

Jaromír Beláček, GGÚ ČSAV

Stejnomený příspěvek ze sborníku ROBUST 86 poukázal na možnost aplikace nej-
jednodušších transformací typu jackknife na koeficientech asociace (matematicky hlad-
kých funkcích frekvenčního multinomického vektoru), které poskytují méně vychýlené od-
hady empirických koeficientů a alternativní (neparametrické) odhady rozptylu. Smyslem
toboto příspěvku je poskytnout informaci o univerzálních zobecněních těchto metod kon-
struovaných pro eliminaci vychylujících členů vyšších řádů nežli $O(n^{-1})$, a to i pro
situaci několika nezávislých náhodných výběrů. Na příkladech vybraných koeficientů aso-
ciace jsou referovány některé zobecnitelné praktické závěry.

V návaznosti na terminologickou diskusi odehrávající se po odeznění prvotního
příspěvku na ZŠ ROBUST 86 si dovoluji upozornit, že termín "jackknife" znamená podle
doslovného překladu (viz. [5], str. 1174) "velký kapesní zavírací nůž, zavírák". V tomto
smyslu interpretoval "jackknife" i profesor P.K.Sen na své loňské předvánoční přednášce
na KPMS přiřazující jednotlivým čepelím svého "nože" rozličné speciální interpretace.
Jak mne ovšem upozornil můj pozorný kolega dr.P.Hrál ([6]), lze na základě věrných his-
torických pramenů přisuzovat anglické terminologické přiřazení dané metodě pracovním
aktivitám českých národnostních menšin před druhou světovou válkou. Jak známo Valaši
(jinak též profesionální klestiči dobytka) se proslavili právě výrobou a aktivním pou-
žíváním kapesních nožů, jejichž prostřednictvím účinně potlačovali nežádoucí aktivity
domácího skotu po celé Evropě. Ilužno podotknout, že kategorická informace o tom, že
adekvátní český ekvivalent anglického "jackknife" je tudíž "oklešťování", způsobila po-
sitivní nezáměr mého nejbližšího nadřízeného dr.J.Řeháka o další terminologické nuance.

1. JEDNOVÝBĚROVÁ JACKKNIFEOVÁ TRANSFORMACE K-TÉHO ŘÁDU

Pro libovolný odhad Θ_n pořízený na základě náhodného výběru o rozsahu n a
 $k < n$ zavedeme veličiny

$$(1) \quad \bar{\Theta}_{n-v} \begin{cases} \binom{n}{v}^{-1} \sum_{1 \leq i_1 < \dots < i_v \leq n} \Theta_{-i_1, \dots, -i_v} & \text{pro } v=1, \dots, k \\ \Theta_n & \text{pro } v=0 \end{cases}$$

jako průměry všech $\binom{n}{v}$ odhadů získaných z původního výběru při postupném vynechávání
 v -tic pozorování. Jednovýběrový jackknifeový odhad k -tého řádu byl definován v [10]
předpisem

$$(2) \quad J_k(\Theta_n) := \begin{vmatrix} \bar{\Theta}_{n-0} & \bar{\Theta}_{n-1} & \cdot & \cdot & \cdot & \bar{\Theta}_{n-k} \\ n^{-1} & (n-1)^{-1} & \cdot & \cdot & \cdot & (n-k)^{-1} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ n^{-k} & (n-1)^{-k} & \cdot & \cdot & \cdot & (n-k)^{-k} \\ \hline 1 & 1 & \cdot & \cdot & \cdot & 1 \\ n^{-1} & (n-1)^{-1} & \cdot & \cdot & \cdot & (n-k)^{-1} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ n^{-k} & (n-1)^{-k} & \cdot & \cdot & \cdot & (n-k)^{-k} \end{vmatrix}$$

jakožto řešení pro $\hat{\Theta}$ soustavy lineárních rovnic

$$(3) \quad \hat{\Theta}_{n-\hat{v}} = \hat{\Theta} + \sum_{v=1}^k \hat{\alpha}_v \cdot (n-\hat{v})^{-v}, \quad \hat{v}=0,1,\dots,k$$

v proměnných $\hat{\Theta}, \hat{\alpha}_1, \dots, \hat{\alpha}_k$. Vzhledem k tomu, že Vandermondův determinant matice soustavy (3) - tj. jmenovatel zlomku v (2) - je různý od nuly, je $J_k(\hat{\Theta}_n)$ definován jednoznačně jako lineární kombinace $\hat{\Theta}_{n-0}, \dots, \hat{\Theta}_{n-k}$ s pevně určenými koeficienty. Rozvojem čitatele v (2) podle prvního řádku lze získat explicitní vyjádření

$$(4) \quad J_k(\hat{\Theta}_n) = \frac{1}{k!} \left[\sum_{v=0}^k (-1)^v \binom{k}{v} (n-v)^k \hat{\Theta}_{n-v} \right]$$

snad poprvé uvedené v [12].

Základní vlastnost transformace $J_k(\cdot)$ vyplývá z níže uvedeného tvrzení - Lemma 1: Necht' odhad $\hat{\Theta}_n$ má konečný rozvoj střední hodnoty typu

$$(5) \quad E \hat{\Theta}_n = \Theta + \sum_{v=1}^{\infty} \alpha_v \cdot n^{-v},$$

kde parametry $\Theta, \alpha_1, \alpha_2, \dots$ nezávisí na n . Pak platí

$$(6) \quad E J_k(\hat{\Theta}_n) = \Theta + O(n^{-k-1}), \quad n \rightarrow \infty.$$

Důkaz viz [10], str. 527 nebo v [3] na str. 54-5.

Je důležité poznamenat, že význam předpokladu (5) nespočívá v explicitní znalosti parametrů $\Theta, \alpha_1, \alpha_2, \dots$, nýbrž v pouhé specifikaci struktury vychylujícího rozvoje odhadu $\hat{\Theta}_n$. Platnost (5) se zbytkem řady nahrazeným symbolem $O(n^{-v})$ pro některé $v > k$ lze ověřit pro většinu standardních odhadů $\hat{\Theta}_n$ založených na posloupnosti nezávislých stejně rozdělených náhodných veličin (viz také [11], §6).

Zajímavá je interpretace vlastnosti (6) z hlediska explicitního vyjádření (4). Koeficienty $(-1)^v \binom{k}{v}$ pro $v=0,1,\dots,k$ charakterizují k-tou zpětnou diferencí běžně zaváděnou pro libovolnou posloupnost $\{x_n\}$ reálných čísel rekurentními vztahy

$$(7) \quad \begin{aligned} \nabla x_n &:= x_n - x_{n-1} \\ \nabla^2 x_n &:= \nabla x_n - \nabla x_{n-1} = x_n - 2x_{n-1} + x_{n-2} \\ &\vdots \\ \nabla^k x_n &:= \nabla^{k-1} x_n - \nabla^{k-1} x_{n-1} = \sum_{v=0}^k (-1)^v \binom{k}{v} x_{n-v} \end{aligned}$$

Není těžké ukázat, že je-li x_n nezáporná mocnina n řádu $j \geq 0$, pak $\nabla^k x_n$ je polynom stupně $j-k$ při $j < k$ identicky rovný nule; dokonce každý pól v nule typu n^{-j} pro kladné j převádí k -tá diference na pól typu $O(n^{-j-k})$. V situaci odpovídající splnění (5) definujeme

$$(8) \quad x_n := n^k \cdot E \hat{\Theta}_n = \Theta \cdot n^k + \sum_{v=1}^{\infty} \alpha_v \cdot n^{k-v}, \quad n > k,$$

takže k -tá diference aplikovaná na x_n (tj. $\nabla^k x_n = (k!) E [J_k(\hat{\Theta}_n)]$) musí být nutně výraz typu (6).

Explicitní znalost transformačního vzorce (4) umožňuje vyšetřovat otázku asymptotického rozdělení $J_k(\hat{\Theta}_n)$ ve vztahu k původnímu rozdělení odhadu $\hat{\Theta}_n$. Platí-li asymptotická normalita pro n.v. $n^{1/2}(\hat{\Theta}_n - \Theta)$, lze stejné normální aproximace použít obvykle i v případě n.v. $n^{1/2}(J_k(\hat{\Theta}_n) - \Theta)$. Odhadem společného asymptotického rozptylu je statistika

$$(9) \quad S_1^2(\hat{\Theta}_n) := (n-1) \sum_{i=1}^n (\hat{\Theta}_{-i} - \hat{\Theta}_{n-1})^2$$

2. S-VÝBĚROVÁ JACKKNIFEOVÁ TRANSFORMACE K-TÉHO ŘÁDU

Definici odhadu (2) lze analogicky rozšířit na situace, kdy parametr Θ je odhadován na základě $S(s+1)$ nezávislých náhodných výběrů s rozsahy n_1, \dots, n_s ($n_1 + \dots + n_s = n$). Pokud zanedbáme členy řádů $O(n^{-s-1})$, lze typickou strukturu vychylujícího rozvoje odpovídajícího odhadu Θ_n^s charakterizovat vyjádřením

$$(10) \quad E \Theta_n^s = \Theta + \sum_{r=1}^{s+k} \sum_{1 \leq a_1 < \dots < a_r \leq s} \sum_{r \leq v_1 + \dots + v_r \leq k} \alpha_{v_1 \dots v_r}^{a_1 \dots a_r} \cdot n_{a_1}^{-v_1} \dots n_{a_r}^{-v_r}$$

kde Θ a všechna $\alpha_{v_1 \dots v_r}^{a_1 \dots a_r}$ jsou parametry nezávislé na n_1, \dots, n_s a poslední sumace sčítá přes systém $\{[v_1, \dots, v_r]; 1 \leq v_1, \dots, v_r \leq k \text{ a } v_1 + \dots + v_r \leq k\}$. V nejjednodušším případě $k=1$ lze jackknifeovou transformaci vyjádřit jako řešení soustavy lineárních rovnic definované vzorcem (viz [7], str. 1013)

$$(11) \quad J_1^s(\Theta_n^s) := \begin{vmatrix} \Theta_n^s & \bar{\Theta}_1^1 & \dots & \dots & \bar{\Theta}_1^s \\ n_1^{-1} & (n_1-1)^{-1} & \dots & \dots & n_1^{-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ n_s^{-1} & n_s^{-1} & \dots & \dots & (n_s-1)^{-1} \end{vmatrix} = (n-s+1)\Theta_n^s - \sum_{a=1}^s (n_a-1)\bar{\Theta}_1^a$$

$$\begin{vmatrix} 1 & 1 & \dots & \dots & 1 \\ n_1^{-1} & (n_1-1)^{-1} & \dots & \dots & n_1^{-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ n_s^{-1} & n_s^{-1} & \dots & \dots & (n_s-1)^{-1} \end{vmatrix}$$

kde $\bar{\Theta}_1^a$ pro $a=1, \dots, s$ je průměr všech n_a odhadů získaných z Θ_n^s vynecháním právě jednoho pozorování v rámci stratifikační proměnné a .

Pro obecné $k < \min\{n_1, \dots, n_s\}$ máme k dispozici systém pomocných odhadů

$$(12) \quad \hat{\Theta}_{\hat{v}_1 \dots \hat{v}_r}^{\hat{a}_1 \dots \hat{a}_r} := \left(\hat{v}_1^{\hat{a}_1}\right)^{-1} \dots \left(\hat{v}_r^{\hat{a}_r}\right)^{-1} \sum_{1 \leq \hat{a}_1 < \dots < \hat{a}_r \leq n_{\hat{a}_1}} \dots \sum_{1 \leq \hat{a}_r < \dots < \hat{a}_r \leq n_{\hat{a}_r}} \Theta_{-1_{\hat{a}_1}^1, \dots, -1_{\hat{a}_1}^1, \dots, -1_{\hat{a}_r}^1, \dots, -1_{\hat{a}_r}^1}$$

při

$$(13) \quad 1 \leq r \leq k, \quad 1 \leq \hat{a}_1 < \dots < \hat{a}_r \leq s, \quad \hat{r} = \hat{v}_1 + \dots + \hat{v}_r \leq k$$

Pro libovolnou sestavu indexů z (13) zřejmě platí

$$(14) \quad E \hat{\Theta}_{\hat{v}_1 \dots \hat{v}_r}^{\hat{a}_1 \dots \hat{a}_r} = \Theta + \sum_{r=1}^{s+k} \sum_{1 \leq a_1 < \dots < a_r \leq s} \sum_{r \leq v_1 + \dots + v_r \leq k} \alpha_{v_1 \dots v_r}^{a_1 \dots a_r} \cdot \tilde{n}_{a_1}^{-v_1} \dots \tilde{n}_{a_r}^{-v_r}$$

když

$$(15) \quad \tilde{n}_a = \begin{cases} n_a & \text{pro } a \in \{\hat{a}_1, \dots, \hat{a}_r\} \\ (n_a - \hat{v}_o) & \text{kdykoli } a \text{ je rovno některému indexu } \hat{a}_o \\ & \text{z množiny } \{\hat{a}_1, \dots, \hat{a}_r\} \end{cases}$$

Nahradíme-li střední hodnoty na levých stranách rovnic (10) a (14) jejich přirozenými odhady Θ_n^s a $\hat{\Theta}_{\hat{v}_1 \dots \hat{v}_r}^{\hat{a}_1 \dots \hat{a}_r}$ a všechny teoretické parametry doplníme symbolem $\hat{\cdot}$, definuje řešení pro $\hat{\Theta}$ takto vzniklé soustavy $\binom{s+k}{s+k}$ lineárních rovnic o $\binom{s+k}{s+k}$ neznámých

odhad $J_k^s(\theta_n^s)$ zobecnující jednovýberový odhad ze stati 1 (speciálně $J_k^s(\theta_n^s) = J_k^1(\theta_n^1)$).

Jednoznačnost definice i přesnější představu o struktuře odhadu $J_k^s(\theta_n^s)$ poskytuje následující technické tvrzení -

Lemma 2: Zavedme symbol $\nabla^{a_1 \dots a_r}$ pro $\binom{k}{r}$ složkový vektor sestavený z prvků

$$(16) \quad \nabla_{v_1 \dots v_r}^{a_1 \dots a_r} := \hat{\theta}_{v_1 \dots v_r}^{a_1 \dots a_r} - [\hat{\theta}_{v_1 \dots v_{r-1}}^{a_1 \dots a_{r-1}} + \dots + \hat{\theta}_{v_2 \dots v_r}^{a_2 \dots a_r}] + \dots + (-1)^{r-1} [\hat{\theta}_{v_1}^{a_1} + \dots + \hat{\theta}_{v_r}^{a_r}] + (-1)^r \hat{\theta}_n^s$$

pro systém indexů korespondující s (13).

Platí

$$(17) \quad J_k^s(\theta_n^s) = \theta_n^s + \sum_{r=1}^{s \wedge k} (-1)^r \left\{ \sum_{1 \leq a_1 < \dots < a_r \leq s} \hat{\alpha}_{(r)}^s(\nabla^{a_1 \dots a_r}) \right\},$$

kde $\hat{\alpha}_{(r)}^s(\nabla^{a_1 \dots a_r})$ je jediné řešení pro $\hat{\alpha}_{1 \dots 1}^{a_1 \dots a_r}$ soustavy lineárních rovnic

$$(18) \quad \frac{\nabla_{v_1 \dots v_r}^{a_1 \dots a_r}}{\hat{v}_1 \dots \hat{v}_r} = \sum_{r \wedge v_1 + \dots + v_r \leq k} \hat{v}_{v_1 \dots v_r}^{a_1 \dots a_r} \cdot (n_1 - \hat{v}_1)^{-v_1} \dots (n_r - \hat{v}_r)^{-v_r}, \quad 1 \leq r \leq s \wedge k$$

Důkaz je podrobně rozepsán v [3].

Při $k < \min\{n_1, \dots, n_s\}$ lze ověřit jednoznačnost řešení každé soustavy typu (18) a za předpokladu (10) i požadovanou vlastnost

$$(19) \quad E J_k^s(\theta_n^s) = \theta_n^s.$$

Explicitní vyjádření (17) pro $k=1, 2, 3$ udávají vzorce (11) a

$$(20) \quad J_2^s(\theta_n^s) = \theta_n^s - \left\{ \sum_{a=1}^s [(n_a-1)^2 \nabla_1^a - \frac{1}{2}(n_a-2)^2 \nabla_2^a] \right\} + \left\{ \sum_{1 \leq a_1 < a_2 \leq s} [(n_{a_1}-1)(n_{a_2}-2) \nabla_{11}^{a_1 a_2}] \right\},$$

$$(21) \quad J_3^s(\theta_n^s) = \theta_n^s - \left\{ \frac{1}{6} \sum_{a=1}^s [3(n_a-1)^3 \nabla_1^a - 3(n_a-2)^3 \nabla_2^a + (n_a-3)^3 \nabla_3^a] \right\} +$$

$$+ \left\{ \frac{1}{2} \sum_{1 \leq a_1 < a_2 \leq s} [(n_{a_1}-2)^2 (n_{a_2}-1) \nabla_{21}^{a_1 a_2} - 2(n_{a_1}-1)(n_{a_2}-1) \begin{vmatrix} (n_{a_1}-1) & (n_{a_2}-2) \\ (n_{a_2}-2) & (n_{a_1}-1) \end{vmatrix} \nabla_{11}^{a_1 a_2} + \right.$$

$$\left. + (n_{a_1}-1)(n_{a_2}-2)^2 \nabla_{12}^{a_1 a_2}] \right\} - \left\{ \sum_{1 \leq a_1 < a_2 < a_3 \leq s} (n_{a_1}-1)(n_{a_2}-1)(n_{a_3}-1) \nabla_{111}^{a_1 a_2 a_3} \right\},$$

kde každá sumace přes systém $1 \leq a_1 < \dots < a_r \leq s$ pro $r > s$ je definována jako nula.

3. APLIKACE NA KOEFICIENTY ASOCIACE

Uvažujme koeficienty asociace typu $\hat{\theta}(\underline{f}_n^s)$, kde $\hat{\theta}(\cdot): \mathbb{R}^m \rightarrow \mathbb{R}$ je funkce a proměnných omezená na svém definičním oboru a $\underline{f}_n^s := (f_{n_1}^1, \dots, f_{n_s}^s)$ sestává z s nezávislých vektorů odpovídajících po řadě c -rozměrným multinomickým rozdělením $M(n_1, p^1), \dots, M(n_s, p^s)$. Příklad $s=1$ pokrývá jednovýberovou situaci z příspěvku [2], při $s>1$ odpovídají nezávislé výběry přirozené stratifikaci vzhledem k hodnotám některé proměnné kontingenční tabulky. Korektní aplikace jackknifeové transformace $J_k^s(\cdot)$ ze stati 2 požaduje pro $\theta_n^s = \hat{\theta}(\underline{f}_n^s)$ strukturu vychylujícího rozvoje typu (10). Za předpokladu dostatečné hladkosti funkce $\hat{\theta}(\cdot)$ na okolí teoretického bodu $\underline{p} := (p^1, \dots, p^s)$ lze tuto vlastnost odvodit ze standardní aproximace Taylorovým rozvojem vyjádřené ve tvaru

$$(22) \quad [E \hat{\theta}(\underline{f}_n^s) - \hat{\theta}(\underline{p})] = \sum_{j=1}^{2k} \frac{1}{j!} D_j \cdot [E(d\underline{f}_n^s)^j] + o(n^{-k}), \quad n \rightarrow \infty,$$

kde D_j je matice j -tých parciálních derivací $\hat{\theta}(\cdot)$ v bodě \underline{p} blokově strukturovaná jakožto $D_j := [D_j^{a_1 \dots a_j}]_{1 \leq a_1, \dots, a_j \leq s}$, $[E(d\underline{f}_n^s)^j]$ symbolizuje matici j -tých centrálních momentů $m = (s \cdot c)$ -složkového vektoru $(d\underline{f}_n^s) := (f_{n_1}^s - p)$ a tečka značí skalární součin.

Elementární rozpis každého součinu z (22) pro matice L_j symetrické vzhledem k permutacím indexů

$$(23) D_j \cdot [E(df_{\underline{n}}^a)^j] = \sum_{r=1}^s \sum_{1 \leq a_1 < \dots < a_r \leq s} \sum_{\substack{j_1, \dots, j_r \\ j_1 + \dots + j_r = j}} \frac{j!}{j_1! \dots j_r!} D_j^{a_1 \dots a_r} [E(df_{n_{a_1}}^{a_1})^{j_1} \dots E(df_{n_{a_r}}^{a_r})^{j_r}]$$

umožňuje využít explicitních vyjádření (7) resp. (8) z článku [2] pro výpočet každého dílčího součinu z (23) typu $D_j^{a_1 \dots a_r} \cdot [E(df_{\underline{n}}^a)^j]$. S ohledem na (22) a (23) lze snadno nahlédnout explicitní závislost každého koeficientu $\alpha_{v_1 \dots v_r}^{a_1 \dots a_r}$ z (10) pouze na parametrech p^{a_1}, \dots, p^{a_r} a dané sestavě v_1, \dots, v_r . Speciálně při vynechání všech nulových sčítanců v (23) - odpovídajících situacím $j_a = 1$ pro některé $a=1, \dots, r$ - dostaneme vyjádření koeficientů u řádů n_a^{-1}, n_a^{-2} a $n_{a_1}^{-1} n_{a_2}^{-1}$ tvaru

$$(24) \begin{aligned} \alpha_1^a &= \frac{1}{2} D_2^a \cdot [E(df_{\underline{n}}^a)^2], & a=1, \dots, s, \\ \alpha_2^a &= \frac{1}{6} D_3^a \cdot [E(df_{\underline{n}}^a)^3] + \frac{1}{8} D_4^a \cdot [E(df_{\underline{n}}^a)^2 \times E(df_{\underline{n}}^a)^2], \\ \alpha_{11}^{a_1 a_2} &= \frac{1}{4} D_4^{a_1 a_2} \cdot [E(df_{n_{a_1}}^{a_1})^2 \times E(df_{n_{a_2}}^{a_2})^2], & 1 \leq a_1 < a_2 \leq s, \end{aligned}$$

kde e^a, e^{a_1}, e^{a_2} jsou nezávislé multinomické vektory s rozdělením $M(1, p^a), M(1, p^{a_1})$ a $M(1, p^{a_2})$. Analogické výpočtové vzorce se komplikují s přechodem k vyšším řádům $O(n^{-k})$ a k nesymetrickým maticím L_j .

Naznačené možnosti teoretického výpočtu vychylujících parametrů pro každý zvolený koeficient $\delta(\underline{f}_{\underline{n}}^a)$ a jejich následné využití pro účely eliminace vychýlení jsou ovšem nivelizovány při aplikaci transformací jackknife, které mají univerzální schopnost eliminace těchto členů bez ohledu na analytický předpis $\delta(\cdot)$. Výpočtové vzorce (12) pro $\delta_{v_1 \dots v_r}^{a_1 \dots a_r} := \delta(\underline{f}_{v_1 \dots v_r}^{a_1 \dots a_r})$ budou mít sice složitější strukturu typu

$$(25) \delta(\underline{f}_{v_1 \dots v_r}^{1 \dots r}) = \left(\binom{n_1}{v_1}^{-1} \dots \binom{n_r}{v_r}^{-1} \sum_{\substack{0 \leq m_1^1, \dots, m_c^1 \leq v_1 \\ m_1^1 + \dots + m_c^1 = v_1 \\ \vdots \\ 0 \leq m_1^r, \dots, m_c^r \leq v_r \\ m_1^r + \dots + m_c^r = v_r}} \left\{ \binom{n_1^1}{m_1^1} \dots \binom{n_c^1}{m_c^1} \dots \binom{n_1^r}{m_1^r} \dots \binom{n_c^r}{m_c^r} \delta(\underline{f}_{-(m_1^1 \dots m_c^1), \dots, -(m_1^r \dots m_c^r)}) \right\} \right)$$

kde sumace v (25) sčítá přes počty $m_1^1, \dots, m_c^1, \dots, m_1^r, \dots, m_c^r$ všech sestav indexů vyhovujících po řadě z 1. až c-té třídy 1. až r-té stratifikační proměnné a $n_1^1, \dots, n_c^1, \dots, n_1^r, \dots, n_c^r$ jsou odpovídající absolutní četnosti (pokud $n_b^a < m_b^a$ pro některou dvojici $[a, b] \in \{1, \dots, r\} \times \{1, \dots, c\}$, je celý výraz ve složené závorce definován jako nula), ale jsou formálně stejné pro libovolný koeficient $\delta(\underline{f}_{\underline{n}}^a)$.

4. ZÁVĚRY ZE SIMULAČNÍHO EXPERIMENTU

Statistické chování jackknifeových transformací bylo zkoumáno na jednovýběrové $J_k(\cdot)$ pro $k=1, \dots, 10$ aplikované na korelačním poměru η^2 (viz [9], str. 251) počítaném z kontingenční tabulky 4×5 a na logaritmu součinného poměru v tabulce 2×2 , který lze chápat jako speciální případ logaritmické interakce z [4]. V obou případech lze teoreticky zdůvodnit volbu řádu k jackknifeové transformace odpovídající faktické numerické významnosti několika prvních členů v (5) pro dostatečně velké n (kupř. $n \geq 100$).

V případě korelačního poměru η^2 je možno ověřit nárůst významnosti teoretických koeficientů z (24) při snižující se hodnotě některé z marginálních řádkových (resp. stratifikačních) pravděpodobností - z teoretických i simulačních výsledků však vyplývá pro všechny běžné volby teoretických parametrů (a rozsahů n) potřeba použití $J_k(\cdot)$ (resp. $J_k^{\#}(\cdot)$) pro k nejvýše rovné dvěma. Méně jednoznačná je volba "správného řádu k " v situaci logaritmické interakce, jejíž teoretickou hodnotu pro jednovýběrový případ lze definovat vztahem (viz [1], str.101)

$$(26) \quad \delta_{\underline{y}}(p) := \sum_{m=1}^c \gamma_m \cdot \log p_m \quad ,$$

kde $p := (p_1, \dots, p_c)$ je pravděpodobnostní rozdělení obecně vícerozměrné kontingenční tabulky s nenulovými složkami a $\underline{\gamma} := (\gamma_1, \dots, \gamma_c)$ nenulový vektor konstant vyhovující podmínce $\gamma_1 + \dots + \gamma_c = 0$. Problémy s definicí (26) na hranici jednotkového simplexu se v případě výběrového koeficientu $\delta_{\underline{y}}(\underline{f}_n^1)$ kompenzují připočtením jisté korekční konstanty ke všem absolutním četnostem analyzované tabulky. Dá se ukázat, že tato úprava nic nezmění na platnosti asymptotického vztahu

$$(27) \quad \begin{aligned} \delta_{\underline{y}}(\underline{f}_n^1) &= \delta_{\underline{y}}(p) - n^{-1} \left\{ \frac{1}{2} \Gamma_{-1} \right\} + n^{-2} \left\{ \frac{1}{2} \Gamma_{-1} - \frac{5}{12} \Gamma_{-2} \right\} - n^{-3} \left\{ \frac{1}{2} \Gamma_{-1} - \frac{1}{4} \Gamma_{-2} + \frac{3}{4} \Gamma_{-3} \right\} + \\ &+ n^{-4} \left\{ \frac{1}{2} \Gamma_{-1} - 2 \frac{11}{12} \Gamma_{-2} + 4 \frac{1}{2} \Gamma_{-3} - 2 \frac{11}{120} \Gamma_{-4} \right\} - \\ &- n^{-5} \left\{ \frac{1}{2} \Gamma_{-1} - 6 \frac{1}{4} \Gamma_{-2} + 18 \frac{3}{4} \Gamma_{-3} - 20 \frac{11}{12} \Gamma_{-4} + 7 \frac{11}{12} \Gamma_{-5} \right\} + \\ &+ n^{-6} \left\{ \frac{1}{2} \Gamma_{-1} - 12 \frac{11}{12} \Gamma_{-2} + 67 \frac{1}{2} \Gamma_{-3} - 135 \frac{23}{24} \Gamma_{-4} + 118 \frac{3}{4} \Gamma_{-5} - 37 \frac{439}{504} \Gamma_{-6} \right\} + o(n^{-6}), n \rightarrow \infty, \end{aligned}$$

$$(28) \quad \text{pro} \quad \Gamma_{-v} := \sum_{m=1}^c (\gamma_m) (p_m)^{-v}, \quad v=1, \dots, 6 \quad ,$$

který lze odvodit z aproximací typu (2c) při $k=6$ s ohledem na jednoduchou (diagonální) strukturu všech matic parciálních derivací a na znalost centrálních momentů elementárního multinomického rozdělení (viz [3], str.35-7). Faktická významnost koeficientů stojících v (27) u jednotlivých řádů n^{-v} tedy bezprostředně závisí na rychlosti konvergence hodnoty $1/(n \cdot \min\{p_1, \dots, p_c\})$ k nule. Pro dosti malá n ($n \leq 50$) a konkrétní volby vektorů $\underline{\gamma}$ a \underline{p} je tudíž oprávněno použití jacksonifových transformací vyššího řádu k .

Formální simulační aplikace $J_k(\cdot)$ na oba uvažované koeficienty přinesla tyto zobecněné závěry:

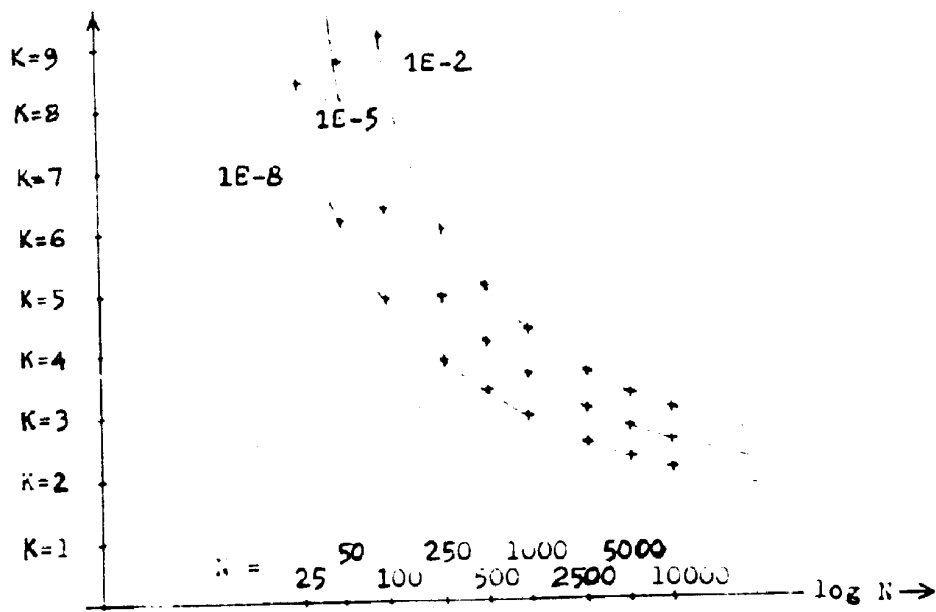
i/ Obecný vzorec (4) nelze použít pro neomezeně velká n a k , neboť sčítání a odčítání příslušných členů typu $O(n^k)$ klade zvýšené nároky na numerickou přesnost použité výpočetní techniky. Z tohoto důvodu je nutné počítat hodnoty koeficientů asociace, na něž se aplikuje transformace $J_k(\cdot)$ při $k > 2$, s vysokou přesností (double precision). Volba maximálního řádu k v závislosti na rozsahu n a zvolené přesnosti zaokrouhlovací chyby vyplývá z empirického grafu na obr.1.

ii/ Počet $\binom{n}{k}$ vynechávaných pozorování pro účely následného přepočítávání hodnot původního odhadu podle (1) lze v případě koeficientů asociace úspěšně redukovat na počet $\binom{c+k-1}{k}$ pro c kategoriálních tříd. Jak ovšem vyplývá z grafu na obr.2, vzrůstá při zvyšující se c a k i tento počet nad časově zvládnutelnou hranici.

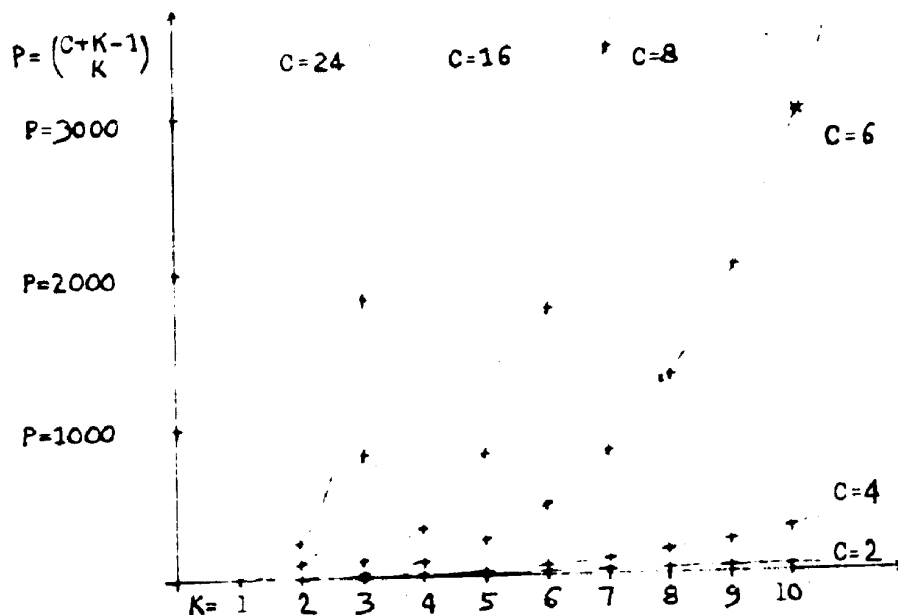
iii/ Pokud zvolený koeficient $\delta(p)$ není definován na celém jednotkovém simplexu, může dojít k situaci, že od jistého k není (na rozdíl od empirické hodnoty $O(\underline{f}_n^1)$) definován odhad $J_k(\delta(\underline{f}_n^1))$. Bodové vlastnosti jacksonifových odhadů jsou potom při malých n velmi citlivé na případná predefinování-korekce jako za vzorcem (26).

iv/ Respektujeme-li omezení vyplývající z i/-iii/, vykazují všechny numericky dostupné odhady značnou stabilitu projevující se vysokou párovou korelací mezi odhady $J_1(\delta(\underline{f}_n^1)), \dots, J_k(\delta(\underline{f}_n^1))$, z nichž všechny se statisticky významně liší od původní-vychýlené hodnoty $\delta(\underline{f}_n^1)$. Pro malé rozsahy n je ovšem třeba mít na paměti, že rozptyl jacksonifových odhadů je obvykle téhož řádu jako původní odhad, takže jednotlivé pozorované hodnoty -přestože jsou fakticky nevychýlené- mohou být značně variabilní!!

Obr.1: Závislost maximálního výpočetního řádu K jackknifeové transformace $J_K(\cdot)$ na rozsahu N a zvolené přesnosti E zokrouhlené větší chyby.



Obr.2: Závislost počtu P znovu přepočítávaných hodnot koeficientu asociace na řádu K jackknifeové transformace $J_K(\cdot)$ a počtu C kategoriálních tříd.



LITERATURA:

- [1] Anděl J.(1973): Interaction in contingency tables. Aplikace matematiky 18, str.99-109
- [2] Běláček J.(1986): Jackknifeování koeficientů asociace. Sborník ROBUST 86, str.21-6
- [3] Běláček J.(1988): Eliminace vychýlení obecných měr asociace. Kand.práce MPP UK, Praha
- [4] Goodman L.(1964): Interactions in multidimensional contingency tables. Ann.Math.Statist.35, str.632-46
- [5] Hais K.-Hodek B.(1984): Velký anglicko-český slovník. Academia, Praha
- [6] Hrala P.(1987) - osobní sdělení
- [7] Krewski L.-Rao J.N.(1981): Inference from stratified samples - properties of the linearization, jackknife and balanced repeated replication methods. Ann.Statist.5, str.1010-19
- [8] Parr W.C.-Tolley H.D.(1982): Jackknifing in categorical data analysis. Austral.J.Statist.24, str.67-79
- [9] Řehák J.-Řeháková B.(1986): Analýza kategorizovaných dat v sociologii. Academia, Praha
- [10] Schucany W.R.-Gray H.L.-Owen L.B.(1971): On bias reduction in estimation. JASA 66, str.524-33
- [11] Serfling R.J.(1980): Approximation theorems of mathematical statistics. J.Wiley and sons, New York-Chichester-Brisbane-Toronto, str.210-42
- [12] Smith E.P.-van Belle G.(1984): Nonparametric estimation of species richness. Biometrics 40, str.119-29