

# Change point (bod zlomu)

Julie Váňová

12. dubna 2023

Statistický seminář 1

Uvažujeme regresní model

$$Y_{in} = \mathbf{x}_{in}^T \beta + \mathbf{x}_{in}^T \delta_n \cdot \mathbf{1}_{[i > m_n]} + \epsilon_i, \quad i = 1, \dots, n.$$

- $\mathbf{x}_{in} = (\mathbf{x}_{i1,n}, \dots, \mathbf{x}_{ip,n})^T$ ,  $i = 1, \dots, n$ , - předem dané body
- $\epsilon_1, \dots, \epsilon_n$  i.i.d. - chybové členy
- $\beta = (\beta_1, \dots, \beta_p)^T$  a  $\delta_n = (\delta_{1n}, \dots, \delta_{pn})^T \neq \mathbf{0}$  - (neznámé) parametry
- $m_n \leq n$  - neznámý parametr (change point)

- chceme testovat hypotézu

$$H_0 : m_n = n \quad \text{proti} \quad H_1 : m_n < n$$

- zajímá nás odhad  $m_n$

- budeme uvažovat následující testové statistiky

$$T_n = \max_{p < k < n-p} \frac{S_{kn}^T C_{kn}^{-1} C_{nn} C_{kn}^{o-1} S_{kn}}{\hat{\sigma}_n^2}$$

$$T_n(q) = \sup_{0 < t < 1} \frac{S_{\lfloor (n+1)t \rfloor n}^T C_{nn}^{-1} S_{\lfloor (n+1)t \rfloor n}}{q^2(t) \hat{\sigma}_n^2}$$

$$T_n^* = \max_{1 \leq k < n} \sqrt{\frac{n}{k(n-k)}} \cdot \frac{|S_{kn}^*|}{\hat{\sigma}_n}$$

$$T_n^*(q) = \sup_{0 < t < 1} \frac{|S_{\lfloor (n+1)t \rfloor n}^*|}{\sqrt{nq(t)} \hat{\sigma}_n}$$

# Odhad rozptylu

- $\hat{\sigma}_n^2$  je odhad  $\sigma^2$  splňující

$$\hat{\sigma}_n^2 - \sigma^2 = o_p((\log(\log n))^{-1/2}), \quad n \rightarrow \infty$$

- 

$$\hat{\sigma}_n^2 = \frac{1}{n-p} \min_{p < k < n-p} \left( \sum_{i=1}^k (Y_{in} - \mathbf{x}_{in}^T \hat{\beta}_k)^2 + \sum_{i=k+1}^n (Y_{in} - \mathbf{x}_{in}^T \hat{\beta}_k^0)^2 \right),$$

kde  $\hat{\beta}_k$  a  $\hat{\beta}_k^0$  jsou LSE založené na  $Y_1, \dots, Y_k$  a  $Y_{k+1}, \dots, Y_n$

- dá se dokázat, že

$$\hat{\sigma}_n^2 = \frac{1}{n-p} \left( \sum_{i=1}^n (Y_{in} - \mathbf{x}_{in}^T \hat{\beta}_n)^2 - \max_{p < k < n-p} (S_{kn}^T C_{kn}^{-1} C_{nn} C_{kn}^{-1} S_{kn}) \right)$$

- u  $T_n^*$  a  $T_n^*(q)$  se rovněž používá

$$\tilde{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (Y_{in} - \mathbf{x}_{in}^T \hat{\beta}_n)^2$$

Věta 1:

- Nechť platí předpoklady (A1)-(A3) a (B1). Potom za  $H_0$

$$\lim_{n \rightarrow \infty} P(a(\log n)\sqrt{T_n} \leq t + b_p(\log n)) = \exp(-2e^{-t}),$$

$$\lim_{n \rightarrow \infty} P(a(\log n)T_n^* \leq t + b_1(\log n)) = \exp(-2e^{-t}),$$

kde

$$a(y) = \sqrt{2 \log y}, \quad b_p(y) = 2 \log y + \frac{p}{2} \log \log y - \log(\Gamma(\frac{p}{2})), \quad y > 1,$$

$$\text{a } \Gamma(p) = \int_0^\infty t^{p-1} e^{-t} dt.$$

Věta 2:

- Nechť platí předpoklady (A1)-(A3), (B1) a (C1). Potom za  $H_0$

$$\sqrt{T_n(q)} \xrightarrow[n \rightarrow \infty]{D} \sup_{0 < t < 1} \frac{\sqrt{\sum_{i=1}^p B_i^2(t)}}{q(t)},$$

$$T_n^*(q) \xrightarrow[n \rightarrow \infty]{D} \sup_{0 < t < 1} \frac{|B_1(t)|}{q(t)},$$

kde  $\{B_i(t); t \in (0, 1)\}_{i=1}^p$  jsou nezávislé Brownovy mosty.

Věta 3:

- Necht' platí předpoklady (A1)-(A3) a (B1). Potom

$$\lim_{n \rightarrow \infty} P\left(a(\log n)\sqrt{T_n(\mathbf{R})} \leq t + b_p(\log n) | \mathbf{Y}_n\right) = \exp(-2e^{-t})$$

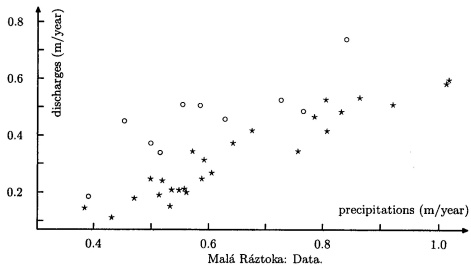
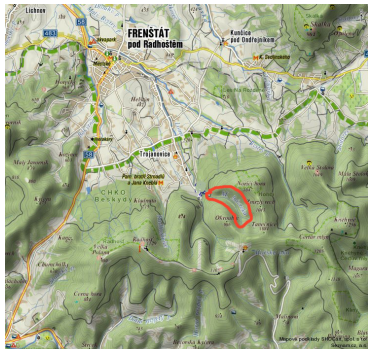
v pravděpodobnosti, a

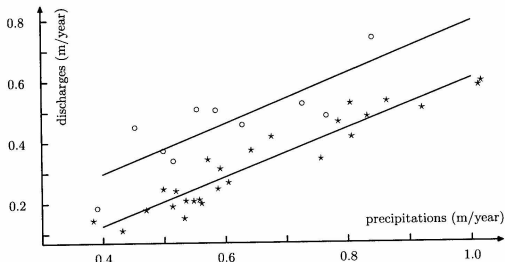
$$\lim_{n \rightarrow \infty} P\left(a(\log n)T_n^*(\mathbf{R}) \leq t + b_1(\log n) | \mathbf{Y}_n\right) = \exp(-2e^{-t})$$

v pravděpodobnosti, kde  $t \in \mathcal{R}_1$ , a  $a(y)$ ,  $b_p(y)$  jsou jako ve větě 1.

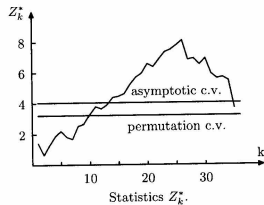
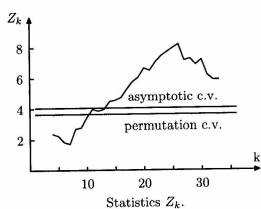


# Příklad

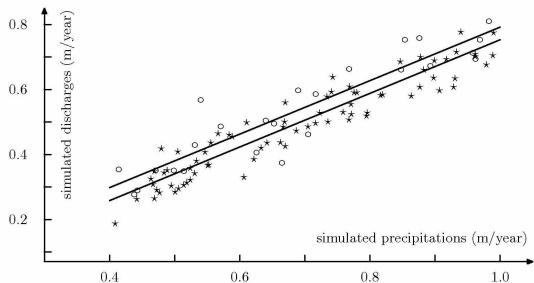




Malá Ráztoka: Data and model.



# Příklad - semi-simulační studie



Simulated data and model.

