

## Pořadové testy s lineárním regresním modelem

a) model regresní přímky:  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$  ( $i = 1, \dots, n$ )

chyby modelu  $\varepsilon_i$  jsou nezávislé, stejné rozdělení s distribuční funkcí  $F$  a hustotou  $g$ .  
 $X_1, \dots, X_n$  jsou pevná čísla

$H_0: \beta_1 = 0$  (je tam vůbec nějaká regrese?)

$H_1: \beta_1 \neq 0$  (případně  $\beta_1 > 0$  nebo  $\beta_1 < 0$ )

### Konstrukce pořadových testů

Pokud bychom znali  $F, \text{ resp. } g$ , pak lokálně nejvhodnější pořadový test  $H_0$  proti  $H_1$  je založen na testové statistice  $S_m = \sum_{i=1}^m (X_i - \bar{x}) A_m(R_i, g)$ , kde

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

Člony  $A_m(i, g)$  můžeme považovat buďto za  $E\psi(U_{(i)}, g)$ , nebo aproximací  $\psi(\frac{i}{m+1}, g)$ .

$$\psi(u, g) = -\frac{g'(\bar{F}(u))}{g(\bar{F}(u))}, \quad 0 < u < 1.$$

V praxi ale  $g$  neznáme, proto zvolíme funkci  $\psi: [0, 1] \rightarrow \mathbb{R}$  neklesající, nekonzantní a integrovatelná se čtením a definujeme  $S_m = \sum_{i=1}^m (X_i - \bar{x}) A_m(R_i) = \sum_{i=1}^m (X_i - \bar{x}) \psi(\frac{R_i}{m+1})$ .

$$\text{kde } A_m(i) = \psi\left(\frac{i}{m+1}\right).$$

Ordné rozdělení  $S_m$  na platnosti  $H_0$ : Vektor pořadí  $(R_1, \dots, R_n)^T$  má rovnoměrné rozdělení na množině všech permutací čísel  $\{1, \dots, m\}$   
(podmíněně při daných  $X_1, \dots, X_m$ )

permutační princip (síta a transformaci)  
společné hodnoty  $S_m$  pro všechny

permutace čísel  $\{1, 2, \dots, m\}$ .

Pořadí jsou invariantní vůči posunům  $\Rightarrow$  rušivý parametr  $\beta_0$  nám nevadí (za  $H_0$  máme model:  $Y_i = \beta_0 + \varepsilon_i$ )

Pro  $n$  střední výpočet márové  $\rightarrow$  asymptotická aproximace

Poznámka

Či platnosti  $H_0$  platí:  $ES_m = 0$  a  $DS_m = \frac{1}{m-1} \left( \sum_{i=1}^m (x_i - \bar{x})^2 \right) \cdot \left( \sum_{j=1}^m (a_{(j)} - \bar{a})^2 \right)$ .

člen  $\left( \sum_{j=1}^m (a_{(j)} - \bar{a})^2 \right)$  můžeme opět pro  $m$  velice aproximovat pomocí  $(m-1) \int_0^1 (\varphi(u) - \bar{\varphi})^2 du$ .

Težka

Či platnosti  $H_0$  a mírných předpokladů na  $x_{11}, \dots, x_{1n}$  platí  $S_m \stackrel{H_0}{\approx} N(0, DS_m)$  pro  $m \rightarrow \infty$ .

b) model vícenásobné regrese:  $Y_i = \beta_0 + x_i^T \beta + \varepsilon_i, i=1, \dots, m$

$\beta \in \mathbb{R}^p$  je  $p$ -rozměrný parametr

$H_0: \beta = 0$   
(je tam vůbec nějaká regrese?)

$H_1: \beta \neq 0$

Pro  $j=1, \dots, p$  definujme  $S_{mj} = \sum_{i=1}^m (x_{ij} - \bar{x}_j) a_m(R_i)$ , kde  $a_m(i) = \varphi\left(\frac{i}{m+1}\right)$  pro nějakou nekombanční, neklesající a integrovatelnou se čtvercem

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

$$S_m = (S_{m1}, \dots, S_{mp})^T = \sum_{i=1}^m (x_i - \bar{x}) \cdot a_m(R_i) = \sum_{i=1}^m (x_i - \bar{x}) \varphi\left(\frac{R_i}{m+1}\right).$$

za predpokladu  $H_0$  platí:  $ES_m = \mathbf{0}$ ,  $DS_m = ES_m S_m^T = \frac{1}{n-1} \underbrace{\left( \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \right)}_{Q_m} \cdot \underbrace{\left( \sum_{j=1}^m (a_{mj} - \bar{a})^2 \right)}_{A_m^2}$ .

a za určitých predpokladov na  $x_1, \dots, x_n$  platí  $S_m \stackrel{H_0}{\approx} N_p(\mathbf{0}, DS_m)$ .

Poznámka

$A_m^2$  môžeme pomocou určitej aproximácie  $(n-1) \int_0^1 (\varphi(u) - \bar{\varphi})^2 du$ .

Testová štatistika:  $T_m = S_m^T (DS_m)^{-1} S_m = \frac{1}{A_m^2} S_m^T Q_m^{-1} S_m$ .

Opäť musíme mať presné rozdelenie  $T_m$  (podmienečne pri  $x_1, \dots, x_n$ ) za predpokladu  $H_0$  pomocou permutačného princípu.

Pro  $n$  strednú pomocou aproximácie: za  $H_0$  a určitých predpokladov na  $x_1, \dots, x_n$  platí  $T_m \stackrel{H_0}{\approx} \chi_p^2$ .

Poznámka

Konkrétne pri hladinách testu je to analogické tým o predchádzajúcich kapitolách.

- Wilcoxon, von der Neumann apod.

# R-odhady v lineárním regresním modelu

a) model regresní přímky:  $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ ,  $i=1, \dots, n$   $\varepsilon_i$  jsou i.i.d. chyby modelu s hustotou  $f$  a hustotou  $f$

Zajímá nás odhad měřnice  $\beta_1$ .

vycházíme z testu hypotéz  $H_0: \beta_1 = b$ , kde  $b$  je reálné číslo

za platnosti  $H_0$  máme model:  $Y_i = \beta_0 + b x_i + \varepsilon_i$

$$\underbrace{Y_i - b x_i}_{Y_i(b) = Y_i^*} = \beta_0 + \varepsilon_i \quad \text{v tomto modelu testujeme } H_0: \beta = 0.$$

Lehmannův statistika:  $S_m(b) = \sum_{i=1}^m (x_i - \bar{x}) a_m(R_i(b))$ , kde  $a_m(u) = \varphi\left(\frac{u}{m+1}\right)$  jsou státní generované nezávislé, nekonstantní, integrovatelné se čtením a antisymetrické kolem  $\frac{1}{2}$  (tj.  $\varphi(u) = -\varphi(1-u)$ ) funkce  $\varphi: [0, 1] \rightarrow \mathbb{R}$ .

$R_i(b)$  je pořadí  $\underbrace{Y_i - b x_i}_{Y_i^*}$  mezi  $\underbrace{Y_1 - b x_1}_{Y_1^*}, \dots, \underbrace{Y_m - b x_m}_{Y_m^*}$ .

$S_m(b)$  jakožto funkce  $b \in \mathbb{R}$  je nerostající, počátečně konstantní funkce

$y$ -li skutečná hodnota parametru  $\beta_1$  (za platnosti  $H_0$ ), pak  $ES_m(b) = 0 \Leftrightarrow$  hledáme  $b$  takové, pro které  $S_m(b) \stackrel{!}{=} 0$ . ← obecně nemáme mil řešení

R-odhad parametru  $\beta_1$  definujeme jako  $\hat{\beta}_1 = \frac{1}{2} (\hat{\beta}_1^+ + \hat{\beta}_1^-)$ , kde  $\hat{\beta}_1^+ = \inf \{ b : S_m(b) < 0 \}$  a  $\hat{\beta}_1^- = \sup \{ b : S_m(b) > 0 \}$ .

## Poznámka

$\hat{\beta}_1$  je asymptoticky nekorelovaný, konzistentní a asymptoticky normální odhad  $\beta_1$  (as. rozptyl závisí na neznámé hustotě  $f$ ).

1) model lineárního regrese:  $Y_i = \beta_0 + x_i^T \beta + \varepsilon_i$ ,  $i=1, \dots, m$

zájímá nás odhad  $p$ -rozměrného parametru  $\beta$ .

Vycházíme z testu hypotézy  $H_0: \beta = b$ , kde  $b$  je známý  $p$ -rozměrný vektor

za předpokladu  $H_0$  máme model:  $Y_i = \beta_0 + x_i^T b + \varepsilon_i$

$$\underbrace{Y_i - x_i^T b}_{Y_i(b) = Y_i^*} = \beta_0 + \varepsilon_i \quad \text{v tomto modelu testujeme } H_0^*: \beta = 0.$$

testová statistika:  $S_m(b) = \sum_{i=1}^m (x_i - \bar{x}) a_m(R_i(b))$ , kde  $a_m(i) = \varphi\left(\frac{i}{m+1}\right)$  je funkce generovaná neklesající měkčím, ind. se čte od 1 a antisymetrickou funkcí  $\varphi: [0,1] \rightarrow \mathbb{R}$ .

$R_i(b)$  je pořadí  $Y_i - x_i^T b$  mezi  $Y_1 - x_1^T b, \dots, Y_m - x_m^T b$ .

$\beta$ -li skutečná hodnota parametru rovená  $\beta$ , pak  $ES_m(b) = 0$ . ( $S_m: \mathbb{R}^p \rightarrow \mathbb{R}^p$  je  $p$ -rozměrná funkce)

$T_m(b) = \|S_m(b)\|$ , kde  $\|\cdot\|$  je libovolná norma v  $\mathbb{R}^p$  (volba normy neovlivní slabsími výsledného odhadu).

$R$ -odhad parametru  $\beta$  definujeme jako  $\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} T_m(b) = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|S_m(b)\|$ .

## Poznámka

$\hat{\beta}$  je as. nekorelovaný, konzistentní a as. normální odhad  $\beta$  (as. rozptyl závisí na neznámé hustotě  $f$ ).

jiná definice R-odhadů: Jaeckel (1972)

rank measure of dispersion:  $D_m(\beta) = \sum_{i=1}^m (y_i - (x_i - \bar{x})^T \beta) a_m(R_i(\beta))$ , kde  $a_m(i)$  a  $R_i(\beta)$  jsou definované výše.

$D_m(\beta)$  je konvexní, po částech lineární, se subgradientem  $-\sum_{i=1}^m (x_i - \bar{x}) a_m(R_i(\beta))$ .

$$\hat{\beta}^{\square} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} D_m(\beta).$$

Tvrzení

Odhady  $\hat{\beta}$  a  $\hat{\beta}^{\square}$  jsou asymptoticky ekvivalentní.

Testy hypotéz o části parametru  $\beta$

$$y_i = \beta_0 + x_i^T \beta + z_i^T \eta + \varepsilon_i, \quad i=1, \dots, m$$

$\begin{matrix} \in \mathbb{R}^p & \in \mathbb{R}^q \\ (x_i^T, z_i^T)^T & \dots \text{vektor regresorů dimenze } p+q \\ (\beta_0, \beta^T, \eta^T) & \dots \text{vektor neznámých parametrů dimenze } 1+p+q \\ \in \mathbb{R}^1 & \in \mathbb{R}^p \end{matrix}$

$$H_0: \eta = \mathbf{0}$$

$\beta_0, \beta$  jsou rušivé parametry

$$H_1: \eta \neq \mathbf{0}$$

a) aligned pořadový test o parametru  $\beta$

- pořadový test je invariantní vůči posunům  $\Rightarrow$  parametru  $\beta$  o nám "neradí".

- odhademe parametru  $\beta$  pomocí R-odhadu, máme rezidua a na ně aplikujeme pořadový test o parametru  $\beta$ .

(i) za předpokladu  $H_0$  máme model:  $Y_i = \beta_0 + x_i^T \beta + \varepsilon_i$   $\leftarrow$  o tomto modelu odhademe  $p$ -rozměrný parametru  $\beta$

$$S_m(\beta) = \sum_{i=1}^m (x_i - \bar{x}) \Psi\left(\frac{R_i(\beta)}{m+1}\right), \text{ kde } \Psi \text{ je netles., netond. a int. nečíslem a } \Psi(t) = -\Psi(1-t)$$

$R_i(\beta)$  je pořadí  $Y_i - x_i^T \beta$  mezi  $Y_1 - x_1^T \beta, \dots, Y_m - x_m^T \beta$ .

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \|S_m(\beta)\|, \text{ kde } \|\cdot\| \text{ je nějaká norma v } \mathbb{R}^p.$$

(ii) definujeme rezidua  $\hat{\varepsilon}_i = Y_i - x_i^T \hat{\beta}$ ,  $i = 1, \dots, m$  by ale nejsou ani nezávislé, ani nějak rozdělené (přesto "na ně" aplikujeme pořadový test o par.  $\beta$ )

novou statistiku  $S_m = \sum_{i=1}^m (z_i - \bar{z}) a_m(\hat{R}_i) = \sum_{i=1}^m (z_i - \bar{z}) \Psi\left(\frac{\hat{R}_i}{m+1}\right)$ , kde  $\Psi$  je netles., netond. a int. nečíslem

$\hat{R}_i$  je pořadí  $\hat{\varepsilon}_i$  mezi  $\hat{\varepsilon}_1, \dots, \hat{\varepsilon}_m$ .

$$D_m = \frac{1}{m-1} \sum_{i=1}^m (z_i - \bar{z})(z_i - \bar{z})^T, A_m^2 = \sum_{j=1}^m (a_m(j) - \bar{a})^2$$

$$T_m = \frac{1}{A_m^2} S_m^T D_m^{-1} S_m.$$

Nyní již nelze jednoduše vyjádřit průměrné rozdělení  $T_m$  na platnosti  $H_0$  ( $\hat{e}_i$  jsou závislé)  $\Rightarrow$  nutnost asymptotické aproximace

za určitých předpokladů  $T_m \stackrel{H_0}{\approx} \chi^2_q$  při  $n \rightarrow \infty$ .

b) test založený na regresečních pořádkových statistikách

na platnosti  $H_0$  máme model:  $Y_i = \beta_0 + x_i^\top \beta + \varepsilon_i$   $\leftarrow$  v tomto modelu najdeme regreseční pořádkové statistiky

označme je  $\hat{a}_i(\alpha)$   $i=1, \dots, m$  pro  $0 \leq \alpha \leq 1$  ty jsou invariantní vůči rušivé regresi v  $x_i$  (analogie pořadí)  $\leftarrow$  rušivé parametry  $\beta_0$  a  $\beta$  nevadí.

označíme:  $X = \begin{pmatrix} 1 & x_1^\top \\ \vdots & \vdots \\ 1 & x_m^\top \end{pmatrix}$  a  $Z = \begin{pmatrix} z_1^\top \\ \vdots \\ z_m^\top \end{pmatrix}$ .

Zvolíme neklesající, nekondantní a int. rektivovanou funkci  $\varphi: [0,1] \rightarrow \mathbb{R}$  a definujeme statistiku  $\hat{b}_i = - \int_0^1 \varphi(t) d\hat{a}_i(t)$  pro  $i=1, \dots, m \rightarrow \hat{b} = (\hat{b}_1, \dots, \hat{b}_m)^\top$ .

testová statistika:  $S_m = \sum_{i=1}^m (z_i - \hat{z}_i) \cdot \hat{b}_i$ , kde  $\hat{z} = \begin{pmatrix} \hat{z}_1 \\ \vdots \\ \hat{z}_m \end{pmatrix} = X(X^\top X)^{-1} X^\top Z$  je projekce  $Z$  do podprostoru sloupců matice  $X$ .

$$\hat{D}_m = \frac{1}{n-1} \sum_{i=1}^m (z_i - \hat{z}_i)(z_i - \hat{z}_i)^\top, \quad A_m^2 = (n-1) \int_0^1 (\varphi(t) - \bar{\varphi})^2 dt.$$

$$\tilde{T}_m = \frac{1}{A_m^2} \cdot S_m^\top \hat{D}_m^{-1} S_m \quad \text{má za určitých předpokladů asymptoticky } \chi^2_q.$$

Poznámky: Oba testy mají stejné asymptotické rozdělení (i za  $H_1$ ), tedy i asymptoticky stejnou sílu.