

NMAI059 Pravděpodobnost a statistika 1

12. přednáška

Robert Šámal

Přehled

Bayesovská statistika

Srovnání dvou přístupů ke statistice

Frekventistický/klasický přístup

- ▶ Pravděpodobnost je dlouhodobá frekvence (z 6000 hodů kostkou padla šestka 1026-krát). Je to objektivní vlastnost reálného světa.
- ▶ Parametry jsou pevné, neznámé konstanty. Nelze o nich říkat smysluplné pravděpodobnostní výroky.
- ▶ Navrhujeme statistické procedury tak, aby měly žádané dlouhodobé vlastnosti. Např. 95 % z našich intervalových odhadů pokryje neznámý parametr.

*približně $\frac{1}{6} \cdot 6000 \times$
 $P = \frac{1026}{6000}$ *

Bayesovský přístup

- ▶ Pravděpodobnost popisuje, jak moc věříme nějakému jevu, jak moc jsme ochotní se vsadit. (Pravděpodobnost, že Thomas Bayes měl 18. prosince 1760 šálek čaje, je 90 %.)
- ▶ Můžeme vyslovovat pravděpodobnostní výroky i o parametrech (třebaže jsou to pevné konstanty).
- ▶ Spočítáme distribuci ϑ a z ní tvoříme bodové a intervalové odhady, atd.

Bayesovská metoda – základní popis

"flat prior" ... konst. disto.
 $\Theta \sim U(0,1)$

- ▶ neznámý parametr považujeme za náhodnou veličinu Θ
- ▶ zvolíme *apriorní distribuci* (prior distribution), neboli hustotu pravděpodobnosti $f_{\Theta}(\vartheta)$ nezávislou na datech.
- ▶ zvolíme statistický model $f_{X|\Theta}(x|\vartheta)$, který popisuje, co naměříme (s jakou pravděpodobností), v závislosti na hodnotě parametru
- ▶ poté, co pozorujeme hodnotu $X = x$, spočítáme *posteriorní distribuci* (posterior distribution) $f_{\Theta|X}(\vartheta|x)$
- ▶ z té pak odvodíme, co potřebujeme např. najdeme a, b , aby $\int_a^b f_{\Theta|X}(\vartheta|x) d\vartheta \geq 1 - \alpha$

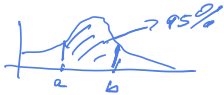
$P_{\Theta}(\vartheta)$
pro diskret.
n.v.

evolu.
podm. pravd. fce

podle Bayesovy
věty

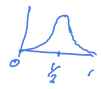
$$P_{X|\Theta}(x|\vartheta)$$

$$P_{\Theta|X}(\vartheta|x)$$



$$P(\Theta \in (a,b) | X=x) \geq 1 - \alpha$$

zabrany předchozí
zobraz



- ▶ $\vartheta = \theta$ malá théta, Θ je velká théta

Bayesova věta

$P(X=x | \Theta=\vartheta) \dots$ podmíněná pravděp. dle realizace příjtu
v klas. statist. $P(X=x; \vartheta) \dots$ normální pravděp.,
 ϑ je parametr

Věta (Bayesova pro diskrétní náhodné veličiny)

X, Θ jsou diskrétní n.v.

$$P(X=x | \Theta=\vartheta) P(\Theta=\vartheta)$$

$$P(\Theta=\vartheta | X=x)$$

$$p_{\Theta|X}(\vartheta|x) = \frac{p_{X|\Theta}(x|\vartheta)p_{\Theta}(\vartheta)}{\sum_{\vartheta' \in I_{m\Theta}} p_{X|\Theta}(x|\vartheta')p_{\Theta}(\vartheta')}$$

(sčítance s $p_{\Theta}(\vartheta') = 0$ považujeme za 0).

Věta (Bayesova pro spojité náhodné veličiny)

X, Θ jsou spojité n.v., které mají hustotu f_X, f_{Θ} i sdruženou hustotu $f_{X,\Theta}$

$$f_{\Theta|X}(\vartheta|x) = \frac{f_{X|\Theta}(x|\vartheta)f_{\Theta}(\vartheta)}{\int_{\vartheta' \in \mathbb{R}} f_{X|\Theta}(x|\vartheta')f_{\Theta}(\vartheta')d\vartheta'}$$

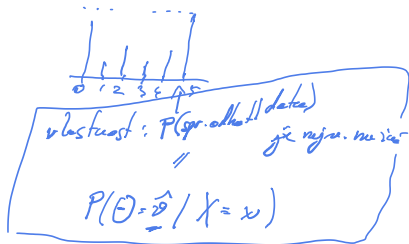
(sčítance s $f_{\Theta}(\vartheta') = 0$ považujeme za 0).

Bayesovské bodové odhady – MAP a LMS

MAP – Maximum A-Posteriori

Volíme $\hat{\vartheta}$ tak, aby maximalizovalo

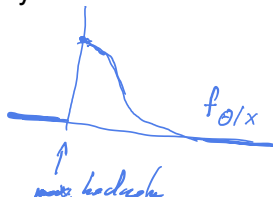
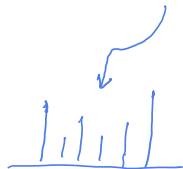
- ▶ $p_{\Theta|X}(\vartheta|x)$ v diskrétním případě
- ▶ $f_{\Theta|X}(\vartheta|x)$ ve spojitém případě



LMS – Least Mean Square

Těž metoda podmíněné střední hodnoty.

- ▶ Volíme $\hat{\vartheta} = \mathbb{E}(\Theta \mid X = x)$



Příklad 1

naivní

money money ...

Bayesovský klasifikátor spamů:

- ▶ vytvoříme seznam podezřelých slov (money, win, pharmacy, ...)
- ▶ N.v. X_i popisuje, zda email obsahuje podezřelé slovo w_i .
- ▶ N.v. Θ popisuje, zda email je spam $\Theta = 1$ nebo ne $\Theta = 0$.
- ▶ Z předchozích emailů získáme odhady $p_{X|\Theta}$ a p_{Θ}
 - 1 → podíl spamů 20%
 - 0 → 20%
- ▶ Použijeme Bayesovu větu na výpočet $p_{\Theta|X}$
 - $x = (x_1, \dots, x_n)$

$$P_{\Theta|X}(1 | X=x) > P_{\Theta|X}(0 | X=x)$$

→ email spam

$$P_{\Theta|X}(0 | X=x) > \text{krit. hod.}$$

→ email ne-spam

$$P_{X_i|\Theta}(x_i | \Theta)$$

$$= P_{X_1|\Theta}(x_1 | \Theta) \cdot P_{X_2|\Theta}(x_2 | \Theta) \dots$$

Příklad 2

Romeo a Julie se mají sejít přesně v poledne. Julie ale přijde pozdě o dobu popsanou náhodnou veličinou $X \sim U(0, \vartheta)$.
 Parametr ϑ modelujeme náhodnou veličinou $\Theta \sim U(0, 1)$. Co z naměřené hodnoty $X = x$ usoudíme o ϑ ? — Juste: $\vartheta \geq x$

aprioras hustota $f_{\Theta}(\vartheta) = 1$ pro $\vartheta \in (0, 1)$

$f_{X|\Theta}(x|\vartheta) = \frac{1}{\vartheta}$ pro $x \in (0, \vartheta)$ $\vartheta \in (x, 1)$

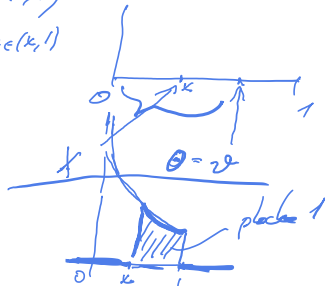
$$f_{\Theta|X}(\vartheta|x) = \frac{f_{X|\Theta}(x|\vartheta) \cdot f_{\Theta}(\vartheta)}{\int_0^1 f_{X|\Theta}(x|\vartheta') \cdot f_{\Theta}(\vartheta') d\vartheta'}$$

$\vartheta \in (0, 1)$

požadav $\vartheta \geq x$

$$\frac{\frac{1}{\vartheta}}{\int_x^1 \frac{1}{\vartheta'} d\vartheta'} = \frac{\frac{1}{\vartheta}}{[\log]_x^1} = \frac{1}{\vartheta \log x}$$

posteriora distribuce



pro $\vartheta < x$... $f_{X|\Theta}(x|\vartheta) = 0$

$f_{\Theta|X}(\vartheta|x) = 0$
 MAP ... $\hat{\vartheta} = x$

podľa št. lematy.

$$\hat{\theta} := E(\Theta | X=x) = \int_{-\infty}^{\infty} \theta f_{\Theta|X}(\theta|x) d\theta$$

$$= \int_x^1 \theta \frac{1}{2 \ln x} d\theta = \int_x^1 \frac{\theta}{\ln x} = \frac{1-x}{\ln x}$$

variance: máme nezávislé a pokusy X_1, \dots, X_n

Příklad 3

nezávislé

Pozorujeme náhodné veličiny $X = (X_1, \dots, X_n)$,
 předpokládáme $X_i \sim N(\vartheta, \sigma_i^2)$ a ϑ je hodnota náhodné veličiny
 $\Theta \sim N(x_0, \sigma_0^2)$. Co z naměřených hodnot $X = x = (x_1, \dots, x_n)$
 usoudíme o ϑ ?

$$f_{\Theta}(\vartheta) = c_1 \exp\left(-\frac{(\vartheta - x_0)^2}{2\sigma_0^2}\right)$$

*postev. disto. vyjít opět stejného
 typu (normální disto)*

$$f_{X_i|\Theta}(x_i|\vartheta) = c_2 \exp\left(-\frac{(x_i - \vartheta)^2}{2\sigma_i^2}\right) \dots \dots$$

$$\exp\left(-\frac{(x_n - \vartheta)^2}{2\sigma_n^2}\right)$$

konst $\frac{1}{\sigma_1 \sigma_1} \cdot \frac{1}{2\sigma_2} \dots \sim$

$$= c_3 \cdot \exp\left(-\frac{(\vartheta - m)^2}{2\sigma^2}\right) \dots \dots \sim N(m, \sigma)$$

Bayes:

$$f_{\Theta|X}(\vartheta|x) = \frac{f_{\Theta}(\vartheta) \cdot f_{X|\Theta}(x|\vartheta)}{\int \dots d\vartheta'}$$

$$= c_3 \cdot c_1 \exp\left(-\frac{(\vartheta - x_0)^2}{2\sigma_0^2}\right) \cdot c_2 \exp \dots$$

nezávislé σ_0^2

$$= c \cdot \exp\left(-\sum_{i=1}^n \frac{(x_i - \vartheta)^2}{2\sigma_i^2}\right)$$

$$\sigma^2 = \frac{1}{\frac{1}{\sigma_0^2} + \sum_{i=1}^n \frac{1}{\sigma_i^2}}$$

$$m = \frac{x_0 \frac{1}{\sigma_0^2} + \sum_{i=1}^n x_i \frac{1}{\sigma_i^2}}{\frac{1}{\sigma_0^2} + \sum_{i=1}^n \frac{1}{\sigma_i^2}}$$

$$-\frac{1}{2} \sum_{i=1}^n \frac{x_i^2}{\sigma_i^2} + \sigma_0^2 \sum_{i=1}^n \frac{x_i}{\sigma_i^2} \Rightarrow \sum_{i=1}^n \frac{x_i}{2\sigma_i^2} = -\frac{(\vartheta - m)^2}{2\sigma^2} + \text{const.}$$

Příklad 4

Házíme mincí, pravděpodobnost, že padne panna je ϑ . Z n hodů padla panna v $X = k$ případech. Pokud naše apriorní distribuce byla $U(0, 1)$, jaká bude distribuce posteriorní?

$$f_{\theta}(\vartheta) = 1 \quad (\vartheta \in (0, 1))$$

$$P(X=k | \theta, \vartheta) = \binom{n}{k} \cdot \vartheta^k (1-\vartheta)^{n-k}$$

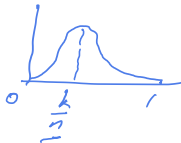
$$P_{X|\theta}(k, \vartheta)$$

$$f_{\theta|X}(\vartheta | k) = \frac{f_{\theta}(\vartheta) P_{X|\theta}(k | \vartheta)}{\int \dots} = \frac{c \cdot \vartheta^k (1-\vartheta)^{n-k}}{\int_0^1 \vartheta^k (1-\vartheta)^{n-k} d\vartheta}$$

post.m. $f_{\theta|X}$: $E(\theta | X=k) = \int_0^1 c \vartheta \cdot \vartheta^k (1-\vartheta)^{n-k} d\vartheta = \frac{k+1}{n+2}$

nice met $\frac{k}{n}$ a $\frac{1}{2}$

MAP \Leftrightarrow ML
 $\hat{\vartheta} = \frac{k}{n}$



Loe postulat s posterior
 distri. $f_{\theta}(\vartheta) = c \vartheta^k (1-\vartheta)^{n-k}$

Střední hodnota a součet čtverců

Věta

Pro libovolnou n.v. Θ je hodnota $\mathbb{E}(\Theta - \hat{\vartheta})^2$ nejmenší pro $\hat{\vartheta} = \mathbb{E}(\Theta)$.

✗

$$\text{var } Z = \mathbb{E}Z^2 - (\mathbb{E}Z)^2$$

$$\mathbb{E}Z^2 = \text{var } Z + (\mathbb{E}Z)^2$$

$$\mathbb{E}(\Theta - \hat{\vartheta})^2 = \text{var } \Theta + (\mathbb{E}(\Theta) - \hat{\vartheta})^2 \geq 0$$

minimální pro $\mathbb{E}(\Theta) = \hat{\vartheta}$.

$\hat{\vartheta}$ [↑] reálné číslo

$$Z := \Theta - \hat{\vartheta}$$

$$\text{var } Z = \text{var } \Theta$$

$$\mathbb{E}Z = \mathbb{E}\Theta - \hat{\vartheta}$$

Podmíněná hodnota dává nejmenší součet čtverců

Věta

Bodový odhad $\hat{\vartheta} = \mathbb{E}(\Theta \mid X = x)$ je nestranný a má nejmenší možnou hodnotu $\mathbb{E}(\Theta - \hat{\vartheta})^2 \mid X = x$