

Vlivné body v lineární regresi
 Jiří Militký, VÚTZ, Dvůr Králové n.L.

1. Úvod

Je známo, že klasická metoda nejmenších čtverců je silně nerobustní. Je tedy velmi důležité identifikovat vlivné body a omezit jejich působení na odhady parametrů resp. jejich statistické charakteristiky.

V této práci jsou na jednoduchém simulačním experimentu demonstrovány nerobustní vlastnosti metody nejmenších čtverců a porovnány různé postupy pro identifikaci resp. eliminaci vlivných měření (při odhadu regresních parametrů).

2. Základní pojmy

Vyjděme ze standardní situace, kdy máme k dispozici n -tici bodů, $\{y_i, \underline{x}_i\}_{i=1, \dots, n}$ kde vektory deterministických vysvětlujících proměnných mají složky $\underline{x}_i^T = (x_{i1}, \dots, x_{im})$. Nechť dále platí aditivní model měření

$$y_i = \underline{x}_i^T \underline{a} + \varepsilon_i \quad (1)$$

kde $\underline{a}^T = (a_1, \dots, a_m)$ jsou regresní parametry. Pak za předpokladu, že náhodné chyby $\underline{\varepsilon}^T = (\varepsilon_1, \dots, \varepsilon_n)$ mají

- a) nulovou střední hodnotu $E(\underline{\varepsilon}) = \underline{0}$
- b) diagonální kovarianční matici $E(\underline{\varepsilon} \cdot \underline{\varepsilon}^T) = \sigma^2 E_n$ ($\sigma^2 < \infty$)
- c) stejné rozdělení, takže $f(\underline{\varepsilon}) = \prod_{i=1}^n f(\varepsilon_i)$

a navíc, že

- d) hodnost matice X je právě n
- e) na parametry \underline{a} nejsou kladena omezení
- f) platí model (2)

kde nejlepší nestranné lineární odhady $\hat{\underline{a}}$ parametrů \underline{a} získat minimalizací kritéria nejmenších čtverců odchylek

$$S(\underline{a}) = \| \underline{y} - X \underline{a} \|_2^2 \quad (2)$$

kde $X^T = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n)$ je $(m \times n)$ matice. Z rov. (2) je zřejmé, že vektor

$$\hat{\underline{y}} = X \hat{\underline{a}} \quad (3)$$

kde $\hat{\underline{a}}$ minimalizuje rov. (2) je kolmá projekce vektoru \underline{y} do roviny vymezené sloupci matice X . Z toho také plyne, že vektor reziduí

$$\hat{\underline{e}} = \underline{y} - \hat{\underline{y}} \quad (4)$$

musí být kolmý na sloupce matice X a tedy

$$X^T \hat{\underline{e}} = \underline{0} \quad \text{resp.} \quad X^T (\underline{y} - X \hat{\underline{a}}) = \underline{0} \quad (5)$$

Z rov. (5) již přímo plyne, že $\hat{\underline{a}}$ minimalizující kritérium $S(\underline{a})$ lze vyjádřit ve tvaru

$$\hat{\underline{a}} = (X^T X)^{-1} X^T \underline{y} \quad (5)$$

Dosadíme-li z rov. (5) do rov. (3) dostaneme

$$\hat{\underline{y}} = X (X^T X)^{-1} X^T \underline{y} = H \underline{y} \quad (6)$$

kde H ($n \times n$) je zřejmě projekční matice. Pokud je $\underline{\varepsilon} \in N(\underline{0}, \sigma^2 E_n)$, jsou odhady $\hat{\underline{a}}$ z rov. (5) maximálně věrohodné.

Snadno lze dokázat, že kovarianční matice odhadů $\hat{\underline{a}}$ má tvar

$$C(\hat{\underline{a}}) = \sigma^2 (X^T X)^{-1} \quad (7)$$

kde reziduální rozptyl σ^2 lze nahradit odhadem $\hat{\sigma}^2 = \| \hat{\underline{e}} \|^2 / (n - m)$.
 Z rov. (6) plyne, že kovarianční matice predikce je rovna

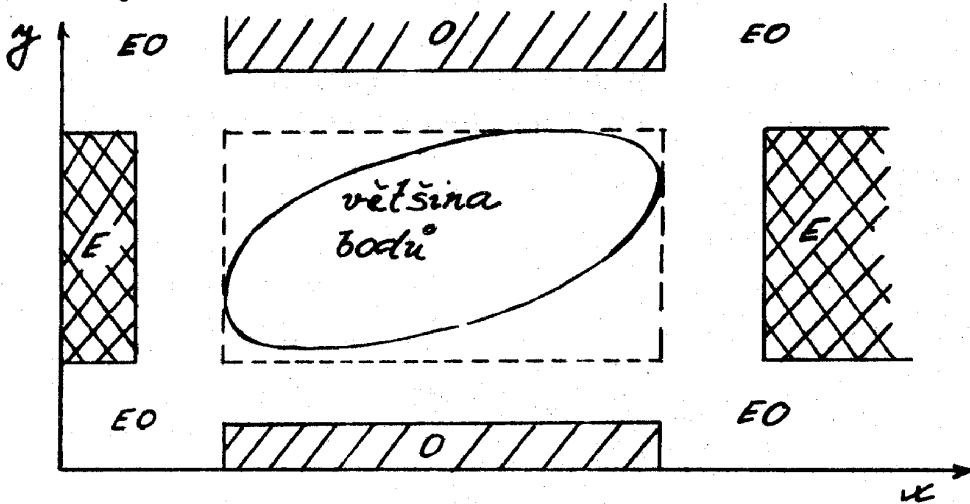
$$C(\hat{\underline{y}}) = \sigma^2 H \quad (8)$$

a z rov. (4) při dosazení z rov. (6) lze určit kovarianční matici reziduí

$$C(\hat{\underline{e}}) = \sigma^2 (I_n - H) \quad (9)$$

Platí, že odhady $\hat{\underline{a}}, \hat{\underline{y}}$ a jejich statistické charakteristiky jsou citlivé na
 - vybočující body (outliers) - označení O
 - extrémní body (remote points) - označení E

Na obr. 1 jsou tyto pojmy graficky znázorněny



Obr. 1. Rozdíl mezi oblastmi kde leží vybočující pozorování // // // a extrémní body XXXX. EO označující vlivné body, které jsou jak vybočující, tak i extrémní.

3. Simulační experiment

Pro účely této práce je použito simulovaných dat $(y_i, x_i)_{i=1, \dots, 50}$. Vysvětlující proměnné x_i byla generována z rovnoměrného rozdělení $R(0, 1)$ a transformována do intervalu $[10, 20]$. Odpovídající hodnoty y_i byly určeny ze vztahu

$$y_i = 1 \cdot x_i + 1 + N(0, 1)$$

kde $N(0, 1)$ je realizace náhodné veličiny s normalizovaným normálním rozdělením. Do takto získané n-tice bodů se zaváděly vybočující (O) resp. extrémní (E) hodnoty o různých velikostech A pro bod s indexem $j = 37$ resp. body s indexy $j = 12, 25, 37$.

Pro ilustraci se dále v tabulkách většinou uvádějí odhady směrnice \hat{a}_1 regrese přímkou určené různými postupy z takto generovaných dat.

V tab. 1 jsou uvedeny odhady směrnice \hat{a}_1 určené metodou nejmenších čtverců pro případ, že se měnily souřadnice pouze bodu $j = 37$. (Uveďme, že původní hodnoty tohoto bodu jsou $x_{37} = 19,2$ a $y_{37} = 20,8$.)

Tab. 1 : Odhady \hat{a}_1 získané MNC pro data s jedním vlivným měřením (O resp. E o velikostech A)

A	E	O
20	0,968	0,967
30	0,71	1,07
40	0,46	1,18
50	0,31	1,28
60	0,22	1,39
70	0,17	1,49
80	0,13	1,602

Z tab. 1 je patrné, že přítomnost extrémů E vede k výrazně horším výsledkům než přítomnost vybočujících pozorování. (Klasická analýza reziduí obyčejně E nezjisti).

4. Identifikace vlivných bodů

Jako vlivné se obecně označují takové body, které výrazně ovlivňují výsledky regrese.

Poznámky:

Ve většině postupů identifikace vlivných měření se využívá diagonálních prvků h_{ii} projekční matice H . Tyto prvky mají řadu zajímavých vlastností plynoucích z faktu, že H je idempotentní symetrická matice. Platí, že /1/

- $0 \leq h_{ii} \leq 1$
- $\sum_{i=1}^n h_{ii} = m$ (pokud je hodnota matice rovna m)
- průměrná hodnota $h_{ii} \approx m/n$
- lze psát, že $h_{ii} = h_{ii} + \sum_{j \neq i} h_{ij} h_{ji}$. Z toho plyne, že pro $h_{ii} = 1$ resp. $h_{ii} = 0$ jsou vždy všechna h_{ij} také nulová.
- pokud mají vysvětlující proměnné vícerozměrné normální rozdělení, má veličina Fisherova rozdělení s $(m-1)$ a $(n-m)$ stupni volnosti.
 $F = (n-m)[h_{ii} - (1/n)] / [(1-h_{ii})(m-1)]$
- Čím je h_{ii} větší, tím více ovlivňuje bod $[y_i, x_i^T]$ predikci \hat{y}_i . To plyne přímo z rov. (6), kterou lze pro \hat{y}_i vyjádřit ve tvaru

$$\hat{y}_i = y_i h_{ii} + \sum_{j \neq i} h_{ij} y_j \quad (10)$$

Je zřejmé, že pro $h_{ii} \rightarrow 1$ jsou $h_{ij} \rightarrow 0$ a $\hat{y}_i \approx y_i$ (tedy veškerá variabilita v místě x_i je objasněna regresním modelem).

- V práci /2/ je dokázáno, že pro centrované proměnné x_i^T je h_{ii} velké, pokud je $x_i^T x_i$ velké (to je případ extrémních bodů) resp. pokud je x_i^T ve směru vlastního vektoru odpovídajícího malému vlastnímu číslu kovarianční matice vysvětlujících proměnných. Malé h_{ii} znamená vždy, že $x_i^T x_i$ je malé (x_i^T leží blízko "těžiště" vysvětlujících proměnných).
- Z rov. (8) a rov. (9) přímo plyne, že $\text{var}(\hat{y}_i) = h_{ii} \sigma^2$ a $\text{var}(\hat{\epsilon}_i) = \sigma^2(1-h_{ii})$. Pro extrémní body (s velkým h_{ii}) bude tedy rozptyl reziduí malý, takže bude obtížné na základě analýzy reziduí identifikovat vybočující měření.
- Pokud $h_{ii} \rightarrow 1$, nejsou klasické metody robustní regrese efektivní. Z těchto vlastností je patrné, že h_{ii} mohou přímo identifikovat extrémní body.

Speciálně pro případy, kdy je účelem identifikovat vybočující body, využívá se různých typů reziduí.

Klasická nenormalizovaná rezidua \hat{a}_i jsou definována vztahy

$$\hat{a}_i = y_i - x_i^T \hat{\beta} \quad i = 1, \dots, n \quad (11)$$

Z rov. (4) a (6) přímo plyne, že

$$\hat{\underline{a}} = (E_n - H)\underline{y} = (E_n - H)\underline{\epsilon} \quad (12)$$

Platí tedy, že každé reziduum je lineární kombinací náhodných chyb $\epsilon_1, \dots, \epsilon_n$. (Proto se u malých výběrů projevuje efekt supernormality, t.j. rozdělení $\hat{\underline{a}}$ je přibližně normální, i když rozdělení $\underline{\epsilon}$ normální není).

Z rov. (9) plyne, že rezidua $\hat{\underline{a}}$ jsou vzájemně korelovaná s robustním rozptylem. Také tzv. normalizovaná rezidua

$$\hat{\underline{a}}_N = (n-m)\hat{\underline{a}} / \|\hat{\underline{a}}\|_2 \quad (13)$$

sčítavěji navzájem korelovaná a mají nekonstantní rozptyl $\text{var}(\hat{\underline{a}}_N) = (E_n - H)$. Zkonstantnění rozptylu lze docílit podělením $\hat{\underline{a}}$ odpovídajícími rozptyly z rov. (8).

Wyjścisław tzw. Studentizované rezidua

$$\hat{t}_i = \frac{\hat{a}_i}{\hat{\sigma} \sqrt{1-h_{ii}}} \quad i = 1, \dots, n \quad (14)$$

Veličiny \hat{t}_i již mají nulovou střední hodnotu a jednotkový rozptyl. Navíc platí, že $\hat{t}_i^2 / (n-m)$ má beta rozdělení.

Nanotřní transformací reziduí \hat{t}_i jsou tzv. Jackknife rezidua

$$\hat{\epsilon}_i = \hat{r}_i \sqrt{\frac{n-m-1}{n-m-\hat{r}_i^2}} \quad i = 1, \dots, n \quad (15)$$

Platí, že $\hat{\epsilon}_i$ jsou testovací statistiky pro test hypotézy $H_0: d=0$ v modelu jednoduchého modelu vybočujícího pozorování vychýleného co do polohy

$$y = xa + \underline{x}_i d + \underline{\epsilon} \quad (16)$$

kde vektor \underline{x}_i má 1-tou souřadnici rovnou jedné a ostatní nulové. Pokud platí H_0 , má $\hat{\epsilon}_i$ Studentovo rozdělení s $n - m - 1$ stupni volnosti.

Poznámka :

Pokud se zavede reziduální rozptyl $\hat{\sigma}_{(-i)}^2$ počítaný vždy bez 1-tého bodu, pro který platí

$$\hat{\sigma}_{(-i)}^2 = \frac{(n-m)\hat{\sigma}^2 - a_i^2/(1-h_{ii})}{n-m-1} = \hat{\sigma}^2 \left(\frac{n-m-\hat{r}_i^2}{n-m-1} \right) \quad (17)$$

můžeme vyjádřit Jackknife rezidua analogicky jako Studentizovaná rezidua. Platí

$$\hat{\epsilon}_i = \frac{\hat{a}_i}{\hat{\sigma}_{(-i)} \sqrt{1-h_{ii}}} \quad (17a)$$

V mnoha aplikacích je výhodné použít predikovaná rezidua

$$\hat{a}_{(-i)} = y_i - x_i^T \hat{a}_{(-i)} \quad i = 1, \dots, n \quad (18)$$

kde $\hat{a}_{(-i)}$ jsou odhady získané metodou nejmenších čtverců při vynechání 1-tého bodu (y_i, x_i^T). Lze snadno ukázat, že $\hat{a}_{(-i)}$ jsou odhady \hat{d} parametru d v modelu (16) získané metodou nejmenších čtverců. Z výpočetního hlediska je výhodné použít pro jejich určení vztahu

$$\hat{a}_{(-i)} = \hat{a}_i / (1-h_{ii}) = \hat{d} \quad (19)$$

Platí, že pokud ϵ_i mají normální rozdělení $N(0, \sigma^2)$, mají $\hat{a}_{(-i)}$ také normální rozdělení $N(0, \sigma^2/(1-h_{ii}))$ a jsou stejně korelovaná jako ϵ_i .

V tabulce 2 jsou pro data ze simulačního experimentu (kap.3) s jedním vlivným měřením (0 resp. E) o velikosti A v bodě $i = 37$ určeny $h_{ii}, \hat{a}_{Ni}, \hat{r}_i, \hat{\epsilon}_i, \hat{a}_{(-i)}$.

Tabulka 2 : Různé typy reziduí a h_{ii} pro data s jedním vlivným měřením (0,E) o velikostech A.

	0					E				
A	\hat{a}_{Ni}	\hat{r}_i	$\hat{\epsilon}_i$	$\hat{a}_{(-i)}$	h_{ii}	\hat{a}_{Ni}	\hat{r}_i	$\hat{\epsilon}_i$	$\hat{a}_{(-i)}$	h_{ii}
20	-0,41	-0,43	-0,42	-0,41	0,06	-0,384	-0,4	-0,396	-0,39	0,078
40	6,36	6,56	20,1	19,6	0,06	-3,64	-6,08	-12,6	-19,8	0,641
60	6,62	6,83	40,6	39,6	0,06	-2,44	-6,36	-15,9	-39,3	0,853
80	6,67	6,88	61,1	59,6	0,06	-1,77	-6,43	-17,1	-58,7	0,924

Z tab. 2 plyne, že :

- \hat{a}_{Ni} vyhovují (alespoň částečně) pro případ vybočujících pozorování. Pro případ extrémních pozorování nejsou vhodné pro identifikaci vlivných bodů (s růstem A jejich hodnota klesá)
- ostatní typy reziduí $\hat{r}_i, \hat{\epsilon}_i, \hat{a}_{(-i)}$ jsou vhodné pro identifikaci obou typů vlivných bodů. Vyjdeme-li ze skutečnosti, že původní hodnota $Y_{37} = 20,8$ jsou $\hat{a}_{(-i)}$ poměrně přesné odhady ($A - Y_{37}$).

- veličina h_{ii} dobře identifikuje extrémní body. (Za extrémní se orientačně považuje bod, pro který je $h_{ii} > \frac{2m}{n}$, což je pro tento případ 0,08).

Poznámka :

Všechny výše uvedené druhy reziduí jsou stejně korelované. Nekorelovaných reziduí je pouze $(n - m)$. Patří mezi ně např. BLUS rezidua resp. rekursivní rezidua /2.3/.

Pro identifikaci vlivných bodů je možno použít různých modifikací vlivové křivky, která je obecně definována např. v práci /4/.

Empirická vlivová křivka (vektor) EIC_i je definována vztahem /2/

$$EIC_i = n(X^T X)^{-1} x_i \hat{\epsilon}_i \quad (20)$$

Pro Jackknife vlivovou křivku (počítanou bez i-tého bodu) platí /2/

$$EIF_i = (n-1)(X^T X)^{-1} x_i \hat{\epsilon}_i / (1-h_{ii})^2 \quad (20a)$$

Pro konečné výběry je možno definovat také výběrovou vlivovou křivku /5/

$$SIC_i = (n-1)(X^T X)^{-1} x_i \hat{\epsilon}_i / (1-h_{ii}) \quad (20b)$$

Je zřejmé, že rozdíl mezi těmito vlivovými křivkami je pouze v mocnině členu $(1-h_{ii})$ ve jmenovateli (která určuje citlivost dané vlivové křivky na přítomnost extrémních bodů).

S využitím známého vztahu

$$\hat{a}_{(-i)} = \hat{a} - (X^T X)^{-1} x_i \hat{\epsilon}_i / (1-h_{ii}) \quad (21)$$

lze snadno ukázat, že výběrová vlivová křivka SIC_i je (až na konstantu $n - 1$) rovna rozdílu $\hat{a} - \hat{a}_{(-i)}$.

Poznámka :

Také řada dalších diagnostických veličin pro identifikaci vlivných bodů porovnávané statistické charakteristiky pro všechna data a pro data bez i-tého bodu.

Pro charakterizaci vzdálenosti mezi \hat{a} a $\hat{a}_{(-i)}$ je vhodné použít Cookovu statistiku

$$D_i = \frac{(\hat{a}_{(-i)} - \hat{a})(X^T X)(\hat{a}_{(-i)} - \hat{a})}{m \hat{\sigma}^2} = \frac{1}{m} \cdot \frac{1}{h_{ii}} \cdot \frac{h_{ii}}{1-h_{ii}} \quad (22)$$

která má F - rozdělení s m a $n-m$ stupni volnosti. (Z rov. (22) plyne, že jde o analogii známého F-testu pro lineární regresní modely).

Atkinson /6/ navrhuje místo D_i použít modifikovanou funkci

$$T_i = \sqrt{\frac{(n-m) h_{ii}}{m (1-h_{ii})}} |\hat{\epsilon}_i| \quad (23)$$

která je citlivější na vybočující měření i extrémní body (Velká T_i identifikují vlivné body).

Pro vyjádření vlivu i-tého bodu na predikci lze použít rozdíl

$$F_i = \hat{y}_i - \hat{y}_{(-i)} = h_{ii} \hat{\epsilon}_i / (1-h_{ii}) \quad (24a)$$

resp. jako standardizované verze

$$DF_i = \frac{\hat{y}_i - \hat{y}_{(-i)}}{\hat{\sigma}_{(-i)} h_{ii}} = \sqrt{\frac{h_{ii}}{1-h_{ii}}} \cdot \hat{\epsilon}_i \quad (24)$$

Orientačně platí, že vlivný bod má DF_i vyšší než $2\sqrt{m/n}$.

Pro vyjádření vlivu i-tého bodu na změnu kovarianční matice odhadů lze využít determinantů (objemů konfidenčních oblastí).

Tedy

$$COV_i = \det \left[\frac{\hat{\sigma}^2}{\hat{\sigma}^2} (\hat{\sigma}_{(-i)}^2) \right]^{-1} / \det \left[\frac{\hat{\sigma}^2}{\hat{\sigma}^2} (\hat{\sigma}^2) \right]^{-1} = \left[\frac{\hat{\sigma}_{(-i)}^2}{\hat{\sigma}^2} \right]^m \cdot \frac{1}{1-h_{ii}} \quad (25)$$

V rov. (25) je $X_{(i)}$ matice, ve které chybí 1-tý řádek x_i^T . Pokud vyjde $|COV_i - 1| > 3m/n$, považuje se 1-tý bod za významný.

Jednoduchou možností indikace vlivných bodů je uvažovat rozšířenou matici $X = (X | \frac{y}{\hat{\sigma}})$, její řádky lze (pro stochastické vysvětlující proměnné) chápat jako $m+1$ rozměrné pozorování. Vybočující pozorování x_i^T se indikuje na základě rozdílu mezi tímto pozorováním a aritmetickým průměrem ostatním pomocí Wilksovy statistiky

$$\lambda_i = [n/(n-1)] \cdot (1-h_{ii}) \left[1 + \frac{\hat{\epsilon}_i^2}{n-m-1} \right]^{-1} \quad (26)$$

Pokud je matice X výběrem z vícerozměrného Gaussova rozdělení, má veličina

$$w_i = (1-\lambda_i) \sqrt{n-m-1} [m \cdot \lambda_i]^{-1} \quad (27)$$

Fisherovo rozdělení s m a $n-m$ stupni volnosti. Další druhy charakteristik vlivných bodů lze nalézt v [1-3].

V tabulce 3 jsou pro stejné data jako v tab. 2 určeny charakteristiky D_i , T_i , DF_i , COV_i a w_i .

Tabulka 3 : Různé charakteristiky vlivných bodů pro data s jedním vlivným měřením (O, E) o velikostech A

A	O					E				
	D_i	T_i	DF_i	COV_i	w_i	D_i	T_i	DF_i	COV_i	w_i
20	$6 \cdot 10^{-3}$	0,52	-0,11	1,1	0,78	$6,7 \cdot 10^{-3}$	0,562	-0,12	1,12	1,25
40	1,38	24,9	5,08	0,012	11,5	33	82,3	-16,8	0,15	23,5
60	1,5	50,3	10,3	$8,5 \cdot 10^{-4}$	22,2	118	188	-38,4	0,17	80,2
80	1,52	75,8	15,5	$1,7 \cdot 10^{-4}$	32,9	251	292	-59,7	0,26	169

Z tabulky 3 plyne, že :

- D_i nejsou vhodné pro identifikaci vybočujících měření (pro $n = 50$ $m = 2$ je $F_{95}(50,48) = 1,61$). Na extrémní body jsou D_i dostatečně citlivé.
 - statistiky T_i a w_i mohou identifikovat jak vybočující, tak i extrémní body
 - kritická hodnota pro DF_i je $\pm 0,4$. Tedy i tato statistika indikuje oba typy vlivných bodů
 - kritická hodnota pro $|COV_i - 1|$ je $\sim 0,12$.
- Tedy pouze pro $A = 20$ (kde nedošlo k prakticky žádnému zkreslení) nejsou nalezeny vlivné body (obou typů).

Pro rychlé identifikační účely je tedy výhodné použít místo běžně doporučených D_i modifikované statistiky T_i .

5. Odhady parametrů v přítomnosti vlivných bodů

Pro odhady parametrů, které jsou méně citlivé na přítomnost vlivných bodů se používá různých variant robustní regrese.

Klasické metody robustní regrese vedoucí na M - odhady parametrů nahrazují čtverce odchylek méně rychle rostoucí funkcí. Odpovídající kritérium regrese má tvar

$$S_M(a) = \sum_{i=1}^n \rho(d_i) = \sum_{i=1}^n \rho \left[\frac{(y_i - x_i^T a)^2}{\sigma^2} \right] \quad (28)$$

Minimalizace $S_M(\underline{a})$ je ekvivalentní řešení soustavy rovnic

$$\sum_{i=1}^n x_{ij} \psi(d_i) = 0 \quad (29)$$

kde $\psi(d_i) = \rho'(d_i)$.

Poznámka :

Pro odhady parametrů $\hat{\underline{a}}_M$ je možno využít faktu, že rov. (29) lze při zavedení vah $w_i = \psi(d_i)/d_i$ převést na soustavu normálních rovnic váženě metody nejmenších čtverců. Postačuje tedy iterativně počítat odhady parametrů $\hat{\underline{a}}_M$ pomocí vážené MNC s vahami určenými vždy z předchozí iterace.

V této práci je zvolena robustní Krasherova funkce

$$\rho(d) = \frac{w^2}{2} \left[1 - \exp\left(-\left(\frac{d}{w}\right)^2\right) \right] \quad (30)$$

Při volbě $w = 2,985$ má v případě, že data mají normální rozdělení, robustní regrese s funkcí (30) 95 %ní efektivnost.

Ve třetím a čtvrtém sloupci tab. 4 jsou pro simulační model s jedním vlivným bodem o velikosti A (pro $i = 37$) uvedeny odhady směrnice regreseční přímky \hat{a}_{M1} minimalizující rov. (28) při $\rho(d)$ definované rov. (30).

Tabulka 4 : Směrnice regreseční přímky \hat{a}_1 určené různými robustními metodami (Krasher pro $\rho(d)$ z rov. (30) resp. Krasher - Welch se speciálními vahami).

Metoda	Krasher		Krasher - Welch	
	A		E	O
40	0,926	0,973	0,958	0,962
70	0,2033	0,989	0,965	0,959
80	0,1401	0,969	0,965	0,959

Z tab. 4 je patrné, že klasické postupy robustní regrese jsou málo robustní vůči extrémním bodům (viz 2. sl. tab. 4)

Robustnost vůči vybočujícím měřením je velmi dobrá.

Eliminace extrémních bodů se dá provést pomocí speciálních vah w_{Si} , které souvisejí s velkými hustotami vlivných bodů.

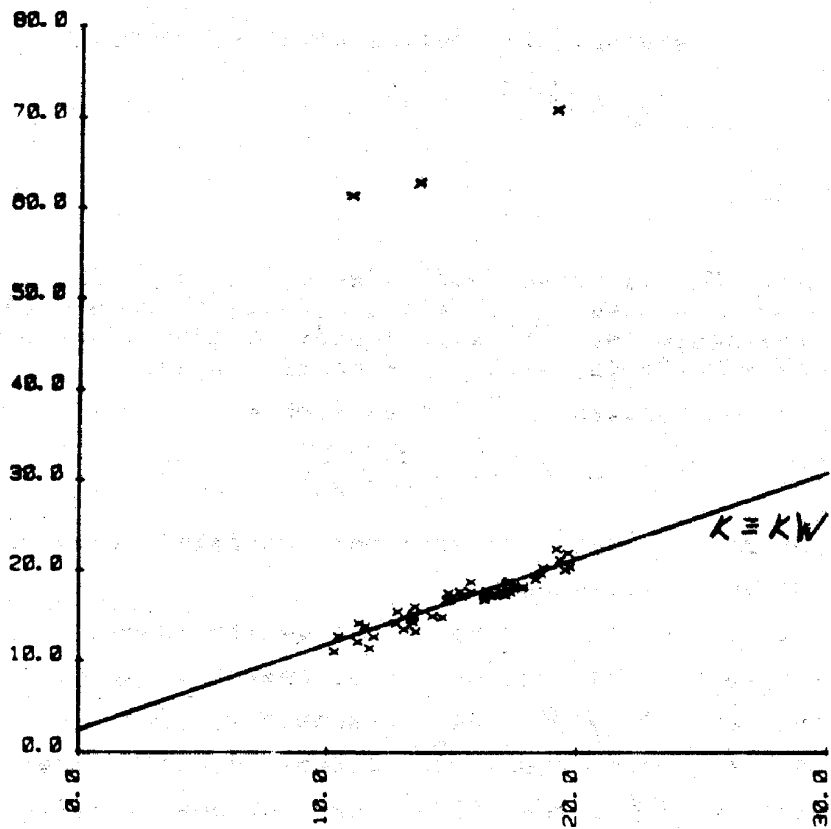
V této práci se využívá jednoduchých vah doporučených Krasherem a Welchem /8/.

$$w_{Si} = (1 - h_{ii}) / \sqrt{h_{ii}} \quad (31)$$

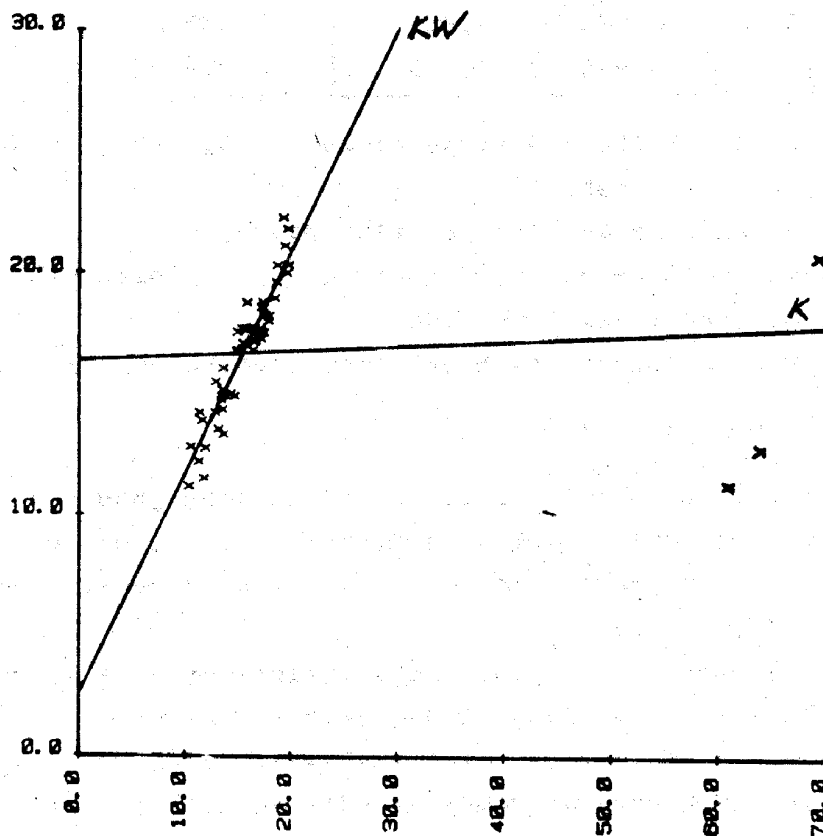
Ostatní způsoby volby vah eliminujících extrémní body jsou uvedeny např. v /7 - 9/. Obecně zavedení vah w_{Si} nijak nekomplikuje výpočty, protože stačí zavést globální váhy $w_i = w_{Si} \psi(d_i) / d_i$ a lze opět použít iterativně váženou metodu nejmenších čtverců /9/.

V pátém a šestém sloupci tab. 4 jsou pro simulační model s jedním vlivným bodem o velikosti A (pro $i = 37$) uvedeny odhady směrnice regreseční přímky \hat{a}_{M1} určené při použití vah dle (31) pro korekci extrémních bodů. Je zřejmé, že teprve pro tento případ je výsledek regrese robustní vůči oběma typům vlivných bodů.

Pro ilustraci rozdílů mezi klasickou robustní regresí a regresí robustní i vůči extrémním bodům jsou na obr. 2 odhady regreseční přímky určené oběma způsoby pro případ, kdy jsou v simulačním modelu tři vlivné body o velikosti $A = 50$ pro $i = 18, 25$ a 37.



Obr. 2 Regresní přímky (pro data se třemi vybočujícími body) počítané dle Krashera (K) resp. Krashera - Welsche (KW).



Obr. 3 Regresní přímky (pro data se třemi extrémními body) počítané dle Krashera (K) resp. Krashera - Welsche (KW)

Opět je jasně patrné, že klasické robustní techniky (Krašerova) jsou robustní vůči vybočujícím měřením a nikoliv extrémním bodům.

Pro automatickou eliminaci vybočujících měření byla navržena také technika "uřezané" regrese, kde se uřezání provádí s využitím regresních kvantilů. To znamená, že se nejdříve vyloučí body ležící mimo pás vymezený regresními kvantily pro $\alpha = 0,1$ a $\alpha = 0,9$ a pak se pro zbylé body použije standardní metoda nejmenších čtverců.

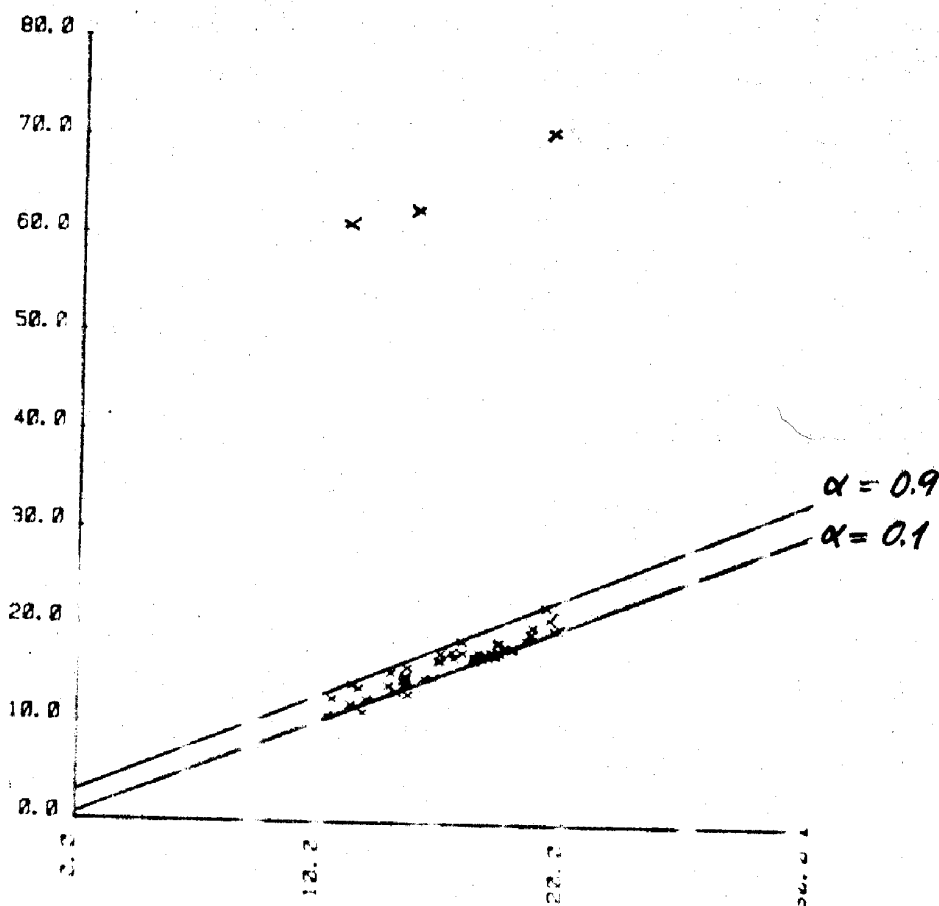
Poznámka :

Uvidíme, že α -regresní kvantil ($0 < \alpha < 1$) je regresní model $\hat{y} = x \hat{a}_\alpha$, kde pro odhad \hat{a}_α se minimalizuje vztah

$$\hat{a}_\alpha = \min_a [\alpha \sum_{i \in (y_i \geq x_i^T a)} |a_i| + (1-\alpha) \sum_{i \in (y_i < x_i^T a)} |a_i|] \quad (32)$$

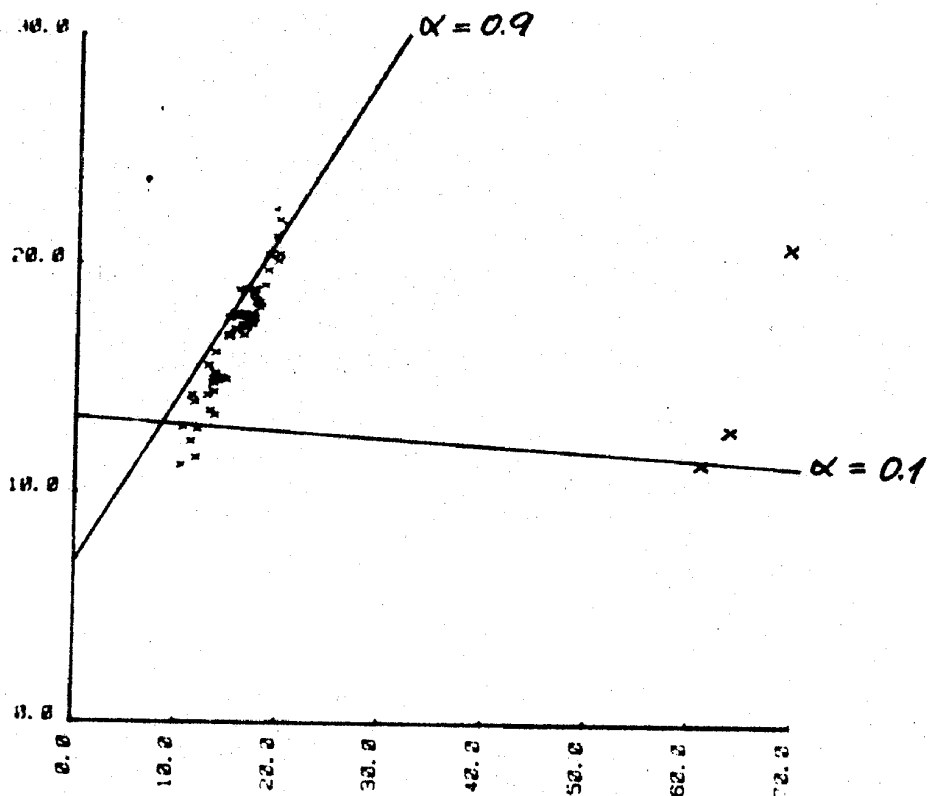
kde $a_i = y_i - x_i^T a$. Pro α -regresní kvantil platí, že
 - nejvíce $n \cdot \alpha$ bodů je pod a nejméně $(n-m) \cdot \alpha$ bodů je nad modelem $x \hat{a}_\alpha$
 - pro homoskedastické normální rozdělené data jsou kvantilové čáry (pro α nepříliš blízké 0 resp. 1) přibližně lineární.

Na obr. 4 jsou pro data z obr. 2 vybočující měření uvedeny α -kvantilové regresní přímky.



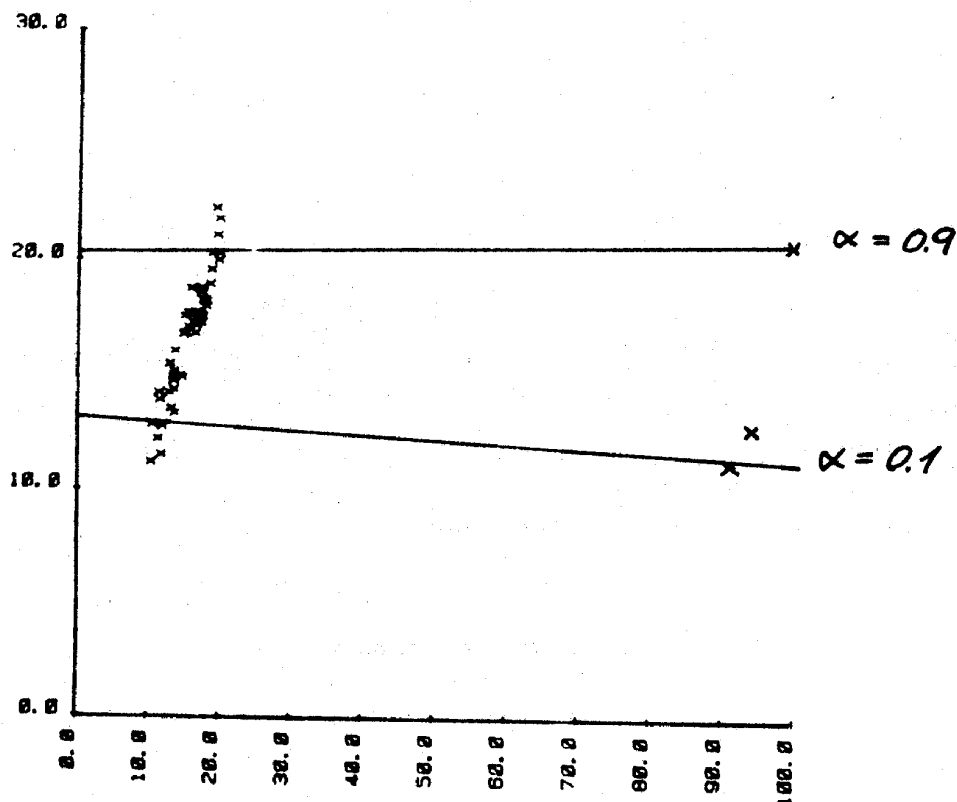
Obr. 4 Regresní α -kvantilové přímky pro data z obr. 2.

Na obr. 5 jsou pro data z obr. 3 (extrémní body) uvedeny α -kvantilové regresní přímky.



Obr. 5 Regresní α -kvantilové přímky pro data z obr. 3.

Z těchto obrázků je zřejmé, že pro případ, kdy se v datech vyskytnou extrémní body, nejsou regresní α -kvantily robustní a nebude možno použít postup z práce /10/. To, že pro dostatečně "extrémní" body jsou oba regresní α -kvantily nerobustní, je patrné z obr. 6 (simulovaná data se třemi extrémními body při $A = 80$ a $i = 12, 25, 37$)



Obr. 6 Regresní α -kvantilové přímky pro simulovaná data se třemi extrémními body ($A = 80$, $i = 12, 25, 37$).

Z uvedeného je zřejmé, že pro vyloučení vlivných bodů bude nutné nejdříve provést eliminaci extrémních bodů (např. přes Cookovu statistiku D_i z rov. (22)). Pak již bude možné využít postup eliminace vybočujících bodů s využitím regresních kvantilů.

6. Závěr

V příspěvku byly zmíněny některé postupy identifikace vlivných bodů a jejich eliminace resp. omezení jejich vlivu při odhadu parametrů lineárních modelů. Nebyly zde zmíněny postupy využívající různých typů grafů (ty jsou sumárně diskutovány např. v práci /11/) ani metody pro testaci většího počtu vlivných bodů /2/. Při konstrukci regresních α -kvantilů bylo využito programu, který sestavil J. Antoch z MFF UK.

Literatura

- /1/ Belsey D.A., Kuk E., Welsch R.E.: Regression Diagnostics, J.Wiley, New York 1980
- /2/ Cook R.D., Weisberg S.: Residuals and Influence in Regression, Chapman and Hall, New York 1982
- /3/ Militký J.: Tvorba matematických modelů - I, skripta pro kurs DT ČSVTS Ostrava 1983
- /4/ Huber P.J.: Robust Statistics, J.Wiley, New York 1981 (překlad v ruštině 1984)
- /5/ Cook R.D., Weisberg S.: Technometrics 22, 495 (1980)
- /6/ Atkinson A.C.: Biometrika 68, 13 (1981)
- /7/ Huber P.J.: J.A.S.A. 78, 66 (1983)
- /8/ Krashinsky W.S., Welsch R.E.: J.A.S.A. 77, 595 (1982)
- /9/ Hill R.W., Commun. Statist. A 11, 849 (1982)
- /10/ Antoch J. a kol., Proc. Compstat '84, Praha, září 1984
- /11/ Militký J.: Sborník přednášek z konference, Numerické metody ve fyzikální metalurgii, Blansko, říjen 1984