

NMAI059 Pravděpodobnost a statistika 1

12. přednáška

Robert Šámal

Statistika – Co už víme

- ▶ základní nastavení: uvažujeme náhodný výběr X_1, \dots, X_n z distribuce F_ϑ — popisuje proces měření, jak mohlo měření proběhnout
 - ▶ naměříme data – konkrétní čísla, tzv. realizaci náhodného výběru x_1, \dots, x_n — jak naše měření skutečně proběhlo
1. bodové odhady: máme určit co nejlepší číslo, odhad pro parametr ϑ , nebo nějakou jeho funkci $g(\vartheta)$.
 2. intervalové odhady: máme určit interval, ve kterém parametr ϑ pravděpodobně leží
 3. testování hypotéz

$$\bar{x}_n - \sigma \quad \hat{\vartheta}_n = \bar{x}_n \quad \bar{x}_n + \sigma$$

Přehled

$$\begin{aligned} n.u. z_2 \quad X &\sim N(\mu, \sigma^2) \\ Y &\sim N(\mu', \sigma'^2) \end{aligned} \Rightarrow -Y \sim N(-\mu', \sigma'^2)$$

↓

$$\underline{X+Z} \sim N(\mu+\mu', \sigma^2+\sigma'^2)$$

Testování hypotéz

$$f_Z(z) = \int f_X(x) f_Y(z-x)$$

Testy dobré shody

Lineární regrese

Testování hypotéz – ilustrace

- ▶ Chceme testovat, zda je mince spravedlivá.
- ▶ H_0 : je spravedlivá *očekávaný stav světa*
- ▶ H_1 : není spravedlivá *překvapivé zjištění* („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

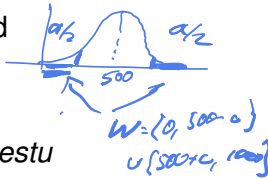
$$H_{\text{muce}}: X_1, \dots, X_n \sim \text{Ber}(p)$$
$$H_0: p = \frac{1}{2} \quad H_1: p \neq \frac{1}{2}$$

- ▶ Vybereme vhodný statistický model.
- ▶ Volíme *hladinu významnosti (significance level)* α : pravd. chybného zamítnutí H_0 . Typicky $\alpha = 0.05$.
- ▶ Určíme *testovou statistiku* $T = h(X_1, \dots, X_n)$, 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 $h(x_1, \dots, x_n) \in W$.
- ▶ $\alpha = P(h(X) \in W; H_0)$
- ▶ $\beta = P(h(X) \notin W; H_1)$... $1 - \beta$ je tzv. *síla testu*
- ▶ často α nevolíme předem, ale spočítáme tzv. *p-hodnotu*: minimální α , pro které bychom H_0 zamítlí.

F_T značka

$$T = \sum_{i=1}^n X_i \sim \text{Bin}(n, p)$$

$n=1000$



volíme c : $P(T < 500-c) = \alpha/2$

$$F_T(500-c) = \frac{\alpha}{2}$$
$$\left\lfloor F_T^{-1}\left(\frac{\alpha}{2}\right) \right\rfloor = 500-c$$

Testování hypotéz – příklad

mírná teplota
okresu $\mu = 5^\circ\text{C}$

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

$$\underline{T} = \frac{x_1 + \dots + x_n}{n} = \bar{X}_n \sim N(\vartheta, \sigma^2/n)$$

$$\underline{S} = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$



$$W = \{s \in \mathbb{R} : |s| > z_{\alpha/2}\}$$

$$z_{\alpha/2} = \Phi^{-1}(1 - \frac{\alpha}{2}) = \underline{1.96}$$

$$\alpha = 0.05$$

... μ známe

1. z důvěryže známe
známe pro vzpěl μ

2. podání μ ... $\vartheta = \mu = 5$
a σ^2 také známe \Rightarrow

známe μ a σ^2

ZNAMÉ!

Testování hypotéz – příklad dvojběrového testu

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

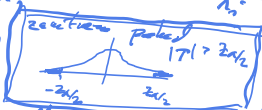
↪ dvě různé léčby
 X_i, Y_j ... zda byla
 léčba úspěšná

$$\hat{\vartheta}_X = \frac{X_{11} + \dots + X_{1n_1}}{n_1}$$

... odhad ϑ_X příbl. norm. rozd. CLV

$$\hat{\vartheta}_Y = \frac{Y_{11} + \dots + Y_{1n_2}}{n_2}$$

... odhad ϑ_Y



$X_i =$ { i-tý pacient, který dostal I. léčbu se uzdravil }

$$Z := \hat{\vartheta}_X - \hat{\vartheta}_Y$$

pohlád 121 je větší, tak H_0 zamítáme

$Y_j =$ { j-tý pac. ... II. léčba ... se uzdravil }

Pohlád pohlád H_0

$$E\hat{\vartheta}_X = E\hat{\vartheta}_Y \Rightarrow EZ = 0$$

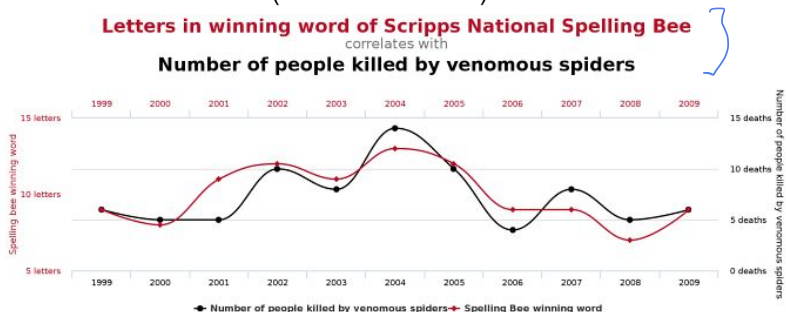
$\vartheta = \vartheta_X = \vartheta_Y$ Z je příbl. $N(0, \sigma^2)$ neznačím $\vartheta(1-\vartheta)$

$$\hat{\sigma}^2 = \text{var}(Z) = \text{var}(\hat{\vartheta}_X) + \text{var}(\hat{\vartheta}_Y) = \frac{\text{var } X_i}{n_1} + \frac{\text{var } Y_j}{n_2} = \frac{\vartheta(1-\vartheta)}{n_1} + \frac{\vartheta(1-\vartheta)}{n_2} = \frac{\vartheta(1-\vartheta)}{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$\hat{\sigma} = \frac{\sum X_i + \sum Y_j}{n_1 + n_2} \dots \text{odhad } \vartheta \Rightarrow \hat{\sigma}^2 = \left(\frac{1}{n_1} + \frac{1}{n_2}\right) \hat{\sigma}(\hat{\sigma}) \Rightarrow T := \frac{\hat{\vartheta}_X - \hat{\vartheta}_Y}{\hat{\sigma}} = \frac{\hat{\vartheta}_X - \vartheta}{\sqrt{\dots}} \sim N(0,1)$$

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)



Přehled

Testování hypotéz

Testy dobré shody

Lineární regrese

χ_k^2 – rozdělení χ -kvadrát

$\rightarrow \text{var}(Z_i) = 1 \Rightarrow E Z_i^2 = 1$

$$EQ = E Z_1^2 + \dots + E Z_k^2 = k$$
$$\text{var } Q = \text{var } Z_1^2 + \dots + \text{var } Z_k^2$$

Definice

$Z_1, \dots, Z_k \sim N(0, 1)$ n.n.v. Rozdělení náhodné veličiny

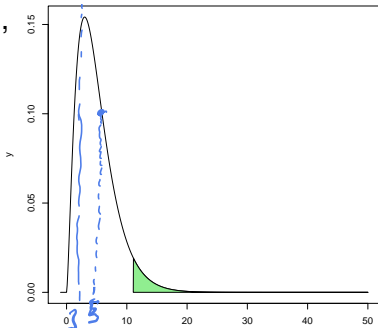
$$\text{var } Z_i^2 = E Z_i^4 - (E Z_i^2)^2$$

||
1

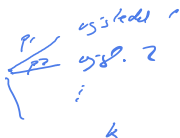
$$Q = Z_1^2 + \dots + Z_k^2$$

se nazývá χ -kvadrát s k stupni volnosti. (Opravdu k !)

- ▶ $E(Q) = k$ (lehké)
- ▶ $\text{var}(Q) = 2k$ (pro info, netřeba pamatovat)
- ▶ hustota jde napsat vzorcem, jde najít např. na Wikipedii
- ▶ $Q \doteq N(k, 2k)$
pro velká k (CLV)



Multinomické a kategoriální rozdělení



Definice

Dána $p_1, \dots, p_k \geq 0$ tak, že $p_1 + p_2 + \dots + p_k = 1$.

n -krát zopakuj pokus, kde může nastat jedna z k možností, i -tá má pravděpodobnost p_i .

$X_i :=$ kolikrát nastala i -tá možnost (X_1, \dots, X_k) má multinomické rozdělení s parametry $n, (p_1, \dots, p_k)$.

- ▶ triviální případ: $X_i =$ počet hodů kostkou, kdy padlo i
- ▶ důležitý případ: $X_i =$ počet výskytů i -tého písmene, i -tého slovního druhu, ...

▶
$$P(X_1 = x_1, \dots, X_k = x_k) = \underline{\binom{n}{x_1, \dots, x_k} p_1^{x_1} \dots p_k^{x_k}}$$

Pearsonova χ^2 statistika

- ▶ (X_1, \dots, X_k) – multinomické rozdělení s parametry $n, (p_1, \dots, p_k)$ jako minule
- ▶ $E_i := \mathbb{E}(X_i) = np_i$
- ▶ Pearsonova χ^2 statistika je funkce O_i ... observed

$$T := \sum_{i=1}^k \frac{(X_i - E_i)^2}{E_i}$$

k pevné, $n \rightarrow \infty$

▶ **Věta** $T \xrightarrow{d} \chi_{k-1}^2$

Dle pro $k=2$ $X_1 + X_2 = n, p_1 + p_2 = 1, E_1 = np_1$

$$T = \frac{(X_1 - E_1)^2}{E_1} + \dots = \frac{(X_1 - np_1)^2 p_2}{np_1 p_2} + \frac{(n - X_1 - n(1-p_1))^2 p_1}{np_2 p_1} = \frac{(X_1 - np_1)^2 \cdot (p_1 p_2)}{np_1 p_2}$$

$$= \left(\frac{X_1 - np_1}{\sqrt{np_1(1-p_1)}} \right)^2 \xrightarrow{d} N(0,1) \text{ podle CLT}$$

Test dobré shody (goodness of fit)

- ▶ (X_1, \dots, X_k) – multinomické rozdělení s parametry $n, \vartheta = (\vartheta_1, \dots, \vartheta_k)$ jako minule
- ▶ n známe, ϑ neznáme.
- ▶ Hypotéza $H_0: \vartheta = \vartheta^*$
- ▶ $E_i := n\vartheta_i^*$ pro všechna i
- ▶ Použijeme statistiku $\chi^2 = T := \sum_{i=1}^k \frac{(X_i - E_i)^2}{E_i}$ \xrightarrow{d} χ_{k-1}^2
- ▶ Hypotézu H_0 zamítneme, pokud $T > \gamma$
- ▶ $\gamma := F_Q^{-1}(1 - \alpha)$, kde $Q \sim \chi_{k-1}^2$ $n \rightarrow \infty$
- ▶ $P(\text{chyba prvního druhu}) = P(T > \gamma; H_0) \rightarrow P(Q > \gamma) = \alpha$



Test dobré shody – příklad

- ▶ Házíme opakovaně kostkou. Jednotlivá čísla padla s četností 92, 120, 88, 98, 95, 107.
- ▶ Je kostka spravedlivá?

$$n = 92 + 120 + \dots = 600 \quad \mathcal{P}^* = \left(\frac{1}{6}, \dots, \frac{1}{6}\right), \quad E_i = n \frac{1}{6} = 100$$
$$T = \sum_{i=1}^6 \frac{(X_i - 100)^2}{100} = \frac{(92-100)^2}{100} + \frac{(120-100)^2}{100} + \frac{12^2}{100} + \frac{2^2}{100} + \frac{5^2}{100} + \frac{7^2}{100}$$
$$= \frac{(8^2 + 20^2 + 12^2 + 2^2 + 5^2 + 7^2)}{100} = \underline{\underline{6.86}}$$

$$Q \sim \chi_5^2 \quad F_Q^{-1}(0.95) = 11.1 =: j$$

$$p\text{-hodnota: } 1 - F_Q(6.86) = 1 - 0.72 = 0.28$$

G-test, ...

Další rozšíření

- ▶ Pro zkoumání rozdělení libovolné n.v. Y můžeme vybrat „příhrádky“ B_1, \dots, B_k (rozklad \mathbb{R}) a zkoumat, kolikrát je $Y \in B_i$
- ▶ Obdobný test pro nezávislost (diskrétních) náhodných veličin



Přehled

Testování hypotéz

Testy dobré shody

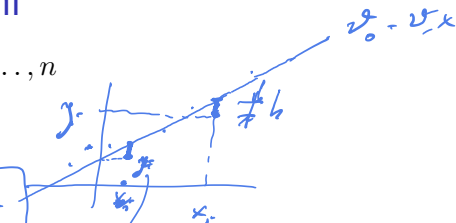
Lineární regrese

Lineární regrese – zadání

- ▶ data: (x_i, y_i) pro $i = 1, \dots, n$
- ▶ cíl: $y = \vartheta_0 + \vartheta_1 x$

skládkou:

$$y = \vartheta_0 + \vartheta_1 x + \text{neh. chyba}$$



- ▶ měříme pomocí kvadratické odchylky

$$\sum_{i=1}^n (y_i - (\vartheta_0 + \vartheta_1 x_i))^2$$

Lineární regrese – řešení

- ▶ Minimalizujeme výraz

$$\sum_{i=1}^n (y_i - (\vartheta_0 + \vartheta_1 x_i))^2$$

- ▶ řešení: Optimální parametry jsou

$$\hat{\vartheta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{\vartheta}_0 = \bar{y} - \vartheta_1 \bar{x},$$

kde $\bar{x} := (x_1 + \dots + x_n)/n$, $\bar{y} := (y_1 + \dots + y_n)/n$.

$\frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})$ ——— y-levová kovariance

$\frac{1}{n-1} \sum (x_i - \bar{x})^2$ ———> y-levová variace

Lineární regrese – proč součet čtverců?

- ▶ Předpokládejme, že x_1, \dots, x_n jsou pevná, y_i je zvoleno jako hodnota náhodné veličiny

$$Y_i = \vartheta_0 + \vartheta_1 x_i + W_i$$

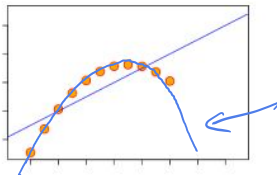
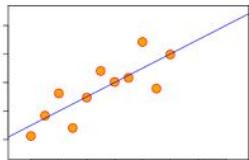
- ▶ $W_i \sim N(0, \sigma^2)$ pro všechna i ; W_1, \dots, W_k nezávislé.
- ▶ metoda maximální věrohodnosti:

$$L(y; \vartheta) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i - \vartheta_0 - \vartheta_1 x_i)^2}{2\sigma^2}}$$

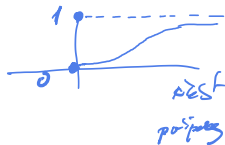
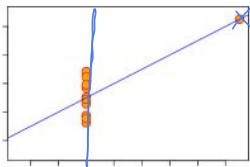
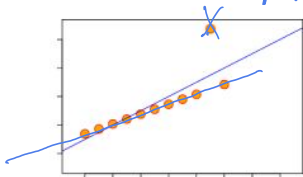
- ▶ $\ell(y; \vartheta) = \log L(y; \vartheta) = a + b \sum_{i=1}^n \underline{\underline{(y_i - \vartheta_0 - \vartheta_1 x_i)^2}}$

Limity regrese

bodie



čas na
připravu



(data: Francis Anscombe 1973, obrázek: wikieditor Schutz)

► nelineární regrese

→ logistická regrese

..... ~~konstante~~

x_i -- reálné číslo

y_i -- 0/1

Simpson's paradox

Treatment Stone size	Treatment A	Treatment B
Small stones	Group 1 93% (81/87)	Group 2 87% (234/270)
Large stones	Group 3 73% (192/263)	Group 4 69% (55/80)
Both	78% (273/350)	83% (289/350)

