

Použití technik Jackknife a Bootstrap při odhadu modelových parametrů metodami nelineární regrese.

Jiří Militký, Výzkumný ústav zušlechťovací
544 28 Dvůr Králové n.L.

1. Úvod

Odhad parametrů v nelineárních regresních modelech patří mezi základní úlohy řešené v mnoha oblastech vědy a techniky. V naprosté většině případů se používá metody nejmenších čtverců (MNČ) v nejjednodušším tvaru (bez statistických vah měření). Tato metoda však často vede, zejména pro malé výběry (malý počet experimentálních bodů), k vychýleným odhadům. Důvody jsou v tom, že:

- a) regresní model nevystihuje přesně chování sledovaného systému (je více či méně zdařilou approximací jisté neznámé funkční závislosti)
- b) jednotlivá data ve výběru se neřídí předpoklady, za nichž byla MNČ odvozena (jde zejména heteroskedasticitu a porušení podmínek normality)
- c) kromě náhodných chyb měření se v datech vyskytuje také hrubé odchyly (jejichž zdrojem může být i vnitřní variabilita sledovaného systému).

Pro odhad parametrů lze v těchto případech použít robustní metody (numerické postupy pro nelineární modely uvádí např. Dennis a Welch /1/).

V poslední době se v literatuře objevují práce pojednávající o použití technik snižujících vychýlenost bodových odhadů (metody typu Jackknife a Bootstrap) pro lineární i nelineární regresní modely /14-16/.

V tomto příspěvku je uveden přehled základních variant technik Jackknife a Bootstrap. Je popsán program umožňující stanovení regresních odhadů metodou nejmenších čtverců a pomocí těchto technik. Na jednom modelu (se simulovanými daty) je provedeno porovnání odhadů získaných různými postupy s ohledem na jejich vychýlení.

2. Snižování vychýlenosti bodových odhadů.

Mějme náhodnou veličinu X , která má frekvenční funkci $f(x) = f(x, \Theta)$. Účelem je nalézt odhad $\hat{\Theta}$ neznámého parametru Θ na základě výběru sestaveného z N realizací (x_1, \dots, x_N) této náhodné veličiny. Odhad $\hat{\Theta}$ se počítá podle jisté funkce $\hat{\Theta} = S(x_1, \dots, x_N)$. Pokud je $\hat{\Theta}$ vychýlený, platí, že

$$E(\hat{\Theta} - \Theta) = h(\hat{\Theta}) \neq 0 \quad (1)$$

kde $E(\cdot)$ je operátor matematického očekávání a $h(\hat{\Theta})$ je vychýlení, které je funkci také velikosti výběru N . Obecně lze jednotlivým pozorováním ve výběru přiřadit kladné váhy (w_1, \dots, w_N) . Základní problém je v tom, že často neznáme přesný tvar $f(x)$. Na základě výběru však můžeme definovat empirickou frekvenční funkci

$$g(x) = \left(1/\sum_{(i)} w_i\right) \cdot \sum_{(i)} w_i \cdot \delta(x - x_i) \quad i = 1, \dots, N \quad (2)$$

kde $\delta(x - x_i)$ je Diracova funkce.

Parametr $\hat{\Theta} = S(f(x))$ pak odhadujeme pomocí funkce

$$S(g(x)) = \hat{\Theta} = S(x_1, \dots, x_N, w_1, \dots, w_N)$$

Je zřejmé, že čím bude $g(x)$ blíže k $f(x)$ (např. ve smyslu Prochorovovy metriky /6/), tím se bude více $\hat{\Theta} \rightarrow \Theta$.

Pro stanovení vychýlení $h(\hat{\Theta})$ uvažujme frekvenční funkci $h(x)$, která je dána vztahem

$$h(x) = (1-t)f(x) + t g(x) \quad (3)$$

Odhad $S(h(x))$ založený na této frekvenční funkci lze vyjádřit Taylorovým rozvojem /6/

$$S(h(x)) = S(f(x)) + t A_1 + t^2/2 A_2 + \dots \quad (4)$$

kde A_1, A_2 jsou derivace $S(h(x))$. Např. parametr $A_1 = \int f(x) g'(x) dx$ obsahuje tzv. von Miessovu derivaci prvního řádu $Y(x)$.

Pokud je $h(x)$ blízké $f(x)$, lze pro $t=1$ vyjádřit rozvoj (1) ve tvaru

$$S(g(x)) = \Theta + \int Y(x) g(x) dx + \frac{1}{2} \iint Y(x, y) g(x) g(y) dx dy + \dots \quad (5)$$

Dosazením za $g(x)$ z rov. (2) a úpravami vyjde

$$\hat{\Theta} = \Theta + \frac{\sum_i w_i Y(x_i)}{\sum_i w_i} + \frac{1}{2} \cdot \frac{\sum_i \sum_j w_i w_j Y(x_i, x_j)}{\left(\sum_i w_i\right)^2} + \dots \quad (6)$$

Poznámka: Funkce $\psi(x_i)$ vyjadřuje vliv hodnoty x_i na odhad $\hat{\theta}$. Proto ji Hampel [9] nazval vlivová funkce. Jde vlastně o první von Misesovu derivaci ve směru Diracovy funkce. Lze dokázat, že $\psi(x_i) = \lim_{\epsilon \rightarrow 0} \{[S(g(x)) - S(f(x))] / \epsilon\}$.

Z rov. (6) je patrné, že vychýlení $b(\hat{\theta})$ lze pro případ, kdy $x_i = 1, i = 1, \dots, N$ vyjádřit vztahem

$$b(\hat{\theta}) = \frac{\alpha_1}{N} + \frac{\alpha_2}{N^2} + \dots \quad (7)$$

V roce 1956 navrhl Quenouille [2] jednoduchý způsob snížení možné vychýlenosti bodových odhadů, který vychází z rozvoje (7).

Vyšel přitom z odhadu $\hat{\theta}_i = S(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N)$, který získal po vynechání i-tého pozorování x_i . Jeho vychýlení lze vyjádřit tvarem

$$b(\hat{\theta}_i) = \frac{\alpha_1}{(N-1)} + \frac{\alpha_2}{(N-1)^2} + \dots \quad (8)$$

Snadno pak dokázal, že odhad

$$\tilde{\theta}_i = N\hat{\theta} - (N-1)\hat{\theta}_i \quad (9)$$

je méně vychýlený, protože již neobsahuje první člen rozvoje.

Vychýlení $b(\tilde{\theta}_i)$ je dán vztahem

$$b(\tilde{\theta}_i) = -\frac{\alpha_2}{N(N-1)} + \dots \approx -\frac{\alpha_2}{N^2} + O(N^{-3}) \quad (10)$$

Postupným opakováním této techniky lze eliminovat i další členy rozvoje (7). Aby nedošlo ke ztrátě efistence odhadu, je nutné nahradit $\tilde{\theta}_i$ průměrnou hodnotou $\tilde{\theta}_J$.

$$\tilde{\theta}_J = \frac{1}{N} \sum_{i=1}^N \tilde{\theta}_i \quad (11)$$

Tukey označil $\tilde{\theta}_i$ jako pseudohodnoty a $\tilde{\theta}_J$ jako Jackknife odhad. V abstraktu [4] předpokládal, že $\tilde{\theta}_i$ jsou přibližně nezávisle náhodné veličiny, pro které lze stanovit výběrový rozptyl

$$\text{var}(\tilde{\theta}_J) = \frac{1}{N(N-1)} \sum_{i=1}^N (\tilde{\theta}_i - \tilde{\theta}_J)^2 \quad (12)$$

Pro odhad konfidenčních intervalů lze použít statistiky

$$\sqrt{N}(\tilde{\theta}_J - \theta) \left\{ \frac{1}{N-1} \sum_{i=1}^N (\tilde{\theta}_i - \tilde{\theta}_J)^2 \right\}^{-1/2} \quad (13)$$

která má přibližně t-rozdělení s N-1 stupni volnosti.

Poznámka: Miller /8/ uvádí některé příklady, které ukazují, že odhad konfidenčních mezi podle rov. (13) nemusí být vždy lepší než obvyklé odhady (založené např. na předpokladu normality).

Rey /6/ navrhl zobecnění techniky Jackknife na případ, kdy váhy měření w_i jsou nekonstantní. Místo vynechávání bodů provádí úpravu w_i pomocí faktoru $(1+t)$. Odhad $\hat{\theta}_i$ se pak určí ze vztahu

$$\hat{\theta}_i(t) = S(x_1, \dots, x_N, w_1, \dots, w_{i-1}, (1+t)w_i, w_{i+1} \dots, w_N) \quad (14)$$

Odpovídající pseudohodnoty $\tilde{\theta}_{w_i}$ jsou definovány vztahem

$$\tilde{\theta}_{w_i} = [(\bar{t}w_i + \sum_{j=1}^N w_j)\hat{\theta}_i(t) - (\sum_{j=1}^N w_j)\hat{\theta}] / t \quad (15)$$

Jackknife odhad $\tilde{\theta}_{w_j}$ je pak vážený průměr

$$\tilde{\theta}_{w_j} = (1/\sum w_i) \sum_{(i)} \tilde{\theta}_{w_i} \quad (15a)$$

Odpovídající vychýlení $b(\tilde{\theta}_{w_j})$ má tvar

$$b(\tilde{\theta}_{w_j}) = \frac{\sum w_i y_i}{\sum w_i} + \frac{1}{2} (1+t) \frac{\sum \sum x_i w_j y_i (x_i, w_j)}{(\sum w_i)^2} + \dots \quad (16)$$

Volbou $t=-1$ vypadne v rov. (16) druhý člen. Odhad $\tilde{\theta}_{w_j}$ je pak váženou variantou původní Jackknife techniky. Pro odhad rozptylu platí

$$\text{var}(\tilde{\theta}_{w_j}) = \sum_{(i)} (\tilde{\theta}_{w_i} - \bar{w}_i \tilde{\theta}_{w_j})^2 / [(\sum_{(i)} w_i)^2 - \sum_{(i)} w_i^2] \quad (17)$$

Pro malá záporná t lze vyjádřit odhad vychýlení \hat{b} vůči $\tilde{\theta}_{w_j}$ vztahem /7/

$$\hat{b}(\hat{\theta}, t) = (\hat{\theta} - \tilde{\theta}_{w_j}) = \frac{1-t}{Nt^2} \left\{ \frac{1}{N} \sum \hat{\theta}_i(t) - \hat{\theta} \right\} \quad (18)$$

Jaeckel (viz /7/) navrhl tzv. infinitesimální Jackknife odhad $\hat{\theta}_{IJ}$ ve tvaru

$$\hat{\theta}_{IJ} = \hat{\theta} - \lim_{t \rightarrow 0} \hat{b}(\hat{\theta}, t) \quad (18a)$$

Rozšíření techniky váženého Jackknife odhadu na případ, kdy

odhadované parametry tvoří vektor, je publikováno v práci /6/.

Schucany, Gray a Owen /3/ stanovili, že ze dvou vychýlených odhadů $\hat{\theta}_1$ a $\hat{\theta}_2$ lze vypočítat nevychýlený odhad $\tilde{\theta}_3$ podle vztahu

$$\tilde{\theta}_3 = (\hat{\theta}_1 - R \hat{\theta}_2)(1-R)^{-1} \quad (19)$$

kde $R = b(\hat{\theta}_1)/b(\hat{\theta}_2)$ je poměr vychýlení odhadů. Lze snadno dokázat, že pro $\hat{\theta}_1 = \hat{\theta}_2$, $\hat{\theta}_2 = \sum \hat{\theta}_i/N$ a $R = (N-1)/N$ je $\tilde{\theta}_3 = \hat{\theta}_1$. V téže práci je navržen postup pro tvorbu odhadů, jejichž vychýlení $b(\tilde{\theta}_j)$ neobsahuje prvních k členů v mocninném rozvoji /7/. Dalšího snížení vychýlenosti odhadů lze docílit rozdělením výběru do J skupin o velikosti $n (N=nJ)$ a vylučováním celých skupin při výpočtu pseudohodnot. Pak

$$\hat{\theta}_j = S(x_1, \dots, x_{j-1}, m, x_{j+1}, \dots, x_N) \quad j=1, \dots, J \quad (20)$$

$$\tilde{\theta}_{GJ} = J\hat{\theta} - \frac{J-1}{J} \sum_{j=1}^J \hat{\theta}_j$$

Aziz /5/ navrhl vylučovat všechny skupiny kromě l -té. Pak

$$\hat{\theta}_l = S(x_{l-1}, n+1, \dots, x_n) \quad l=1, \dots, J \quad (21)$$

$$\tilde{\theta}_{AJ} = \frac{J}{J-1} \hat{\theta} - \frac{1}{J(J-1)} \sum_{l=1}^J \hat{\theta}_l$$

Jak je zřejmé, je odhad $\tilde{\theta}_{AJ}$ výpočetně složitý (zejména pro velké N).

V práci /5/ je dokázáno, že pro všechna reálná $p > 0$ taková, že

$$0 < \lim_{N \rightarrow \infty} N^p b(\hat{\theta}) < \infty$$

je $b(\tilde{\theta}_{AJ}) < b(\hat{\theta})$ jen v případě, že $p < \ln(2g-1)/\ln g$.

Za těchž podmínek je však $b(\tilde{\theta}_{AJ}) > b(\tilde{\theta}_{GJ})$ (pro $p \neq 1$).

Aziz /5/ také navrhl, jak určit optimální J s ohledem na minimalizaci střední kvadratické odchyly odhadů.

Další možnosti je odstranění n prvků ve výběru všemi ($\binom{N}{n}$) možnými způsoby.

Výhody a nevýhody jednotlivých způsobů dělení výběru do skupin diskutuje např. Rey /6/.

Jak je zřejmé z výše uvedeného, umožňuje technika Jackknife snížení vychýlenosti odhadů a stanovení jejich rozptylů bez detailních znalostí o výběrové frekvenční funkci. Podmínkou je však splnění předpokladů, za nichž byla odvozena rov. (7).

velmi blízká technikám Jackknife je metoda Bootstrap /10/.

Postup při jejím použití lze vyjádřit posloupností těchto kroků /11/

- I. Vytvoření empirické výběrové frekvenční funkce $g(x)$, např. podle rov. (2) položením $u_i = 1 (i = 1 \dots N)$
- II. Sestavení náhodného výběru o velikosti N z $g(x)$, tedy tzv. Bootstrap výběru $(x_{ij}^*, x_{ij}^*, x_{Nj}^*)$
- III. Stanovení odhadu $\hat{\theta}_{Bj} = S(x_{1j}^*, \dots, x_{Nj}^*)$
- IV. Opakování kroků II a III pro $j \neq 1 \dots M$ a pak výpočet Bootstrap odhadu

$$\tilde{\theta}_B = \frac{1}{M} \sum_{j=1}^M \hat{\theta}_{Bj}. \quad (22)$$

Na rozdíl od Jackknife se náhodné výběry tvorí z množiny (x_1, \dots, x_N) s vracením (tj. některá x_i se v Bootstrap výběru vyskytuje víckrát a některá vůbec ne). Vztah mezi $\tilde{\theta}_B$ a $\hat{\theta}_j$ je diskutován v práci /10/.

3. Odhad parametrů v regresních modelech.

Mějme modelovou funkci $\varphi(\bar{x}, \bar{\theta})$ a výběrová data (experimentální body) $(x_i, y_i)_{i=1 \dots N}$. Účelem je nalézt odhady $\hat{\theta}$ modelových parametrů $\bar{\theta}$. V nejjednodušším případě lze uvažovat model měření ve tvaru

$$y_i = \varphi(\bar{x}_i, \bar{\theta}) + \varepsilon_i \quad (23)$$

kde ε_i jsou náhodné chyby, o nichž se obvykle předpokládá, že mají stejné rozdělení s frekvenční funkcí $f(\varepsilon)$, a že

$$E(\varepsilon_i) = 0, \quad E(\varepsilon_i \varepsilon_j) = 0 \quad i \neq j = 1 \dots N$$

Odhad parametrů $\bar{\theta}$ lze pro známé $f(\varepsilon)$ získat maximalizací věrohodnostní funkce

$$\hat{\theta} = \max_{\bar{\theta} \in E^P} \prod_{i=1}^N f(\varepsilon_i) \quad (24)$$

Pokud $\varepsilon_i \in N(0, \sigma^2)$, vede řešení problému (24) na metodu nejménších čtverců odchylek, kdy

$$\hat{\theta} = \min_{\bar{\theta} \in E^P} \sum_{i=1}^N [y_i - \varphi(x_i, \bar{\theta})]^2 \quad (25)$$

Lze dokázat, že normální rozdělení je stabilní (definice stabilitnosti viz /12/) ve třídě všech hustot s ohraničeným rozptylem.

Na druhé straně však bylo stanoveno, že MNČ poskytuje v případě, že $f(\epsilon)$ je Laplaceovo oboustranné rozdělení, odhadu s dvojnásobným rozptylem než metoda nejmenších absolutních odchylek (L_1 aproximace) /13/.

Je také známo, že určené MNČ jsou citlivé na výskyt hruživých chyb a porušení předpokladu homoskedasticity.

Obvykle používané robustní metody regrese nahrazují součet čtverců v rov. (25) jinou, méně rychle rostoucí funkci. To vede zejména u nelineárních modelů ke značným numerickým obtížím.

Při použití technik Jackknife a Bootstrap se vychází z opakování řešení rov. (25) pro různě modifikovaná data. První práci týkající se aplikace Jackknife metody při odhadu parametrů v lineárních regresních modelech publikoval Miller /14/. Ten vyšel z modelu

$$\bar{Y} = X \bar{\theta} + \bar{\epsilon} \quad (25a)$$

kde $\bar{Y} = (y_1, \dots, y_N)^T$, $\bar{\theta} = (\theta_1, \dots, \theta_p)^T$, $\bar{\epsilon}$ je chybový vektor a X je matice $(N \times p)$ obsahující jako řádky vektory $\bar{x}_i = (x_{1,i}, \dots, x_{p,i})^T$. Metodou nejmenších čtverců lze určit podle výrazu

$$\hat{\theta} = (X^T X)^{-1} X^T \bar{Y} \quad (26)$$

Odhad $\hat{\theta}_i$ vzniklý vyloučením bodu (y_i, \bar{x}_i^T) je možno vypočítat podle relace (viz /14/)

$$\hat{\theta}_i = \hat{\theta} - \frac{(X^T X)^{-1} \bar{x}_i^T (y_i - \bar{x}_i^T \hat{\theta})}{1 - \bar{x}_i^T (X^T X)^{-1} \bar{x}_i^T} \quad (27)$$

Rov. (27) tedy umožňuje snadný výpočet $\hat{\theta}_i$ při znalosti $\hat{\theta}$ z MNČ. Z odhadu $\hat{\theta}_i$ se dle rov. (9) snadno určí pseudohodnoty $\tilde{\theta}_{LJ}$ a z nich pomocí rov. (11) linearizovaný Jackknife odhad $\tilde{\theta}_{LJ}$ (pro lineární regresní modely je $\tilde{\theta}_{LJ} = \hat{\theta}_J$). Hinkley /15/ nazývá tento postup jako rovnovážný Jackknife.

Pro případ, že jsou "váhy" $\bar{x}_i (\mathbf{X}^T \mathbf{X})^{-1} \bar{x}_i^T$ v jednotlivých bodech "regresního plánu" \mathbf{X} nestejně, navrhuje tzn. nerovnovážný Jackknife, pro který lze pseudohodnoty $\tilde{\theta}_{Ni}$ vyjádřit tvarem

$$\tilde{\theta}_{Ni} = \hat{\theta} + N(\mathbf{X}^T \mathbf{X})^{-1} \bar{x}_i (y_i - \bar{x}_i \hat{\theta}) \quad (28)$$

Odpovídající Jackknife odhad $\tilde{\theta}_{Nj}$ však nesnižuje pro lineární modely žádné vychýlení, protože platí

$$\tilde{\theta}_{Nj} = \frac{1}{N} \sum_{i=1}^N \tilde{\theta}_{Ni} = \hat{\theta} \quad (29)$$

Použití nerovnovážného Jackknife odhadu je tedy vhodné pouze při výpočtu rozptylu odhadů podle rov. (12).

Miller /14/ dokázal, že Jackknife odhad $\tilde{\theta}_{Nj}$ funkce $\theta = \varphi(\beta)$ regresních parametrů β v lineárním modelu $\bar{Y} = \mathbf{X}\beta + \bar{E}$ je asymptoticky normální i pro případ, že E_i jsou negaussovské náhodné veličiny.

Pro nelineární regresní modely $\varphi(\bar{x}, \hat{\theta})$ se obyčejně využívá linearizace pomocí Taylorova rozvoje /14, 16/

$$\bar{\theta} - \hat{\theta} = (\mathcal{Z}^T \mathcal{Z})^{-1} \mathcal{Z}^T \bar{r} \quad (30)$$

kde $\bar{r} = (y_1 - \varphi(\bar{x}_1, \hat{\theta}), \dots, y_N - \varphi(\bar{x}_N, \hat{\theta}))$

a \mathcal{Z} je Jakobián, jehož řádky tvoří vektory

$$\bar{x}_i^T = \left(\frac{\partial}{\partial \theta_1} \varphi(\bar{x}_i, \hat{\theta}), \dots, \frac{\partial}{\partial \theta_p} \varphi(\bar{x}_i, \hat{\theta}) \right)$$

Nahradíme-li v rov. (27) a (28) matici \mathbf{X} maticí \mathcal{Z} , a vektory \bar{x}_i pomocí vektorů $\bar{x}_{(i)}$, lze snadno určit odhadu $\tilde{\theta}_{Lj}$ i $\tilde{\theta}_{Nj}$, aniž je nutno řešit $(n+1)$ krát iterativně rov. (25) (pro výpočet $\tilde{\theta}_{Nj}$). I pro odhad parametrů nelineárních modelů lze určit odhad rozptylů podle rov. (12) /17/. Rozdíl mezi $\tilde{\theta}_{Nj}$ a $\tilde{\theta}_{Lj}$ je měřítkem toho, nakolik je linearizace (30) dobrou approximaci $\varphi(\bar{x}, \hat{\theta})$.

Chambers /19/ navrhl pro approximaci pseudohodnot $\tilde{\theta}_{ci}$ v nelineárních modelech vztah

$$\tilde{\theta}_{ci} = \hat{\theta} + (N-1)(\mathcal{Z}_{(i)}^T \mathcal{Z}_{(i)})^{-1} \bar{r}_{(i)} \quad (31)$$

Zde $\mathcal{Z}_{(i)}$ resp. $\bar{r}_{(i)}$ jsou matice \mathcal{Z} resp. vektor \bar{r} s vyněchaným i-tým bodem $(\bar{x}_{(i)}^T, y_i)$. Dosazením $\tilde{\theta}_{ci}$ do rov. (11)

rezultuje approximativní Jackknife odhad $\tilde{\hat{\theta}}_{\text{ej}}$.

Pro stanovení zrobustněných Jackknife odhadů se doporučuje nahradit aritmetický průměr v rov. (11) např. 10%ním useknutým průměrem nebo jiným robustním odhadem polohy /15/.

Při odhadu regresních parametrů technikou Bootstrap se postupuje takto /10/.

- I. Na základě původního výběru (\bar{x}_i^T, \bar{y}_i) $i=1..N$ se metodou MNČ (podle rov. (25)) stanoví odhad $\hat{\theta}$.
- II. Určí se rezidua $\hat{\epsilon}_i = \bar{y}_i - f(\bar{x}_i, \hat{\theta})$ a jejich empirická frekvenční funkce $g(\hat{\epsilon})$ podle rov. (2) (při $w_i = 1$, $i=1..N$).
- III. Z empirické frekvenční funkce $g(\hat{\epsilon})$ se pomocí generátoru náhodných čísel sestaví Bootstrap výběr o velikosti N ($\hat{\epsilon}_1^{(j)} \dots \hat{\epsilon}_N^{(j)}$). Jde o výběr s vracením.
- IV. Určí se nové "náhodné" proměnné $\bar{y}_i^{(j)} = f(\bar{x}_i, \hat{\theta}) + \hat{\epsilon}_i^{(j)}$ $i=1..N$
- V. Provede se odhad parametrů $\tilde{\hat{\theta}}$ z rov. (25) pro výběr $(\bar{x}_i^T, \bar{y}_i^{(j)})$ $i=1..N$. Získaný odhad označme $\tilde{\hat{\theta}}^{(j)}$.
- VI. Provede se M opakování ($j=1..M$) kroků III až V.

Bootstrap odhad $\tilde{\hat{\theta}}_B$ je pak aritmetickým průměrem

$$\tilde{\hat{\theta}}_B = \frac{1}{M} \sum_{j=1}^M \tilde{\hat{\theta}}^{(j)} \quad (32)$$

a odpovídající rozptyl lze vyčíslit z výrazu

$$\text{var}(\tilde{\hat{\theta}}_B) = \frac{RSC}{N-p} (\mathbf{Z}^T \mathbf{Z})^{-1} \quad (33)$$

kde RSC je reziduální součet čtverců.

4. Program pro odhad modelových parametrů.

Na základě výše uvedených vztahů byl sestaven program v jazyce HPL pro stolní kalkulátor HP 9825 A. Ten umožňuje výpočet odhadů $\hat{\theta}$ (MNČ z rov. (25)), $\tilde{\hat{\theta}}$ (z rov. (9) a (11)), $\tilde{\hat{\theta}}_{\text{ej}}$ (z rov. (27) a (11)), $\tilde{\hat{\theta}}_{NJ}$ (z rov. (29) a (11)), $\tilde{\hat{\theta}}_{\text{ej}}$ (z rov. (31) a (11)) a $\tilde{\hat{\theta}}_B$ (z rov. (32)). Při výpočtu odhadu $\tilde{\hat{\theta}}_B$ je použito $M=N$ Bootstrap výběrů. Pro průměrování pseudohodnot je použito jak aritmetického tak i 10%ního useknutého průměru. Pro odhadы $\hat{\theta}$ a $\tilde{\hat{\theta}}_B$ jsou kovarianční matice určeny z rov. (33) pro ostatní z rov. (12).

Rov. (25) (MNČ) je řešena Marquardtovou metodou s modifikacemi navrženými Nagelem a Wolfem /18/ a adaptivní korekci maximální délky přirůstek v každé iteraci.

5. Příklad

V literatuře existuje zatím velmi málo prací, ve kterých byly techniky Jackknife a Bootstrap prakticky použity při odhadu parametrů v nelineárních modelech $\hat{Y}(\hat{\alpha}, \hat{\theta})$ v případě malého počtu měření N . Zde je pro porovnání výše uvedených odhadů použit jednoduchý Arrheniův model vyjadřující závislost rychlostní konstanty K na teplotě T vztahem

$$K = \exp(K^* - E/RT) \quad (34)$$

kde K^* a E jsou modelové parametry (E má význam aktivační energie daného děje) a R je univerzální plynová konstanta. Protože má E známý fyzikální význam, je přirozené požadovat, aby její odhad byl co nejméně vychýlený. Z praktického hlediska však často nelze experimentálně stanovit dostatečně velký výběr $(K_i, T_i)_{i=1..N}$ ani provést opakování měření v bodech T_i .

Aby bylo možno porovnat chování jednotlivých odhadů i za těchto podmínek, byly připraveny simulované výběry o velikosti 10. Bylo zvoleno $K^* = \ln(251 \cdot 10^9)$ a $E = 59$ a rovnoramné dělení teploty $T_i = 323,15 + 10(i-1) \quad i = 1 \dots 10$

Jednotlivé rychlostní konstanty byly určeny ze vztahu

$$K_i = \exp(251 \cdot 10^9 - 59/RT_i) + \varepsilon_i$$

Náhodná proměnná ε_i byla volena tak, aby bylo možno stanovit vliv heteroskedasticity i pro případy, kdy jde o součet dvou gaussovský rozdelených veličin. Jednotlivá ε se počítala ze vztahu

$$\varepsilon_i = (1-G) \cdot N(0, \tilde{\sigma}_i^2) + G \cdot N(0, \tilde{\sigma}^2) \quad (35)$$

Pomocí parametru G (bylo voleno $G = 0, 0.25, 0.5, 0.75$ a 1) lze získat ε_i , která se více či méně liší od předpokladu $\varepsilon_i \in N(0, \tilde{\sigma}^2)$.

První člen v rov. (35) odpovídá heteroskedastickým datům. Rozptyl $\hat{\sigma}_i^2$ se určoval z výrazu $\hat{\sigma}_i^2 = \hat{G}_i^2 K_i$ a bylo voleno $\hat{G}^2 = 15$. Celý simulační experiment není ještě ukončen. Zatím byly pro každé G generovány čtyři náhodné výběry. Vzhledem k tomu, že tento počet nepostačuje pro statistické zpracování výsledků, bylo postupováno tak, že pro každý výběr bylo stanovano pořadí podle vzdálenosti odhadů aktivační energie E od zadané hodnoty 59. V tab. I jsou pro jednotlivá G uvedena vždy průměrná pořadí (čím *následně* pořadí, tím je odhad méně vychýlený).

Tab. I Pořadí odhadů aktivační energie (průměry ze čtyř náhodných výběrů).

G	0	0.25	0.5	0.75	1
Odhad					
$\hat{\theta}_{MN}$	3-4	5-6	6-7	3-4	2-3
$\hat{\theta}_{CJ}$	5-6	8	2-3	6	10
$\hat{\theta}_{CJ}^*$	7	10	5	5	5-6
$\hat{\theta}_{LJ}$	9	9	9-10	10	5-6
$\hat{\theta}_{LJ}^*$	1-2	1	4	8	7
$\hat{\theta}_{NJ}$	3-4	5-6	6-7	3-4	2-3
$\hat{\theta}_{NJ}^*$	1-2	4	2-3	1-2	8
$\hat{\theta}_B$	8	2	8	1-2	9
$\hat{\theta}_J$	10	7	9-10	9	1
$\hat{\theta}_J^*$	5-6	3	1	7	4

*) je použito 10%ního useknutého průměru

Vzhledem k malému počtu simulací zde nejsou ani porovnávány jednotlivé kovarianční matice odhadů. I když pořadí není kvantitativní měrou vychýlenosti jednotlivých odhadů, lze i z těchto údajů zjistit zajímavé závěry:

- a) Odhad $\hat{\theta}$ a $\hat{\theta}_{NJ}$ jsou prakticky stejné. To svědčí o *ustálenosti* experimentálního "plánu" /15/.
- b) Odhad $\hat{\theta}_J$ a $\hat{\theta}_{LJ}$ resp. $\hat{\theta}_{CJ}$ se liší. To svědčí o faktu, že linearizace vyjádřená rov. (30) není příliš dobrou approximací rov. (34).
- c) Náhradou aritmetického průměru useknutým průměrem dochází většinou ke snížení vychýlení. Jsou však případy (zejména u $\hat{\theta}_{CJ}$), kdy dojde k malému zhoršení.

- d) Počet výběru $M = N = 10$ zřejmě nestačí pro získání méně vychýlených odhadů $\hat{\theta}_B$.
- e) Chování odhadů je silně ovlivněno typem chyb ε_i .

Vě většině případů se však ukázalo, že pokud byla $\hat{\theta}$ podstatněji vychýlená, došlo u všech Jackknife i Bootstrap techniky buď k malému snížení vychýlenosti nebo dokonce i ke zvýšení $\hat{\theta}(\hat{\theta})$. Výhodnější by zřejmě bylo vypuštění "bodů", jimž odpovídající pseudohodnoty nebyly zařazeny do výpočtu useknutého průměru.

Definitivnější závěry však bude možno stanovit až po provedení celého zamýšleného simulačního experimentu.

6. Závěr:

V příspěvku byly uvedeny postupy pro snižování vychýlenosti bodových odhadů. Bylo ukázáno, jak lze využít při odhadu parametrů v nelineárních regresních modelech.

Literatura:

- /1/ Dennis J., Welch R.E.: Commun. Statist. B7, 345 (1978)
- /2/ Quenouille M.H.: Biometrika 43, 353 (1956)
- /3/ Schucany W.R., Gray H.L., Owen D.B.: J. Amer. Stat. Assoc. 66, 524 (1971)
- /4/ Tukey J.W.: Ann. Math. Statist. 29, 614 (1958)
- /5/ Aziz E.S.: On Jackknife Estimators Based on Dependent Samples, Rept. Rouen University No 4, July 1979
- /6/ Rey W.J.J.: Robust Statistical Methods, Lect. Notes in Mathematics, Springer Verlag, Berlin 1978
- /7/ Miller R.G.: Biometrika 61, 1 (1974)
- /8/ Miller R.G.: Ann. Math. Statist. 35, 1594 (1964)
- /9/ Hampel F.R.: J. Amer. Stat. Assoc. 69, 383 (1974)
- /10/ Efron B.: Ann. Statist. 7, 1 (1979)
- /11/ Efron B.: SIAM Rev. 21, 460 (1979)
- /12/ Vapnik V.N.: Vosstanovlenije zavisimostěj po empiričeskim dannym, Nauka Moskva 1979
- /13/ Mudrov V.I., Kuško V.L.: Metody obrabotki izmerenij, Sov. Radio, Moskva 1976
- /14/ Miller R.G.: Ann. Statist. 2, 880 (1974)
- /15/ Hinkley D.V.: Technometrics 19, 285 (1977)
- /16/ Fox T., Hinkley D.V., Larntz K.: Technometrics 22, 29 (1980)
- /17/ Duncan G.T.: Technometrics 20, 123 (1978)
- /18/ Nagel G., Wolf O.: Biometrische Ztsch. 16, 431 (1974)
- /19/ Chambers J.M.: Biometrika 60, 1 (1973)