

13. cvičení z PSt

Připomenutí teorie

- **Čebyševova nerovnost:** Nechť X má konečnou střední hodnotu μ a rozptyl σ^2 . Pak

$$P(|X - \mu| \geq a \cdot \sigma) \leq \frac{1}{a^2}.$$

- **Centrální limitní věta:** Nechť X_1, \dots, X_n jsou stejně rozdělené n.n.v. se střední hodnotou μ a rozptylem σ^2 . Označme $Y_n = ((X_1 + \dots + X_n) - n\mu)/(\sqrt{n} \cdot \sigma)$. Pak $Y_n \xrightarrow{d} N(0, 1)$. Neboli

$$\lim_{n \rightarrow \infty} F_{Y_n}(x) = \Phi(x) \quad \text{pro každé } x \in \mathbb{R}.$$

- **De Moivre–Laplaceho věta** říká o něco silnější věc: v okolí pn je i pravděpodobnostní funkce $Bin(n, p)$ dobře aproximována hustotou $N(np, np(1-p))$, tedy vhodně přeškálovanou Gaussovou funkcí φ . Přesněji:

$$\binom{n}{k} p^k (1-p)^{n-k} \simeq \frac{1}{\sqrt{2\pi np(1-p)}} e^{-\frac{(k-np)^2}{2np(1-p)}}$$

pro k blízké pn . Ještě přesněji: pokud pro nějaké c je $|k - pn| < c\sqrt{np(1-p)}$ a n se blíží nekonečnu, tak poměr dvou výrazů nahoře se blíží k jedné.

Aplikace nerovností a Centrální Limitní Věty

1. Statistik chce odhadnout průměrnou výšku h (v metrech) lidí v nějaké populaci, pomocí n nezávislých vzorků X_1, \dots, X_n , které vybíráme uniformně náhodně ze všech možných lidí. Pro odhad použije výběrový průměr $\bar{X}_n = (X_1 + \dots + X_n)/n$. Odhaduje, že směrodatná odchylka jednoho měření je nejvýše 1 metr.

(a) Jak velké n má volit, aby směrodatná odchylka \bar{X}_n byla nejvýše 1 cm?

(b) Pro jaké n zajistí Čebyševova nerovnost, že \bar{X}_n se liší od h nejvýše o 5 cm s pravděpodobností alespoň 99 %?

(c) Statistik si všimne, že všichni měření lidí mají výšku v intervalu (1.4, 2.1). Jak má upravit odhad směrodatné odchylky? Jak se změní odpovědi na předchozí otázky?

2. Odhadněte $\binom{100}{30}$ pomocí CLV. Náповěda: použijte CLV pro odhad $P(29.5 < X < 30.5)$ pro vhodnou n.v. X . Na druhou stranu pro $P(X = 30)$ máme vzorec $\binom{100}{30}/2^{100}$ z binomického rozdělení. Alternativně, můžete použít Moivre–Laplaceho větu. Pokud máte po ruce vhodný stroj, vyčíslete.

Bodové odhady

- Zkoumáme posloupnost n.n.v. se stejným rozdělením, např. $Geom(\theta)$, $U(0, \theta)$, kde θ je parametr.
- Zapisujeme $X_1, \dots, X_n \sim F_\theta$, tzv. **náhodný výběr** z F_θ (model s parametrem).
- Naměříme $X_1 = x_1, \dots$, chceme odhadnout θ .
- $\hat{\theta}$... nějaká metoda jak odhadnout θ pomocí naměřených dat (hodnot X_1, \dots, X_n), angl. *estimator*
- $m_r(\theta) = \mathbb{E}(X^r)$ pro $X \sim F_\theta$... **r -tý moment**, ideální vlastnost rozdělení
- $\hat{m}_r(\theta) = \frac{1}{n} \sum_{i=1}^n X_i^r$... **r -tý výběrový moment**, náhodná veličina, funkce našeho naměřeného vzorku (tj. statistika)
- **Odhad metodou momentů** vyřešíme rovnici $m_1(\theta) = \hat{m}_1(\theta)$ pro neznámou θ .
- event. soustavu rovnice $m_r(\theta) = \hat{m}_r(\theta)$ pro $r = 1, 2, \dots$ podle potřeby.

- $L(\theta; x_1, \dots, x_n) = P(X_1 = x_1 \& \dots \& X_n = x_n) \dots$ pravd. pozorovaných dat závislá na parametru θ .
- nebo $L(\dots) = f_{X_1, \dots, X_n}(x_1, \dots, x_n) \dots$ hustota pravděpodobnosti \dots
- $\ell(\theta; x_1, \dots, x_n) = \log L(\dots) \dots$ pro snazší výpočty.
- **Odhad metodou maximální věrohodnosti (Maximal Likelihood)** hledáme θ , pro které je maximální $L(\theta; x_1, \dots, x_n)$, resp. $\ell(\dots)$. Obvykle pomocí derivací funkce L , resp. ℓ .
- **bias (vychýlení):** $\mathbb{E}(\hat{\theta} - \theta) \dots$ θ skutečný parametr, $\hat{\theta}$ náš odhad (náhodná veličina, protože závisí na naměřených datech)
- odhad je **nevychýlený/nestranný/unbiased:** $bias = 0$
- odhad je **asymptoticky nevychýlený:** bias konverguje k 0, neboli $\mathbb{E}(\hat{\theta}) \rightarrow \theta$
- odhad je **konzistentní:** $\hat{\theta} \xrightarrow{P} \theta$: pro všechna $\varepsilon > 0$ $P(|\hat{\theta} - \theta| > \varepsilon) \rightarrow 0$
- **MSE (mean square error, střední kvadratická odchylka):** $\mathbb{E}((\hat{\theta} - \theta)^2)$
- Věta: $MSE(\hat{\theta}) = bias(\hat{\theta})^2 + var(\hat{\theta})$.

Pro praktickou ukázkou, viz pythonový notebook na webu cvičení.

- Máme náhodný výběr $X_1, \dots, X_n \sim U(0, \theta)$.
 - Navrhněte bodový odhad θ momentovou metodou.
 - Navrhněte bodový odhad θ metodou maximální věrohodnosti.
 - Pro každý z nich zjistěte, zda je nestranný a konzistentní. (Stačí experimentálně na počítači.)
 - Pro každý z nich spočítejte střední kvadratickou odchylku (MSE). (Stačí experimentálně na počítači.)
 - Který odhad je lepší? Napadá vás nějaký ještě lepší?
- Pro náhodný výběr $X_1, \dots, X_n \sim Geom(p)$ řešte části (a)–(c) jako výše.
- Pro náhodný výběr $X_1, \dots, X_n \sim Exp(1/\theta)$ řešte části (a)–(d) jako výše.

K procvičení

- Označme $S = \sum_{k=0}^{30} \binom{100}{k}$. Označme dále $X = \sum_{i=1}^{100} X_i$, kde X_i je 0 nebo 1, obojí s pravděpodobností 1/2 a veličiny X_1, \dots, X_n jsou nezávislé. Je tedy $X \sim Bin(100, 1/2)$.
 - Vyjádřete S pomocí distribuční funkce F_X .
 - Použijte CLV na odhad této pravděpodobnosti.
 - Případně vyčíslíte S vhodným softwarem a srovnejte.
- Chceme odhadnout, zda naše mince (a způsob jak s ní házíme) je spravedlivá. Pokud ze sta hodů padne orel více než 55-krát, řekneme, že spravedlivá není. Jaká je pravděpodobnost, že se zmýlíme?
- Nechť $X \sim Exp(\lambda)$ popisuje dráhu, kterou uletí radioaktivní částice, necht' se rozpadne. Náš přístroj její rozpad (a polohu rozpadu, tj. hodnotu X) zachytí, ale jen pokud $1 \leq X \leq 2$. Formálně, budeme zkoumat náhodný výběr $X_1, \dots, X_n \sim F_{X|B}$ pro jev $B = 1 \leq X \leq 2$.
 - Navrhněte bodový odhad λ momentovou metodou.
 - Navrhněte bodový odhad λ metodou maximální věrohodnosti.
 - Pro každý z nich zjistěte, zda je nestranný a konzistentní.