

Jan ŘEHÁK (Ústav pro filosofii a sociologii ČSAV, Praha)

Blanka ŘEHÁKOVÁ (Ústav pro výzkum veřejného mínění při FSÚ, Praha)

Problém shody zahrnuje v analytické praxi širší třídu úloh, než jakje tradičně uváděn.

1. Úloha a motivace

Uveďme příklad čtyř kategorií $A = (A_1, A_2, A_3, A_4)$ s rovnoměrným rozdělením předpokládaných pravděpodobností výskytu

A_1	A_2	A_3	A_4	
.25	.25	.25	.25	teoretické rozdělení
.25	.25	.35	.15	první empirické rozdělení
.35	.25	.25	.15	druhé " " " "

V případě, že jde o prostou klasifikaci (nominální proměnnou), se první a druhé empirické rozdělení od teoretického liší ve zcela stejném stupni, pouze kvalita odchýlení (tj. kategorie, do níž se přesouvají jednotky z A_4) se mění. Mpřípadě, že jde o uspořádanou klasifikaci (v pořadí indexů), je už rozdíl velmi značný. Jestliže mezi A_k uvážujeme jinou relaci vzhledem ke které shodu hodnotíme, může být situace neshody ještě více diferencována. (A_k mohou být např. místa na mapě - Bratislava, Praha, Most, Litvínov tvoří názorný extrémní případ; kategorie mohou reprezentovat také určité typy, které se od sebe liší v různém stupni). Obecně složitější relace mezi A_k se vyskytují často ve společenských, biologických a lékařských vědách, v ekologické, diagnostické, typologické, dialektologické a jiné problematice).

Z příkladu plynou tři hlavní problémy:

- Jak odlišit různé situace (relace mezi A_k) při testování shody?
- Jak měřit neshodu?
- Jak specifikovat neshodu?

Poslední problém se řeší pomocí testování reziduí (případ prosté klasifikace a případy, kdy klasifikace vzniká jako kombinace dvou, přičemž marginální distribuce jedné nebo obou klasifikací je známa viz B.Řeháková [1979]). Měření neshody bylo pro nominální proměnnou navrženo H. Weilerem [1966].

II. Weilerův koeficient neshody a jeho diskuse

Označíme

A :	A_1, A_2, \dots, A_k	kategorie,
π :	$\pi_1, \pi_2, \dots, \pi_k$	hypotetické pravděpodobnosti,
f :	f_1, f_2, \dots, f_k	empirické relativní četnosti,
p :	p_1, p_2, \dots, p_k	skutečné pravděpodobnosti.

Test hypotézy $H_0 : \pi = p$ za předpokladu multinomického rozdělení provádíme např. pomocí Pearsonova χ^2

$$(1) \quad \chi^2 = n \frac{\sum_1 (\pi_1 - f_1)^2}{\pi_1} \sim \chi_{k-1}^2 \quad \text{za } H_0.$$

H. Weiler [1966] definoval koeficient neshody $\bar{\varphi}$ jako normalizaci $\sqrt{\chi^2}$ mezi 0 a 1

$$(2) \quad \bar{\varphi} = \sqrt{\frac{n_{\min}}{1-n_{\min}} \sum_{i=1}^K \frac{(n_i - f_i)^2}{n_i}}$$

a) podal jeho diskusi. Koeficient však má dvě základní nevýhody:

- Je vhodný pro prostou klasifikaci a není zřejmé jeho rozšíření na β -decnější typy proměnné.
- Platí

$$(3) \quad \bar{\varphi}^2 = \frac{n_{\min}}{1-n_{\min}} \left[\frac{(n_{\min} - f_m)^2}{n_{\min}} + \sum_{i \neq m} \frac{(n_i - f_i)^2}{n_i} \right]$$

kde m je index příslušný k n_{\min} (je-li jediné). Z (3) plyne, že $\lim_{n_{\min} \rightarrow 0} \bar{\varphi}^2 = f_m^2$ a tudíž příspěvek druhého členu v (3) je pro malá n_{\min} mizivý oproti příspěvku v kategorii f_m . (Obdobně výsledky platí i v případě, kdy je více kategorií s minimálním n_{\min}). V tomto případě tedy $\bar{\varphi}$ neodráží rovnoměrně neshodu ve všech kategoriích (obdobný nedostatek má i χ^2).

III.3. Distanční model pro analýzu zobecněných kategorizovaných proměnných (D-model).

Výše uvedené aspekty se snažíme překonat pomocí distanční analýzy distribucí (J. Řehák, B. Řeháková [1979]). Označíme $A = \{A_1, \dots, A_k\}$ množinu hodnot kategorizace A , $f = \{f_1, \dots, f_k\}$ distribucí z K -rozměrného simplexu $Q_k = \{f: \sum f_i = 1, f_i \geq 0\}$. Dále označíme $D = \|d_{ij}\|$, typu $K \times K$ matrici skóru vytvářejících typ proměnné, d_{ij} = skóry určující relace (A_i, A_j) a splňující podmínky $d_{ij} = d_{ji}$ a $d_{ii} = 0$. Zobecněnou proměnnou definujeme jako $A = \{A_1, \dots, A_k, D\}$.

Příklady:

- $d_{ij} = 1$ pro $i \neq j$, pak A se nazývá nominální (prostá kategorizace),
- $d_{ij} = |i-j|$, pak A se nazývá ordinální (uspořádaná kategorizace).
- $d_{ij} = (x_i - x_j)^2$, kde x_i jsou čísla přiřazená kategoriím ($x_i = x(A_i)$), pak proměnnou nazýváme kardinální (číselná kategorizace).
- Jestliže existuje K vektorů s M složkami $x_k = (x_{k1}, \dots, x_{kM})$, $k=1, \dots, K$ tak, že $d_{ij} = \sum_{m=1}^M (x_{im} - x_{jm})^2$, pak hovoříme o M -dimenzionální metrické proměnné; pro $M=2$ se tato proměnná nazývá ergální.

Základní charakteristiky (definice) a vlastnosti

- Zobecněná míra rozptylu

$$(4) \quad \text{Senvar } f = f' D f$$

- Míra polohy

$$(5) \quad c = A_k \geq d_k^* = \min_j d_j, \text{ kde } d_j = E d_{ij} = \sum f_i d_{ij}$$

- Věta (existence symetrie)

Pro $f, g \in Q_k$ platí

$$(6) \quad D(f, g) = \sqrt{(f-g)' D (g-f)}$$

je symetrická na Q_k , jestliže $D^* = \|d_{ij}^*\|$, $d_{ij}^* = d_{ik} + d_{jk} - d_{ij}$, typu $(K-1) \times (K-1)$ je pozitivně semidefinitní. Je-li D^* pozitivně definitní, pak $D(f, g)$ je metrika na Q_k . $D(f, g)$ nazýváme D -symetrikou, resp. D -metrikou na Q_k .

- Věta (rozklad zobecněné variance)

Buď $f = \sum p_r f(r)$, kde $p_r \in Q_R$, $f(r) \in Q_k$, pak pro D tvořící symetrickou platí

$$(7) \quad \text{Senvar } f = \sum p_r \text{Senvar } f(r) + \frac{1}{2} \sum_r \sum_s p_r p_s D^2(f(r), f(s)) = \sum p_r \text{Senvar } f(r) + \sum p_r D^2(f(r), f)$$

Příklady:

	nominální	ordinální	kardinální
Genvar 1	$1 - \sum f_i^2$	$2 \sum F_i(1-F_i)$	Zvar X
c	modální kat.	mediánová kat.	kategorie průměru
$D(f, g)$	$\sqrt{\sum (f_i - g_i)^2}$ (metrika)	$\sqrt{2 \sum (F_i - G_i)^2}$ (metrika)	$\sqrt{2} \bar{x}_f - \bar{x}_g $ (semimetrika)

kde F_i, G_i jsou distribuční funkce, \bar{x}_f, \bar{x}_g jsou průměry.

Důsledky:

- Nominální analýza odpovídá např. metodě CATANOVA (analýza rozptylu pro kategorizovaná data) zavedené v práci Light, Margolin [1971], aplikujeme-li (7) na nominální proměnnou.
- Analýza kardinálních proměnných přechází k lineárnímu modelu a metodě nejmenších čtverců.
- Model poskytuje přístup k analýze ordinálních dat, která je analogická k přístupu Cramér-von Misesovy statistiky pro spojitá data (viz J. Řehák [1976]).

IV. Aplikace D-modelu na testování dobré shody (asymptotická teorie)

U testů dobré shody tvoří D přirozenou ztrátovou funkci odvozenou z typu kategorizace a tudíž poskytuje váhy pro odchylky od předpokladu. Testy dobré shody můžeme založit na statistice typu (8):

$$(8) \quad d^2 = D^2(\mathcal{N}, f)$$

pro hypotézu $\mathcal{N} = p$. Její významnost lze určit podle algoritmu AS 106 (viz Sheil, O'Muir-cheartaigh [1977]).

Druhou možností je převedení na χ^2 distribuci pomocí následujícího postupu. Položme $f^* = (f_1, \dots, f_{k-1})$, obdobně \mathcal{N}^*, g^* . Existuje a takové, že $D^* = a a'$, $a' a = I_M$, kde I_M je jednotková matice typu $M \times M$, $M =$ hodnost D^* . Platí $(f^* - \mathcal{N}^*) = g^* \sim N(0, \Sigma)$ za platnosti H_0 . Pro všechna $\mathcal{N}_i > 0$ je rozložení regulární ($\sigma_{ij} = n^{-1}[\delta_{ij} \mathcal{N}_i - \mathcal{N}_i \mathcal{N}_j]$). Odtud $h = (h_1, \dots, h_M) = g^* a \sim N(0, a' \Sigma a)$ a tudíž

$$(9) \quad x_D^2 = h(a' \Sigma a)^{-1} h' = (f^* - \mathcal{N}^*) a (a' \Sigma a)^{-1} a' (f^* - \mathcal{N}^*)'$$

je za platnosti H_0 rozdělena jako χ_{M}^2 .

Pro ordinální D je $h = (F_1, \dots, F_{k-1})$, $\Sigma_0 = a' \Sigma a$ je určena prvky $\sigma_{0,ij} = \prod_i (1 - \prod_j)$ pro $i \leq j$, kde \prod_i je distribuční funkce k \mathcal{N} .

Pro kardinální D je problém převeden na obvyklý z-test $x_D^2 = z^2 = n(\bar{x}_f - \bar{x}_g)^2 / \text{var}(X/\mathcal{N})$. Nominální případ vede na x^2 z (1).

V. Aplikace na měření shody

Pro $A = \{A_1, A_2, \dots, A_k\} \cup \emptyset$, $\mathcal{N}, p, f \in \mathcal{Q}_k$ zavedla B. Řeháková [1980] koeficient shody

$$(10) \quad \rho = \frac{D(\mathcal{N}, p)}{\max_p D(\mathcal{N}, p)}$$

V praxi používáme výběrový analog (odhad metodou maximální věrohodnosti)

$$(11) \quad \hat{\rho} = \frac{D(\mathcal{N}, f)}{\max_p D(\mathcal{N}, p)}$$

Maximum ve jmenovateli (11) pro běžné typy proměnných (nominální, ordinální, kardinální) je

$$\sqrt{\frac{1 - 2x_{\min} + \sum_{k=1}^K x_k^2}{2 \max \left(\sum_{k=1}^{K-1} \pi_k^2, \sum_{k=1}^{K-1} (1 - \pi_k)^2 \right)}, \quad \pi_k = \sum_{i=1}^k x_i}$$

$$\sqrt{2 \max (x_{\max} - x_{\min}, x_{\max} - \bar{x}_x)},$$

kde x_{\min} , x_{\max} jsou minimální a maximální skóre přiřazené kategoriím, $\bar{x}_x = \sum_k x_k$.

Vlastnosti R :

- a) $0 \leq R \leq 1$; $R = 0 \iff x = f$ pro B - metriku
 $x = f \implies R = 0$ pro B - semimetriku
- b) $R \approx N \cdot K \left(\rho, \frac{v(R)}{N} \right)$, $v(R)$ jsou uvedeny v citované práci.

Asymptotická normalita umožňuje

- a) konstrukci intervalů spolehlivosti $R \pm z_{\alpha} \sqrt{\frac{v(R)}{N}}$,
- b) test $\rho = \rho_0$,
- c) testy $\rho_1 = \rho_2 = \dots = \rho_N$ a párová porovnání $\rho_i = \rho_j$ pro nezávislé výběry.

Poznámka : Koeficienty neshody i jejich testy našly praktickou aplikaci při vyhodnocování reprezentativity výzkumů veřejného mínění a při porovnávání distribucí s nezávislým standardem (např. shoda reality a záměru u působení časopisu, televize, rozhlasu ap.).

VI. Další aplikace B - modelu

1. Koeficienty asociace - speciální případ je Mallisovo η i korelační poměr η^2 aplikovaný na kontingenční tabulku, též koeficient β pro ordinální proměnné (viz Řehák [1976]).
2. Koeficienty parciální asociace.
3. Analýza kategorizovaných dat ve schématech ANOVA (včetně testů homogenity pro tabulku R x S).
4. Seskupování řádků tabulky (např. vážená centroidní metoda - viz Řehák, Řeháková [1982]).
5. Klasické multidimenzionální škálování řádků tabulky.

Výhody B-modelu spočívají

- a) v překonání obtíží, které byly diskutovány v úvodních částech referátu, při řešení všech komparačních úloh spojených s obecnými typy kategorizací;
- b) v poskytnutí jednotného pohledu a jednotné interpretace odvozených měr a statistických testů;
- c) v možnosti práce s obecnými typy kategorizací, ale i v jednoduché metodice pro analýzu ordinálních dat;
- d) v jednoduchém heuristickém základu a v možnosti jednoduchých grafických reprezentací výsledků (eukleidovské reprezentace simplexu Q_K).

Literatura

1. Light, R.J., Margolin, B.H. [1971] : An Analysis of Variance for Categorical Data. Journal of American Statistical Association 66, 534-544.
2. Rao, R.C. [1978] : Lineární metody statistické indukce a jejich aplikace. Academia, Praha.
3. Řehák, J. [1976] : Základní deskriptivní míry pro rozložení ordinálních dat. Sociologický časopis 4, 416-431.

4. Řehák, J., Řeháková B. [1979] : Základní charakteristiky proměnných s konečným počtem hodnot a distanční analýza jejich rozložení. Sociologický časopis 2, 214-233.
5. Řehák, J. Řeháková B. [1981] : Analýza kategorizovaných dat v sociologii (nepublikováno).
6. Řeháková, B. [1979] : Statistické ověřování reprezentativity : testy dobré shody. Sociologický časopis 6, 615-629.
7. Řeháková, B. [1980] : Statistické ověřování reprezentativity : koeficienty neshody. Sociologický časopis 6, 612-627.
8. Weiler, H. [1966] : A Coefficient Measuring the Goodness of Fit. Technometrics 2, 327-334.
9. Sheil J., O'Muircheartaigh I. [1977] : Algorithm AS 106, The Distribution of Non-negative Quadratic Forms in Normal Variables. Applied Statistics, 92-98.