / lze ovšem zavést analýzu pomocí distribučních funkcí - vzorec (4) /.

Tyto vzdálenosti umožňují zavést pro komparační závěry celou řadu výhodných měr:

a) vzdálenost $f_{(r)}$ a směsi ostatních distribucí $\bar{f}_{(r)} = \frac{1}{(1 - w_r)} \sum_{i \neq r} w_i f_{(i)}$,
 $w = (w_1, \dots, w_R)'$ je marginální sloupcové rozložení v tabulce nebo jinak určené váhy:
 $D(f_{(r)}, \bar{f}_{(r)})$ resp. $D(f_{(r)}, \bar{f}_{(r)}) / \max D(f, g)$;

b) vzdálenost $f_{(r)}$ a marginální distribuce $f = \sum_{i=1}^R w_i f_{(i)}$, $D(f_{(r)}, f)$ resp. $D(f_{(r)}, f) / \max D(f, g)$;

c) míru neshody distribucí $\frac{\sum_r \sum_{r'} D^2(f_{(r)}, f_{(r')})}{\max \sum_r \sum_{r'} D^2(f_{(r)}, f_{(r')})}$ a vzdálenost dvou skupin distribucí a pod.

2. Zobrazení $f_{(r)} \xrightarrow{\quad} X_{(r)}$

Předpokládejme, že pro nějaké M existují X_r tak, že platí vztah (1):

$$D(f_{(r)}, f_{(r')})^2 = d^2(X_r, X_{r'}), \quad d \text{ je metrika v } E_M.$$

Dvojitým centrováním této matice dostaneme matici C o rozměru $R \times R$ a o prvcích

$$(6) \quad c_{rr'} = -\frac{1}{2} (D^2(f_{(r)}, f_{(r')}) - \frac{1}{R} \sum_j D^2(f_{(j)}, f_{(r')}) - \frac{1}{R} \sum_j D^2(f_{(r)}, f_{(j)}) + \frac{1}{R^2} \sum_j \sum_{j'} D^2(f_{(j)}, f_{(j')})) = \sum_m (x_{rm} - \bar{x}_{.m})(x_{r'm} - \bar{x}_{.m}) = \langle X_r, X_{r'} \rangle;$$

specificky: $c_{rr} = \sum_m \bar{x}_{rm}^2 = d^2(X_r, 0) = D^2(f_r, f)$,
 $c_{rr'} = \sum_m \bar{x}_{rm} \bar{x}_{r'm} = (f_{(r)} - f)' D (f - f_{(r')})$, kde \bar{x}_r je vektor X_r centrovány do těžiště bodů X_1, \dots, X_R , t.j. $\bar{x}_r = \frac{1}{R} \sum_j X_j$, $f = \frac{1}{R} \sum_j f_{(j)}$.

(7) Platí, že $C = \bar{X}' \bar{X}$, kde X je matice $R \times M$ tvořená řádky \bar{x}_r .

Dále $c_{kk} = \|\bar{x}_k\|^2$, stopa $C = \sum \|\bar{x}_k\|^2 = \sum_r \sum_{r'} d^2(\bar{x}_r, \bar{x}_{r'}) = \sum_r \sum_{r'} d^2(X_r, X_{r'}) = \sum_r \sum_{r'} D^2(f_{(r)}, f_{(r')})$.

Pro kanonický /spektrální/ rozklad matice C /viz C.R. Rao, [1978]/

(8) $C = \sum_{m=1}^M \lambda_m P_m P_m'$, kde $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M > 0$ jsou charakteristická čísla matice C , P_1, P_2, \dots, P_M jsou charakteristické sloupcové vektory k nim postupně příslušné a platí, že

a) P_m jsou normalizované ortogonální vektory, $P_i' P_j = \delta_{ij}$,

b) pro matice $U_L = [P_1, P_2, \dots, P_L]$, $\Lambda_L = \text{diag} \{ \lambda_1, \dots, \lambda_L \}$ je

$U_L \Lambda_L U_L'$ nejlepší aproximace matice C mající hodnotu L ve smyslu euklidovské normy $\min_{B: h(B)=L} \|C - B\|^2 = \|C - U_L \Lambda_L U_L'\|^2 = \lambda_{L+1} + \dots + \lambda_M$,

c) stopa $C = \sum_{\nu} \left(\sum_{m=1}^M x_{\nu m}^2 \right) = \frac{1}{2R} \sum_{\nu} \sum_{\beta} D^2(x_{\nu}, x_{\beta}) = \sum_{\nu} D^2(t_{(\nu)}, t) = \frac{1}{2R} \sum_{\nu} \sum_{\beta} D^2(t_{(\nu)}, t_{(\beta)})$

d) stopa $U_L \Lambda_L U_L' = \lambda_1 + \lambda_2 + \dots + \lambda_L$.

Je tedy vidět, že matice

(9) $X = U_M \Lambda_M \frac{1}{2}$ je řešením rovnice (1) a

(10) $X = U_L \Lambda_L \frac{1}{2}$ je řešením úlohy (2), a že řádky x_{ν} matice X jsou souřadnice orthonormální báze v E_M resp. v E_L . Podíl m -té souřadnice na rozptýlení bodů x_{ν} je možno měřit koeficientem determinace m -té škály

$$(11) \quad r_m^2 = \frac{\lambda_m}{\sum_{m=1}^M \lambda_m}$$

Podíl m -té souřadnice na L -rozměrné škále hodnotu

$$(12) \quad r_m^2(L) = \frac{\lambda_m}{\sum_{m=1}^L \lambda_m}$$

Podíl L -rozměrné škály /prvních L dimenzi/ na úplném rozptýlení bodů

$$(13) \quad R_L^2 = \frac{\sum_{m=1}^L \lambda_m}{\sum_{m=1}^M \lambda_m}$$

(Obdobně platí (13) pro jakýkoliv soubor L dimenzi: $\sum_{m=1}^L$ se zamění sumou příslušných charakter. čísel λ_j).

Řešení (9) je, až na násobení vektorů P_m čísly ± 1 jednoznačné v tom smyslu, že posloupnost $U_L' \Lambda_L U_L$ jsou postupně nejlepší aproximace matice C , $L=1, 2, \dots, M$ vzhledem k euklidovské normě.

Výslednou bázi je možno rotovat, tj. přejít orthonormální transformací T k jiné bázi s vyjádřením vektorů

$$(14) \quad x^* = XT$$

3. Numerický příklad

Níže uvedená kontingenční tabulka vyjadřuje komparaci 4 skupin pokusných zvířat, které byly vystaveny experimentálnímu působení v různých podmínkách.

Kategorie závislé proměnné odpovídají osmi kombinacím třech zjišťovaných symptomů A, B, C : ABC, AB \bar{C} , A $\bar{B}C$, $\bar{A}BC$, A $\bar{B}\bar{C}$, $\bar{A}B\bar{C}$, $\bar{A}\bar{B}C$, $\bar{A}\bar{B}\bar{C}$. Komparační analýzu pomocí programu DISTAN provedeme jednak pro běžný nominální typ a jednak pro zobecněný typ daný maticí skóru $D = [d_{ij}]$, d_{ij} = počet symptomů u nichž se vyskytuje rozdíl u kategorie "i" a kategorie "j". /např.: $d(ABC, A\bar{B}\bar{C}) = 2$, $d(A\bar{B}\bar{C}, \bar{A}\bar{B}\bar{C}) = 1$ a p/

Rozložení četnosti

skupina	Vyskyty symptomů								velikost souboru
	ABC	AB \bar{C}	A $\bar{B}C$	$\bar{A}BC$	A $\bar{B}\bar{C}$	$\bar{A}B\bar{C}$	$\bar{A}\bar{B}C$	$\bar{A}\bar{B}\bar{C}$	
1	.390	.339	.085	.025	.102	.034	.008	.017	118
2	.055	.407	.251	.090	.085	.050	.040	.020	199
3	.374	.122	.033	.016	.057	.098	.033	.268	123
4	.060	.080	.120	.060	.080	.090	.300	.210	100

Maticе D, vytvářející typ zobecněné proměnné

	ABC	AB \bar{C}	A $\bar{B}C$	$\bar{A}BC$	A $\bar{B}\bar{C}$	$\bar{A}B\bar{C}$	$\bar{A}\bar{B}C$	$\bar{A}\bar{B}\bar{C}$
ABC	0.000	1.000	1.000	1.000	2.000	2.000	2.000	3.000
AB \bar{C}	1.000	0.000	2.000	2.000	1.000	1.000	3.000	2.000
A $\bar{B}C$	1.000	2.000	0.000	2.000	1.000	3.000	1.000	2.000
$\bar{A}BC$	1.000	2.000	2.000	0.000	3.000	1.000	1.000	2.000
A $\bar{B}\bar{C}$	2.000	1.000	1.000	3.000	0.000	2.000	2.000	1.000
$\bar{A}B\bar{C}$	2.000	1.000	3.000	1.000	2.000	0.000	2.000	1.000
$\bar{A}\bar{B}C$	2.000	3.000	1.000	1.000	2.000	2.000	0.000	1.000
$\bar{A}\bar{B}\bar{C}$	3.000	2.000	2.000	2.000	1.000	1.000	1.000	0.000

Reprezentace distribucí v euklidovském prostoru

nominální případ

$$M = 3$$

$$\lambda_1 = .161$$

$$\lambda_2 = .148 \quad \sum \lambda = .323$$

$$\lambda_3 = .014$$

zobecněný případ

$$M = 3$$

$$\lambda_1 = .625$$

$$\lambda_2 = .020 \quad \sum \lambda = .656$$

$$\lambda_3 = .011$$

Souřadnice

skupina	1	2	3
1	-.191	-.145	-.074
2	-.194	.222	.051
3	.114	-.231	.066
4	.271	.154	-.043
r^2 (v%)	49.9	45.8	4.4

Souřadnice

skupina	1	2	3
1	.469	.009	-.065
2	.184	-.089	.057
3	-.046	.108	.044
4	-.607	-.028	-.036
r^2 (v%)	95.3	3.1	1.6

Pro popis a komparaci distribucí je zapotřebí dvourozměrná reprezentace ($L=2$). První souřadnice odlišuje (1,2) proti (3,4), druhá (1,3) proti (2,4).

Distribuce jsou umístěny na jednorozměrné škále; první souřadnice je řadí do pořadí a charakterizuje podíly jejich vzdálenosti.

V obou případech lze nakreslit obrázek: v nominálním případě dvou-
rozměrný, v zobecněném případě stačí ke zobrazení umístit skupiny na přímku.

ZÁVĚR

Vzhledem k tomu, že pro podmíněné distribuce kategorizovaných dat
je možno zavést smysluplný pojem vzdálenosti, vystačíme po geometrickém zobrazo-
vání do euklidovského prostoru s klasickým mnohorozměrným škálováním.

Vyhodnost metody se projeví ovšem především při velkém počtu pod-
míněných distribucí.

LITERATURA

- Lingoes J.C. [1973] : The Guttman - Lingoes Nonmetric
Program Series. Ann Arbor -
Michigan, Mathesis Press
- Rao R.C. [1978] : Lineární metody statistické indukce
a její aplikace, Praha, Academia
- Řehák J., Loučková I. [1983] : Klasické mnohorozměrné
škálování /Aplikace metody DISTAN/,
Sociologický časopis XIX, s. 535-554
- Řehák J., Řeháková B. [1979] : Základní charakteristiky
proměnných s konečným počtem hodnot
a distanční analýza jejich rozložení.
Sociologický časopis XV, s. 214-233
- Řehák J., Řeháková B. [1982] : Distanční přístup k ana-
lyze kategorizovaných dat a jeho
aplikace na problem shody, Robust,
Podkost
- Řehák J., Řeháková B. [1984] : Analýza kategorizovaných
dat v sociologii, Praha, Academia
/v tisku/