

NMAI059 Pravděpodobnost a statistika 1

9. přednáška

Robert Šámal

Přehled

Statistika – úvod

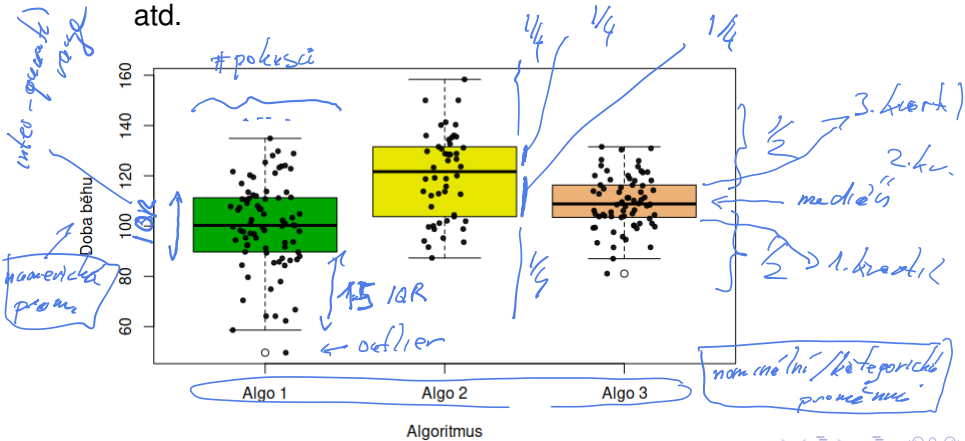
Statistika – bodové odhady

Statistika – intervalové odhady

Intro – explorační analýza dat (exploratory data analysis)

* / ** / ***
25 / 50 / 75
ordinální - prom.

- ▶ posbíráme data (a dáme pozor na systémové chyby – nezávislost, nezaujatost, ...)
- ▶ různé tabulky (třeba v Excelu a spol.)
- ▶ vhodné obrázky: histogram, krabicový diagram (boxplot), atd.



Náhodný výběr

▶ s vracením

▶ bez vracení

$\Omega = \{\text{všechny } n\text{-tice obyvatel ČR}\}$

Pro $\omega = (\omega_1, \dots, \omega_n)$ zvolíme $X_i = I(\omega_i \text{ je levák})$.

$$|\Omega| = N^n$$

$$|\Omega| = N(N-1)\dots(N-k+1)$$

$$n=17, N = \# \text{obyv. ČR}$$

$$= 10.6 \text{ mil}$$

$\rightarrow X_1, \dots, X_n$ nejsou nezávislé

$\rightarrow X_1, \dots, X_n$ jsou nezávislé

(ale rozděl je už zanedbáme)

Statistika – přehled

- ▶ nezávislá měření – hodnoty n.n.v. $X_1, \dots, X_n \sim F$
náhodný výběr s distribuční funkcí F s rozsahem n
..... (n = 20)
- ▶ neparametrické modely: povolujeme velkou třídu F
- ▶ 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.

Zkoumané úlohy – cíle konfirmační analýzy (confirmatory data analysis)

- ▶ bodové odhady
 - ▶ intervalové odhady
 - ▶ testování hypotéz
 - ▶ (lineární) regrese
- ▶ *statistika* – libovolná funkce náhodného výběru, tj. např. aritmetický průměr, medián, maximum, atd.
Tj. $T = T(\underbrace{X_1, \dots, X_n})$.

$$\bar{T}(x_1, \dots, x_n) = \frac{1}{n} \sum x_i$$

$$\bar{T}_2(x_1, \dots, x_n) = \max(x_1, \dots, x_n)$$

Přehled

Statistika – úvod

Statistika – bodové odhady

Statistika – intervalové odhady

Výběrový průměr a rozptyl

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

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

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

Besselova korekce

$\frac{n}{n-1}$

Cíle

některá statistika / $\hat{\Theta}_n$ je konzistentní $\Leftrightarrow \hat{\Theta}_n \rightarrow \mu$ s.j.
(nebo aspoň v pravd.)

Definice

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

- ▶ **neustranný (unbiased)** – pokud $\vartheta = \mathbb{E}(\hat{\Theta}_n)$
- ▶ **asymptoticky neustranný (asymptotically unbiased)** – pokud $\vartheta = \lim_{n \rightarrow \infty} \mathbb{E}(\hat{\Theta}_n)$
- ▶ **vychýlení (bias)** $bias_{\vartheta}(\hat{\Theta}_n) = \mathbb{E}(\hat{\Theta}_n) - \vartheta$
- ▶ **střední kvadratická chyba (mean squared error, MSE)** je $\mathbb{E}((\hat{\Theta}_n - \vartheta)^2)$

$$\mathbb{E}(\hat{\Theta}_n - \vartheta) = \text{bias}$$

Věta

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

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

Parametry výběrového momentu a rozptylu

Vkl 1 $E\bar{X}_n = \frac{1}{n} E(X_1 + \dots + X_n) = \frac{1}{n} E n\mu = \mu$ $\text{zvc} \Rightarrow \bar{X}_n$ je konzistentní.

Věta

- \bar{X}_n je konzistentní nestranný odhad $\mu = E X_i$
- \bar{S}_n je konzistentní asymptoticky nestranný odhad $\sigma^2 = \text{var } X_i$
- \hat{S}_n je konzistentní nestranný odhad σ^2

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

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

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

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

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

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

$= \frac{1}{n} \sum_{i=1}^n [E(X_i - \mu)^2 - (E(X_i - \mu))^2]$

$= \sigma^2 - \frac{\sigma^2}{n} \cdot \frac{n-1}{n} \sigma^2 \xrightarrow{n \rightarrow \infty} \sigma^2$

$\text{var } \bar{X}_n = \frac{\text{var } X_1 + \dots + X_n}{n}$

$= \frac{1}{n^2} \sum \text{var } X_i = \frac{\sigma^2}{n}$

Metoda momentů

- ▶ $\widehat{m}_r(\vartheta) := \mathbb{E}(X^r)$ 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 $m_r(\vartheta)$

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

ZVČ pro X^r : $\widehat{m}_r(\vartheta) \rightarrow m_r(\vartheta)$ S.J.

- ▶ Odhad metodou momentů je řešení soustavy rovnic

\leftarrow specifické \leftarrow specifické & dost

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

\leftarrow z obf. F_ϑ

Metoda momentů – příklady

1) X_1, \dots, X_n n.v. z $Be(p)$

$$m_1(x) = \mathbb{E}X_1 = p$$

$$\widehat{m_1(x)} = \frac{1}{n} (X_{11} + \dots + X_{1n}) = \overline{X_{1i}}$$

matrice
náhodná fronta
dozvíme se nešé
pořadí

2) X_1, \dots, X_n n.v. z $U(0, a)$

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

$$\widehat{m_1(x)} = \overline{X_{1i}}$$

$$m_2(x) = m_2(x)$$

$$\frac{1}{2}a = \overline{X_{1i}}$$

$$a = 2\overline{X_{1i}}$$

statistika

$$\widehat{a} = 2 \cdot \frac{1}{n} \sum_{i=1}^n X_{1i}$$

je asi odhad
prav. a

Metoda maximální věrohodnosti (maximal likelihood, ML)

- ▶ náh. výběr $X = (X_1, \dots, X_n)$ z modelu s parametrem ϑ
- ▶ možný výsledek $x = (x_1, \dots, x_n)$ *p = nej: $\omega \in \mathcal{R}$ $x_1 = X_1(\omega)$ $x_2 = X_2(\omega)$...*
- ▶ ... 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

$$L(x; \vartheta) = P_X(x; \vartheta) = P_X(x_1, \dots, x_n; \vartheta) = \prod_{i=1}^n p(x_i; \vartheta)$$

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

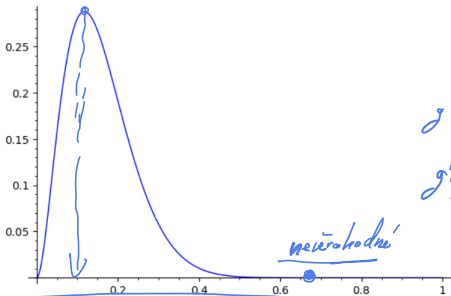
Handwritten annotations: A long horizontal arrow under the sum points to the word "max" written below it. A vertical arrow points from the word "max" to the term $p(x_i; \vartheta)$ in the sum. A double-headed arrow connects the two "max" labels.

ML – leváci

$n = 17, 2$ leváci

$$P_x(\sum X_i = 2) = \binom{17}{2} p^2 (1-p)^{15}$$

```
plot(binomial(17,2)*p^2*(1-p)^15, [0,1])
```



max. $\log p^2 (1-p)^{15}$

$$g(p) = 2 \log p + 15 \log (1-p)$$

$$g'(p) = \frac{2}{p} - \frac{15}{1-p} = 0$$

$$\frac{1-p}{p} = \frac{15}{2}$$

$$\frac{1}{p} - 1 = \frac{15}{2}$$

$$\frac{1}{p} = \frac{17}{2} \quad p = \frac{2}{17}$$

$$x_1, \dots, x_{17} = 1, 0, 0, 1, 0, \dots, 0$$

$$P(X_1 = x_1) = p, P(X_2 = x_2) = 1-p, \dots$$

$$L(x, p) = p (1-p) (1-p) \dots p (1-p) \dots$$

$$l(x, p) = \log p + \log (1-p) + \dots$$

$$X_1 - X_n \sim U(0, a)$$

$$X_1, \dots, X_n \in \mathbb{R}^+$$



$$L(x; a) = f_X(x; a) = \prod_{i=1}^n f_{X_i}(x_i; a) = \frac{1}{a^n} \quad \text{p.kid } x_1 - x_n \in (0, a]$$

$$= \begin{cases} 0, & \text{jezich} \end{cases}$$

met. max. ver.: volme a , abzgl $0 \leq X_1 - X_n \leq a$
 $a \geq \max(X_1 - X_n)$

$$2) \frac{1}{a^n} \text{ co nejiv.}$$

$$a \text{ co nejiv.}$$

$$\hat{a} = \max(X_1 - X_n)$$

Přehled

Statistika – úvod

Statistika – bodové odhady

Statistika – intervalové odhady

Intervalové odhady

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

g ... neznaný param., konstanta

Definice

něk. veličiny

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

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

"
A

"
B

~~$P(\vartheta \in I) = 1 - \alpha$~~

$P(I \ni \vartheta) = 1 - \alpha$

*sp. 95 %
něk. interval vždy s
sp. hodnotou*

Intervalové odhady normální náhodné veličiny