

## ODHAD BODU VZNIKU KVADRATICKÉHO TRENDU

DANIELA JARUŠKOVÁ

ABSTRAKT. The problem of least squares method estimation of the parameter  $\tau^*$  in the regression model

$$Y_i = \beta^*(i/n - \tau^*)_+ + \gamma^*(i/n - \tau^*)_+^2 + e_i$$

is considered. Supposing  $\gamma^* \neq 0$  and  $\tau^* \in (0, 1)$  it is shown that the distribution of  $(\hat{\tau} - \tau^*)$  is largely affected by the value of  $\beta^*$ . In the case  $\beta^* \neq 0$  the variable  $\sqrt{n}(\hat{\tau} - \tau^*)$  is asymptotically normally distributed whereas in the case  $\beta^* = 0$  the variable  $\sqrt{n}(\hat{\tau} - \tau^*)^2$  has the same distribution as  $\max(0, Z)$  where  $Z$  has a zero mean normal distribution.

The problem of least squares method estimation of the parameter  $\tau^*$  in the regression model

$$Y_i = \beta^*(i/n - \tau^*)_+ + \gamma^*(i/n - \tau^*)_+^2 + e_i$$

is considered. Supposing  $\gamma^* \neq 0$  and  $\tau^* \in (0, 1)$  it is shown that the distribution of  $(\hat{\tau} - \tau^*)$  is largely affected by the value of  $\beta^*$ . In the case  $\beta^* \neq 0$  the variable  $\sqrt{n}(\hat{\tau} - \tau^*)$  is asymptotically normally distributed whereas in the case  $\beta^* = 0$  the variable  $\sqrt{n}(\hat{\tau} - \tau^*)^2$  has the same distribution as  $\max(0, Z)$  where  $Z$  has a zero mean normal distribution.

**Резюме:** Уважается проблема оценивания параметра  $\tau^*$  в модели линейной регрессии

$$Y_i = \beta^*(i/n - \tau^*)_+ + \gamma^*(i/n - \tau^*)_+^2 + e_i$$

при использовании метода наименьших квадратов. Показывается, что предполагая  $\gamma^* \neq 0$  и  $\tau^* \in (0, 1)$ , значение параметра  $\beta^*$  имеет большое влияние на распределение  $(\hat{\tau} - \tau^*)$ . В слуае когда  $\beta^* \neq 0$ , потом  $(\hat{\tau} - \tau^*)$  имеет асимптотически нормальное распределение. Наоборот, когда  $\beta^* = 0$ , потом  $(\hat{\tau} - \tau^*)$  обладает темже паспределением как  $(0, Z)$ , где  $Z$  следует нормальное распределение с нуловым средним.

## 1. ÚVOD

V praxi občas narazíme na problém, kde se lineární závislost v neznámém časovém okamžiku změni v závislost kvadratickou. Sledujeme-li například u jistých typů slitin závislost napětí  $Y$  na zatížení  $X$ , ukazuje se, že při malém zatížení je zvyšování napětí lineární (elastická oblast), zatímco po překročení určitého kritického zatížení se měni kvadraticky (quasielastická oblast), tj.  $Y = f(X)$ , kde

$$f(x) = p + qx + \beta^*(x - \tau^*)_+ + \gamma^*(x - \tau^*)_+^2,$$

při označení  $a_+ = \max(0, a)$ .

V našem článku si poněkud zjednodušíme situaci tím, že budeme předpokládat, že parametry  $p$  a  $q$  jsou známé, a tedy je bez újmy na obecnosti můžeme položit rovny nule. Dále budeme předpokládat, že veličina  $X$  je měfena v equidistantních vzdálenostech  $\Delta, 2\Delta, \dots, n\Delta$ , a tudíž ji lze transformovat na veličinu nabývající

2000 *Mathematics Subject Classification.* Primary 62F12.

*Klíčová slova.* Detekce bodu změny (change-point problem), nelineární regrese, odhady.

Práce byla částečně podporována granty GAČR 201/00/0769 a MSM 210000001.

hodnot  $1/n, 2/n, \dots, 1$ , to jest veličinu, která nabývá hodnot pouze v intervalu  $[0, 1]$ . Jestliže předpokládáme aditivní vliv náhodných chyb  $\{e_i\}$  na naměřené hodnoty nezávisle proměnné  $Y$ , dospíváme k regresnímu modelu

$$(1) \quad Y_i = \beta^* (i/n - \tau^*)_+ + \gamma^* (i/n - \tau^*)_+^2 + e_i,$$

kde  $\beta^*$ ,  $\gamma^*$  a  $\tau^*$  jsou neznámé parametry.

Pro jednoduchost předpokládejme, že náhodné chyby  $\{e_i\}$  jsou nezávislé stejně rozdělené s rozdělením  $N(0, \sigma^2)$ , kde  $\sigma^2$  je známé, a tudíž opět bez újmy na obecnosti můžeme položit  $\sigma^2 = 1$ .

Úlohou, kterou se budeme zabývat, je odhad neznámých parametrů  $\beta^*$ ,  $\gamma^*$  a  $\tau^*$ . Především nás bude zajímat odhad parametru  $\tau^*$ , který se nazývá bod změny. Budeme předpokládat, že  $\tau^* \in (0, 1)$ . Parametr  $\beta^*$  odpovídá první derivaci zprava a  $2\gamma^*$  druhé derivaci zprava regresní funkce v bodě  $\tau^*$ .

Vzhledem k normalitě chyb  $\{e_i\}$  lze maximálně věrohodné odhady získat metodou nejmenších čtverců.

## 2. BODOVÝ ODHAD PARAMETRŮ

Model (1) je speciálním modelem nelineární regrese, kterému se někdy říká semilineární model, viz Knowles et al. (1991). Kdybychom totiž znali hodnotu  $\tau$  parametru  $\tau^*$ , tj.  $\tau^* = \tau$ , pak by model (1) byl modelem lineární regrese s dvěma vysvětlujícími proměnnými, tj. lineární model s maticí plánu experimentu

$$\mathbf{X}_n(\tau) = \begin{pmatrix} 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ \frac{[n\tau]+1}{n} - \tau & \left(\frac{[n\tau]+1}{n} - \tau\right)^2 \\ \frac{[n\tau]+2}{n} - \tau & \left(\frac{[n\tau]+1}{n} - \tau\right)^2 \\ \vdots & \vdots \\ 1 - \tau & (1 - \tau)^2 \end{pmatrix}.$$

Označme

$$(\tilde{\beta}_\tau, \tilde{\gamma}_\tau)^T = (\mathbf{X}_n^T(\tau)\mathbf{X}_n(\tau))^{-1}\mathbf{X}_n^T(\tau)\mathbf{Y},$$

kde  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ , a

$$S_n(\tau, \beta, \gamma) = \sum_{i=1}^n \left( Y_i - \left( \beta(i/n - \tau)_+ + \gamma(i/n - \tau)_+^2 \right) \right)^2.$$

Pak

$$S_n(\hat{\tau}, \hat{\beta}, \hat{\gamma}) = \min_{\beta, \gamma, \tau \in (0,1)} S_n(\tau, \beta, \gamma) = \min_{\tau \in (0,1)} S_n(\tau, \tilde{\beta}_\tau, \tilde{\gamma}_\tau).$$

Nejčastějším numerickým postupem pro nalezení přibližných hodnot odhadů  $\hat{\tau}$ ,  $\hat{\beta}$ ,  $\hat{\gamma}$  spočívá v tom, že pro hodnoty  $\tau$  z dosti husté mříže bodů  $T$  v intervalu  $(0,1)$  počítáme residuální součet čtverců  $S_n(\tau, \tilde{\beta}_\tau, \tilde{\gamma}_\tau)$  a pak hledáme  $\min_{\tau \in T} S_n(\tau, \tilde{\beta}_\tau, \tilde{\gamma}_\tau)$ .

Statistik se však obvykle nespokojí s bodovým odhadem neznámých parametrů, ale zajímá se též o intervaly spolehlivosti, respektive oblasti spolehlivosti.

## 3. OBLASTI SPOLEHLIVOSTI V NELINEÁRNÍ REGRESI

Uvažujme obecný model nelineární regrese s  $k$  - rozměrným vektorem parametrů  $\boldsymbol{\theta}^* = (\theta_1^*, \dots, \theta_k^*)^T$

$$Y_i = f_i(\boldsymbol{\theta}^*) + e_i, \quad i = 1, \dots, n,$$

kde  $\{f_i(\cdot)\}$  jsou známé funkce a  $\{e_i\}$  jsou nezávislé stejně rozdělené náhodné veličiny s rozdělením  $N(0, 1)$ .

Označme  $S_n(\boldsymbol{\theta}) = \sum (Y_i - f_i(\boldsymbol{\theta}))^2$  a  $\hat{\boldsymbol{\theta}}$  odhad  $\boldsymbol{\theta}^*$  metodou nejmenších čtverců, tj.  $S_n(\hat{\boldsymbol{\theta}}) = \min_{\theta_1, \dots, \theta_k} \sum (Y_i - f_i(\theta_1, \theta_2, \dots, \theta_k))^2$ . Dále označme  $(\tilde{\theta}_2(\theta_1), \dots, \tilde{\theta}_k(\theta_1))$  odhad metodou nejmenších čtverců při pevné hodnotě prvního parametru rovnajícího se  $\theta_1$ , tj.  $S_n(\theta_1, \tilde{\theta}_2(\theta_1), \dots, \tilde{\theta}_k(\theta_1)) = \min_{\theta_2, \dots, \theta_k} \sum (Y_i - f_i(\theta_1, \theta_2, \dots, \theta_k))^2$ .

Užívané konfidenční oblasti pro celý vektor parametrů  $\boldsymbol{\theta}$  jsou obvykle dvou typů:

- 1 a)  $\{\boldsymbol{\theta}, (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T \mathbf{A} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) < C_1\}$ , kde  $\mathbf{A}$  je nějaká symetrická matice,
- 2 a)  $\{\boldsymbol{\theta}, S_n(\boldsymbol{\theta}) < S_n(\hat{\boldsymbol{\theta}}) + C_2\}$ .

První oblast je eliptická, zatímco druhá může mít naprosto obecný tvar a může být dokonce nesouvislá. Druhé oblasti se někdy říká exaktní, neboť je odvozena přímo z poměru věrohodnosti.

Analogické konfidenční oblasti pro jeden parametr, řekněme  $\theta_1$ , mají tvar:

- 1 b)  $\{\theta_1, |\theta_1 - \hat{\theta}_1| < K_1\}$ ,
- 2 b)  $\{\theta_1, S_n(\theta_1, \tilde{\theta}_2(\theta_1), \dots, \tilde{\theta}_k(\theta_1)) < S_n(\hat{\boldsymbol{\theta}}) + K_2\}$ .

Konstanty  $C_1, C_2$ , resp.  $K_1, K_2$ , jsou obvykle odvozeny z asymptotického rozdělení  $\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$  a  $S_n(\hat{\boldsymbol{\theta}}) - S_n(\boldsymbol{\theta}^*)$ , resp.  $\hat{\theta}_1 - \theta_1^*$  a  $S_n(\hat{\boldsymbol{\theta}}) - S_n(\theta_1^*, \tilde{\theta}_2(\theta_1^*), \dots, \tilde{\theta}_k(\theta_1^*))$ .

Je-li splněna celá řada podmínek, viz například podmínky A(1), ..., A(9) nebo B(1), ..., B(8) v knize Seber & Wild (1989), kapitola 12, pak  $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$  má asymptoticky normální rozdělení. Většinou mezi tyto podmínky patří existence spojitých prvních i druhých parciálních derivací  $\frac{\partial f_i(\boldsymbol{\theta})}{\partial \theta_r}$  a  $\frac{\partial^2 f_i(\boldsymbol{\theta})}{\partial \theta_r \partial \theta_s}$ ,  $i = 1, \dots, n$ ,  $r, s = 1, \dots, k$ , na nějakém okolí správné hodnoty  $\boldsymbol{\theta}^*$ . Navíc se předpokládá, že matice  $\frac{1}{n}(\mathbf{F}_n(\boldsymbol{\theta}))^T (\mathbf{F}_n(\boldsymbol{\theta}))$ , kde

$$\mathbf{F}_n(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial f_1}{\partial \theta_1} & \cdots & \frac{\partial f_1}{\partial \theta_k} \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial \theta_1} & \cdots & \frac{\partial f_n}{\partial \theta_k} \end{pmatrix},$$

konverguje stejnoměrně na nějakém okolí  $\boldsymbol{\theta}^*$  k nesingulární matici  $\mathbf{G}(\boldsymbol{\theta})$ . Matice  $\mathbf{G}^{-1}(\boldsymbol{\theta}^*)$  je pak limitní varianční maticí vektoru  $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$ .

## 4. OBLASTI SPOLEHLIVOSTI V MODELU (1)

Je patrné, že v modelu (1) nemá regresní funkce vzhledem k parametru  $\tau^*$  již ani první derivaci. Přesto se však dá ukázat, že v případě, že  $\beta^* \neq 0$ ,  $\gamma^* \neq 0$  a  $\tau^* \in (0, 1)$ , má  $\sqrt{n}(\hat{\tau} - \tau^*, \hat{\beta} - \beta^*, \hat{\gamma} - \gamma^*)$  asymptoticky normální rozdělení se symetrickou varianční maticí  $\mathbf{G}^{-1}$ , kde

$$\mathbf{G} = \begin{pmatrix} \beta^{*2}(1-\tau^*) + 4\beta^*\gamma^* \frac{(1-\tau^*)^2}{2} + 4\gamma^{*2} \frac{(1-\tau^*)^3}{3} & \dots & \dots \\ -\beta^* \frac{(1-\tau^*)^2}{2} - 2\gamma^* \frac{(1-\tau^*)^3}{3} & \frac{(1-\tau^*)^3}{3} & \dots \\ -\beta^* \frac{(1-\tau^*)^3}{3} - 2\gamma^* \frac{(1-\tau^*)^4}{4} & \frac{(1-\tau^*)^4}{4} & \frac{(1-\tau^*)^5}{5} \end{pmatrix}$$

Zřejmě  $\mathbf{G} = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{F}_n^{*T} \mathbf{F}_n^*$ , kde

$$\mathbf{F}_n^* = \begin{pmatrix} 0 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \\ -\beta^* - 2\gamma^* \left( \frac{[n\tau^*]+1}{n} - \tau^* \right) & \frac{[n\tau^*]+1}{n} - \tau^* & \left( \frac{[n\tau^*]+1}{n} - \tau^* \right)^2 \\ \vdots & \vdots & \vdots \\ -\beta^* - 2\gamma^*(1-\tau^*) & 1-\tau^* & (1-\tau^*)^2 \end{pmatrix}$$

K důkazu lze použít stejný postup, který použila Hušková (1998) pro jednodušší model typu

$$(2) \quad Y_i = \beta^* (i/n - \tau^*)_+ + e_i,$$

a který spočívá v tom, že na okolí správné hodnoty aproximujeme funkci nejmenších čtverců kvadratickou funkcí. Speciálně ukázala, že na okolí bodu  $\tau^*$  lze aproximovat

$$S_n(\tau, \tilde{\beta}_\tau) - S_n(\tau^*, \tilde{\beta}_{\tau^*})$$

funkcí

$$-C n (\tau - \tau^*)^2 + 2 X \sqrt{n} (\tau - \tau^*),$$

kde  $X/C$  má normální rozdělení  $N(0, 1/C)$ .

Jiný postup použil Feder (1975), který modely (1) i (2) považuje za speciální případy regresní funkce

$$\begin{aligned} f(x, \boldsymbol{\theta}) &= f_1(x, \boldsymbol{\theta}_1) \quad \text{pro } 0 \leq x \leq \tau^* \\ &= f_2(x, \boldsymbol{\theta}_2) \quad \text{pro } \tau^* \leq x \leq 1, \end{aligned}$$

kde  $f_i(x, \boldsymbol{\theta}) = \sum_{j=1}^{K(i)} \theta_{ij} f_{ij}(x)$ . Za přípustnou množinu parametrů  $\Theta$  pak uvažuje jen takové parametry, při kterých se funkce  $f_1(x, \boldsymbol{\theta}_1)$  a  $f_2(x, \boldsymbol{\theta}_2)$  protínají uvnitř intervalu  $(0,1)$ . Ukazuje, že limitní rozdělení  $\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$  závisí na tom, zda  $\boldsymbol{\theta}^*$  je vnitřním nebo krajním bodem množiny  $\Theta$ . Parametr  $\tau^*$  pak odhaduje jako průsečík funkcí  $f_1(x, \hat{\boldsymbol{\theta}}_1)$  a  $f_2(x, \hat{\boldsymbol{\theta}}_2)$ . Řád konvergence  $\hat{\tau}$  k  $\tau^*$  závisí na tom, v kolika derivacích se shodují funkce  $f_1(x, \boldsymbol{\theta}^*)$  a  $f_2(x, \boldsymbol{\theta}^*)$  v bodě  $\tau^*$ .

V případě modelu (1) s  $\beta^* \neq 0$ ,  $\gamma^* \neq 0$  a  $\tau^* \in (0, 1)$  lze z přístupů Huškové (1998) i Federa (1975) odvodit asymptotickou normalitu  $(\hat{\tau} - \tau^*, \hat{\beta} - \beta^*, \hat{\gamma} - \gamma^*)$ . Speciálně platí

$$(3) \quad \sqrt{n} (\hat{\tau} - \tau^*) \sim N\left(0, \frac{9}{\beta^{*2}(1-\tau^*)}\right)$$

a  $S_n(\tau^*, \tilde{\beta}_{\tau^*}, \tilde{\gamma}_{\tau^*}) - S_n(\hat{\tau}, \hat{\beta}, \hat{\gamma})$  má asymptoticky  $\chi^2$  rozdělení o 1 stupni volnosti.

Všimněme si, že asymptotický rozptyl odhadu  $\hat{\tau}$ , viz (3), závisí na poloze  $\tau^*$ , což je jakási globální vlastnost regresní funkce. Je přirozené, že parametr  $\tau^*$  lépe

odhadneme, máme-li možnost pozorovat kvadratický trend delší dobu. Dále však je asymptotický rozptyl  $\hat{\tau}$ , a tedy i délka intervalu spolehlivosti typu 2 a), ovlivněna první derivací kvadratické funkce v bodě  $\tau^*$ , což je lokální vlastnost regresní funkce.

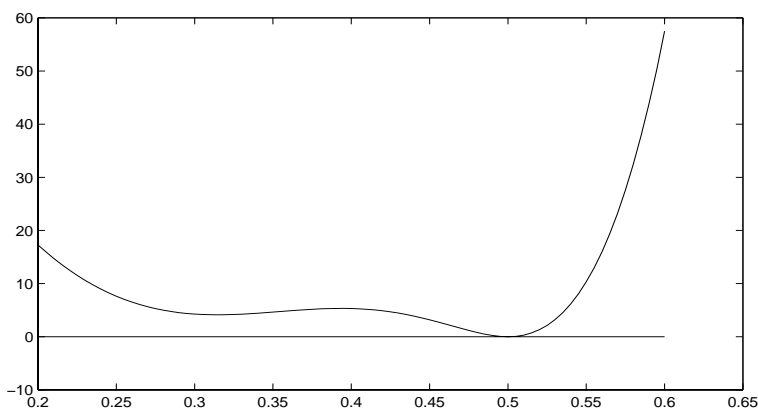
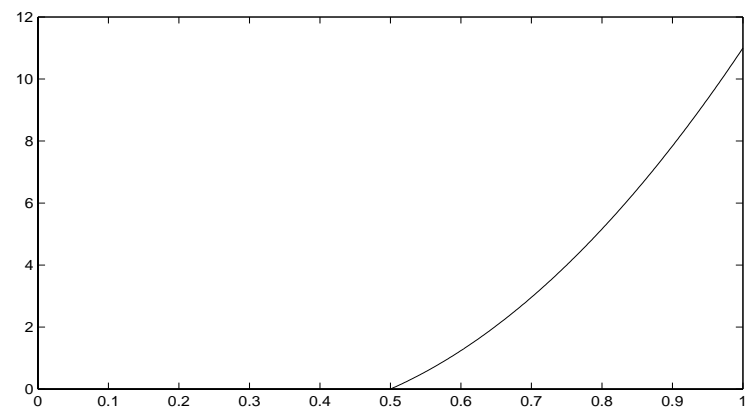
Je-li  $\beta^*$  malé, můžeme (použijeme-li postup 1 b)) dostat pro konečná  $n$  dokonce interval spolehlivosti, jehož krajní body leží mimo interval  $(0, 1)$ . V teorii nelineární regrese se tento jev nazývá špatná podmíněnost. Špatná podmíněnost je způsobena tvarem regresní funkce.

Horní část obrázku 1 představuje regresní funkci

$$(4) \quad f(x, \tau^*, \beta^*, \gamma^*) = \beta^*(x - \tau^*)_+ + \gamma^*(x - \tau^*)_+^2$$

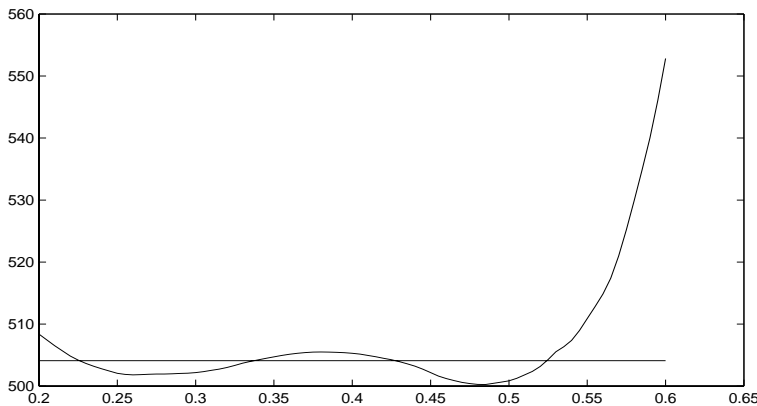
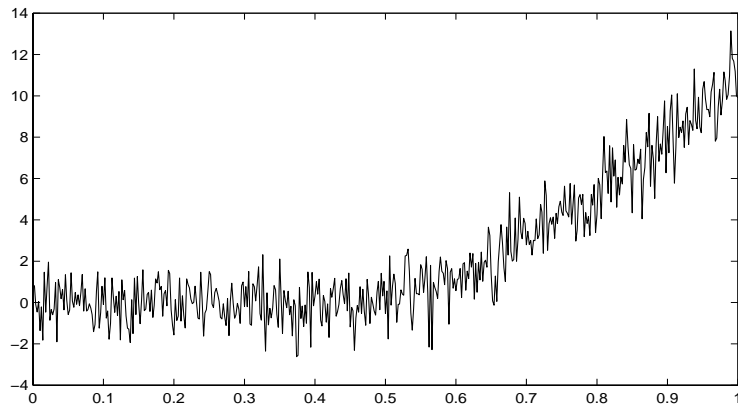
pro  $\tau^* = 0.5$ ,  $\beta^* = 10$  a  $\gamma^* = 24$  a dolní část obrázku zobrazuje odpovídající funkci nejmenších čtverců  $B_n(\tau, \tau^*)$  pro  $n = 500$  a  $\tau \in (0.2, 0.6)$  za předpokladu, že model neobsahuje žádné chyby, tj.

$$B_n(\tau, \tau^*) = \frac{1}{n} \sum_{i=1}^n \left( f\left(\frac{i}{n}, \tau, \tilde{\beta}_\tau, \tilde{\gamma}_\tau\right) - f\left(\frac{i}{n}, \tau^*, \tilde{\beta}_{\tau^*}, \tilde{\gamma}_{\tau^*}\right) \right)^2.$$



Obrázek 1

Jestliže uvažujeme model (1) s regresní funkcí (4) včetně náhodných chyb, pak může funkce nejmenších čtverců a 95% oblast spolehlivosti typu 2 b) vypadat například jako na obrázku 2.



Obrázek 2

Ačkoliv k tomu nemáme žádný teoretický důvod doporučujeme ze zkušeností používat spíše oblast spolehlivosti typu 2 b), kde  $K_2$  je příslušný kvantil  $\chi^2$  - rozdělení o 1 stupni volnosti. V každém případě doporučujeme vždy při statistické analýze vykreslit průběh funkce  $S_n(\tau, \tilde{\beta}_\tau, \tilde{\gamma}_\tau)$ .

V krajním případě, kde  $\beta^* = 0$ ,  $\gamma^* \neq 0$  a  $\tau^* \in (0, 1)$ , lze ukázat, že v okolí bodu  $\tau^*$  lze rozdíl funkcí čtverců

$$S_n(\tau, \tilde{\beta}_\tau, \tilde{\gamma}_\tau) - S_n(\tau^*, \tilde{\beta}_{\tau^*}, \tilde{\gamma}_{\tau^*})$$

aproximovat funkcí

$$-C n (\tau - \tau^*)^4 + 2 X \sqrt{n} (\tau - \tau^*)^2,$$

kde  $X/C \sim N(0, 1/C)$  a  $C = \gamma^{*2} (1 - \tau^*)/9$ . Je zřejmé, že funkce

$$-C x^2 (x^2 - 2 X/C)$$

nabývá maxima v nule, pokud  $X$  je záporné, a hodnoty  $X/C$ , pokud  $X$  je kladné. Odtud vyplývá, že  $\sqrt{n} (\hat{\tau} - \tau^*)^2$  má asymptoticky stejné rozdělení jako  $\max(0, Z)$ , kde veličina  $Z$  má normální rozdělení

$$(5) \quad Z \sim N\left(0, \frac{9}{\gamma^{*2}(1 - \tau^*)}\right).$$

Výsledek (5) lze použít pro konstrukci symetrického intervalu spolehlivosti 2 a). Můžeme též zkonstruovat oblast typu 2 b) s využitím toho, že  $S_n(\tau^*, \tilde{\beta}_{\tau^*}, \tilde{\gamma}_{\tau^*}) - S_n(\hat{\tau}, \hat{\beta}, \hat{\gamma})$  má asymptoticky stejné rozdělení jako  $(\max(0, Y))^2$ , kde  $Y$  má standardní normální rozdělení.

Pro zajímavost uveďme, že pokud odhadujeme parametr  $\tau^*$  v modelu

$$Y_i = \gamma^*(i/n - \tau^*)_+^2 + e_i, \quad i = 1, \dots, n,$$

pak  $\sqrt{n}(\hat{\tau} - \tau^*)$  má asymptoticky normální rozdělení  $N\left(0, 12/(\gamma^{*2}(1 - \tau^*)^3)\right)$ . Odtud je zřejmé, že informace, zda-li je první derivace v bodě změny nulová, je při odhadování tohoto bodu velmi velmi důležitá.

#### LITERATURA

- Feder P. J. (1975). *On asymptotic distribution theory in segmented regression problems - identified case*. The Annals of Statistics **3**, 49–63.
- Hušková M. (1998). *Estimators in the location model with gradual changes*. Comment. Math. Univ. Carolinae **39**, 147–157.
- Knowles M., Siegmund D. a Zhang H. (1991). *Confidence regions in semilinear regression*. Biometrika **78**, 1991, 15–31.
- Seber G. A. F. a Wild C. J. (1989). *Nonlinear regression*. J. Wiley, New York.

ČVUT V PRAZE, STAVEBNÍ FAKULTA, KATEDRA MATEMATIKY, THÁKUROVA 7, 166 29 PRAHA 6  
E-MAIL: jarus@mat.fsv.cvut.cz