

NMAI059 Pravděpodobnost a statistika 1

11. přednáška

Robert Šámal

Přehled

Testy dobré shody

Lineární regrese

Bayesovská statistika

χ_k^2 – rozdělení χ -kvadrát

z_i^2 vel. se sk. b. 1
rozdy. 2 (2)

$$z_i: \Omega \rightarrow \mathbb{R}$$

Definice

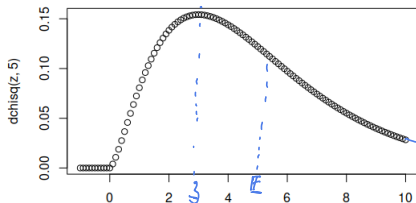
$Z_1, \dots, Z_k \sim N(0, 1)$ n.n.v. Rozdělení náhodné veličiny

$$\text{var } z_i = \mathbb{E}z_i^2 - (\mathbb{E}z_i)^2 = \mathbb{E}z_i^2 \quad Q = \underbrace{Z_1^2}_{\dots} + \dots + \underbrace{Z_k^2}_{\dots} : \Omega \rightarrow \mathbb{R}$$

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

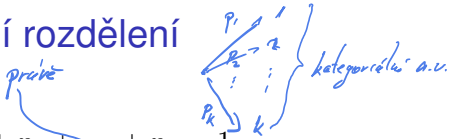
- ▶ $\mathbb{E}(Q) = k$ (lehké) $\mathbb{E}Q = \mathbb{E} \sum_i z_i^2 = \sum_i \mathbb{E}z_i^2 = \sum_i 1 = k$
- ▶ $\text{var}(Q) = 2k$ (pro info, netřeba pamatovat) $\mathbb{E}Q^2 = \mathbb{E}(\sum_i z_i^2)(\sum_j z_j^2)$
- ▶ hustota jde napsat vzorcem, jde najít např. na Wikipedii
- ▶ Q je pro velké n blízke $N(k, 2k)$ (CLV)

$$\begin{aligned} &= \mathbb{E} \sum_{i,j} z_i^2 \cdot z_j^2 \\ &= \sum_{i,j} \mathbb{E} z_i^2 \cdot z_j^2 \\ &\leftarrow \begin{cases} i=j & \mathbb{E} z_i^2 \cdot \mathbb{E} z_i^2 = 1 \cdot 1 \\ i \neq j & \mathbb{E} z_i^2 \cdot \mathbb{E} z_j^2 = 1 \cdot 1 \end{cases} \end{aligned}$$



χ_5^2

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 $n, (p_1, \dots, p_k)$.

$$n = \sum_{i=1}^k X_i$$

$$0 \leq X_i \leq n$$

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

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

↓ dostáváme množinovou rozsl. $n, (p_1, p_2, p_3, \dots, p_k)$

multinac. koef. $\frac{n!}{x_1! \dots x_k!} = \# \text{ uspoř. } x_1 \times 1, x_2 \times 2, \dots, x_k \times k \text{ do řady}$

Pro kontrolu

$$\sum_{x_1, \dots, x_k} P(X_1 = x_1, \dots, X_k = x_k) = \sum_{x_1, \dots, x_k} \binom{n}{x_1, \dots, x_k} p_1^{x_1} \dots p_k^{x_k} = (p_1 + \dots + p_k)^n = 1$$

→ $X_i \sim \mathcal{B}_m(n, p_i)$
→ slovíčím ta 2 do jedné

kategoriální

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

$\mu = \mathbb{E} X_i$
 $\sigma^2 = \text{Var } X_i = np_i(1-p_i)$

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

$T \xrightarrow{d} Z_i^2$
 $Z_i \sim N(0,1)$

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

$\otimes = \frac{(X_1 - np_1)^2}{np_1 p_2} (p_1 + p_2)$

$T = \left(\frac{X_1 - np_1}{\sqrt{np_1(1-p_1)}} \right)^2$

$k=2$

$X_1 + X_2 = n$

$p_1 + p_2 = 1$

$T = \frac{(X_1 - E_1)^2}{E_1} + \frac{(X_2 - E_2)^2}{E_2}$

$= \frac{(X_1 - np_1)^2}{np_1} + \frac{(n - X_1 - n(1-p_1))^2}{np_2} = \frac{(X_1 - np_1)^2}{np_1 p_2} + \frac{(X_1 - np_1)^2}{np_2 p_1} = \otimes$

$\frac{X_i - \mu}{\sigma} \xrightarrow{d} N(0,1)$ [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 $T := \sum_{i=1}^k \frac{(X_i - E_i)^2}{E_i}$

▶ Hypotézu H_0 zamítneme, pokud $T > \gamma$

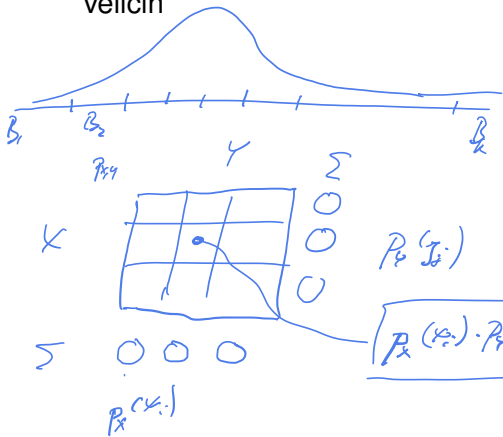
▶ $\gamma := F_Q^{-1}(1 - \alpha)$, kde $Q \sim \chi_{k-1}^2$

▶ $P(\text{chyba prvního druhu}) = P(T > \gamma; H_0) \xrightarrow{\alpha \rightarrow \infty} P(Q > \gamma) = \alpha$



Další rozšíření

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



\rightsquigarrow m. Hranou. rozd.

$$P_i = P(Y \in \underline{B}_i)$$

probl. počítka: $E_{ij} > 5$,
jinak nefunguje

$$\rightarrow T = \sum_{i,j} \frac{(X_{ij} - E_{ij})^2}{E_{ij}}$$

Přehled

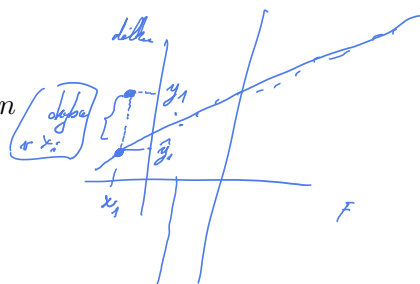
Testy dobré shody

Lineární regrese

Bayesovská statistika

Lineární regrese – zadání

- ▶ data: (x_i, y_i) pro $i = 1, \dots, n$
- ▶ cíl: $y = \underline{\vartheta}_0 + \underline{\vartheta}_1 x$



- ▶ měříme pomocí kvadratické odchylky

minimálně číselně

$$\sum_{i=1}^n (y_i - (\underline{\vartheta}_0 + \underline{\vartheta}_1 x_i))^2$$

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

(bez dk - viz PCA/LA)

↑
y korese průměry

→ rozptyl

ne~~z~~ korr. kovariance

Lineární regrese – proč součet čtverců?

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

bodové odhady pro ϑ_0, ϑ_1

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

*ideální svět → chyba
 $W_i \sim \mathcal{N}(\vartheta_0, \vartheta_1, \sigma^2)$*

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

*ϑ_0 – hodnota W_i
 $N(0, \sigma^2)$*

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

$\varphi(y_i; (\vartheta_0 + \vartheta_1 x_i, \sigma^2))$

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

max

$$= n \log \frac{1}{\sigma \sqrt{2\pi}}$$

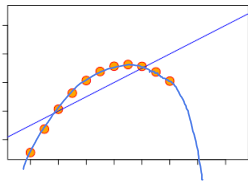
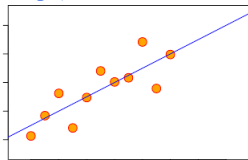
$$- \frac{1}{2\sigma^2} \sum ()^2$$

min

Limity regrese

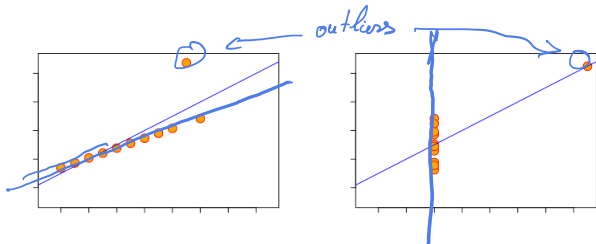
stejně právný, koefice

OK



$$y = g(x)$$

↑
kvadratická!

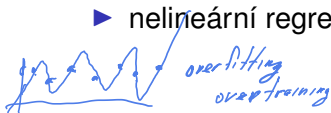


(data: Francis Anscombe 1973, obrázek: wikieditor Schutz)

► nelineární regrese

$$\sum (y_i - g(x_i))^2$$

hledáme g , které minimum součet čtv.



Přehled

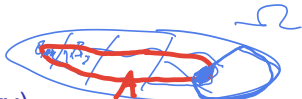
Testy dobré shody

Lineární regrese

Bayesovská statistika

klasická/frekventistická statistika

Bayesova věta



Věta (Bayesova pro jevy)

Pokud B_1, B_2, \dots je rozklad Ω , $A \in \mathcal{F}$ a $P(A), P(B_j) > 0$, tak

$$P(B_j | A) = \frac{P(A | B_j)P(B_j)}{\sum_i P(A | B_i)P(B_i)} \cdot \frac{P(A \cap B_j)}{P(A)}$$

(sčítance s $P(B_i) = 0$ považujeme za 0).

- ▶ apriorní, posteriorní pravděpodobnost (prior, posterior)

Věta (Bayesova pro diskrétní náhodné veličiny)

X, Y jsou diskrétní n.v.

$$P(Y=y | X=x) = P_{Y|X}(y|x) = \frac{p_{X|Y}(x|y)p_Y(y)}{\sum_{y'} p_{X|Y}(x|y')p_Y(y')}$$

(sčítance s $p_Y(y') = 0$ považujeme za 0).

Bayesova věta

Věta (Bayesova pro spojité náhodné veličiny)

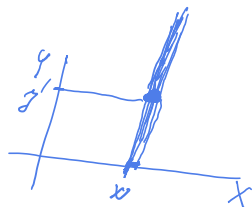
X, Y jsou spojité n.v., které mají hustotu f_X , f_Y i sdruženou hustotu $f_{X,Y}$

$$f_{Y|X}(y|x) = \frac{f_{X|Y}(x|y)f_Y(y)}{\int_{y' \in \mathbb{R}} f_{X|Y}(x|y')f_Y(y')dy'}$$

(sčítance s $f_Y(y') = 0$ považujeme za 0).

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X|Y}(x|y) \cdot f_Y(y)}{\int_{y'} f_{X|Y}(x|y') \cdot f_Y(y') dy'}$$

$$f_X(x) = \int_{y'} f_{X,Y}(x,y') dy' = \int_{y'} f_{X|Y}(x|y') \cdot f_Y(y') dy'$$



Bayesova věta – zbylé varianty

- ▶ X je spojitá, Y diskrétní

$$p_{Y|X}(y|x) = \frac{f_{X|Y}(x|y)p_Y(y)}{\sum_{y' \in \text{Im}Y} f_{X|Y}(x|y')p_Y(y')}.$$

- ▶ X je diskrétní, Y spojitá

$$f_{Y|X}(y|x) = \frac{p_{X|Y}(x|y)f_Y(y)}{\int_{y' \in \mathbb{R}} p_{X|Y}(x|y')f_Y(y')dy'}.$$