

BOD SELHÁNÍ NEKTERÝCH ODHADŮ PARAMETRU POLOHY

Jan Hanousek, MFF UK Praha

1. ÚVOD

Většina výsledků je formulována pro Lévyho okolí distribuční funkce F_0

$$P_\epsilon(F_0) = \left[F | (\forall t) F_0(t-\epsilon) - \epsilon \leq F(t) \leq F_0(t+\epsilon) + \epsilon \right].$$

Tento se tato okolí vyjadřuje pomocí Lévyho metriky d_F

$$d_F(F, G) = \inf \left[\epsilon | (\forall t) F(t-\epsilon) - \epsilon \leq G(t) \leq F(t+\epsilon) + \epsilon \right].$$

Potom

$$P_\epsilon(F_0) = \left[F | d_F(F, F_0) \leq \epsilon \right].$$

Snadno odvodíme, že stochasticky největším prvkem okolí $P_\epsilon(F_0)$ je funkce F_1 , tvaru

$$F_1(x) = \begin{cases} 0 & x \leq x_0 + \epsilon \\ F_0(x-\epsilon) - \epsilon & x > x_0 + \epsilon \end{cases}$$

kde $x_0 = F_0^{-1}(\epsilon)$. Analogicky najdeme stochasticky nejménší funkci F_2

$$F_2(x) = \min \left[F_0(x+\epsilon) + \epsilon; 1 \right] = \begin{cases} F_0(x+\epsilon) + \epsilon & x \leq x_0 - \epsilon \\ 1 & x > x_0 - \epsilon \end{cases}$$

kde $x_0 = F_0^{-1}(1-\epsilon)$.

(Připomínáme si, že F je stochasticky větší než G , pokud pro všechna x platí, že $F(x) \leq G(x)$).

Vzhledem k častému používání uvedeme ještě tvar $F_1^{-1} \circ F_2^{-1}$:

Zde jako obvykle

$$F^{-1}(y) = \sup \left[x : F(x) \leq y \right].$$

Tedy

$$F_1^{-1}(x) = \begin{cases} \epsilon + F_0^{-1} & 0 \leq x \leq 1-\epsilon \\ +\infty & 1-\epsilon < x \leq 1 \end{cases}$$

$$F_2^{-1}(x) = \begin{cases} -\epsilon + F_0^{-1}(x-\epsilon) & \epsilon \leq x \leq 1 \\ -\infty & 0 \leq x < \epsilon \end{cases}$$

2. DEFINICE A PŘÍKLADY

2.1. Definice bodu selhání

Bodem selhání odhadu T pro rozdělení F_0 nazveme

$$\sigma^* = \sigma^*(F_0, T) = \sup \left[\epsilon | b(\epsilon) < b(1) \right],$$

kde

$$b(\epsilon) = \lim_{n \rightarrow \infty} \sup_{F_n} \left| E(F, T_n) \right|.$$

$M(F, T_n)$ je medián rozdělení $L_F[T_n - T(F_0)]$.
Puto obecnou definici budeme modifikovat na naše případy.

Zavedeme veličiny

$$\begin{aligned} b^+(\epsilon) &= \sup [T(F) \mid d_F(F, F_0) \leq \epsilon], \\ b^-(\epsilon) &= \inf [T(F) \mid d_F(F, F_0) \geq \epsilon] \\ b_1(\epsilon) &= \max [b^+(\epsilon); -b^-(\epsilon)]. \end{aligned}$$

a využijeme vět 1.4.1 a 1.4.2 Huber (1981), ze kterých plyne, že $b_1(\epsilon) = b(\epsilon)$ v bodech spojnosti b_1 .

Další úvahy budou patřit vyběrovému bodu selhání.

Nechť $X = (x_1, \dots, x_n)$ je libovolný soubor o pevném rozsahu Δ . Soubor zahrnuje znehodnotit několika způsoby:

a) ϵ -znečištění

Připojime m libovolných hodnot $Y = (y_1, \dots, y_m)$ k souboru. Znehodnocený soubor $X' = X \cup Y$ je třídy $n+m$ a relativní obsah špatných hodnot je roven $\epsilon = \frac{m}{n+m}$.

b) ϵ -vrácení

Do souboru znova zahrime libovolnou podmnožinu m prvků základního souboru s hodnotami y_1, \dots, y_m . Znehodnocený soubor je třídy n a relativní obsah špatných hodnot je $\epsilon = \frac{m}{n}$.

c) ϵ -modifikace

Nechť $p(\cdot)$ je libovolná metrika definovaná na prostoru empirických měr. Nechť F_n je empirická míra odpovídající danému souboru X . Nechť X je nějaký jiný soubor s empirickou mírou G_n takovou, že $p(F_n, G_n) \leq \epsilon$.

V našem případě je nejhodnější užívat případ a). Nyní definujme vyběrový bod selhání odhadu T .

Nechť je dán soubor $X = (x_1, \dots, x_n)$ a k němu přísluší ϵ -znečištěný soubor X' . Maximální spád odhadu T při ϵ -znečištění je

$$b(\epsilon, X, T) = \sup |T(X') - T(X)|,$$

kde supremum se bere přes množinu všech ϵ -znečištěných souborů X' . Potom vyběrový bod selhání odhadu T definujeme jako

$$g(X, T) = \inf [\epsilon \mid b(\epsilon, X, T) = +\infty].$$

Bod selhání je tedy, nepřesně řečeno, nejménší ϵ -znečištění, na kterém může odhad nabýt libovolně velké hodnoty.

2.2. JEDNOROZMÍTRÉ M-ODHADY

VĚTA 2.2

Nechť $p(\cdot)$ je neklesající, ale ne mutně spojitá omezená funkce nabývající kladných i záporných hodnot. Uvažujme M-odhad polohy $T(F)$ definovaný rovností

$$\int p(x - T(F)) dF(x) = 0.$$

Bod selhání odhadu T se definuje vztahy

$$\epsilon^* = \epsilon^*(T, p) = \frac{\gamma}{1-\gamma}, \quad \text{kde}$$

$$\gamma = \min \left[-\frac{p(-\infty)}{p(+\infty)}, \frac{-p(+\infty)}{p(-\infty)} \right]$$

ϵ dosahuje své maximální hodnoty $\epsilon^* = 1/2$ za podmínky $p(-\infty) = -p(+\infty)$.

Důkaz

Označme $\lambda(t, F) = \int p(x-t) dF(x)$. Vážíme si, že $\lambda(\cdot, \cdot)$ je nerestoucí funkce v parametrech t ; v F je tato funkce neklesající, jestliže F je stochasticky větší.

Platí tedy nerovnost

$$\lambda(t, F) \leq \lambda(t, F_1) = \int_R p(x-t) dF_1(x) = \int_{x_0}^{\infty} (x-t+\epsilon) dF_0(x) + \epsilon \cdot p(+\infty),$$
$$x_0 = F_0^{-1}(\epsilon).$$

Z vlastnosti funkce λ a F_1 plyne, že

$$b^+(\epsilon) = \inf [t \mid \lambda(t, F_1) < 0].$$

Je vidět, že

$$b^+(\epsilon) \leq -\infty \Leftrightarrow \lim_{t \rightarrow \infty} \lambda(t, F_1) < 0.$$

$$\lim_{t \rightarrow \infty} \lambda(t, F_1) = (1-\epsilon) \cdot p(-\infty) + \epsilon \cdot p(+\infty)$$

Odhad neselže, pokud

$$(1-\epsilon) \cdot p(-\infty) + \epsilon \cdot p(+\infty) < 0$$

tedy

$$\frac{\epsilon}{1-\epsilon} < -\frac{p(-\infty)}{p(+\infty)}.$$

Obdobně, bereme-li funkci F_2 -stochasticky nejmenší, že

$$b^-(\epsilon) = \sup [t \mid \lambda(t, F_2) > 0]$$

$$b^-(\epsilon) \geq -\infty \Leftrightarrow \lim_{t \rightarrow -\infty} \lambda(t, F_2) < 0.$$

Výpočtem zjistíme, že

$$\lambda(t, F_2) = \int_{-\infty}^{x_0} p(x-t-\epsilon) dF_0(x) + \epsilon \cdot p(-\infty)$$

$$\lim_{t \rightarrow -\infty} \lambda(t, F_2) = (1-\epsilon) \cdot p(+\infty) + \epsilon \cdot p(-\infty).$$

Z podmínky

$$(1-\epsilon) \cdot p(+\infty) + \epsilon \cdot p(-\infty) < 0$$

vyjde

$$\frac{\epsilon}{1-\epsilon} < -\frac{p(+\infty)}{p(-\infty)}.$$

Tím už snadno dojdeme k tvrzení věty 2.2.

2.3. L-ODMODY

VĚTA 2.3

Nechť $M = M^+ - M^-$ je konečná znaménková míra na $(0,1)$.

(Z Jordanova rozkladu dostaneme, že M^+ a M^- jsou konečné, kladné míry na $(0,1)$) a $T(p) = \int p^{-1}(s) dM(s)$. Nechť d je největší reálné číslo, při kterém interval $(0,1-d)$ obsahuje nosič mř M^+ a M^- . Dále nechť v bodech nespojitosti F_0^{-1} má M maleoucí míru. Potom bod selhání $\epsilon^* < d$. Jestliže je M kladná míra, potom $\epsilon^* = d$.

Důkaz

Vzhledem k rozkladu $M = M^+ - M^-$ dostaneme $T = T^+ - T^-$, kde

$$T^+(F) = \int_{\mathbb{R}} F^{-1}(s) dM^+(s)$$

$$T^-(F) = \int_{\mathbb{R}} F^{-1}(s) dM^-(s).$$

Jestliže body selhání T^+ i T^- jsou větší nebo rovny ϵ , pak bod selhání T je zřejmě také větší nebo roven ϵ . Hledejme zatím podmínky, které platí pro bod selhání odhadu $T = \int_{\mathbb{R}} F^{-1}(s) dM(s)$, kde M je konečná, kladná míra na $(0,1)$. Nejprve předpokládejme, že $0 < \epsilon < \epsilon$. Znovu využijeme veličin

$$b^+(\epsilon) = \sup [T(F) \mid d_F(F_0, F) \leq \epsilon]$$

$$b^-(\epsilon) = \inf [T(F) \mid d_F(F_0, F) \leq \epsilon].$$

Z vlastnosti funkcí F_1 a F_2 můžeme spočítat, že

$$b^+(\epsilon) = \int_{F_1^{-1}}(s) dM(s) = \epsilon + \int_{\epsilon}^{1-\epsilon} F_0^{-1}(s+\epsilon) dM(s)$$

$$b^-(\epsilon) = \int_{F_2^{-1}}(s) dM(s) = -\epsilon + \int_{\epsilon}^{1-\epsilon} F_0^{-1}(s-\epsilon) dM(s).$$

Víme, že

$$b_1(\epsilon) = \max [b^+(\epsilon); -b^-(\epsilon)] \leq b^+(\epsilon) - b^-(\epsilon) \rightarrow 0$$

pro $\epsilon \rightarrow 0$ neboť body nespojitosti F_0^{-1} mají při M nulovou míru (funkcionál T je spojitý). Tedy veličina $b_1(\epsilon)$ je konečná při $\epsilon \ll \epsilon$. Odtud plyne, že $\epsilon \geq \epsilon$. Pokud M je konečná kladná míra, je vidět, že bod selhání $\epsilon \ll \epsilon$. Tím je věta 2.3 dokázána.

2.4. R-ODHADY

VĚTA 2.4

Nechť funkce J , sloužící k výpočtu vah u R-odhadů, je neklesající, integrovatelná a nechť J splňuje podmíinku

$$J(1-t) = -J(t), \quad 0 < t < 1.$$

R-odhad $T(F)$ definuje rovností

$$\int J \left[1/2 \left(s+1 - F(2s/F) - F^{-1}(s) \right) \right] ds = 0.$$

Potom bod selhání $\epsilon^*(T, F_0)$ je řešením rovnice

$$\int_{1/2}^{1-\epsilon^*/2} J(s) ds = \int_{1-\epsilon^*/2}^1 J(s) ds.$$

Důkaz

Platí, že funkce

$$\lambda(t, F) = \int J \left[1/2 \left(s+1 - F(2t - F^{-1}(s)) \right) \right] ds$$

je nerestoucí v t a neklesající v F , pokud F je stochasticky větší. Znovu využijeme vlastnosti funkcí F_1 a F_2 , respektive F_1^{-1} a F_2^{-1} . Odtud plyne, že

$$\lambda(t; F_2) \leq \lambda(t; F) \leq \lambda(t; F_1).$$

Spočteme tedy $\lambda(t; F_1) = \lambda(t; F_2)$. Jestliže jsou splněny podmínky

$$0 \leq s \leq 1-\varepsilon \quad \text{a} \quad 2t - F_1^{-1}(s) \geq x_0 + \varepsilon$$

$$\text{takže} \quad 0 \leq s \leq 1-\varepsilon \quad \text{a} \quad s \leq F_0(2(t-\varepsilon) - x_0) - \varepsilon,$$

dostaneme, že

$$F_1[2t - F_1^{-1}(s)] = F_0[2(t-\varepsilon) - F_0^{-1}(s+\varepsilon)] - \varepsilon.$$

Oznáčme-li bod, ve kterém nastává zlom

$$s_0 = [F_0(2(t-\varepsilon) - x_0) - \varepsilon]^+,$$

můžeme vyjádřit

$$\lambda(t; F_1) = \int_0^{s_0} J\left[\frac{1}{2}\left(s + \varepsilon + 1 - F_0(2(t-\varepsilon) - F_0^{-1}(s+\varepsilon))\right)\right] ds + \int_{s_0}^1 J\left(\frac{1}{2}(s+1)\right) ds.$$

Znadno se ukáže, že

$$b^+(\varepsilon) = \inf\left[t \mid \lambda(t; F_1) < 0\right]$$

a že

$$b^+(\varepsilon) < \infty \Leftrightarrow \lim_{t \rightarrow \infty} \lambda(t; F_1) < 0.$$

Limitním přechodem zjistíme

$$\lim_{t \rightarrow \infty} \lambda(t; F_1) = \int_0^{1-\varepsilon} J\left[\frac{1}{2}(s+\varepsilon)\right] ds + \int_{1-\varepsilon}^1 J\left[\frac{1}{2}(s+1)\right] ds.$$

Jednoduchou substitucí a s využitím vlastnosti J máme

$$\lim_{t \rightarrow \infty} \lambda(t; F_1) = 2 \cdot \left(\int_{1-\varepsilon/2}^1 J(s) ds - \int_{1/2}^{1-\varepsilon/2} J(s) ds \right).$$

V tomto případě odhad vidíme, že bude-li

$$\int_{1-\varepsilon/2}^1 J(s) ds = \int_{1/2}^{1-\varepsilon/2} J(s) ds.$$

Podobně budeme postupovat pro stočnosticky nejménší funkci F_2 . Vyjde tedy výsledky

$$\lambda(t; F_2) = \int_{s_0}^1 J\left[\frac{1}{2}\left(s + 1 - F_0(2(t+\varepsilon) - F_0^{-1}(x-\varepsilon))\right)\right] ds + \int_0^{s_0} J\left(\frac{1}{2}s\right) ds$$

kde

$$s_0 = F_0(2(t+\varepsilon) - x_0) + \varepsilon \vee 1.$$

Dále spočteme

$$\begin{aligned} \lim_{t \rightarrow -\infty} \lambda(t; F_2) &= \int_0^1 J\left(\frac{1}{2}(s+1-\varepsilon)\right) ds + \int_0^{\varepsilon} J\left(\frac{1}{2}s\right) ds = \\ &= 2 \cdot \left(\int_{1/2}^{1-\varepsilon/2} J(s) ds - \int_{1-\varepsilon/2}^1 J(s) ds \right). \end{aligned}$$

Opět $b^-(\epsilon)$ je konečné, pokud $\lim_{t \rightarrow \infty} \lambda(t; F_2) > 0$.
 Tedy i v tomto případě dostaneme pro bod selhání c podmínu

$$\int_{1/2}^{1-\epsilon/2} J(s) ds = \int_{1-\epsilon/2}^1 J(s) ds,$$

což je tvrzení věty 2.4.

2.5. REDESENEDNÍ M-ODHADY S OMEZENOU FUNKCI

Předpokládejme, že $\rho(\cdot)$ má minimum v 0 , $\rho(0) = -1$ a že ρ je neklesající na obou stranách ($k \rightarrow \infty$ a $k \rightarrow -\infty$) a $\lim_{x \rightarrow \infty} \rho(x) = 0$.

VĚTA 2.5

Položí-li

$$\sum_x \rho(x - T(X)) = -A,$$

pak (výběrový) bod selhání T v modelu ϵ -značitelní je roven

$$s(X, T) = \frac{m'}{n + m'},$$

kde m' je přirozené číslo, vyhovující vztahu $|A| \leq m' < |A| + 1$.

Když existují nějaké $c < \infty$ takové, že $\rho(x) = 0$ pro $x \geq c$, potom $m' = |A|$.

2.6. REDESENEDNÍ M-ODHADY S NEOMEZENOU FUNKCI

Předpokládejme, že ρ je sudá, $\rho(0) = 0$ a ρ je neklesající na $(0, \infty)$.

Dále je

$$\lim_{x \rightarrow \pm\infty} \rho(x) = \pm\infty$$

$$\lim_{x \rightarrow \infty} \frac{\rho(x)}{x} = 0.$$

Předpokládejme, že $\gamma = \rho'$ je spojitá a že existují x_0 tak, že γ je neklesající pro $0 < x < x_0$ a nerostoucí pro $x_0 < x < \infty$.

VĚTA 1.6

Za výše uvedených podmínek je bod selhání M-odhadu v modelu ϵ -značitelní roven $\epsilon^2 = 1/2$.

Důkaz

Viz Huber (1984).

2.7. P-ODHADY

Odhad Pitmanova typu, neboli p-odhad polohy se definuje

$$T_p = \frac{\int \exp \left[-\sum \rho(x_i - \theta) \right] \phi d\theta}{\int \exp \left[-\sum \rho(x_i - \theta) \right] d\theta},$$

pokud má výraz smysl.

VĚTA 2.7

Předpokládejme, že ρ je sudá a konvexní a $0 < \lim_{x \rightarrow \pm\infty} \rho(x) \Leftrightarrow x = c < \infty$.

Potom bod selhání P-odhadu je $\epsilon^* = 1/2$.

Důkaz

Viz Huber (1984).

2.8. PRÍKLADY

1. Hodges - Lehmannův odhad

H-L odhad $T = \operatorname{med}_{i>j} [(x_i + x_j)/2]$ máme dostat také jako R-odhad s funkcí $J(t) = t - 1/2$.

Spočteme jeho bod selhání.

a) pomocí věty 2.4

Médáme řešení rovnice

$$\int_{1/2}^{1-\delta/2} (t - 1/2) dt = \int_{1-\delta/2}^1 (t - 1/2) dt.$$

Spočteme, že δ musí splňovat rovnici

$$2 \cdot \frac{\delta^2}{2} - 4\delta + 1 = 0$$

$$\text{a z ní plyne } \epsilon^* = 1 - \frac{\sqrt{2}}{2} \approx 0,239$$

b) z definice bodu selhání

Připojíme-li k výběru g hodnot, odhad selže, bude-li

$$\binom{n}{2} < \frac{1}{2} \binom{n+g}{2}$$

odtud vypočteme

$$g = \frac{1 - 2n + \sqrt{(2n^2 - 2n + 1)}}{2}$$

$$\epsilon^* = \frac{g}{n+g} = \frac{1/2 - n + \sqrt{(2n^2 - 2n + 1/4)}}{1/2 + \sqrt{(2n^2 - 2n + 1/4)}} = \frac{\sqrt{2} - 1}{\sqrt{2}} = \frac{1 - \sqrt{2}}{2}.$$

2. Výběrový průměr

Odhad selže už připojením jediného čpatného porovnání.

$$\epsilon^* = \frac{1}{n+1}$$

3. g - určitý průměr

Tj. zanedbáme g největších a nejméněch pozorování. Snažme spočteme

$$\epsilon^* = \frac{g+1}{n+g+1}$$

4. median

Pomocí věty 2.2 - $\gamma = \operatorname{med}\nolimits_{i>j} x_i$, že $\epsilon^* = 1/2$

5. ϵ -uříšnuttí průměr

Pomocí věty 2.3 lze snadno ukázat, že $\epsilon^* = \epsilon$.

6. R-odhad s normálními vahami

Tj. $J(t) = \phi^{-1}(t)$. Opět hledáme řešení rovnice

$$\int_{1/2}^{1-1/2\delta} \phi^{-1}(s) ds = \int_{1-\sqrt{2}/2}^1 \phi^{-1}(s) ds.$$

Nejprve použijeme substituci $s = \phi(x)$, potom $\frac{ds}{dx} = t$ a dejeme k rovnici

$$\ln 2 - \frac{[\phi^{-1}(1 - 1/2)]^2}{2} = 0.$$

Dopodíleme a máme

$$\epsilon^* = 2(1 - \phi(\sqrt{\ln 2})) = 2\phi(-\sqrt{\ln 2})$$

$$\underline{\epsilon^* \approx 0,239}$$

POUŽITÁ LITERATURA

Bonche, D. L. and Huber, P. J. (1982). The notion of breakdown point. In: Postscript in Honor of Erich Lehmann, Ed. by K. Doksum and J. L. Hodges, Wadsworth, Belmont, CA.

Huber, P. J. (1981). Robust Statistics. Wiley, New York.

Huber, P. J. (1984). Finite sample breakdown of M-and P-estimators. The Annals of Statistics, 1984, Vol. 12, No. 1, 119-126.