

Robustní odhady v lineárním regresním modelu

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_{p-1} X_{i,p-1} + \varepsilon_i, \quad i = 1, \dots, m$$

$$Y_i = X_i^T \beta + \varepsilon_i$$

$$X_i = (1, X_{i1}, \dots, X_{i,p-1})^T \dots p\text{-rozměrný vektor regresorů (převážně čísla)}$$

$$\beta = (\beta_0, \beta_1, \dots, \beta_{p-1})^T \dots p\text{-rozměrný vektor neznámých parametrů}$$

ε_i jsou chyby modelu, i.i.d. s distribuční funkcí F (obecně neznáma)

$$E\varepsilon_i = 0, \quad D\varepsilon_i = \sigma^2$$

Odhad parametrů β metodou nejmenších čtverců: $\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^m (Y_i - X_i^T \beta)^2 = (X^T X)^{-1} X^T Y$ maximálně robustní odhad β , jsou-li chyby modelu ε_i normální

M-odhady

M-odhad parametrů β je řešením minimalizace: $\hat{\beta}^M = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^m \rho(Y_i - X_i^T \beta)$, kde ρ je nějaká vhodná funkce (symetrická)

ex. - li použít derivace $\psi(x) = \rho'(x)$, pak $\hat{\beta}^M$ řeší soustavu rovnic $\sum_{i=1}^m X_i \cdot \psi(Y_i - X_i^T \hat{\beta}^M) = \mathbf{0}$.

Poznámky

- volba ρ , resp. ψ děje jako v modelu polohy

- tento definovaný M-odhad nemá ekvivalentní vůči měřítku \rightarrow studentizace

$$X = \begin{pmatrix} 1 & X_{11} & \dots & X_{1,p-1} \\ 1 & X_{21} & \dots & X_{2,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{m1} & \dots & X_{m,p-1} \end{pmatrix}$$

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix}$$

• odhademe σ^2 pomocí vhodné skalové statistiky S_m (klasický odhad $\sqrt{\frac{1}{m-p} \sum_{i=1}^m (\hat{y}_i - y_i)^2}$).

• studentizovaný M-odhad: $\hat{\beta}^{SM} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^m \rho\left(\frac{y_i - x_i^T \beta}{S_m}\right)$, příp. řeší rovnici $\sum_{i=1}^m x_i \cdot \psi\left(\frac{y_i - x_i^T \hat{\beta}^{SM}}{S_m}\right) = \mathbf{0}$.

Influenciční funkce M-odhadu je $IF(\underbrace{x_i, y_i}_{(p+1) \text{ proměnných}}, \mathbf{T}, P) = \bar{B}^{-1} \times \underbrace{\psi(y - x^T \mathbf{T}(P))}_{\text{lineární } \psi \text{ volil omezenou}}$, kde $B = \int_{\mathbb{R}^{p+1}} x x^T \psi'(y - x^T \mathbf{T}(P)) dP(x, y)$
monotonní funkce

Bod selhání je $\mathbf{0}$.

GM-odhady

$\hat{\beta}^{GM} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^m \rho(y_i - x_i^T \beta) \cdot w(x_i)$, kde $w: \mathbb{R}^p \rightarrow \mathbb{R}$ je vhodná váhová funkce mizející u silně odlehlých hodnot x_i .

$\hat{\beta}^{GM}$ řeší rovnici: $\sum_{i=1}^m x_i \cdot w(x_i) \psi(y_i - x_i^T \beta) = \mathbf{0}$.

$IF(x_i, y_i, \mathbf{T}, P) = \bar{B}^{-1} \underbrace{x_i \cdot w(x_i)}_{\text{volil } \psi \text{ omezenou}} \psi(y_i - x_i^T \mathbf{T}(P))$

w volíme tak, aby funkce $x \cdot w(x)$ byla omezená

např. $w(x_i) = \frac{1}{h_{ii}}$, kde h_{ii} je diagonální prvek matice $H = X(X^T X)^{-1} X^T$

L-odhady

1978 Koenker, Bassett definovali α -regresní kvantil $\hat{\beta}(\alpha)$ jako řešení minimalizace:

$$\hat{\beta}(\alpha) = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^m \rho_{\alpha}(Y_i - x_i^T \beta), \text{ kde } 0 < \alpha < 1 \text{ je pevné a}$$

$$\rho_{\alpha}(x) = \begin{cases} d \cdot x & , x > 0 \\ (d-1)x & , x \leq 0 \end{cases} = |x| \cdot \left(d \cdot \mathbb{1}_{\{x > 0\}} + (1-d) \cdot \mathbb{1}_{\{x \leq 0\}} \right).$$

Poznámka

$\alpha = 0,5$, $\rho_{\alpha}(x) = 0,5 \cdot |x|$, $\hat{\beta}(0,5) = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^m 0,5 \cdot |Y_i - x_i^T \beta| = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^m |Y_i - x_i^T \beta|$ je regresní median.

regresní kvantily se hledají pomocí lineárního programování - simplexová metoda

řešení duální úlohy $\hat{a}(\alpha) = (a_1(\alpha), \dots, a_m(\alpha))^T \dots$ regresní pořadové skály

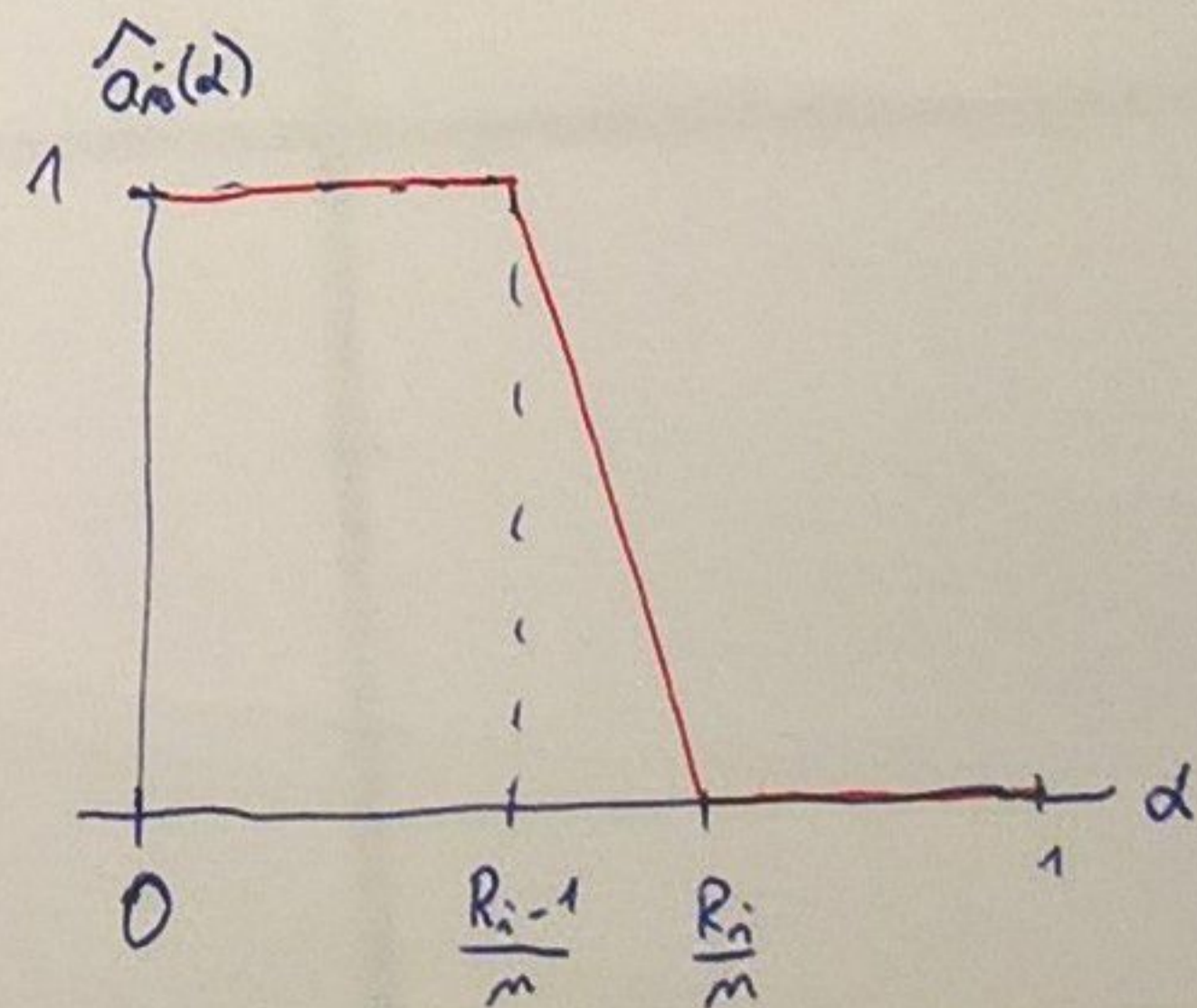
Poznámky

Regresní pořadové skály jsou invariantní vůči regresi X, Y . $\hat{a}(\alpha, Y + X^T b) = \hat{a}(\alpha, Y)$, $\forall b \in \mathbb{R}^p$

speciálně v modelu polohy: $X = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$, $Y_i = \beta_0 + \varepsilon_i$

$$\hat{a}_i(\alpha) = \begin{cases} 1 & , d \leq \frac{R_i - 1}{m} \\ R_i - dm & , \frac{R_i - 1}{m} < d < \frac{R_i}{m} \\ 0 & , d \geq \frac{R_i}{m} \end{cases}$$

R_i je pořadí Y_i mezi Y_1, \dots, Y_m .



o duálním modelu (díky dualitě) platí:

$$\hat{a}_i(d) = \begin{cases} 1 & , Y_i > x_i^T \hat{\beta}(d) \\ 0 & , Y_i < x_i^T \hat{\beta}(d) \end{cases} \quad \text{a } \hat{a}_i(d) \in (0,1), \text{ pokud } Y_i = x_i^T \hat{\beta}(d)$$

$\hat{a}_i(d)$ je spojitá, po částech lineární $\hat{a}_i(0) = 1, \hat{a}_i(1) = 0$.

Vraťme se zpět k regresním kvantilům:

- L_1 -odhad (regresní medián, 0,5-regresní kvantil, LAD)

- bod selhání 0 (robustní jen v odchytkách v Y_i)

- trimmed least squares (TLS) - uzelnuté nejmenších čtverců

$0 < \underline{d}_1 < \underline{d}_2 < 1$ a vybereme jen ta pozorování, která leží mezi d_1 a d_2 -regresními kvantily a s jejich pomocí spočteme odhad metodou nejmenších čtverců

$$a_i = \mathbb{1} \{ x_i^T \hat{\beta}(d_1) \leq Y_i \leq x_i^T \hat{\beta}(d_2) \} \quad \dots i = 1, \dots, n$$

$$A = \text{diag} \{ a_i \}$$

$$\hat{\beta}^{\text{TLS}} = (X^T A X)^{-1} X^T A Y$$

- preliminary estimate (PE)

• preliminary estimate (PE)

- vhodné počáteční odhad $\hat{\beta}^{(0)}$ parametru β . Jako $\hat{\beta}^{(0)}$ zvolí $\left\{ \begin{array}{l} \text{LSE} \\ \text{LAD} \\ \frac{1}{2} (\hat{\beta}^{(2)} + \hat{\beta}^{(1-d)}) \end{array} \right.$

- odhadáme $\lfloor d \cdot n \rfloor$ pozorování s nejmenšími a $\lfloor n - d \rfloor$ pozorování s největšími residui.

- ze zbylých pozorování určíme odhad metodou nejmenších čtverců

Odhady s vysokým bodem reliability

• least median of squares (LMS): $\hat{\beta}^{LMS} = \underset{\beta \in RP}{\operatorname{argmin}} \operatorname{med}_{i=1, \dots, m} \{ (Y_i - x_i^T \beta)^2 \}$ $\epsilon^* = 0,5$, nemá však příliš efektivní

• least trimmed squares (LTS): $\hat{\beta}^{LTS} = \underset{\beta \in RP}{\operatorname{argmin}} \sum_{i=1}^h r_i^2(\beta)$, kde $r_i^2(\beta) = (Y_i - x_i^T \beta)^2$, $r_{(1)}^2 \leq r_{(2)}^2 \leq \dots \leq r_{(m)}^2$ $h = \lfloor \frac{m}{2} \rfloor + \lfloor \frac{p+1}{2} \rfloor \dots \epsilon^* = 0,5$

• least trimmed absolute deviation (LTA): $\hat{\beta}^{LTA} = \underset{\beta \in RP}{\operatorname{argmin}} \sum_{i=1}^{n-2\lfloor d \rfloor} |r(\beta)|_{(i)}$, kde $|r_i(\beta)| = |Y_i - x_i^T \beta|$, $|r(\beta)|_{(1)} \leq |r(\beta)|_{(2)} \leq \dots \leq |r(\beta)|_{(n)}$

model regrese přímky: $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$

• repeated median: $\hat{\beta}_1 = \operatorname{med}_{i=1, \dots, m} \operatorname{med}_{\substack{j=1, \dots, m \\ j \neq i}} \frac{Y_i - Y_j}{x_i - x_j}$ $\epsilon^* = 0,5$

• Theilov-Genov odhad: $\hat{\beta}_1 = \operatorname{med}_{i < j} \frac{Y_i - Y_j}{x_i - x_j}$ $\epsilon^* = 1 - \frac{1}{12} = 0,9167$

Robustní odhady parametru σ

- median absolute deviation (MAD)

$\beta^{(0)}$ počáteční odhad parametru β

$$r_i = y_i - x_i^T \beta^{(0)}$$

$$\tilde{r} = \text{med} \{ r_1, \dots, r_m \}$$

$$S_m = \text{med} \{ |r_1 - \tilde{r}|, |r_2 - \tilde{r}|, \dots, |r_m - \tilde{r}| \}.$$

- L-stabilitý odhad na regresních kvantilech

$$S_m = \|\hat{\beta}(d_2) - \hat{\beta}(d_1)\|, \quad 0 < d_1 < d_2 < 1 \quad \text{a } \|\cdot\| \text{ je eukleidovská norma v } \mathbb{R}^p.$$