

NĚKTERÉ POSTUPY PRO VÝPOČET ROBUSTNÍCH ODHADŮ V LINEÁRNÍM REGRESNÍM MODELU

Jaromír ANTOCH

1. ÚVOD

Uvažujme klasický lineární regresní model

$$(1.1) \quad \tilde{Y} = \tilde{X} \tilde{\beta} + \tilde{\varepsilon},$$

kde $\tilde{X} = (x_{ij})_{i=1, \dots, n}^{j=1, \dots, p}$ je matici známých regresních konstant, $\tilde{Y} = (Y_1, \dots, Y_n)$ vektor pozorování a $\tilde{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$ vektor chyb. O vektoru chyb budeme předpokládat, že jeho složky jsou nezávislé, stejně rozdělené dle distribuční funkce F , jež není obecně specifikována.

Je-li F normální, je všeobecně známo, že optimálním způsobem odhadu vektoru parametrů $\tilde{\beta}$ je klasická metoda nejménších čtverců. Tzn., že odhad $\hat{\beta}$ dostaneme jako

$$(1.2) \quad \min_{\tilde{\beta}} \sum_{i=1}^n \left(Y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2.$$

V praxi se ale občas setkáváme s pozorováními, která určitým způsobem "vybočují z fady", čímž způsobují, že klasické metody ztrácejí svoji účinnost. Proto byly v nedávné minulosti navrženy nové metody odhadu, méně citlivé na výskyt takovýchto pozorování. Jednu z nejdůležitějších tříd mezi nimi zaujmají tzv. robustní metody, jimž je v poslední době věnována velká pozornost ve všech oborech matematické statistiky.

Cílem tohoto příspěvku je zrekapitulovat některé algoritmy a ukázat postupy pro výpočet robustních odhadů v lineárním regresním modelu. Zaměřím se přitom na robustnost vzhledem k rozdělení chyb, tj. proti výskytu odlehlych pozorování v datech.

2. Submodel parametru polohy

Pro submodel parametru polohy, tj. $p=1 \times X_{ij}=1, i=1, \dots, n$, byly navrženy a podrobně studovány v řadě prací v zásadě 3 alternativy vzhledem k (1.2), totiž: M-odhad (odhad maximálně věrohodného typu), L-odhad (odhad založené na lineárních kombinacích pořádkových statistik) a R-odhad (odhad založené na pořadí pozorování). Zájemci o praktické použití patrně najdou nejlepší poučení v knize Andrenwall [2], která shrnuje nejenom četné teoretické vlastnosti odhadů vhodných pro tento model, ale především výsledky rozsáhlých simulací a z nich vyplývající hodnocení a porovnání 62 různých odhadů výše uvedených tříd. Pro praktika, jakož i osoby nezabývající se programováním, je více než přijemným faktem to, že v knize jsou uvedeny programy pro výpočet všech 62 uvažovaných odhadů, napsané ve Fortranu IV., jež byly optimalizovány jak z hlediska spotřeby strojového času, tak nároků na paměť.

Pro ty, jež nechtějí ztrácat mnoho času hledáním v literatuře, připomeneme, že jsou to 10% usazený průměr, tj. průměr počítaný z pozorování, jež zbudou po odstranění cca 10% minimálních a 10% maximálních pozorování výběru, spolu s M-odhadu značenými v [2] H12, resp. H15, jež lze doporučit pro převážnou většinu praktických situací.

Vedle těchto tříd tzv. robustních odhadů je třeba se ještě v krátkosti zmínit i o existenci dalších typů odhadů, které za jistých podmínek mohou dát neméně dobré výsledky jako odhady předchozí. Jedná se především o tzv. SA-odhady, zavedené Martinem v [36] a adaptivní odhady, viz např. [18], [22] či [23].

Vypracováno během studijního pobytu v Laboratoire de Statistique et Probabilités, Université Paul Sabatier, Toulouse.

3. Postupy robustní vzhledem k rozdělení chyb

Společným rysem všech metod robustních vzhledem k rozdělení chyb (distributional robustness) je omezit možný vliv velkých reziduí

$$(3.1) \quad \delta_i(\beta) = Y_i - \sum_{j=1}^p X_{ij}\beta_j, \quad i=1, \dots, n,$$

odpovídajících většinou odlehlym pozorováním či chybám v datech. Zpravidla se tak děje volbou funkce méně citlivé na extrémní hodnoty reziduí než je tomu v (1.2)

V následujících paragrafech se zaměříme na některé způsoby výpočtu robustních odhadů v lineárním modelu, a to především na R, L a M odhad, včetně krátké zmínky o SA-odhadech a adaptivních odhadech.

3.1. R - odhady

R-odhad vektoru parametrů β dostaneme, zhruba řečeno, jako řešení úlohy minimalizovat

$$(3.2) \quad \sum_{i=1}^n a_n(R_i) \delta_i(\beta)$$

vzhledem k $\beta = (\beta_1, \dots, \beta_p)$, kde R_i je pořadí i-tého rezidua $\delta_i(\beta)$ v $(\delta_1(\beta), \dots, \delta_n(\beta))$ a $a_n(\cdot)$ je některá monotóní skórová funkce [většinou pro jednoduchost normovaná tak, aby $\sum_{i=1}^n a_n(i) = 0$]. Diferencujeme-li v (3.2) podle β_j , $j=1, \dots, p$, lze převést řešení minimalizační úlohy (3.2) na řešení soustavy přibližných rovnic

$$(3.3) \quad \sum_{i=1}^n a_n(R_i) X_{ij} \approx 0, \quad j=1, \dots, p.$$

Tato varianta byla zkoumána např. Jaceckem v [25].

Jurečková v [26] naopak místo řešení (3.3) uvažovala minimalizaci výrazu

$$(3.4) \quad \sum_{j=1}^p \left| \sum_{i=1}^n a_n(R_i) X_{ij} \right| .$$

Lze ukázat, že oba výše uvedené přístupy jsou (za jistých podmínek regularity) asymptoticky ekvivalentní a mají mnoho dobrých asymptotických vlastností. Bohužel, ani pro jeden z nich se mi nepodařilo nalézt v literatuře podrobně rozebraný fungující algoritmus pro jejich výpočet v obecném lineárním modelu. Metody numerické matematiky jsou však stále větší, což může v dohledné době přivodit zvrát i v této oblasti.

Jednokrokovým R-odhadem se v [33] zabývali Kraft a van Eeden. Jejich idea spočívá ve snaze vylepšit některý dobrý počáteční odhad (reasonable preliminary estimator) dodatečným členem počítaným za použití pořadí reziduí. Přesněji řečeno:

(a) Nechť skórová funkce $a_n(\cdot)$ je generována nekonstantní a neklesající funkcí $\Psi(t) \in L_2(0,1)$ tak, že $a_n(i) = \Psi(\frac{i}{n+1})$. Označme $K_\Psi = \int_0^1 \Psi(u) du$.

(b) Nechť $\hat{\beta}^1$ je některý počáteční odhad β a $\hat{\sigma}$ je některý konzistentní odhad parametru měřítka chyb σ . Autoři doporučují vzít za $\hat{\beta}^1$ odhad získaný metodou nejmenších čtverců či LAE odhad (tj. least absolute error estimator).

Potom odhad $\hat{\beta}$ vektoru parametrů β je definován vztahem

$$(3.5) \quad \hat{\beta} = \hat{\beta}^1 + \frac{\hat{\sigma}}{K_\Psi} (X' X)^{-1} \cdot X' \cdot R(\hat{\beta}^1),$$

kde $R(\hat{\beta}^1) = \left\{ \Psi\left(\frac{(Y-X\hat{\beta}^1)_i}{n+1}\right) \right\}_{i=1}^n$ značí $(nx1)$ rozměrný vektor a $R = (Y-X\hat{\beta}^1)$ je pořadí i-té komponenty reziduí počítaných pro počáteční odhad $\hat{\beta}^1$.

Poznámka 3.1. Základní nevýhodou tohoto přístupu je požadavek některého dostatečně dobrého počátečního odhadu. Autoři článku uvažují především klasický odhad metodou nejmenších čtverců. Jsou-li počáteční data "slušná", je i počáteční odhad

metodou nejmenších čtverců slušný a je sporné, zda "případné nevelké vylepšení" tohoto odhadu odpovídá námaze s jeho výpočtem. Obsahuje-li však počáteční data hrubé chyby, a nemusí jich být tak mnoho, je odhad získaný metodou nejmenších čtverců zpravidla špatný a lze jen těžko věřit, že jej lze výrazně vylepšit jednou opravou. Na základě osobních zkušeností s M-odhady tomu věřit nelze. Autori sice v poznámce pod čarou hovoří o možnosti opakování použití vztahu (3.5), ale bez důkazu o konvergenci takovéto procedury, nemluvě o stanovení stop pravidla.

Poznámka 3.2. Kraft a van Eeden v [33] pracují se známým rozdělením F vektoru chyb a vzhledem k němu volí optimální funkci Ψ generující skóry. Je-li však F neznámá, je třeba doplnit proceduru o odhad funkce Ψ [resp. F], jak je tomu u tzv. adaptivních procedur, viz např. [18], [22] či [23]. Adaptivní postupy mají sice výborné asymptotické vlastnosti, ale vzhledem k jejich pomalé konvergenci je třeba mít k dispozici značně rozsáhlé data pro dosažení dostatečné přesnosti, což poněkud omezuje možnosti jejich nasazení. Blíže se lze o těchto odhadech dočít v příspěvku M. Huškové v tomto sborníku.

3.2. L-odhad

Lineární kombinace pořádkových statistik byly dlouhou dobu středem zájmu jako odhady parametru polohy či měřítka. Důvody jsou nasnadě - snadno se s nimi počítá, jsou dostatečně asymptoticky vydatné a mnohé z nich, např. α -useknutý průměr, se ukázaly být robustní vzhledem k mnoha různým navrženým kritériím. Jejich zobecnění pro lineární model však nebylo zdaleko tak přímočáre jako u M a R odhadů. Naopak, po dlouhou dobu to byl pouze návrh Bickela [8], který se touto problematikou zabýval. Teprve pojem tzv. regresních kvantilů, zavedených Koenkerem a Bassetem v [31], zlomil předchozí bariéru a dovolil definovat L-odhad parametrů v lineárním modelu, které díky snadnosti výpočtu a velmi dobrým teoretickým vlastnostem směle konkuruji ostatním robustním odhadům a často je předčí.

Zmiňme se nejprve krátce o Bickelově návrhu, který je částečně analogický návrhu Krafta a van Eednové. Vychází se opět z některého počátečního odhadu $\hat{\beta}^*$ a nový odhad $\hat{\beta}$ je počítán ze vztahu

$$(3.6) \quad \hat{\beta} = \hat{\beta}^* + (\mathbf{X}'\mathbf{X})^{-1} \cdot \mathcal{L}(\hat{\beta}^*),$$

kde $\mathcal{L}(\hat{\beta}^*)$ je oprava počítaná z uspořádaných reziduí pro $\hat{\beta}^*$. Vzhledem k tomu, že přesný tvar této opravy je dosti složitý a potřebuje dlouhé zavedení, bylo od něj zde upuštěno a zájemci se odkazují na [8]. Pro praktické použití Bickelova postupu nelze než zopakovat poznámku 3.1.

O mnoho zajímavější a schůdnější pro praxi se jeví definice L-odhadů v regresním modelu založená na pojmu tzv. regresních kvantilů, zavedených Koenkerem a Bassetem v [31]. Vzhledem k modelu (1.1) lze α -tý regresní kvantil, $0 < \alpha < 1$, definovat jako řešení úlohy minimalizovat, vzhledem k $\hat{\beta} = (\beta_1, \dots, \beta_p)$

$$(3.7) \quad \sum_{i=1}^n \mathbb{P}_{\alpha}^*(Y_i - \sum_{j=1}^p X_{ij} \beta_j),$$

kde

$$(3.8) \quad \mathbb{P}_{\alpha}(x) = x \cdot (\alpha - I_{[x < 0]}), \quad x \in \mathbb{R}_1.$$

Nevychází se zde přitom z obvyklé definice výběrového kvantilu založeného na pořádkových statistikách, nýbrž z definice následující, viz [16]. Je-li X_1, \dots, X_n náhodný výběr rozdělený podle téže distribuční funkce F , potom výběrový α -kvantil, $0 < \alpha < 1$, budež definován jako řešení minimalizačního problému

$$(3.9) \quad \min_{\beta \in \mathbb{R}^p} \left\{ \sum_{\{i | X_i > \beta\}} \alpha \cdot |X_i - \beta| + \sum_{\{i | X_i < \beta\}} (1 - \alpha) \cdot |X_i - \beta| \right\}.$$

Prohlédneme-li si vzorce (3.7)-(3.9) pečlivěji, vidíme, že se vlastně nejedná o nic jiného, než o jistý druh minimalizace v L_1 normě. Přesněji řečeno, řešením

je vektor minimalizující vážený součet absolutních hodnot residuí, což je problém, o jehož řešení je známo poměrně mnoho. Tato metoda, tj. hledání řešení minimalizujícího L_1 -normu, byla již v "dávné statistické minulosti" doporučována v některých učebnicích na místo metody nejmenších čtverců pro případ, obáváme-li se výskytu divokých pozorování. Podívejme se nyní na 2 základní typy L -odhadů v našem modelu, tož α -useknutý a α -winsorizovaný odhad metodou nejmenších čtverců. O regresním mediánu, tj. $1/2$ -regresním kvantilu, se není třeba blíže zmiňovat, neboť se shoduje s řešením v L_1 -normě, tj. jedná se o LAE (least absolute error) odhad.

$L_n(\alpha)$: α -useknutý odhad metodou nejmenších čtverců

Nechť $0 < \alpha < \frac{1}{2}$ je předem zvolená konstanta odpovídající procentu useknutí.

1⁺). Nechť $\hat{\beta}(\alpha) = (\hat{\beta}_1(\alpha), \dots, \hat{\beta}_p(\alpha))'$ je řešení (3.7) pro α , tj. $\hat{\beta}(\alpha)$ je α -tý regresní kvantil.

2⁺) Nechť $\hat{\beta}(1-\alpha) = (\hat{\beta}_1(1-\alpha), \dots, \hat{\beta}_p(1-\alpha))'$ je řešení (3.7) pro $(1-\alpha)$, tj. $\hat{\beta}(1-\alpha)$ je $(1-\alpha)$ -tý regresní kvantil.

3⁺) Odstraňme z původního výběru všechna pozorování taková, pro něž

$$Y_i - \sum_{j=1}^p X_{ij} \hat{\beta}_j(\alpha) < 0 \quad \text{nebo} \quad Y_i - \sum_{j=1}^p X_{ij} \hat{\beta}_j(1-\alpha) > 0, \quad i=1, \dots, n.$$

4⁺) α -useknutý odhad metodou nejmenších čtverců $L_n(\alpha)$ pak obdržíme jako odhad metodou nejmenších čtverců ze zbyvajících pozorování.

$W_n(\alpha)$: α -winsorizovaný odhad metodou nejmenších čtverců

Nechť $0 < \alpha < \frac{1}{2}$ je předem zvolená konstanta odpovídající procentu winsorizace.

Nechť $\hat{\beta}(\alpha)$, $\hat{\beta}(1-\alpha)$ a $L_n(\alpha)$ jsou postupně α -tý regresní kvantil, $(1-\alpha)$ -tý regresní kvantil a α -useknutý odhad metodou nejmenších čtverců počítané podle 1⁺) - 4⁺).

5⁺) α -winsorizovaný odhad metodou nejmenších čtverců $W_n(\alpha)$ vypočteme ze vztahu

$$(3.10) \quad W_n(\alpha) = n^{-1} \cdot \left\{ [n\alpha] \cdot (\hat{\beta}(\alpha) + \hat{\beta}(1-\alpha)) + (n-2[n\alpha]) \cdot L_n(\alpha) \right\}.$$

Poznámka 3.3. : O některých zkušenostech s témito odhady se lze dočíst např. v [3] či [5]. Naše zkušenosti s nimi jsou více než dobré a lze je, zvláště α -useknutý odhad metodou nejmenších čtverců, jenom doporučit.

Domníváme-li se, že naše data jsou "celkem čistá", jak tomu zpravidla v praxi bývá, pak volbou $\alpha \approx 0.05-0.1$ je chránění mnohem lépe než při přímém použití metody nejmenších čtverců před případným výskytom několika nahodilých chyb. Obáváme-li se výskytu mnoha hrubých chyb, a máme-li k dispozici dost pozorování, pak lze doporučit zvýšení α až do hranice cca 0.2. Obsahují-li data i potom hrubé chyby, je nejlepší je zahodit a provést nová měření, nebo co jiného dělat s daty, v nichž je více než 40% "podivných pozorování". Pokud si takový "luxus" nemůžeme dovolit, použijme regresní medián, tj. LAE odhad.

Programy ve Fortranu IV. spolu s návodem na jejich použití lze obdržet na KPMS MFF UK od autora tohoto příspěvku.

Poznámka 3.4. O odhadech parametrů a jejich vlastnostech je známo poměrně mnoho. Mnohem méně prací však bylo věnováno problémům testování hypotéz o odhadovaných parametrech, tolik potřebných při tvorbě modelu etc. Ve speciálním případě pro model polynomické regrese je k dispozici práce J. Jurečkové [28]. Pro obecný lineární model pak lze vyjít např. z [32] či [21]. Nicméně lze konstatovat, že na tomto poli stále zůstává dostatek místa pro další iniciativu.

3.3. M - odhady

Podíváme-li se do literatury, vidíme, že jsou to M-odhady, jimiž byla věnována největší pozornost. Pro některá shrnutí viz např. [21]. Podívejme se nyní krátce, co nám nabízejí.

Pro omezení vlivu odlehých pozorování bylo navrženo místo (1.2) hledat spíše

$$(3.11) \min_{\beta} \sum_{i=1}^n \zeta \left(Y_i - \sum_{j=1}^p X_{ij} \beta_j \right)$$

vzhledem k $\beta = (\beta_1, \dots, \beta_p)$. Funkce ζ by přitom měla omezit vliv velkých residuí a měla by tudíž být funkcí, jež konverguje do nekonečna pomaleji než kvadratická funkce. U volbě ζ optimální vzhledem k apriorní informaci o rozdělení chyb se lze dočíst např. v [42] či [21]. Naneštěstí však takováto procedura není obecně invariantní vzhledem ke změně měřítka a tento parametr je třeba odhadnout současně s β , viz [21]. Tedy, místo (3.11) musíme hledat

$$(3.12) \min_{\beta} \sum_{i=1}^n \zeta \left(\frac{Y_i - \sum_{j=1}^p X_{ij} \beta_j}{\sigma} \right),$$

např. za podmínky

$$(3.13) \sum_{i=1}^n \chi \left(\frac{Y_i - \sum_{j=1}^p X_{ij} \beta_j}{\sigma} \right) = (n-p) \frac{\gamma}{2}$$

kde

$\chi(t) = t \Psi(t) - \Phi(t)$, $\Psi = \Phi'$, $\gamma = 2 \int x \chi(x) d\Phi(x)$, $\Phi(x)$ je distribuční funkce $N(0,1)$. Důvod k zavedení podmínky (3.13) je snaha získat asymptoticky nestraný odhad σ v případě, jsou-li chyby normálně rozděleny, viz např. [21].

Předchozí úloha se dá převést na hledání minima výrazu

$$(3.14) \sum_{i=1}^n \zeta \left(\frac{Y_i - \sum_{j=1}^p X_{ij} \beta_j}{\sigma} \right) \cdot \sigma + a \sigma$$

vzhledem k σ a β . Ma-li ζ spojitou derivaci Ψ a zvolíme-li $a = (n-p) \frac{\gamma}{2}$, lze se snadno přesvědčit, že minimalizace (3.14) vzhledem k σ a β je ekvivalentní minimalizaci (3.12) vzhledem k β , přičemž (3.13) je spiněna.

Uvedme nyní tzv. jednoduchý algoritmus pro výpočet robustních odhadů, tak jej nazývají jeho autoři Huber a Dutter v [20]. Pro úplnost dodejme, že jeho konvergence byla dokázána za následujících podmínek na funkci $\zeta(x)$: $\zeta(x)$ je 2x differencovatelná a konkavní symetrická funkce, $\zeta'(t) \geq 0$, $\zeta'(0) = 0$, $\zeta''(t)/t$ je konkavní pro $t < 0$ a konkávní pro $t > 0$, $\lim_{t \rightarrow 0} \zeta''(t)/t < \infty$, $0 \leq \zeta'' \leq 1$.

Jednoduchý algoritmus pro M-odhady

Nechť $\hat{\beta}^{(0)}$ a $\hat{\sigma}^{(0)}$ jsou některé počáteční odhady β a σ .

1^o) Nechť $m=0$, kde m je konstanta označující počet iterací.

2^c) Spočteme rezidua $\delta_i^{(m)} = Y_i - \sum_{j=1}^p X_{ij} \hat{\beta}_j^{(m)}$, $i=1, \dots, n$.

3^c) Spočteme nový odhad $\hat{\sigma}^{(m+1)}$ z výrazu

$$(\hat{\sigma}^{(m+1)})^2 = \frac{1}{n} \sum_{i=1}^n \chi \left(\frac{\delta_i^{(m)}}{\hat{\sigma}^{(m)}} \right) \cdot (\hat{\sigma}^{(m)})^2.$$

4^c) Spočteme hodnoty $\Delta_i^{(m)}$ tzv. winsorizovaných residuí, tj.

$$\Delta_i^{(m)} = \Psi \left(\frac{\delta_i^{(m)}}{\hat{\sigma}^{(m+1)}} \right) \cdot \hat{\sigma}^{(m+1)}, \quad i=1, \dots, n.$$

5^c) Spočteme minimum, vzhledem k $\hat{\tau}$, výrazu

$$\sum_{i=1}^n (\Delta_i^{(m)})^2 - \sum_{j=1}^p X_{ij} \hat{\tau}_j)^2$$

a označme jej $\hat{\tau}^{(m)}$.

6^o) Nový odhad $\hat{\beta}^{(m+1)}$ dostaneme jako

$$\hat{\beta}^{(m+1)} = \hat{\beta}^{(m)} + \hat{\zeta}^{(m)} \cdot K,$$

kde $0 < K < 2$ je některá předem zvolená konstanta pomáhající nám zrychlit (zpomalit) proceduru, blíže viz [10].

7^o) Jestliže $\|\hat{\zeta}\| < \varepsilon_0$, kde ε_0 je některá předem předepsaná tolerance, ukončeme proceduru a položme $\hat{\beta} = \hat{\beta}^{(m+1)}$ a $\hat{\sigma} = \hat{\sigma}^{(m+1)}$. V opačném případě položme $m=m+1$ a vraťme se do bodu 2^o.

Podívejme se ještě na 2 modifikace předchozího algoritmu.

Modifikace I. : spořívá v přehození bodů 3^o a 4^o a poněkud odlišném způsobu odhadu σ , tj.

$$3^o) \quad \Delta_i^{(m+1)} = \psi\left(\frac{\delta_i^{(m)}}{\hat{\sigma}^{(m)}}\right) \cdot \hat{\sigma}^{(m)}, \quad i=1, \dots, n,$$

$$4^o) \quad (\hat{\sigma}^{(m+1)})^2 = \frac{1}{2a} \sum_{i=1}^n (\Delta_i^{(m+1)})^2.$$

Ostatní kroky jsou shodné s přehození procedurou. Soudě podle výsledků v [11], nedává nám tato modifikace odhady zásadně odlišné od metody předchozí. Naopak lze říci, že odhady počítané z týchž dat jsou si zcela rovnocenné.

Modifikace II. : vychází také z předchozího algoritmu, ale využívá více vlastností lineárního modelu.

Předpokládejme, že matice X má plnou hodnost, potom řešení bodu 5^o lze psát ve tvaru

$$(3.15) \quad \hat{\zeta}^{(m)} = (X'X)^{-1} X' \Delta^{(m)},$$

kde $\Delta^{(m)}$ je vektor winsorizovaných reziduí z bodu 4^o. Pro další iteracní krok potřebujeme vektor reziduí $\delta^{(m+1)}$ vzhledem $\hat{\beta}^{(m+1)} = \hat{\beta}^{(m)} + K \hat{\zeta}^{(m)}$, který lze ale počítat přímo z $\delta^{(m)}$ a $\Delta^{(m)}$, neboť

$$(3.16) \quad \delta^{(m+1)} = Y - X \hat{\beta}^{(m+1)} = Y - X \hat{\beta}^{(m)} - X \hat{\zeta}^{(m)} \Rightarrow \delta^{(m+1)} = Y - X(X'X)^{-1} X' \Delta^{(m)}.$$

Spočteme-li si tedy na počátku prjekční matici $X(X'X)^{-1} X'$ a schováme její horní trojuhelníkovou část, lze využít předchozí algoritmus, aniž bychom explicitně v každém kroku počítali opravny vektor $\hat{\zeta}^{(m)}$. Proceduru zastavíme tehdy, je-li $\|X(X'X)^{-1} X' \Delta^{(m)}\| < \varepsilon_0$, kde ε_0 je předem zadaná tolerance (různá od ε_0). Výsledný odhad vektoru parametrů $\hat{\beta}$ pak má tvar

$$(3.17) \quad \hat{\beta} = (X'X)^{-1} X' (Y - \delta^{(m+1)}),$$

tj. $\hat{\beta} = \hat{\beta}^{(m+1)}$.

Poznámka 3.5.: Předložený algoritmus (a jeho varianty) není jediný způsob pro hledání M-odhadů, viz např. [10]. Prostudujeme-li si však většinu ostatních algoritmů pečlivěji, zjistíme, že se jedná více méně o variace na předchozí téma, zpravidla využívají apriorní informaci až již o funkci Φ , či tvaru chyb ap. Vyjimkou jsou metody využívající technik vězených nejmenších čtverců, viz [41], [15] či [24].

Poznámka 3.6.: Připomeňme, že v MS BÚ ČSAV Praha je k dispozici rozsáhlá knihovna podprogramů pro výpočet M-odhadů, viz [17]. Od vydání sborníku ROBUST 80 byla mezi tím p. Malinou napsána řada programů umožňujících zájemcům o praktické vyzkoušení pouze přijít, naděrovat data v předepsaném tvaru a odeslat k výpočtu. Více snad již ani nabídnout nelze - jen to děrování, ale na to zase nejsou lidí.

4.4. SA - odhady

Odstavec 3 by nebyl úplný, kdybychom se vedle R, L a M odhadů nezmínili o odhadech založených na algoritmu stochastických approximací (SA) Robins-Monroova (RM) typu, jež jsou dány jednoduchým rekurzivním vztahem a jsou vhodné pro aplikace v reálném čase.

Nechť $\hat{\beta}^{(0)}$ je některý počáteční odhad a $\gamma > 0$ konstanta. Potom odhad $\hat{\beta}_{SA}$ vektoru $\hat{\beta}$ v modelu (1.1), definovaný rekurzivním vztahem

$$(3.18) \quad \hat{\beta}^{(n)} = \hat{\beta}^{(n-1)} + \frac{\gamma}{n} (\mathbf{x}' \mathbf{x})^{-1} \mathbf{x}' \Psi(\mathbf{y} - \mathbf{x}' \hat{\beta}^{(n-1)}),$$

$n=1, 2, \dots$ nazveme SA-odhadem příslušným k funkci Ψ . Funkce Ψ se volí zpravidla též jako u M-odhadů. Blížší informaci o těchto odhadech lze nalézt např. v [36], [37] či [38].

Odhady tohoto typu se v ČSSR zabývá např. dr.J.Novovičová z ÚTIA ČSAV Praha se svými spolupracovníky. Tento tým nejenže odvodil řadu zajímavých teoretických výsledků o odhadech typu SA-RM, ale připravil řadu programů pro jejich výpočet. Zájemce o případné praktické využití proto odkazujeme na výše uvedené pracovníky.

4.5. Závěr

Jak už jsem se zmínil v úvodu a snažil se dokumentovat v textu, četnost a kvalita prací věnovaných výpočetní stránce robustních metod není zdaleka ve shodě s ohromným množstvím prací teoretických. A to jak v oblasti algoritmů, tak ještě ve větší míře v oblasti skutečných realizací.

Nicméně nelze říci, že by na tomto poli nic neexistovalo. Naopak, možnosti existují nejenom ve světě, ale i u nás. K dispozici jsou poměrně velmi dobře fungující knihovny programů na výpočet jak L, tak M odhadů. Tyto programy umožňují experimentovat nejenom na poli statistiky - při zkoumání chování odhadů na datech umělých i z praxe, ale i na poli numeriky a programování - vylepšování programů, jejich optimalizace, hledání nových algoritmů etc.

Zbývá však ještě krok třetí, nejtěžší. Je jím dovenení robustních metod mezi běžné uživatele do každodenní praxe, neboť teprve potom bude všechno to nezměrné úsilí náležitě zhodnoceno. Na tomto poli jsme zatím postoupili nejméně, a to nejenom u nás, ale i ve světě. I zde se však vlastně rozjíždí a akce typu RCBUST se snaží jej co nejvíce zrychlit.

Literatura

- [1] Adelmalek N.N., On the discrete linear l_1 -approximation and l_1 solutions of overdetermined linear equations. Journal of Approximation theory 11(1974), 38-53.
- [2] Andrews D.F., Bickel P.J., Hampel F.R., Huber P.J., Tukey J.W., Robust estimates of location. Survey and advances. Princeton University Press, 1972.
- [3] Antoch J., Collomb G., Hassani S., Robustness in parametric and non-parametric regression estimation: An investigation by computer simulations. Proceedings of COMPSTAT 1984, Physica Verlag, Viena, 1984.
- [4] Antoch J., Soubor programů pro robustní regresní analýzu. Publ.de Laboratoire de Statistique et Probabilités, Toulouse, 1984.
- [5] Antoch J., L -estimates in linear regression model. Publ.de Laboratoire de Statistique et Probabilités, Toulouse, 1984.
- [6] Barrodale I., Roberts F.D.K., An improved algorithm for discrete l_1 -linear approximation. SIAM J.Numer.Anal. 10,5 (1973), 839-848.
- [7] Barrodale I., Roberts F.D.K., Solution of an overdetermined system of equations in the l_1 -norm. CACM 17 (1974), 319-320.
- [8] Bickel P.J., On some analogues to linear combinations of order statistics in the linear model. AS 1 (1973), 597-616.

- [7] Dione L., Efficient nonparametric estimators of parameters in the general linear hypothesis. AS 4 (1970).
- [10] Dutter R., Robust regression: Different approaches to numerical solutions and algorithms. Research report N° 6, ETH Zürich 1975.
- [11] Dutter R., Numerical solution of robust regression problems: Computational aspects and comparison. Research report N° 7, ETH Zürich, 1975.
- [12] Dutter R., Slinwdr: Robust and bounded influence regression. Proc.of COMPSTAT 1982, Caussinus & al.(eds.), Physica-Verlag, Vienna, 1982
- [13] Dutter R.; Computer program BLINWDR for robust and bounded influence regression. Res.Rep.No 6 (Version 2). Inst.für Statistik, Technische Universität Graz, 1983.
- [14] Dutter R., Covinter : A computer program for computing robust covariances and for plotting tolerance ellipses. Res.Rep. No.10, idem.
- [15] Dutter R., Statistical computing: robust regression. Handout at 6th International Summer School on Problems of Model Choise and Parameter Estimation in Regression Analysis, Sellin 1983, NDR.
- [16] Ferguson T.S., Mathematical statistics. A decision theoretic approach. New York, Academic Press, 1967.
- [17] Havránek T. - Antoch J., ROBETH-knihovna programů pro robustní statistické metody. Sborník ROBUST 80, 25-33, JČSMF 1980.
- [18] Hogg R.V., Adaptive robust procedures. JASA 69(1974), 909-923.
- [19] Huber P.J., Robust estimation of regression parameter. AMS 35(1964), 73-101.
- [20] Huber P.J., Dutter R., Numerical solution of robust regression problems. Proc.of COMPSTAT 1974, Physica Verlag, Vienna ,1974.
- [21] Huber P.J., Robust Statistics. J.Wiley & Sons, Inc., 1981.
- [22] Hušková M., Adaptive procedures. Acta Univ.Carolinæ-Math. et Ph.24(1983), 41-48.
- [23] Hušková M., Adaptivní odhady. Sborník ROBUST 84, JČSMF 1984.
- [24] Chambers J.M., Computational methods of data analysis, J.Wiley, 1977.
- [25] Jaeckel L.A., Estimating regression coefficients by minimizing the dispersion of the residuals. AMS 43 (1972), 1449-1458.
- [26] Jurečková J., Nonparametric estimate of regression coefficients. AMS 42 (1971), 1328-1338.
- [27] Jurečková J., Winsorized least-squares estimator and its M-counterpart. Contributions to Statistics. Essays in Honour of N.L.J.Johnson, P.K.Sen(ed.), 237-245, North Holland, 1983.
- [28] Jurečková J., Trimmed polynomial regression. CMUC 24,4(1983), 597-607.
- [29] Jurečková J., Robust estimators of location and regression parameters and their second order asymptotic relations. Trans. 9th Prague Conf.on Inf. Theory etc.,
- [30] Jurečková J., Regression quantiles and trimmed least-squares estimator under a general design. Subm.to Kybernetika, 1984.
- [31] Koenker R., Bassett G., Regression quantiles. Econometrica 46(1978), 33-51.
- [32] Koenker R., Bassett G., Tests of linear hypotheses and l_1 -estimation. Econometrica 50 (1982), 1577-1583.
- [33] Kraft C., van Eeden C., Linearized rank estimates and sign-rank estimates based on ranks in regression. AMS 42 (1972), 42-57.
- [34] Lawson C.L., Hanson R.J., Solving least squares problems. Prentice Hall, Inc., Englewood Cliffs, New Jersey, 1974.
- [35] Marazzi A., Robeth: A subroutine library for robust statistical procedures. Proc.of COMPSTAT 1980, Barriett et al.(eds.), Physica-Verlag, Vienna, 1980.
- [36] Martin R.D., Robust estimation of signal amplitude. IEEE Trans.Inform.Th. II-18, 596-606.
- [37] Novovičová J., M-odhadý a SA-odhadý v lineárním regresním modelu. Sborník ROBUST 80, 119-133, JČSMF, 1980.
- [38] Novovičová J., SA-odhadý regrese. Kandidátská disertační práce, ÚTIA ČSAV Praha, 1983.

3 Ruppert D., Carroll R.J., Robust regression by trimmed least-squares estimation. Inst. of Statist. Nimeo Series No 1186, Univ. of North Carolina at Chapel Hill, 1978.

4 Ruppert D., Carroll R.J., Trimmed least squares estimation in the linear model. JASA 75 (1980), 828-838.

41 Samarov A., Welsch R.E., Computational procedures for bounded-influence regression. Proc. of COMPSTAT 1982, Caussinus et all.(eds.), Physica Verlag, Vienna, 1982.

42 Чепкин Я.З., Поляк Б.Т., Огрубленный метод максимального правдоподобия. Сб. "Динамика систем", 12, 22-46.