

NMAI059 Pravděpodobnost a statistika 1

11. přednáška

Robert Šámal

Přehled

Statistika – bodové odhady (point estimation)

Statistika – intervalové odhady

Testování hypotéz

Výběrový průměr a rozptyl

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

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

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

něk. výběr
 x_1, \dots, x_n

measure MSE
rozptyl S_n^2

nostrovy - odhad

$$\mathbb{E} \hat{S}_n^2 = \text{rozptyl } X_i$$

$$\text{MSE} = \mathbb{E} (\text{odhad} - \text{skutečnost})^2$$

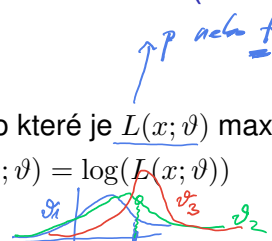
Metoda maximální věrohodnosti (maximal likelihood, ML)

► **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



$$L(x; \vartheta) = L(x_1; \vartheta) \dots L(x_n; \vartheta)$$

$$\ell(x; \vartheta) = \ell(x_1; \vartheta) + \dots + \ell(x_n; \vartheta)$$

derivace... $p = \frac{p}{20}$

hledáme největší max. v řádce

Bin(20,p)	0.2	0.3	0.4	0.45	0.5	0.55	0.6
7	0.0545	0.1643	0.1659	0.1221	0.0739	0.0366	0.0146
8	0.0222	0.1144	0.1797	0.1623	0.1201	0.0727	0.0355
9	0.0074	0.0654	0.1597	0.1771	0.1602	0.1185	0.071
10	0.002	0.0308	0.1171	0.1593	0.1762	0.1593	0.1171
11	0.0005	0.012	0.071	0.1185	0.1602	0.1771	0.1597
12	0.0001	0.0039	0.0355	0.0727	0.1201	0.1623	0.1797
13	0	0.001	0.0146	0.0366	0.0739	0.1221	0.1659
14	0	0.0002	0.0049	0.015	0.037	0.0746	0.1244

k

část tabulky pro p=0.3

ML – další ilustrace

$$X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2) \quad \vartheta = (\mu, \sigma)$$

naměřené n čísel $x_1 = 1.5, x_2 = 2.7, x_3 = \dots \dots \dots x_n$
 $x = (x_1, \dots, x_n)$

chceme: μ, σ

$$L(x_i; \vartheta) = f(x_i; \vartheta) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$l(x_i; \vartheta) = -\frac{1}{2} \left(\frac{x_i - \mu}{\sigma}\right)^2 - \log \sigma - \log \sqrt{2\pi}$$

↑ měření
↑ vstup
↑ stejné μ, σ
↑ $\pi = 3.14$
↑ cíl: μ, σ
↑ abychom $l(x_i; \vartheta)$
bylo max.

$$l(x; \vartheta) = \sum_{i=1}^n l(x_i; \vartheta) = -\frac{1}{2} \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma}\right)^2 - n \log \sigma - n \log \sqrt{2\pi}$$

→ zděruj: palb σ } palb $\mu = 0 \Rightarrow 2$ rovnice pro μ a σ
 → $-\log \sigma$ μ 2 rovnice, $\mu = \sigma$

$$\frac{\partial l}{\partial \mu} = + \frac{1}{\sigma} \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma}\right) \cdot \frac{1}{\sigma} = \frac{1}{\sigma^2} \sum (x_i - \mu) = 0 \Rightarrow \sum x_i = n \cdot \mu$$

$$\frac{\partial l}{\partial \sigma} = + \frac{1}{\sigma} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^3} \cdot (-1) - n \cdot \frac{1}{\sigma} = 0$$

→ odhad $\hat{\mu} = \bar{x}$

$$\sum (x_i - \mu)^2 = n \sigma^2 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

Met. name. $X_1, \dots, X_n \sim N(\mu, \sigma^2)$

parovhane 1, 2, ..., r-1 moment $\left\{ \begin{array}{l} \text{idealar} \\ \text{y-bary} \end{array} \right.$

$$EX_i = \mu = m_1(\vartheta)$$

$$\underline{EX_i^2} = \underbrace{EX_i^2 - (EX_i)^2}_{\text{var}(X_i)} + (EX_i)^2 = \sigma^2 + \mu^2 = m_2(\vartheta)$$

$\stackrel{||}{\sigma^2}$

$$\hat{m}_1(\vartheta) = \bar{X}_n = \frac{X_1 + \dots + X_n}{n} = m_1(\vartheta) = \mu$$

$$\hat{m}_2(\vartheta) = \frac{1}{n} \sum_{i=1}^n X_i^2 = m_2(\vartheta) = \sigma^2 + \mu^2$$

$$\boxed{U(0, \vartheta)}$$

$$\hat{\mu} = \bar{X}_n$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum X_i^2 - (\bar{X}_n)^2$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum (X_i - \bar{X}_n)^2$$

Přehled

Statistika – bodové odhady (point estimation)

Statistika – intervalové odhady

Testování hypotéz

Intervalové odhady (interval estimation)

ϑ - neznámý parameter

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

zvolíme

$\alpha \in (0,1)$

např.: $\alpha = 0.05$

Definice

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

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

ne'hodnota

pevně!

- ▶ tohle jsou tzv. oboustranné odhady
- ▶ jednostranný odhad: $[\hat{\Theta}^-, \infty)$ nebo $(-\infty, \hat{\Theta}^-]$

..... 11% levička
11 ± 3%

(8% , 14%)

"... v tomto interválu"
je správná hodnota
v 95% pokusech

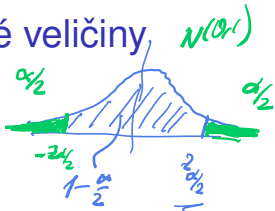
Intervalové odhady normální náhodné veličiny $N(\vartheta, \sigma^2)$

Věta

X_1, \dots, X_n je náhodný výběr z $N(\vartheta, \sigma^2)$.

σ známe, ϑ chceme určit, $\alpha \in (0, 1)$.

Nechť $\Phi(z_{\alpha/2}) = 1 - \alpha/2$. Zvolíme $\hat{\Theta}_n := \bar{X}_n$.



odhadujeme \rightarrow $C_n := [\hat{\Theta}_n - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \hat{\Theta}_n + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$

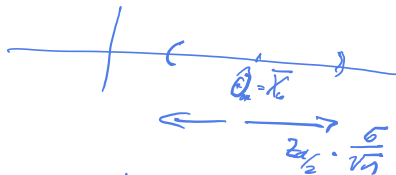
Pak $P(C_n \ni \vartheta) = 1 - \alpha$.

Důkaz.

$$C_n \ni \vartheta \Leftrightarrow |\hat{\Theta}_n - \vartheta| \leq z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

$$\Leftrightarrow \left| \frac{\hat{\Theta}_n - \vartheta}{\sigma/\sqrt{n}} \right| \leq z_{\alpha/2}$$

$Z \sim N(0,1)$



$$\begin{aligned} P(C_n \ni \vartheta) &= P(|Z| \leq z_{\alpha/2}) \\ &= \Phi(z_{\alpha/2}) - \Phi(-z_{\alpha/2}) \\ &= (1 - \alpha/2) - (\alpha/2) = 1 - \alpha. \end{aligned}$$

Intervalové odhady pomocí CLV

Věta

X_1, \dots, X_n je náhodný výběr z rozdělení se střední hodnotou ϑ , rozptylem σ^2 .

σ známe, ϑ chceme určit, $\alpha \in (0, 1)$.

Nechť $\Phi(z_{\alpha/2}) = 1 - \alpha/2$. Zvolíme $\hat{\Theta}_n := \bar{X}_n$.

(ne nutně $N(\vartheta, \sigma^2)$)

$$C_n := \left[\hat{\Theta}_n - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \hat{\Theta}_n + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

Pak $\lim_{n \rightarrow \infty} P(C_n \ni \vartheta) = 1 - \alpha$.

$n \rightarrow \infty$

F

Dle $Z_n = \frac{\hat{\Theta}_n - \vartheta}{\sigma/\sqrt{n}}$ je pribl. $N(0,1)$ dle CLV

$$Z_n \xrightarrow{d} N(0,1)$$

$$\begin{aligned} P(C_n \ni \vartheta) &= P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) \\ &= F_{Z_n}(z_{\alpha/2}) - F_{Z_n}(-z_{\alpha/2}) \rightarrow \Phi(z_{\alpha/2}) - \Phi(-z_{\alpha/2}) \\ &= 1 - \alpha \end{aligned}$$

Studentovo rozdělení

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

$$\hat{\sigma}_n = \sqrt{\hat{S}_n^2} \dots \text{výběr. směro-odch.}$$

- ▶ Nechť $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ ←
▶ Pak $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$ ←
▶ Studentovo t -rozdělení s $n - 1$ stupni volnosti je rozdělení

protože \hat{S}_n^2 má
část σ

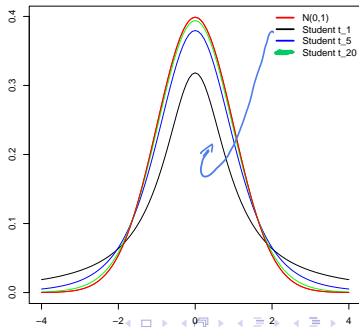
víme z úvodu

$$\text{n.v. } \frac{X_n - \mu}{\hat{S}_n / \sqrt{n}}$$

$\rightarrow N(0,1)$

- ▶ Distribuční funkci budeme značit Ψ_{n-1}
Je v tabulkách,
v R: **pt(x, n-1)**

Důležité → je to stejné
rozdělení pro všechna
 μ, σ



Int. odhady normální n.v. pomocí Studentova t

Věta

$$\alpha = 5\%$$

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

X_1, \dots, X_n je náhodný výběr z $N(\vartheta, \sigma^2)$.

ϑ chceme určit σ neznáme $\alpha \in (0, 1)$. Necht'

$$\Psi_{n-1}(z_{\alpha/2}) = 1 - \alpha/2. \hat{\Theta}_n = \bar{X}_n, \hat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

St. t-odh.

$$C_n := \left[\hat{\Theta}_n - z_{\alpha/2} \frac{\hat{S}_n}{\sqrt{n}}, \hat{\Theta}_n + z_{\alpha/2} \frac{\hat{S}_n}{\sqrt{n}} \right]$$

Pak $P(C_n \ni \vartheta) = 1 - \alpha$.

$$D_k P(C_n \ni \vartheta) = P(|Z| \leq z_{\alpha/2}) = \Psi_{n-1}(z_{\alpha/2}) - \Psi_{n-1}(-z_{\alpha/2})$$

$$Z = \frac{\bar{X}_n - \vartheta}{\hat{S}_n / \sqrt{n}} \quad \text{St. t-odh. s (n-1) stupni.}$$

$$\phi^{-1}(0.975) = 1.96$$

$$\Psi_{19}^{-1}(0.975) = 12.7$$

$$\Psi_{20}^{-1}(0.975) = 2.09$$

Přehled

Statistika – bodové odhady (point estimation)

Statistika – intervalové odhady

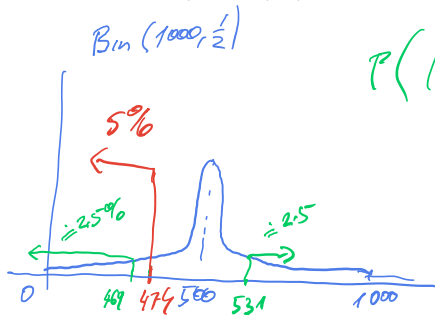
Testování hypotéz

Úvod do testování hypotéz

- ▶ Je naše mince spravedlivá?
 - ▶ Je naše kostka spravedlivá?
 - ▶ Má vylepšený program kratší dobu běhu než původní?
 - ▶ Je léčba nemoci metodou X dobrá? (Lepší než placebo, lepší než metoda Y, ...)
 - ▶ Jsou leváci lepší boxeři?
-
- ▶ dvě hypotézy: H_0 , H_1
 - ▶ H_0 – nulová hypotéza (*null hypothesis*) – značí defaultní, konzervativní model (léčba, mince je spravedlivá)
 - ▶ H_1 – alternativní hypotéza (*alternative hypothesis*) – značí alternativní model „pozoruhodnost“

Testování hypotéz – ilustrace

- ▶ Chceme testovat, zda je mince spravedlivá.
- ▶ Hodíme n -krát mincí, orel padne S -krát.
- ▶ Pokud je $|S - n/2|$ moc velké, tak mince není spravedlivá.



$$P(|S - \frac{n}{2}| \geq 31) = \underline{\underline{5\%}}$$

Neřekneme, že je
mince nespřavedlivá.

(nezamítáme H_0) me chybami
ujiznuvame
5%

psychol. ... 5%

fyzika ... $< 3 \cdot 10^{-3}$ evidence
 $< 3 \cdot 10^{-7}$ objvu.

Testování hypotéz – ilustrace

- ▶ Chceme testovat, zda je mince spravedlivá.
- ▶ H_0 : je spravedlivá
- ▶ H_1 : není spravedlivá („Vědci objevili, že v kasinu byla použita falešná mince.“)
- ▶ Výsledky: zamítneme H_0 /nezamítneme H_0
- ▶ Chyba 1. druhu: chybné zamítnutí. Zamítneme H_0 , i když platí. Trapas.
- ▶ Chyba 2. druhu: chybné přijetí. Nezamítneme H_0 , ale ona neplatí. Promarněná příležitost.
- ▶ Potřebujeme určit k takové, že budeme zamítat H_0 pokud $|S - n/2| > k$.

Testování hypotéz – obecný postup

- ▶ Vybereme vhodný statistický model.
- ▶ Volíme *hladinu významnosti (significance level)* α : pravd. chybného zamítnutí H_0 . Typicky $\alpha = 0.05$.
- ▶ Určíme *testovou statistiku* $S = h(X_1, \dots, X_n)$, kterou budeme určovat z naměřených dat.
- ▶ Určíme *kritický obor (rejection region)* – množinu W .
- ▶ Naměříme hodnoty x_1, \dots, x_n náh. veličin X_1, \dots, X_n .
- ▶ Rozhodovací pravidlo: zamítneme H_0 pokud $h(x_1, \dots, x_n) \in W$.
- ▶ $\alpha = P(h(X) \in W; H_0)$
- ▶ $\beta = P(h(X) \notin W; H_1)$... *síla testu*
- ▶ často α nevolíme předem, ale spočítáme *p-hodnotu*: minimální α , pro které bychom H_0 zamítlí.

Testování hypotéz – příklad

- ▶ X_1, \dots, X_n náhodný výběr z $N(\vartheta, \sigma^2)$
- ▶ σ^2 známe
- ▶ $H_0 : \vartheta = 0$ $H_1 : \vartheta \neq 0$

Testování hypotéz – příklad

- ▶ X_1, \dots, X_{n_1} náhodný výběr z $Ber(\vartheta_X)$
- ▶ Y_1, \dots, Y_{n_2} náhodný výběr z $Ber(\vartheta_Y)$
- ▶ $H_0 : \vartheta_X = \vartheta_Y$ $H_1 : \vartheta_X \neq \vartheta_Y$

p-hacking

- ▶ napřed získáme data, pak v nich hledáme zajímavosti
- ▶ když máme dost dat, tak tam nějaké budou „shodou okolností“
- ▶ *reprodukovatelnost* – po explorační analýze dat uděláme nezávislý sběr dat a ten analyzujeme konfirmačně
- ▶ nebo dopředu náhodně rozdělíme data na část pro tvorbu hypotéz a část pro jejich potvrzení . . . jednoduchý případ křížové validace (cross validation)

