

NMAI059 Pravděpodobnost a statistika 1

10. přednáška

Robert Šámal

Přehled

Statistika – model situace

Statistika – bodové odhady

Statistika – intervalové odhady

Náhodný výběr

① ▶ bez vracení
 $\Omega = \{\text{všechny } n\text{-tice obyvatel ČR}\}$
Pro $\omega = (\omega_1, \dots, \omega_n)$ zvolíme $X_i = I(\omega_i \text{ je levák})$.

② ▶ s vracením
 $\Omega = \{\text{všechny } n\text{-tice obyvatel ČR, mohou se opakovat}\}$
Pro $\omega = (\omega_1, \dots, \omega_n)$ zvolíme $X_i = I(\omega_i \text{ je levák})$.

- ▶ varianty (stratifikovaný výběr)
Chceme adekvátně reprezentovat různé podmnožiny (dané věkem, bydlištěm, ...).
Nebudeme dále zkoumat.

1. zkušeb

x_1, \dots, x_n
je to
nezávislé

n volček
↓

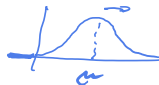
1 = 2

Statistika – model

- ▶ nezávislá měření – hodnoty n.n.v. $X_1, \dots, X_n \sim F$
náhodný výběr s distribuční funkcí F s rozsahem n
z distribuce

- ▶ neparametrické modely: povolujeme velkou třídu F
(všechny možné distribuce)

- ▶ parametrické modely: $F \in \{F_\vartheta : \vartheta \in \Theta\}$



- ▶ příklady:

- ▶ $Pois(\lambda)$ (parametr $\vartheta = \lambda$, $\Theta = \mathbb{R}^+$)
- ▶ $U(a, b)$ (parametr $\vartheta = (a, b)$, $\Theta = \mathbb{R}^2$)
- ▶ $N(\mu, \sigma^2)$ (parametr $\vartheta = (\mu, \sigma)$, $\Theta = \mathbb{R} \times \mathbb{R}^+$)

- ▶ „Všechny modely jsou špatné, ale některé jsou užitečné.“
(George Box)

Zkoumané úlohy – cíle konfirmační analýzy (confirmatory data analysis)

- ▶ bodové odhady
- ▶ intervalové odhady
- ▶ testování hypotéz
- ▶ (lineární) regrese



$$X_1, \dots, X_n$$

- ▶ *statistika* – libovolná funkce náhodného výběru, tj. např. aritmetický průměr, medián, maximum, atd.
Tj. $T = T(\underline{X_1, \dots, X_n})$.

Přehled

Statistika – model situace

Statistika – bodové odhady

Statistika – intervalové odhady

Výběrový průměr a rozptyl

↳ příklady statistek

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\bar{S}_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

$$\hat{S}_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

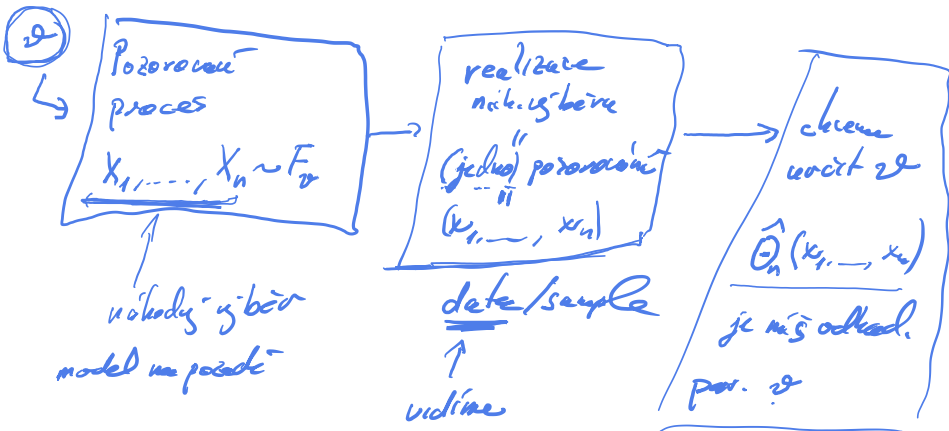
odhad rozptylu

Odhad

θ velká theta
 $\hat{\theta} = \theta$ malá theta

Definice

Odhad je libovolná statistika.



Vlastnosti bodových odhadů

parametry $\theta_1, \dots, \theta_k \rightarrow$ dostává u.v.

obecněji na jekt
fce $g(\vartheta)$

Definice

Odhad $\hat{\Theta}_n = \hat{\Theta}_n(X_1, \dots, X_n)$ parametru ϑ je

- ▶ neustranný (unbiased) – pokud $\vartheta = \mathbb{E}(\hat{\Theta}_n)$ (pro každé ϑ)
- ▶ asymptoticky neustranný (asymptotically unbiased)
– pokud $\vartheta = \lim_{n \rightarrow \infty} \mathbb{E}(\hat{\Theta}_n)$
- ▶ konzistentní (consistent) – pokud $\hat{\Theta}_n \xrightarrow{P} \vartheta$.
- ▶ vychýlení (bias) $bias_{\vartheta}(\hat{\Theta}_n) := \mathbb{E}(\hat{\Theta}_n) - \vartheta$
- ▶ střední kvadratická chyba (mean squared error, MSE) je

$$MSE := \mathbb{E}((\hat{\Theta}_n - \vartheta)^2)$$

popř. chybovost

neuváží konst.

$$\Rightarrow \mathbb{E}(\hat{\Theta}_n - \mu) = 0$$

chyba odhadu

$$\mu = \mathbb{E}\hat{\Theta}_n$$

Věta

$$MSE = bias_{\vartheta}(\hat{\Theta}_n)^2 + var_{\vartheta}(\hat{\Theta}_n)$$

$$\mathbb{E}((\hat{\Theta}_n - \mu) + (\mu - \vartheta))^2 = \mathbb{E}((\hat{\Theta}_n - \mu)^2 - 2(\hat{\Theta}_n - \mu)(\mu - \vartheta) + (\mu - \vartheta)^2)$$

$$= \mathbb{E}((\hat{\Theta}_n - \mu)^2) + 2\mathbb{E}(\hat{\Theta}_n - \mu)(\mu - \vartheta) + (\mu - \vartheta)^2$$

Parametry výběrového momentu a rozptylu

X_1, \dots, X_n náhod. úběr $\sim F_{\mu, \sigma^2}$ např. Norm(0)

chceme určit $\mu = \frac{1}{n} \dots$ nej: fce $\sigma^2 = \text{jiná fce}$

Věta

1. \bar{X}_n je konzistentní nestranný odhad $\mu = EX_1 = EX_2 = \dots$
2. \bar{S}_n je konzistentní asymptoticky nestranný odhad σ^2
3. \hat{S}_n je konzistentní nestranný odhad σ^2

(1)
$$\bar{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n}$$

\bar{X}_n je nestranný, tj. $EX_n = \mu$

\implies i asympt. nestr.

$$\frac{EX_1 + EX_2 + \dots + EX_n}{n} = \frac{\mu + \mu + \dots + \mu}{n} = \mu$$

\bar{X}_n je konsist., tj. $\bar{X}_n \xrightarrow{P} \mu$ (střed. 2.V.C. $(\text{var}(\bar{X}_n) = \frac{\sigma^2}{n})$ & čísl. iev.)

$$(2) \quad \bar{S}_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

$$\begin{aligned} \mu &= EX_i \\ \sigma^2 &= \text{var}(X_i) \\ \mu &= E\bar{X}_n \end{aligned}$$

$$E\bar{S}_n = E \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

$$= E \frac{1}{n} \sum_{i=1}^n \left[(X_i - \mu)^2 - 2(X_i - \mu)(\bar{X}_n - \mu) + (\bar{X}_n - \mu)^2 \right]$$

$$E \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - E \frac{2}{n} \sum_{i=1}^n (X_i - \mu)(\bar{X}_n - \mu) + E(\bar{X}_n - \mu)^2$$

$$E \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - 2 \underbrace{E(\bar{X}_n - \mu)(\bar{X}_n - \mu)}_{=0} + \underbrace{E(\bar{X}_n - \mu)^2}_{=0}$$

$$\frac{1}{n} \sum_{i=1}^n \underbrace{E(X_i - \mu)^2}_{\text{var } X_i = \sigma^2} - \underbrace{E(\bar{X}_n - \mu)^2}_{=0}$$

asympt.
nestor.

$$\sigma^2 - \text{var } \bar{X}_n = \sigma^2 - \frac{1}{n} \sigma^2 = \frac{n-1}{n} \sigma^2 \xrightarrow{n \rightarrow \infty} \sigma^2$$

$$E \bar{S}_n = \frac{n-1}{n} \sigma^2 \rightarrow \sigma^2$$

konzistentni ugredek

Bessel korekce

$$(3) \hat{S}_n = \frac{1}{n-1} \sum (-) ^2 = \frac{n}{n-1} \bar{S}_n$$

$$E \hat{S}_n = \frac{n}{n-1} E \bar{S}_n = \sigma^2 \rightarrow \hat{S}_n \text{ je nestrojenj- odhad}$$

Je boljše \hat{S}_n nebo \bar{S}_n ?

→ \hat{S}_n je nestrojenj, \bar{S}_n ue.

→ Kateri nis boljše M.S.E. ?

Metoda momentů

--- } konstrukce bodové odhadů

- ▶ $m_r(\vartheta) := \mathbb{E}(X^r)$ pro $X \sim F_\vartheta$... r -tý moment
- ▶ $\widehat{m}_r(\vartheta) := \frac{1}{n} \sum_{i=1}^n X_i^r$ pro náhodný výběr X_1, \dots, X_n z F_ϑ
... r -tý výběrový moment

Věta

$\widehat{m}_r(\vartheta)$ je nestranný konzistentní odhad pro $m_r(\vartheta)$

stejný jako u výběr. průměru ($r=1$)

$$\mathbb{E} \widehat{m}_r(\vartheta) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i^r) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i^r) = m_r(\vartheta)$$

- ▶ Odhad metodou momentů je řešení soustavy rovnic

$$m_r(\vartheta) = \widehat{m}_r(\vartheta) \quad r = 1, \dots, k.$$

i.n.v. \leadsto realizace

Metoda momentů – příklady

① $X_1, \dots, X_n \sim \text{Ber}(p)$

$$\vartheta = p \in (0, 1)$$

$$m_1(\vartheta) = \mathbb{E}X_1 = \vartheta$$

$$\widehat{m}_1(\vartheta) = \frac{1}{n} (X_1 + \dots + X_n) = \bar{X}_n$$

--- $X_i = \begin{cases} 1 & \text{je levěh} \\ 0 & \text{je pravěh} \end{cases}$

neměřené data

$$x_1, \dots, x_n$$

$$\frac{x_1 + \dots + x_n}{n} = \vartheta$$

$$\hat{\vartheta}_n(x_1, \dots, x_n)$$

② $X_1, \dots, X_n \sim U(0, \vartheta)$

$$m_1(\vartheta) = \mathbb{E}X_1 = \frac{\vartheta}{2}$$

$$\widehat{m}_1(\vartheta) = \bar{X}_n$$

$$\rightarrow \frac{\vartheta}{2} = \frac{x_1 + \dots + x_n}{n}$$

bed. odhad $\frac{x_1 + \dots + x_n}{n} \cdot 2$

statistika $\hat{\vartheta}_n(x_1, \dots, x_n)$

Metoda maximální věrohodnosti (maximal likelihood, ML)

- ▶ náh. výběr $X = (X_1, \dots, X_n)$ z modelu s parametrem ϑ
- ▶ možný výsledek $x = (x_1, \dots, x_n)$
- ▶ ... sdružená pravděpodobnostní funkce $p_X(x; \vartheta)$
- ▶ ... sdružená hustota $f_X(x; \vartheta)$
- ▶ věrohodnost (likelihood) $L(x; \vartheta)$ značí p_X nebo f_X
- ▶ normálně: máme pevné ϑ , a $L(x; \vartheta)$ je funkce x
- ▶ teď: máme pevné x a $L(x; \vartheta)$ je funkce ϑ

Metoda MV (ML):

volíme takové ϑ , pro které je $L(x; \vartheta)$ maximální

Metoda maximální věrohodnosti (maximal likelihood, ML)

► Metoda MV (ML):

volíme takové ϑ , pro které je $L(x; \vartheta)$ maximální

► definujeme také $\ell(x; \vartheta) = \log(L(x; \vartheta))$

► díky nezávislosti je

$$L(x; \vartheta) = p(x_1; \vartheta) \cdot p(x_2; \vartheta) \cdots p(x_n; \vartheta)$$

$$\ell(x; \vartheta) = \sum_{i=1}^n \log p(x_i; \vartheta)$$

← hledá max.

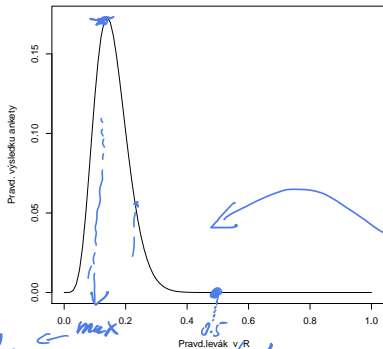
$$0 = \ell'(x; \vartheta) = \sum_{i=1}^n \frac{1}{p(x_i; \vartheta)} \cdot p'(x_i; \vartheta)$$

ML – leváci

6 leváci z 43 měření

X_1, \dots, X_{43}

$X_i =$ "i-tý člověk je levák"



$E X_i = \vartheta =$ procento leváků v ČR

realizace $X_1 = 0, X_2 = 0, X_3 = 1, \dots$
 $\sum_{i=1}^{43} X_i = 6$

$X = X_1 + \dots + X_{43} \sim \text{Bin}(43, \vartheta)$

$$P(X=6) = \binom{43}{6} \vartheta^6 (1-\vartheta)^{37} = f(\vartheta)$$

odkud se pomocí max. věroh.

$X_i =$ i-tý člověk
 kterého se ptáme
 je levák

náh. člověk z ČR

\approx
 \approx
 \approx

$$f'(\vartheta) = e \cdot 6 \cdot \vartheta^5 (1-\vartheta)^{37}$$

i-tý člověk
 na přední straně
 je levák

$$= e \cdot 37 \vartheta^6 (1-\vartheta)^{36}$$

$$= e \vartheta^6 (1-\vartheta)^{37} \left(\frac{6}{\vartheta} - \frac{37}{1-\vartheta} \right) \Rightarrow \vartheta = \frac{6}{43}$$

Přehled

Statistika – model situace

Statistika – bodové odhady

Statistika – intervalové odhady

Intervalové odhady (interval estimation)

- ▶ místo jednoho čísla s nejistým významem vypočítáme z dat interval $[\hat{\Theta}^-, \hat{\Theta}^+]$

Definice

Nechť $\hat{\Theta}^-$, $\hat{\Theta}^+$ jsou n.v. které závisí na náhodném výběru $X = (X_1, \dots, X_n)$. Tyto n.v. určují intervalový odhad, též konfidenční interval o spolehlivosti $1 - \alpha$ ($1 - \alpha$ confidence interval), pokud

$$P(\hat{\Theta}^- \leq \vartheta \leq \hat{\Theta}^+) \geq 1 - \alpha.$$