

# Robustní odhady vícerozměrného parametru polohy

## Příklad (jednorozměrný)

Nechť  $X_1, \dots, X_n$  je náh. vzorek z  $N(\theta, \sigma^2)$ .  $\theta \in \mathbb{R}$  je neznámý parametr,  $\sigma^2 > 0$  je známý parametr  
(parametr polohy) (parametr měřítka)

$$X_i = \theta + \varepsilon_i, \text{ kde } \varepsilon_i \sim N(0, \sigma^2)$$

Nechť  $\hat{\theta}$  je odhad parametru  $\theta$ . Jako měřítko jeho kvality bereme střední čtvercovou chybu:  $E(\hat{\theta} - \theta)^2$ . Chceme najít odhad  $\hat{\theta}$ , který tuto chybu minimalizuje (pro všechny hodnoty parametru  $\theta$ , tj. nejhomogenně).

$\hat{\theta} = \bar{X}$  je hledaný odhad.

vícerozměrný případ:  $\theta \in \mathbb{O} \subset \mathbb{R}^p$  je  $p$ -rozměrný parametr  $\theta = (\theta_1, \dots, \theta_p)^T$ .

$X_1, \dots, X_n$  je náh. vzorek z  $p$ -rozměrného rozdělení

Nechť  $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$  je odhad parametru  $\theta$   
 $(\hat{\theta}_1, \dots, \hat{\theta}_p)^T$

$L: \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  je ztrátová (riziková) funkce, jestliže  $L(x, y) \geq 0$  a  $L(x, x) = 0$

$L(\hat{\theta}, \theta) \dots$  je ztráta, kterou utrpíme, když parametr  $\theta$  odhadneme pomocí  $\hat{\theta}$ .

$$L(\hat{\theta}, \theta) = \|\hat{\theta} - \theta\|^2 = \sum_{i=1}^p (\hat{\theta}_i - \theta_i)^2 \dots \text{ kvadratická ztrátová funkce}$$

$$R(\hat{\theta}, \theta) = EL(\hat{\theta}, \theta) = E\|\hat{\theta} - \theta\|^2 = \sum_{i=1}^p E(\hat{\theta}_i - \theta_i)^2 \dots \text{ riziko}$$

## Poznámka

$p=1$ , pak riziko  $R(\hat{\theta}, \theta) = E(\hat{\theta} - \theta)^2$  ... je střední kvadratická chyba

$\hat{\theta}_1$  a  $\hat{\theta}_2$  dva odhady parametru  $\theta$ . Řekneme, že  $\hat{\theta}_1$  dominuje  $\hat{\theta}_2$ , jestliže  $R(\hat{\theta}_1, \theta) \leq R(\hat{\theta}_2, \theta)$ ,  $\forall \theta \in \Theta$  a existuje  $\theta_0 \in \Theta$ :  
 $R(\hat{\theta}_1, \theta_0) < R(\hat{\theta}_2, \theta_0)$ .

## Definice

Odhad  $\hat{\theta}$  parametru  $\theta$  je přípustný (admissible), jestliže neexistuje žádný jiný odhad parametru  $\theta$ , který by jej dominoval.

$X_1, \dots, X_n$  je náh. vzájemně  $p$ -rozměrného normálního rozdělení se střední hodnotou  $\theta \in \mathbb{R}^p$ .

James, Stein (1961):

•  $p=1, 2$ , pak  $\hat{\theta} = \bar{X}$  je přípustný odhad parametru  $\theta$

•  $p \geq 3$ , pak  $\hat{\theta} = \bar{X}$  není přípustný

## Ústředí (více-rozměrný)

$X_1, \dots, X_n$  je náh. vzájemně  $p$ -rozměrného normálního rozdělení  $N_p(\theta, \sigma^2 \cdot \mathbf{I})$ , kde  $p \geq 3$  a  $\sigma^2 > 0$  je (pro jednoduchost) známé.

$X_i = \theta + \varepsilon_i$ , kde  $\varepsilon_i \sim N_p(\theta, \sigma^2 \cdot \mathbf{I})$

Odhady parametru  $\theta = (\theta_1, \dots, \theta_p)^T$ .

$\hat{\theta} = \bar{X} = (\bar{X}^1, \dots, \bar{X}^p)^T$  je vzájemný průměr

$$R(\hat{\theta}, \theta) = E \|\hat{\theta} - \theta\|^2 = \sum_{i=1}^p E(\hat{\theta}_i - \theta_i)^2 = \sum_{i=1}^p E(\bar{X}^i - \theta_i)^2 = \sum_{i=1}^p D(\bar{X}^i) = p \cdot \frac{\sigma^2}{n}$$

James a Stein:  $\hat{\Sigma}_{JS} = \left( 1 - \frac{(p-2) \frac{\sigma^2}{m}}{\|\bar{X}\|^2} \right) \bar{X}$

ale smí tento odhad nemít 'přípustný'

$$R(\hat{\Sigma}_{JS}, \Sigma) = p \cdot \frac{\sigma^2}{m} - \underbrace{(p-2) \cdot E\left(\frac{1}{\|\bar{X}\|^2}\right)}_{> 0} < p \cdot \frac{\sigma^2}{m}$$

positivně part Jamesův-Steinův odhad:  $\hat{\Sigma}_{JS}^+ = \left( 1 + \frac{(p-2) \cdot \frac{\sigma^2}{m}}{\|\bar{X}\|^2} \right)_+ \bar{X}$

### Ročník

↑ shrnutí odhadů: snížení rozptylu odhadu na úkor vychylenosti.

Jak definovat medián pro p-rozměrná data?

- a) marginální medián - složky jsou příslušné jednorozměrné mediány
- b) geometrický (spatial,  $L_1, \dots$ ) medián

data:  $x_1, \dots, x_m$  (p-rozměrná)

$$\hat{\Theta} = \operatorname{argmin}_{\theta \in \mathbb{R}^p} \sum_{i=1}^m \|x_i - \theta\|, \text{ kde } \|\cdot\| \text{ je eukleidovská norma v } \mathbb{R}^p$$

bod nelháví  $\epsilon^* = 0,5$ .

$p=1$   $\hat{\Theta} = \operatorname{argmin}_{\theta \in \mathbb{R}} \sum_{i=1}^m |x_i - \theta|$  klasický medián

c) median založený na koncepci hloubky dat

hloubková funkce: hloubka bodu  $x \in \mathbb{R}^p$  vzhledem k datovému souboru  $X = \begin{pmatrix} x_1^\top \\ x_2^\top \\ \vdots \\ x_n^\top \end{pmatrix}$  je funkce  $d(x, X)$  splňující následující vlastnosti:

(i) hloubka je afinně invariantní

(ii) nulová v nekonečnu  $d(x, X) \rightarrow 0$   
 $\|x\| \rightarrow \infty$

(iii) je maximální v centru symetrie

(iv) je monotónní od nejhlubšího bodu

Tukeyho hloubka dat  $\text{depth}(x, X) = \frac{1}{n} \cdot \min_{\substack{\|a\|=1 \\ a \in \mathbb{R}^p}} \{n_i : a^\top X_i \geq a^\top x\} = \frac{1}{n} \cdot \left\{ \begin{array}{l} \text{nejmenší počet pozorování } x_i \text{ obsažený v libovolném uzavřeném poloprostoru v } \mathbb{R}^p \\ \text{p hranici procházející bodem } x \end{array} \right\}$

Tukeyho (hloubkový) median:  $\hat{\Phi} = \underset{x \in \mathbb{R}^p}{\text{argmax}} \text{depth}(x, X)$ .

Poznámka

g-li  $X$  v dobré pozici, platí  $\frac{1}{p+1} \leq \text{depth}(\hat{\Phi}, X) \leq \frac{1}{2}$ .

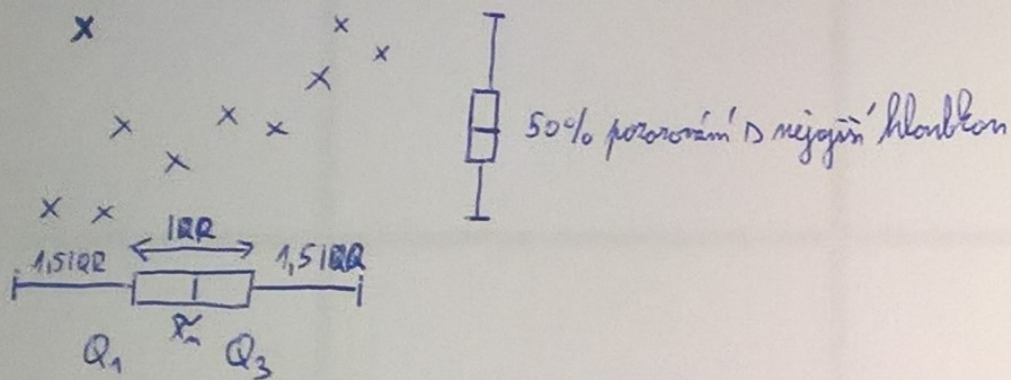
$\text{depth}(\hat{\Phi}, X) = \frac{1}{2} \Leftrightarrow X$  pochází z uhlavě symetrického rozdělení  $\frac{x - \hat{\Phi}}{\|x - \hat{\Phi}\|} \stackrel{d}{=} \frac{\hat{\Phi} - x}{\|x - \hat{\Phi}\|}$ .

$E^* \geq \frac{1}{p+1}$

$\alpha$ -úroveňný príměr:  $\hat{\Theta}_\alpha = \sum_{i: \text{depth}(x_i, X) \geq d} X_i \cdot \frac{1}{|i: \text{depth}(x_i, X)|}$

$\varepsilon^* = \alpha$

Odlehla' porovnaní ve více dimenzích



Bagplot - dvouosměrný boxplot (Rousseeuw, 1999)

- 1) Najdeme oblast s největší hloubkou a definujeme hloubkou (Tukeyho) median  $T^*$  jako ležící tuto oblast.
- 2) Najdeme bag = oblast obsahující 50% porovnaní v největší hloubce

definujeme úrovně množiny  $D_k = \{x: \text{depth}(x, X) \geq \frac{k}{n}\}$  pro  $k=1, 2, \dots$

$D_1 \supset D_2 \supset D_3 \supset \dots$   $|D_1| \geq |D_2| \geq |D_3| \geq \dots$

oznáváme  $|D_k|$  = počet porovnaní  $x_i$  obsažených v  $D_k$

Najdeme  $k$  takové, že  $|D_k| \leq \lfloor \frac{n}{2} \rfloor < |D_{k-1}|$

a nekrajní bag B: konvexní množina mezi  $D_k$  a  $D_{k-1}$  (lineární interpolace relativně vzhledem k  $T^*$ )

3) Uvědomíme si, že bag 3× seřadíme vzhledem ke  $T^*$  (0,5% pozorování z dvourozměrného normálního rozdělení leží mimo fence)

4) Uvědomíme si, že loop - konvexní obal všech bodů mezi fence a bag.

5) Ostatní pozorování jsou označeno jako odlehla (outliers).

outlier  
x

