

# ODHADY PARAMETRU TUKEYOVA A ZOBECNĚNÉHO TUKEYOVA SYSTÉMU ROZDĚLENÍ

Marie Hušková, MFF UK

## 0. ÚVOD

Tukey (1960) zavedl systém jednoparametrických rozdělání s kvantilovou funkcí

$$(1) \quad F^{-1}(u; \lambda) = (u^\lambda - (1-u)^\lambda) / \lambda \quad u \in (0,1), \quad \lambda \in R_1.$$

V literatuře se pro něj ujal název Tukeyův systém rozdělání.

Ramberg a Schmeiser (1974) zavedli systém rozdělání závisející na čtyřech parametrech  $(\lambda_1, \lambda_2, \lambda_3, \lambda_4)$  s kvantilovou funkcí

$$(2) \quad F^{-1}(u; \lambda_1, \lambda_2, \lambda_3, \lambda_4) = \lambda_1 + (u^{\lambda_2} - (1-u)^{\lambda_2}) \lambda_4 \quad u \in (0,1),$$

kde  $\lambda_1 \in R_1$  je parametr polohy,  $|\lambda_4|$  je parametr měřítka a  $\lambda_2, \lambda_3$  jsou parametry charakterizující chování chvostů,  $(\lambda_2, \lambda_3, \lambda_4) \in \mathcal{L} = \{\lambda_2 > 0, i = 2,3,4\} \cup \{\lambda_2 < 0, i = 2,3,4\} \cup \{\lambda_2 < -1, \lambda_3 > 0, \lambda_4 < 0\} \cup \{\lambda_2 > 0, \lambda_3 < -1, \lambda_4 < 0\} \cup \{\lambda_2 = 0, \lambda_3 \lambda_4 > 0\} \cup \{\lambda_3 = 0, \lambda_2 \lambda_4 > 0\}$ .

Systém rozdělání tvořený rozděláními s kvantilovou funkcí (2) s  $\lambda_1 \in R_1$  a  $(\lambda_2, \lambda_3, \lambda_4) \in \mathcal{L}$  se nazývají zobecněný Tukeyův systém rozdělání.

Tyto systémy zahrnují řadu běžně používaných rozdělání popř. jejich aproximace. V práci bude pojednáno o vlastnostech těchto systémů a odhadech jejich parametrů.

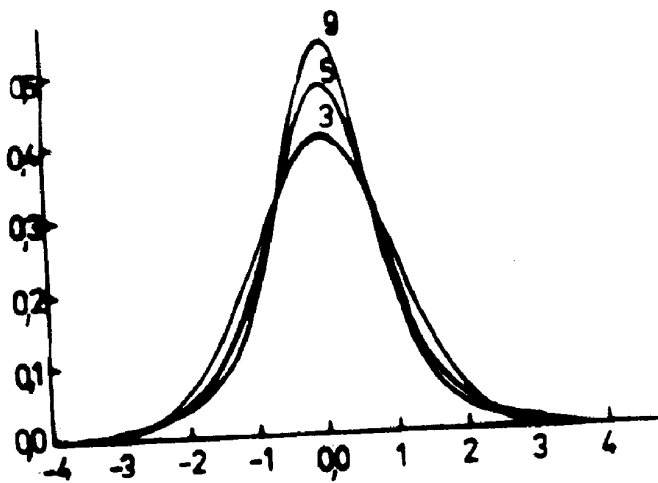
## 1. NĚKTERÉ VLASTNOSTI SYSTÉMU ROZDĚLENÍ

Tukeyův systém rozdělání zahrnuje symetrické rozdělání s tzv. těžkými chvosty ( $\lambda < 0$ ), i lehkými chvosty ( $\lambda > 0$ ). Při  $\lambda < 1$  nebo  $\lambda > 2$  jde o rozdělání unimodální jednovrcholové a při  $1 < \lambda < 2$  má rozdělání tzv. U - tvar, Cauchyovo rozdělání, pro  $\lambda = 0$  logistické a při  $\lambda = 1$  nebo 2 rozdělání rovnoměrné. Při  $\lambda = 0,08$  obdržíme rozumnou aproximaci pro dvojitě exponenciální rozdělání, při  $\lambda = 0,14$  pro normální. Tento systém poskytuje také dobrou aproximaci pro t-rozdělání.

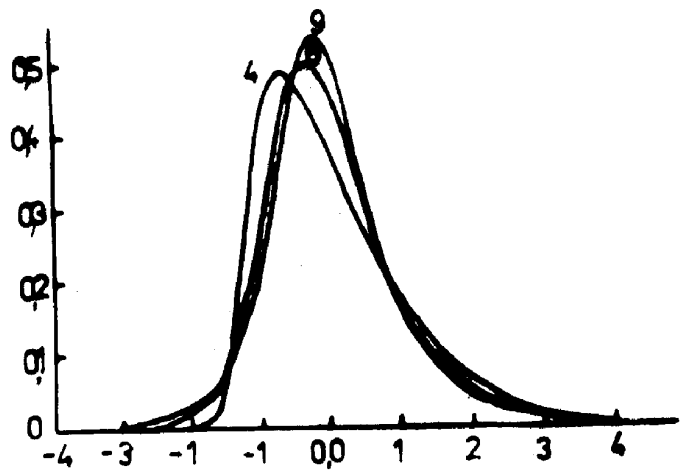
Zobecněný Tukeyův systém rozdělání zahrnuje některé asymetrické rozdělání. Poskytuje dobrou aproximaci např. pro exponenciální rozdělání.

Na obrázcích 1,2,3,4 jsou uvedeny grafy hustot několika rozdělání náležejících do těchto systémů. ( $\alpha_1 = E(X - EX)^4 (E(X - EX)^2)^{-1/2}$   $i = 3,4$ ).

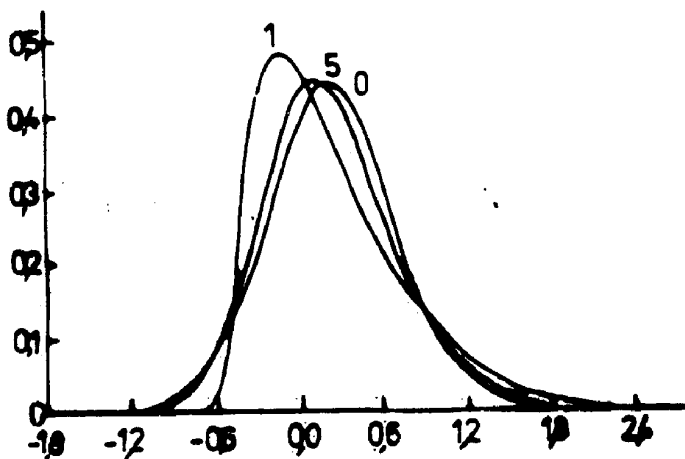
Uvědomíme-li si fakt, že má-li náhodná veličina U rovnoměrné rozdělání na  $(0,1)$  a F je nějaká distribuční funkce, pak  $F^{-1}(U)$  má rozdělání s distribuční funkcí F, lze uvedené systémy rozdělání výhodně využít při různých simulačních studiích. Navíc se oba systémy ukázaly velmi užitečné při reprezentaci dat, jestliže model není znám.



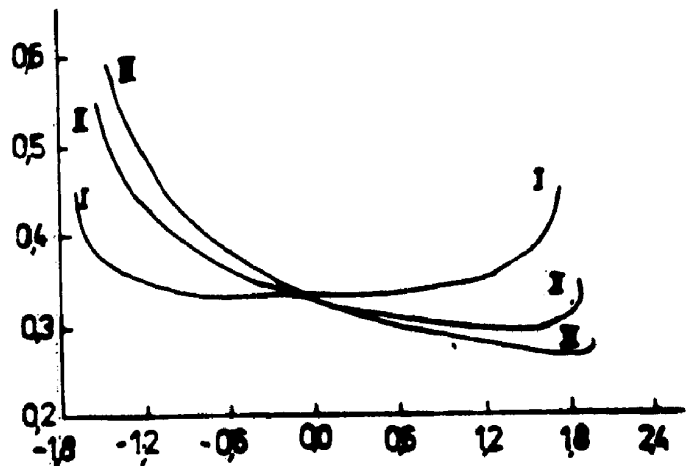
Obr. 8.1 Grafy hustot pro  $\alpha_3 = 0, \alpha_4 = 3, 5, 9$



Obr. 8.2 Grafy hustot pro  $\alpha_3 = 1, \alpha_4 = 4, 6, 9$



Obr. 8.3 Grafy hustot pro  $\alpha_3 = 0, 0.5, 1; \alpha_4 = 4$



Obr. 8.4 Grafy hustot s tzv. U - tvarem

Obecně neexistuje explicitní vyjádření distribuční funkce nebo hustoty. Momenty lze vyjádřit následovně (pro systém (2)):

$$(3) \quad EX = A_1 + A_4 \left( (A_2 + 1)^{-1} - (A_3 + 1)^{-1} \right) \quad \text{pro } A_2 > -1, A_3 > -1,$$

$$(4) \quad E(X - A_1)^k = A_4^k \sum_{j=0}^k \binom{k}{j} (-1)^j B(A_2(k-j) + 1, A_3(j+1)) \quad \text{pro } A_2^k > -1, A_3^k > -1,$$

kde  $B(a, b)$  je beta funkce.

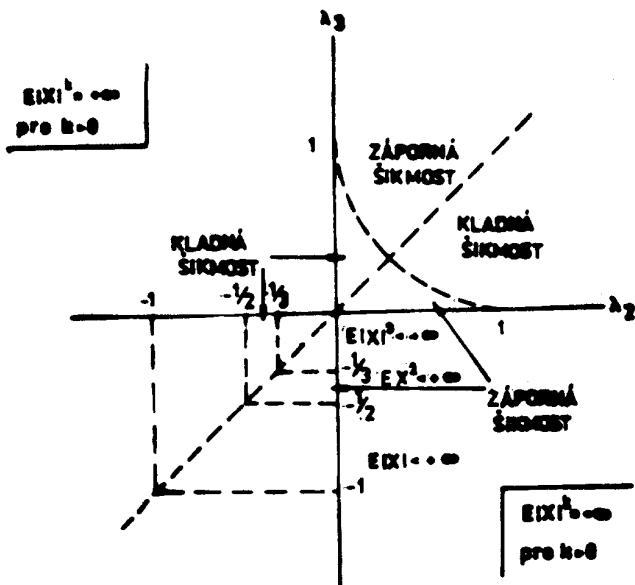
Pro momenty náhodné veličiny  $X$  s kvantilovou funkcí (1) platí

$$(6) \quad EX = 0$$

$$(7) \quad EX^k = A^{-k} \sum_{j=0}^k \binom{k}{j} (-1)^j B(A(k-j) + 1, A_{j+1}) \quad \text{pro } A^k > -1,$$

V slánku Kamberg a kol. (1979) jsou tabulovány hodnoty  $\lambda_1, \lambda_2, \lambda_3, \lambda_4$  pro  $EX = 0$ ,  $EX^2 = 1$ ,  $\alpha_3 = 0(0.05)2$ ,  $\alpha_4 = 1.8(0.2)9$ ,  $\alpha_3 = E(X-EX)^3$ ,  $\alpha_4 = E(X-EX)^4$ .

Tabulka 8.5 obsahuje přehled o konečnosti momentů a kladné či záporné šikmosti ( $E(X-EX)^3 > 0$  resp.  $< 0$ ).



Obr. 5 Některé vlastnosti rozdělení v závislosti na  $\lambda_2, \lambda_3$ .

Poznámeme, že Fisherova míra informace vzhledem k parametru posunutí je konečná pro  $\lambda_2, \lambda_3, \lambda_4 \in (1/2, 2)$ .

## 2. ODHADY PARAMETRŮ $\lambda_1, \lambda_2, \lambda_3, \lambda_4$

Pozornost soustředíme hlavně na systém (2), neboť (1) je speciální případem (2).

Nechť  $X_1, \dots, X_n$  je náhodný výběr z rozdělení s kvantilovou funkcí (2) a  $X_{(1)}, \dots, X_{(n)}$  je příslušný uspořádaný výběr. V dalším bude  $X_{(a)}$  značit  $[a]$ -tou pořádkovou statistiku, kde  $[a]$  je celá část  $a$ .

Za předpokladu konečnosti 4. momentu můžeme odhadnout parametry  $\lambda_1, \lambda_2, \lambda_3, \lambda_4$  metodou momentů, t.j. za odhady vezmeme řešení rovnic:

$$(8) \quad EX^k = \frac{1}{n} \sum_{i=1}^n X_i^k \quad k = 1, 2, 3, 4.$$

Explicitní řešení obecně neexistuje. Můžeme však využít tabulky Kamberga a kol. (79) zmíněné výše. Podobně se odhad parametru  $\lambda$  můžeme vrátit řešení rovnice (dle 7):

$$(9) \quad EX^3 = \frac{1}{n} \sum_{i=1}^n X_i^3,$$

při jejichž řešení můžeme využít tabulky Kamberga a kol. (79).

Můžeme však též použít velice jednoduché odhady založené na pořádkových statistikách. Je-li  $\lambda_2 < 1$ , pak můžeme pro  $\lambda_2$  použít odhad

$$(10) \quad \hat{\lambda}_2(N, a, b, s) = \frac{1}{\log s^{-1}} \log \frac{X_{(aM)} - X_{(bM)}}{X_{(aSM)} - X_{(bSM)}},$$

kde  $M, a, b, s > 0$ ,  $s \neq 1$ ,  $a \neq b$ ,  $\max(a, b, as, bs) < n/2$ . Podobně, je-li  $A_3 < 1$ , lze  $A_3$  odhadnout

$$(11) \quad \hat{A}_3(M, a, b, s) = \frac{1}{\log s^{-1}} \log \frac{X(n - aM) - X(n - bM)}{X(n - asM) - X(n - bsM)}$$

Oba tyto odhady jsou konsistentní, pro  $n \rightarrow \infty$ ,  $M \rightarrow \infty$ ,  $M/n \rightarrow 0$  ( $a, b, s$  jsou pevné) platí

$$(12) \quad \hat{A}_i(M, a, b, s) = A_i + O_p(\max(M^{-1/2}, (M/n)^{1-A_i})) \quad \text{pro } A_i < 1$$

$$= 1 + O_p(1) \quad \text{pro } A_i \geq 1,$$

$i = 2, 3$ . Pro praktické účely je vhodné např. volba  $s = 1/2$ ,  $a = 4$ ,  $b = 2$ , dostaneme pak

$$\hat{A}_2(M, 4, 2, 1/2) = \frac{1}{\log 2} \log \frac{X(4M) - X(2M)}{X(2M) - X(M)}$$

a  $M$  volíme podle rozsahu výběru  $0,02n - 0,05n$ .

Další typ odhadů parametrů ( $A_2, A_3$ ) získáme řešením následujících transcendentních rovnic:

$$(13) \quad \frac{X(a_i n) - X(b_i n)}{X(c_i n) - X(d_i n)} = \frac{a_i^{A_2} (1 - a_i)^{A_3} - b_i^{A_2} (1 - b_i)^{A_3}}{c_i^{A_2} (1 - c_i)^{A_3} - d_i^{A_2} (1 - d_i)^{A_3}} \quad i = 1, 2$$

kde  $1 > a_i, b_i, c_i, d_i > 0$ ,  $(a_i, b_i) \neq (c_i, d_i)$ ,  $(d_i, c_i)$ ,  $i = 1, 2$ ,

$(a_1, b_1, c_1, d_1) \neq (a_2, b_2, c_2, d_2)$ ,  $(c_2, d_2, a_2, b_2)$ ,  $c_i \neq b_i$ ,  $c_i \neq d_i$ .

Získané odhady (ozn.  $(\tilde{A}_2, \tilde{A}_3)$ ) jsou konsistentní; pro  $n \rightarrow \infty$  a  $a_i, b_i, c_i, d_i$ ,  $i = 1, 2$ , pevné platí

$$(14) \quad \tilde{A}_i = A_i + O_p(n^{-1/2}) \quad i = 2, 3.$$

Parametr měřítka  $A_4$  můžeme odhadnout buď

$$(15) \quad \hat{A}_4(M, N) = (X_{(n-M)} - X_{(N)})/2, \text{ je-li } A_2 > 0, A_3 > 0,$$

nebo

$$(16) \quad \tilde{A}_4(a, b) = \frac{X(an) - X(bn)}{a^{A_2} (1-a)^{A_3} - b^{A_2} (1-b)^{A_3}},$$

kde  $M, N < n/2$ ,  $a \neq b \in (0, 1)$ ,  $\tilde{A}_2, \tilde{A}_3$  jsou odhady získané řešením rovnic (13). Oba odhady jsou za jistých předpokladů konsistentní.

Pro  $n \rightarrow \infty$   $a, b$  pevné platí

$$(17) \quad \tilde{A}_4(a, b) = A_4 + O_p(n^{-1/2}).$$

Jestliže navíc  $M/n \rightarrow 0$ ,  $N/n \rightarrow 0$   $M, N$  mohou být pevné nebo konvergovat k nekonečnu,  $A_i > 0$   $i = 2, 3$ , pak lze odhadnout platnost

$$(18) \quad \hat{A}_4(M, N) = A_4 + O_p(\max(\frac{M}{n}, (\frac{M}{n})^{A_2}, \frac{N}{n}, (\frac{N}{n})^{A_3})).$$

Parametr  $A_1$  můžeme odhadnout buď

$$(19) \quad \hat{A}_1(M, N) = (X_{(n-M)} + X_{(N)})/2, \quad \text{je-li } A_2 > 0, A_3 > 0,$$

nebo

$$(20) \quad \tilde{A}_1(a, b) = (X_{(an)} + X_{(bn)})/2 - \tilde{A}_4(a \bar{A}_2 - (1-a) \bar{A}_3 + b \bar{A}_2 - (1-b) \bar{A}_3)/2$$

kde  $a \neq b \in (0, 1)$   $M, N < n/2$ . Odhady mají vlastnosti analogické  $\hat{A}_4(M, N)$  respective,  $\tilde{A}_4(a, b)$ .

Podívejme se nyní na odhad parametru  $A$  Tukeyova systému rozdělení. Nabízejí se dva typy (ve shodě s předchozím) jednak

$$(21) \quad \hat{A}(M, a, b, s) = (\hat{A}_2(M, a, b, s) + \hat{A}_3(M, a, b, s))/2, \quad \text{pro } a < 1$$

kde  $\hat{A}_i(M, a, b, s)$   $i=2, 3$  jsou definovány (10) respective, (11) jednak řešením rovnice

$$\frac{X_{(an)} - X_{(1-a)n}}{X_{(cn)} - X_{(1-c)n}} = \frac{a^A - (1-a)^A}{c^A - (1-a)^A},$$

kde  $a \neq c \in (0, 1)$ .

Odhady parametrů  $A_2, A_3$  popř.  $A$  můžeme velmi výhodně využít např. při konstrukci asymptoticky optimálních odhadů v následujících dvou typických úlohách neparametrických metod.

a) Nechť  $X_1, \dots, X_n$  jsou nezávislé náhodné veličiny,  $X_1$  má distribuční funkci  $F(x - \theta c_1)$ , kde  $(c_1, \dots, c_n)$  jsou známé regresní konstanty,  $\theta$  je neznámý parametr, který chceme odhadnout. Asymptoticky optimální R-odhad je generován skórovou funkcí

$$\psi_f(u) = - \frac{f'(F^{-1}(u))}{f(F^{-1}(u))}, \quad u \in (0, 1)$$

kde  $f$  a  $f'$  jsou hustota a její derivace příslušné  $F$ ,  $F^{-1}$  je kvantilová funkce. Jestliže distribuční funkce  $F$  náleží do zobecněného Tukeyova systému pak platí (pro  $A_2, A_3 \in (1/2, 2)$ ):

$$\psi_f(u) = - \frac{A_2(A_2-1)u^{A_2-2} - A_3(A_3-1)(1-u)^{A_3-2}}{A_4(A_2 u^{A_2-1} + A_3(1-u)^{A_3-1})^2}, \quad u \in (0, 1).$$

Tedy optimální skórovou funkcí známe již na parametry  $A_2, A_3, A_4$ , nahradíme-li je odhady navrženými výše, pak R-odhad generovaný takto vzniklou skórovou funkcí je asymptoticky optimální (ovšem za předpokladu, že  $F$  náleží do zobecněného Tukeyova systému).

b) Nechť  $(X_1, \dots, X_n)$  je náhodný výběr z rozdělení s distribuční  $F(x - \theta)$ , kde  $\theta$  je neznámý parametr, který chceme odhadnout a  $F$  je symetrická distribuční funkce. Asymptoticky optimální L-odhad je generován funkcí

$$J_f(u) = \psi_f'(u) f(F^{-1}(u)) \left( \int \psi_f(y) f(F^{-1}(y)) dy \right)^{-1}, \quad u \in (0, 1).$$

Jestliže  $F$  náleží do Tukeyova systému, známe  $J_p$  až na parametr  $\lambda$ . Nahradíme-li  $\lambda$  některým z odhadů navržených výše, pak  $L$ -odhad, generovaný takto vzniklou funkcí  $J_p$ , je asymptoticky optimální (za předpokladu, že  $F$  náleží do Tukeyova systému). Další podrobnosti může čtenář nalézt v citované literatuře.

## LITERATURA

- Chan, L. K. and Rhodin, L. S. (1980): Robust estimation of location using optimally chosen sample quantiles. *Technometrics* 22, 225-237.
- Filiben, J. J. (1969). Simple and robust estimation of the location parameter of a symmetric distribution. Ph. D. dissertation, Princeton University, Princeton, N. Y.
- Hušková, M. (1986). Simple estimators of the parameters of generalized Tukey's family, vyjde ČMUC.
- Hušková, M. (1986). Adaptive estimators of parameters generalized Tukey's family, zasláno k otištění.
- Joiner, E. L. and Rosenblatt, J. R. (1971). Some properties of the range in samples from Tukey's symmetric  $\lambda$ -distributions. *J. Amer. Statist. Assoc.*, 66, 394-399.
- Jones, D. H. (1979). An efficient adaptive distribution-free test for location, *J. Amer. Statist. Assoc.* 74, 822-828.
- Ranberg, J. S. and Schmeiser, B. W. (1972). An approximate method for generating symmetric random variables, *Comm. of the ACM*, 15, 987-990.
- Ranberg, J. S. and Schmeiser, B. W. (1974). An approximative method for generating asymmetric random variables. *Comm. of the ACM*, 17, 78-82.
- Ranberg, J. S., Tadikamalla, P. R., Dudewicz, E. J. and Mykytka, E. F. (1979). A probability distribution and its uses in fitting data. *Technometrics* 21, 201-214.
- Tukey, J. W. (1960). The practical relationship between the common transformations of percentages or counts and of amounts. Technical Report 36, Statistical Research Group, Princeton University, Princeton, N. Y.