

Asymptotické testy hypotéz pro rozdělení D-proměnné

Blanka Neháková, Ústav pro výzkum veřejného měření při FSÚ, Praha

Model pro obecné kategorizované proměnné, u jejichž hodnot (kategorií) uvažujeme číselně vyjádřitelné relace nepodobnosti, a jeho teorie byly popsány v pracích Neháka [1976], Neháka a Nehákové [1979], [1984], [1986a], [1986b], Nehákové [1982], [1985]. Příspěvek předkládá asymptotické testy dobré shody a homogenity pro kontingenční tabulky, v nichž se obecná proměnná (dále D-proměnná) vyskytuje v roli závisle proměnné.

Testy jsou založeny na předpokladu multinomického výběru vytvářejícího tabulku. Odtud pak plynou asymptotická normalita rozdělení četnosti a z ní použitelnost výsledků teorie rozdělení kvadratických forem normálně rozdělených veličin. Testové statistiky navržené v práci Nehákové [1985] vyplývají přirozeně z modelu, mají tvar kvadratických forem a vycházejí ze (semi)metriky.

Základní definice a vlastnosti

D-proměnnou $A = \{A_1, A_2, \dots, A_K; D\}$ nazveme úplnou soustavu jevů A_1, A_2, \dots, A_K spolu s maticí $D = [d_{jk}]$ typu $K \times K$, pro jejíž prvky $d_{ij} = d(A_i, A_j)$ platí

identita: $d_{ii} = 0$ pro $i = 1, \dots, K$,

symetrie: $d_{ij} = d_{ji}$ pro $i, j = 1, \dots, K$,

nezápornost: $d_{ij} \geq 0$ pro $i, j = 1, \dots, K$ a

$d_{ij} > 0$ pro alespoň jednu dvojici i, j ,

interpretabilita: d_{ij} je tím větší, čím nepodobnější jsou si kategorie A_i, A_j .

Prvky d_{ij} matice D nazveme skóry vzdálenosti kategorií A_i, A_j . D nazveme maticí skóru vytvářejících typ proměnné.

Nominální, ordinální a kardinální kategorizované proměnné jako speciální případy D-proměnné jsou vytvářeny po řadě maticemi s prvky $d_{ij} = 1 - \delta_{ij}$ ($\delta_{ij} = 1$ pro $i=j$, $\delta_{ij} = 0$ pro $i \neq j$), $d_{ij} = |i-j|$, $d_{ij} = (x_i - x_j)^2$ (x_i je číselná hodnota přiřazená kategorii A_i).

Označíme $Q_K = \{\underline{f} : \underline{f} = (f_1, f_2, \dots, f_K)', f_k \geq 0, k = 1, 2, \dots, K, \sum_1^K f_k = 1\}$

K-rozměrný simplex rozdělení. Pro volbu testových statistik jsou důležité následující vlastnosti:

1. Funkce $D(\underline{f}, \underline{g}) = \sqrt{(\underline{f}-\underline{g})' D (\underline{g}-\underline{f})}$ je (semi)metrika, právě když matici $D^{\frac{1}{2}} = [d_{ij}^{\frac{1}{2}}]$ typu $(K-1) \times (K-1)$, pro niž $d_{ij}^{\frac{1}{2}} = d_{ik} + d_{kj} - d_{ij}$ je pozitivně (semi)definitní matici.

2. Pro zobecněnou variaci $Gvar \underline{f} = \underline{f}' D \underline{f}$, $\underline{f} \in Q_K$ platí za předpokladu existence semimetriky $D(\dots)$ rozklad: je-li $\underline{f} = \sum w_r \underline{f}_{(r)}$, $\underline{f}_{(r)} \in Q_K$, $r=1, \dots, R$, $w_r \in Q_R$, pak

$$Gvar \underline{f} = \sum_r w_r Gvar \underline{f}_{(r)} + \frac{1}{2} \sum_r \sum_s w_r w_s D^2(\underline{f}_{(r)}, \underline{f}_{(s)}) = \sum_r w_r Gvar \underline{f}_{(r)} + \sum_r w_r D^2(\underline{f}_{(r)}, \underline{f}).$$

Testy jsou založeny na tomto výsledku (viz např. Ruben [1962]). Nechť $\underline{Z} = (Z_1, \dots, Z_K)'$ má rozdělení $M(Q, \underline{f})$, $h(\underline{f}) = K$ a $Y = \underline{Z}' \underline{A} \underline{Z}$, kde \underline{A} je symetrická pozitivně definitní nebo pozitivně semidefinitní matici. Pak $L[Y] = L[\sum_{i=1}^K \lambda_i Z_i^2]$, kde Z_1^2, \dots, Z_K^2 jsou nezávislé náhodné veličiny, z nichž každá má rozdělení chí-kvadrát o jednom stupni vlnosti a $\lambda_1, \dots, \lambda_K$ jsou charakteristická čísla matice \underline{A} .

Dosazené hladiny významnosti $1 - P(\underline{Z}' \underline{A} \underline{Z} \leq t)$ lze pocítat pomocí algoritmu AS 106 (viz Sheil, O'Muircheartaigh [1977], který kombinuje výsledky Rubena [1962] a Kotze, Johnsona, Boyd [1967]):

$$P(\underline{Z}' \underline{A} \underline{Z} \leq t) = \sum_{k=0}^{\infty} c_k P(K+2k, t/\beta), \quad t > 0$$
$$= 0, \quad t \leq 0$$

kde $F(K', y)$ je distribuční funkce centrálního rozdělení chi-kvadrát o K' stupních volnosti, $K' = h(\sum A_i) = h(A)$. (Volba konstanty β , výpočet c_K a další podrobnosti viz výše uvedená literatura.)
Následující testy A-C zavedla Reháková [1985], test D Rehák (nepublikováno).

A) Test dobré shody s předpokladem

Nechť $\underline{z} \in Q_K$ je vektor relativních četností pro $A = \{A_1, \dots, A_K; B\}$, $n \underline{z} \sim M_K(n, p)$.
Nechť $\underline{g} \in Q_K$ je daný vektor. Uvažujme hypotézu $H_0: p = \underline{g}$ a alternativu $H_A: D(p, \underline{g}) \neq 0$. Testová statistika je $D^2(\underline{z}, \underline{g}) = (\underline{z} - \underline{g})' D^{\frac{1}{2}} (\underline{z} - \underline{g})$, kde D je matici $D(p, \underline{g})$. Testová statistika je $L[n D^2(\underline{z}, \underline{g})] \rightarrow L\left[\sum_{i=1}^{K-1} \lambda_i z_i^2\right]$, kde z_i^2 ($i=1, \dots, K-1$) jsou nezávislé náhodné veličiny s rozdělením chi-kvadrát o jednom stupni volnosti a λ_i ($i=1, \dots, K-1$) jsou charakteristická čísla matice D , $D = D(p, \underline{g}) = D_{K-1} - \underline{g} \underline{g}'$, $D_{K-1} = \text{diag}\{x_1, \dots, x_{K-1}\}$. (Matica $D_{K-1} - \underline{g} \underline{g}'$ je regulařní právě když $x_k > 0$ pro $k=1, \dots, K$).

B) Test shody dvou rozdělení

Dvouvýběrový problém je u D-proměnných řešitelný pomocí statistiky $D^2(\underline{z}, \underline{g}) = (\underline{z} - \underline{g})' D^{\frac{1}{2}} (\underline{z} - \underline{g})$, kde $\underline{z}, \underline{g} \in Q_K$ jsou dvě nezávislá rozdělení proměnné $A = \{A_1, \dots, A_K; B\}$, $n \underline{z} \sim M_K(n, p)$, $m \underline{g} \sim M_K(m, p)$.

$$H_0: p = \underline{z}, \quad H_A: D(p, \underline{z}) \neq 0.$$

Za hypotézy H_0 platí, že když n a m jdou do nekonečna tak, že $\lim n/(m+n) = \omega$, $0 < \omega < 1$, pak

$$L[(n+m) D^2(\underline{z}, \underline{g})] \rightarrow L\left[\sum_{i=1}^{K-1} \lambda_i z_i^2\right],$$

kde z_i^2 jsou nezávislé náhodné veličiny s rozdělením chi kvadrát o jednom stupni volnosti a λ_i ($i=1, \dots, K-1$) jsou charakteristická čísla matice

$$\frac{1}{\omega(1-\omega)} (D_p - \bar{p} \bar{p}') D^{\frac{1}{2}}, \quad D_p = \text{diag}\{p_1, \dots, p_{K-1}\}.$$

V praxi vezmeme za ω číslo $m/(m+n)$ a místo \underline{p} dosadíme vektor se složkami $\hat{p}_i = (n_i + m_i) / (n + m)$ za předpokladu, že $n_i + m_i > 0$ ($i=1, \dots, K$).

Test je speciálním případem obecnější situace, v níž se porovnává $R > 2$ rozdělení (viz C). Jeho důležitost je dána nejen častým výskytem dvouvýběrového problému, ale též proto, že porovnání $R > 2$ rozdělení lze rozložit do $R(R-1)/2$ porovnání dvojic při použití simultánní inference.

C) Test shody R rozdělení

Nechť $\underline{z}(1), \dots, \underline{z}(R)$ jsou nezávislé náhodné vektory (řádková rozdělení četností) proměnné $A = \{A_1, \dots, A_K; B\}$ v tabulce $R \times K$, nechť $\underline{g} = (w_1, \dots, w_K)' \in Q_R$, $w_r > 0$ ($r=1, \dots, R$), speciálně $w_r = n_r/n$, kde n_r je rozsah r-tého výběru, $\underline{z} = \sum n_r \underline{z}(r)$. Rozdělení $n_r \underline{z}(r)$ je $M_K(n_r, p_{rx})$, $r = 1, \dots, R$. Označme $\underline{z} = \sum w_r \underline{z}(r)$.

Testovou statistikou je druhý člen rozkladové formule pro Gvar \underline{z} a to

$$\sum w_r D^2(\underline{z}(r), \underline{z}) = (\underline{z}(1) - \bar{z}, \dots, \underline{z}(R-1) - \bar{z})' (M_{R-1} \otimes D^{\frac{1}{2}}) (\underline{z}(1) - \bar{z}, \dots, \underline{z}(R-1) - \bar{z}).$$

Matice M_{R-1} je symetrická pozitivně definitní matice s prvky

$$w_{rs} = \frac{w_r (w_r + w_s)}{w_R}, \quad r \neq s, \quad r, s = 1, \dots, R-1,$$

\otimes je Kroneckerův součin matic.

$H_0: p_{(1)} = \dots = p_{(R)} (= p)$.

H_A : alespoň pro jednu dvojici $r \neq s$ je $D(p_{(r)}, p_{(s)}) \neq 0$. Nechť n_x jdou do nekonečna tak, že $n_x/n \Rightarrow \omega_x$, $0 < \omega_x < 1$, $(x=1, \dots, R)$, $\sum \omega_x = 1$. Pak za hypotézy H_0

$$L \left[n \sum_{r=1}^R w_r b^2(f_{\omega(r)}, f) \right] \rightarrow L \left[\sum_{i=1}^{(R-1)(K-1)} \lambda_i z_i^2 \right],$$

kde z_i^2 jsou nezávislé náhodné veličiny, z nichž každá má rozdělení chí-kvadrát o jednom stupni volnosti a λ_i ($i=1, \dots, (R-1)(K-1)$) jsou charakteristická čísla matices $(E_{R-1} \otimes E_{K-1}) (E_{R-1} \otimes E^T)$, kde $E_{K-1} = D_p - \bar{P} \bar{E}$, $D_p = \text{diag}\{p_1, \dots, p_{K-1}\}$

a E_{R-1} je symetrická pozitivně definitní matice typu $(R-1) \times (R-1)$ s prvky

$$b_{xx} = \frac{1}{\omega_x} - \frac{2w_x}{\omega_x} + \sum_{i=1}^R \frac{w_i^2}{\omega_i}, \quad b_{xs} = -\frac{w_x}{\omega_x} - \frac{w_s}{\omega_s} + \sum_{i=1}^R \frac{w_i^2}{\omega_i}, \quad x \neq s.$$

V praxi vezmeme za ω_x číslo n_x/n a za p dosadíme větše \hat{p} se složkami

$$\hat{p}_k = \sum_x n_{xk}/n \text{ za předpokladu, že } \hat{p}_k > 0 \text{ pro } k=1, \dots, K.$$

D) Test shody marginálních rozdělení ve čtvercové tabulce

Necké $A = \{A_1, \dots, A_K; B\}$, $B = \{B_1, \dots, B_K; B\}$ jsou dvě D-proměnné o stejném počtu kategorií a stejná matici sloučené D . Čtvercová tabulka $A \times B$ vznikla jako výběr o rozsahu n z multinomického rozdělení $M(n; p_{11}, \dots, p_{KK})$. Označme

n_{ij} četnost v poli (i,j) , $n_{+j} = \sum_i n_{ij}$, $n_{i+} = \sum_j n_{ij}$, $f_{i+} = n_{i+}/n$, $f_{+j} = n_{+j}/n$.

$H_0: E_A = E_B$, $E_A : D(p_A, p_B) \neq 0$, $p_A = (p_{1+}, \dots, p_{K+})^T$, $p_B = (p_{+1}, \dots, p_{+K})^T$.

Testová statistika je $D^2(f_A, f_B) = (\bar{f}_A - \bar{f}_B)^T \bar{E}^T (\bar{f}_A - \bar{f}_B)$, $\bar{f}_A = (f_{1+}, \dots, f_{K-1, +})^T$,

$\bar{f}_B = (f_{+1}, \dots, f_{+, K-1})^T$. Za platnosti H_0

$$L \left[n D^2(f_A, f_B) \right] \rightarrow L \left[\sum_{i=1}^{K-1} \lambda_i z_i^2 \right],$$

kde z_i^2 jsou nezávislé náhodné veličiny, z nichž každá má rozdělení chí-kvadrát s jedním stupněm volnosti a λ_i ($i=1, \dots, K-1$) jsou charakteristická čísla matices \bar{E}^T , kde \bar{E} je symetrická pozitivně definitní matice typu $(K-1) \times (K-1)$ s prvky

$$\sigma_{11} = p_{1+} + p_{+1} - 2p_{11}, \quad \sigma_{ij} = -(p_{ij} + p_{ji}), \quad i \neq j$$

předpokládáme $p_{ij} > 0$ pro $i, j = 1, \dots, K$.

V praxi použijeme odhadu $\hat{p}_{ij} = n_{ij}/n$, $\hat{p}_{1+} = n_{1+}/n$, $\hat{p}_{+1} = n_{+1}/n$.

LITERATURA

Keta S., Johnson N.L., Boyd D.W. [1967] Series Representations of Distributions of Quadratic Forms in Normal Variables I. Central Case. AMS 38, 838-848.

Zuban H. [1962] Probability Content of Regions under Spherical Normal Distributions, IV: The Distribution of Homogeneous and Non-homogeneous Quadratic Functions of Normal Variables. AMS 33, 542-570.

Beňák J. [1976] Základní deskriptivní míry pro rozložení ordinálních dat. Sociologický časopis XII, 416-431.

Beňák J., Beňáková B. [1979] Základní charakteristiky proměnných s konečným počtem hodnot a distanční analýza jejich rozložení. Sociologický časopis XIV, 214-231.

Beňák J., Beňáková B. [1984] Parciální asociacné koeficienty v kontingenčních tabulkách. In: Šantoch J., Jurečková J. (ed.) Robust 84, 105-108.

Beňák J., Beňáková B. [1986a] Classifications with relations: a model for the

- description of distributions and their distances. Kybernetika (v tisku).
- Mehák J., Meháková B. [1986b] Vícenásobná a parciální asociace v kontingenčních tabulkách. Sociologický časopis (v tisku).
- Meháková B. [1982] Koeficienty parciální asociace pro obecní typ kategorizované proměnné. Práce k aspirantskému minimu.
- Meháková B. [1985] Model a metoda pro analýzu kategorizovaných dat s relacemi. Kandidátská disertační práce.
- Sheil J., O'Muircheartaigh I. [1977] The Distribution of Non-negative Quadratic Forms in Normal Variables. Applied Statistics 26, 92-98.