

Přednáška $Ma + b$

Markovova nerovnost

Pro libovolnou nerápornou náhodnou veličinu X platí

$$P(X \geq \varepsilon) \leq \frac{E(X)}{\varepsilon}.$$

Důkaz: Pro spojitou náh. veličinu máme:

$$\begin{aligned} E(X) &= \int_{-\infty}^{\infty} x f(x) dx = \int_0^{\infty} x f(x) dx = \int_0^{\varepsilon} x f(x) dx + \int_{\varepsilon}^{\infty} x f(x) dx \\ &\geq \int_{\varepsilon}^{\infty} \varepsilon f(x) dx = \varepsilon \int_{\varepsilon}^{\infty} f(x) dx = \varepsilon P(X \geq \varepsilon). \quad \blacksquare \end{aligned}$$

Lze navíc pro každou nerápornou náhodnou veličinu, známe-li její střední hodnotu.

Známe-li navíc střední hodnoty i rozptyl lze navíc Čebyševova nerovnost (uvedena dříve).

Příklad ve fotbalové lize padne n průběhem λ branek n zápasů. Odhadněte pravděpodobnost, že padne aspoň 3λ branek.

Rěšení pomocí Čebyševovy Markovovy nerovnosti:

$$P(X \geq 3\lambda) \leq \frac{E(X)}{3\lambda} = \frac{\lambda}{3\lambda} = \frac{1}{3}.$$

Lejn' otkad: Víme, že náhodná veličina X = počet kamek v nápare má Poissonovo rozdělení s parametrem λ

$$f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

Toto rozdělení má střední hodnotu

$$E(X) = \lambda$$

a rozptyl

$$\text{var}(X) = \lambda.$$

Rěšení úlohy pomocí Čebyševovy nerovnosti
ne pak

$$\begin{aligned} P(X \geq 3\lambda) &= P(|X - \lambda| \geq 2\lambda) = \\ &= P(|X - E(X)| \geq 2\lambda) \leq \frac{\text{var}(X)}{(2\lambda)^2} = \frac{\lambda}{4\lambda^2} = \frac{1}{4\lambda}. \end{aligned}$$



Matematická statistika

- zkoumáme výběr z nějakého souboru (populace)
- chceme vědět, do jaké míry jsou spolehlivé výsledky platné pro celý soubor (populaci)
- případně se spolehlivých dat se snažíme upravit vhodný teoretický model pro charakterizaci celého souboru a z něho pak odhadnout pravděpodobnost budoucího jevu

2 přístupy

- frekvenční (klasická) statistika
- bayesovská statistika

Frekvenční statistika

- vychází z mat. abstrakce, se skutečnou pravděpodobností jsou dáány četnosti výsledků jevu a tak velkými množství dat, ^{je} většinou se dobře aproximovat nekonečnými modely a využívat pro odhady teorie středních hodnot a limitní věty.
- má své limity, pokud nejsou data spolehlivá nebo je experiment nevhodný
- nelze používat pro odhad výsledků jednorázových

děju.

Náhodný výběr

Základní soubor (populace) má N jednotek. Pro každou jednotku můžeme určit číselný znak $x_{11}, x_{21}, \dots, x_{N1}$. Ze základního souboru vybereme $n \ll N$ jednotek jednu po druhé a každou jednotku položíme do souboru náh. Tedy každá jednotka má pravděpodobnost výběru $\frac{1}{N}$. Hovoříme o náhodném výběru.

V matematické idealizaci pracujeme s n -licí měřitelných náhodných veličin X_1, X_2, \dots, X_n , které mají vzájemně nezávislé rozdělení pravděpodobnosti, tj. stejnou distribuční funkci $F_X(x)$. Tuto n -licí náhodných veličin nazýváme NÁHODNÝ VÝBĚR ROZSAHU n .

V matematické statistice často pracujeme s transformacemi náhodného výběru

$$Y = f(X_1, X_2, \dots, X_n),$$

takovým náhodným veličinám (případně vektorům) říkáme statistiky.

Zavedeme několik důležitých statistik
a ukážeme jejich souvislost s číselnými
charakteristikami.

Definice Nechtě X_1, X_2, \dots, X_n je náhodný
v'běr. Statistiku

$$M = \frac{1}{n} \sum_{i=1}^n X_i$$

nazýváme v'běrovým průměrem.

Statistiku

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

v'běrovým rozptylem a statistikou $S = \sqrt{S^2}$
v'běrovou směrodatnou odchylkou.

Statistiky jsou opět náhodné veličiny, proto
se ke platí na jejich číselné charakteristiky
- střední hodnota, rozptyl, momenty, rozdělení.

Věta Nechtě X_1, X_2, \dots, X_n je náhodný v'běr rozsa-
hu n a rozdělení x střední hodnotou μ a rozpty-
lem σ^2 . Pak platí

$$(1) \quad EM = \mu$$

$$(2) \quad \text{var } M = \frac{\sigma^2}{n}$$

$$(3) \quad ES^2 = \sigma^2$$

Důkaz: (1) a (2) plyne z vlastností střední hodnoty a rozptylu.

$$(3) \quad \text{Platí} \quad \sum_{i=1}^n (X_i - \mu)^2 = \sum_{i=1}^n (X_i - M)^2 + n(M - \mu)^2$$

Na levé straně se po umocnění vyskytují výrazy typu

$$L: \quad X_i^2, \quad -2X_i\mu, \quad n\mu^2$$

Na pravé straně se vyskytují výrazy

$$P: \quad X_i^2, \quad -2MX_i, \quad 2nM^2, \quad -2nM\mu, \quad n\mu^2$$

První a poslední výrazy se obzvlášť. Dále

$$\sum_{i=1}^n (-2X_i\mu) = -2n \frac{X_1 + \dots + X_n}{n} \mu$$

$$\begin{aligned} \text{a} \quad 2nM^2 - 2M(\sum_{i=1}^n X_i) &= 2nM^2 - 2nM \frac{\sum_{i=1}^n X_i}{n} = \\ &= 2nM^2 - 2nM^2 = 0 \end{aligned}$$

Prova

$$\begin{aligned} E(S^2) &= E\left\{\frac{1}{n-1} \sum_{i=1}^n (X_i - M)^2\right\} = \\ &= E\left\{\frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2 - \frac{n}{n-1} (M - \mu)^2\right\} \\ &= \frac{1}{n-1} \sum_{i=1}^n E(X_i - \mu)^2 - \frac{n}{n-1} E(M - \mu)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \text{var } X_i - \frac{n}{n-1} \text{var } M = \\ &= \frac{n}{n-1} \sigma^2 - \frac{1}{n-1} \sigma^2 = \sigma^2. \quad \blacksquare \end{aligned}$$

Dalšie zali máme, že vyberový súčet M splňuje

$$E(M) = \mu,$$

že to znamená, že hodnota je vždy rovná zodpovedajúceho parametra μ . Pôjde o to, že M je nehromadným odhadom parametra μ .

Obdobne platí, že S^2 je nehromadným odhadom parametra σ^2 .

Poznámka „Přirozenejší“ definovaná statistika

$\frac{1}{n} \sum_{i=1}^n (X_i - M)^2$ není nehromadným odhadem σ^2 , její střední hodnota je $\frac{n-1}{n} \sigma^2$.

Na'hodny' vy'běr z norma'lni'ho rozdělení'

Necht' X_1, X_2, \dots, X_n je na'hodny' vy'běr
a necht' každá z na'hodny'ch veličin
ma' norma'lní rozdělení $N(\mu, \sigma^2)$.

Časem dokážeme a doplníme následující

Věta:

- (1) M a S^2 jsou nesávislé náhodné
veličiny.
- (2) M ma' norma'lní rozdělení se střední
hodnotou μ a rozptylem σ^2/n , tj.

$$M \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

- (3) Náhodná veličina $U = \frac{\sqrt{n}(M - \mu)}{\sigma}$

ma' norma'lní rozdělení se střední
hodnotou 0 a rozptylem 1.

$$U \sim N(0, 1).$$

V následujícím příkladu si ukážeme,
jak tuto větu použít k odhadu μ ,
a náme-li σ^2 .

Příklad V roce 1951 byla rozsáhlým statistickým průzkumem zjištěno, že střední hodnota výšky desetiletých chlapců je 136,1 cm se směrodatnou odchylkou $\sigma = 6,4$ cm. V roce 1961 byla zjištěna výška pouze u 15 náhodně vybraných chlapců. Naměřené hodnoty byly:

130, 140, 136, 141, 139, 133, 149, 151, 139, 136, 138, 142, 127, 139, 147.

Otázkou je, zda se v porovnání s rokem 1951 změnila střední výška chlapců, pokud předpokládáme, že rozptýl výšek se v různých generacích příliš nemění.

Řešení Data lze považovat za náhodný výběr. Zjistíme výběrový průměr $M = 139,133$, $n = 15$. X_1, X_2, \dots, X_{15} je náhodný výběr s normálním rozdělením $N(\mu, \sigma^2)$. Střední hodnotu μ pro rok 1961 neznáme, σ považujeme za stejnou jako v roce 1951, tj. $\sigma = 6,4$. Zajímá nás interval $(M - \varepsilon, M + \varepsilon)$ takový, že

$$P(M - \varepsilon < \mu < M + \varepsilon) = 0,95$$

Neznanosti $M - \varepsilon < \mu < M + \varepsilon$ jsou ekvivalentní
s neznaností

$$|M - \mu| < \varepsilon$$

a to je ekvivalentní s neznaností

$$|U| = \left| \frac{\sqrt{n}(M - \mu)}{\sigma} \right| < \frac{\sqrt{n}}{\sigma} \varepsilon$$

Náhodná veličina U má stand. normální
rozdělení s distribuční funkcí Φ

Proto chceme najít ε tak, aby

$$P\left(-\frac{\sqrt{n}}{\sigma} \varepsilon < U < \frac{\sqrt{n}}{\sigma} \varepsilon\right) = 0,95$$

Pomocí distribuční funkce Φ to je

$$\Phi\left(\frac{\sqrt{n}}{\sigma} \varepsilon\right) - \Phi\left(-\frac{\sqrt{n}}{\sigma} \varepsilon\right) = 0,95$$

$$2\Phi\left(\frac{\sqrt{n}}{\sigma} \varepsilon\right) - 1 = 0,95$$

$$\Phi\left(\frac{\sqrt{n}}{\sigma} \varepsilon\right) = \frac{1,95}{2} = 0,975$$

$$\frac{\sqrt{n}}{\sigma} \varepsilon = \Phi^{-1}(0,975) = 1,96$$

$$\varepsilon = \frac{1,96 \cdot \sigma}{\sqrt{n}} = 3,23$$

Interval $(M - 3,23, M + 3,23)$ je interval s okrají sadany'mi náhodny'mi veličinama-mi. V naší nahladi'me M spočítany'm průměrem 139,13. Počto dostaneme interval

$$(139,13 - 3,23, 139,13 + 3,23) = (135,90; 142,36)$$

Střední hodnota výřek z roku 1951 leží s touto intervalu, nema'me tedy se spolehlivostí 95% spjistěno, že se průměrná výřka změnila.

Pokud si vezmeme spolehlivost pouze 90%, pak je interval

$$(136,41; 141,85).$$

V něm už střední hodnota 136,1 z roku 1951 neleží. Tedy střední hodnota výřky se změnila se spolehlivostí 90%.

Pokud na's raji'ma' pouze dolní odhad, tj interval

$$(M - \varepsilon, \infty)$$

taleny', se $P(M - \varepsilon < u) = 0,95,$

dostaneme interval

$$(136,41, \infty)$$

Tedy s pravděpodobností 95% je střední výška v roce 1961 větší než 136,41 cm a tedy i větší než v roce 1951. ■

Teoretické shrnutí předchozího příkladu

Interval spolehlivosti pro střední hodnotu μ

Nechtě X_1, X_2, \dots, X_n jsou nezávislé veličiny

s rozdělením $N(\mu, \sigma^2)$. Pro Nechtě $\alpha \in (0,1)$

~~Ukážeme~~ dvoustranný interval pro střední hodnotu μ spolehlivosti $1-\alpha$ je

$$(M - \varepsilon, M + \varepsilon)$$

s platností

$$P(M - \varepsilon < \mu < M + \varepsilon) = 1 - \alpha$$

jednostranný interval spolehlivosti je jednostranný

$$(M - \varepsilon, \infty), \text{ kde } P(M - \varepsilon < \mu) = 1 - \alpha$$

a jednostranný $(-\infty, M + \varepsilon)$, kde $P(\mu < M + \varepsilon) = 1 - \alpha$.

Stejně jako v příkladu spočítáme, že

$$\Phi\left(\frac{\sqrt{n}}{\sigma} \varepsilon\right) = 1 - \frac{\alpha}{2}$$

$$\frac{\sqrt{n}}{\sigma} \varepsilon = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$$

Tedy dvoustranný interval spolehlivosti $1-\alpha$ je

$$\left(M - \frac{\sigma}{\sqrt{n}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right), M + \frac{\sigma}{\sqrt{n}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \right)$$

Levostroanný interval spolehlivosti je

$$\left(M - \frac{\sigma}{\sqrt{n}} \Phi^{-1}(\alpha), \infty \right),$$

pravostroanný pak

$$\left(-\infty, M + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1-\alpha) \right).$$



Vraťme se k větě před příkladem

Věta: Nechtě X_1, X_2, \dots, X_n je náhodný výběr s rozložením $N(\mu, \sigma^2)$. Pak M má normální rozložení $N\left(\mu, \frac{\sigma^2}{n}\right)$ a

$$U = \frac{\sqrt{n}(M - \mu)}{\sigma} \sim N(0, 1).$$

Důkaz přechů řešení - druhé již známe.

Víme, že normální rozložení $Z \sim N(0, 1)$ má momentovou funkci $E(e^{tz}) = e^{\frac{t^2}{2}}$.

Normální rozložení $X \sim N(\mu, \sigma^2)$ má momentovou funkci

$$M_X(t) = M_{\mu + \sigma Z}(t) = e^{\mu t} M_Z(\sigma t) = e^{\mu t} \cdot e^{\frac{(\sigma t)^2}{2}}$$

Nechtě $X_1 \sim N(\mu_1, \sigma_1)$ a $X_2 \sim N(\mu_2, \sigma_2)$.

Potom

$$\begin{aligned} M_{X_1+X_2}(t) &= M_{X_1}(t) \cdot M_{X_2}(t) = \\ &= e^{\mu_1 t} \cdot e^{\frac{\sigma_1^2 t^2}{2}} \cdot e^{\mu_2 t} \cdot e^{\frac{\sigma_2^2 t^2}{2}} = \\ &= e^{(\mu_1 + \mu_2)t} e^{\frac{(\sigma_1^2 + \sigma_2^2)t^2}{2}} = M_Y(t) \end{aligned}$$

keďže $Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$. Tedy súčet nezávislých náhodných veličín s normálnym rozdelením je opäť náhodná veličina s normálnym rozdelením.

Preto keďže $X_i \sim N(\mu_i, \sigma_i^2)$, pak

$$X_1 + X_2 + \dots + X_n \sim N(n\mu, n\sigma^2)$$

a $M = \frac{1}{n} (X_1 + X_2 + \dots + X_n) \sim N(\mu, \frac{\sigma^2}{n})$. \square

Budeme sa zaoberať ďalšími rozdeleniami, ktoré sú s normálnym rozdelením úzko súvisiace.

Γ rozdelení $\Gamma(a, b)$

ma' hustotu

$$f(x) = \begin{cases} 0, & x \leq 0 \\ c x^{a-1} e^{-bx}, & x > 0, \end{cases}$$

keďže $c = \frac{b^a}{\Gamma(a)}$, $a \geq 0$, $b \geq 0$,

-15-

$$\Gamma(a) = \int_0^{\infty} t^{a-1} e^{-t} dt.$$

Poznámka: Funkce Γ je zobecněním faktoriálu, neboť $\Gamma(n) = (n-1)!$ pro $n \in \mathbb{N}$.

Dále platí $\Gamma(1/2) = \sqrt{\pi}$

$$\Gamma(a+1) = a \cdot \Gamma(a).$$

Věta Momentová vyjádření funkce rozdělení

$\Gamma(a, b)$ je

$$M(t) = \left(\frac{b}{b-t} \right)^a,$$

střední hodnota je a/b a rozptyl a/b^2 .

Důkaz sa. DU

Rozdělení χ^2 - chi kvadrat $\chi^2(1)$

Nechť Z je náhodná veličina s rozdělením $N(0, 1)$. Spojíme-li me křivku f pro

náhodnou veličinu Z^2 :

zřejmě $f(x) = 0$ pro $x \leq 0$. Pro kladná x máme

$$F(x) = P(Z^2 < x) = P(-\sqrt{x} < Z < \sqrt{x}) =$$

$$= \int_{-\sqrt{x}}^{\sqrt{x}} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz = \int_0^x \frac{1}{\sqrt{2\pi}} t^{-\frac{1}{2}} e^{-\frac{t}{2}} dt$$

Hastota derivace derivaci distribucni funkce

$$f(x) = \frac{d}{dx} F(x) = \frac{1}{\sqrt{2\pi}} x^{-\frac{1}{2}} e^{-\frac{x}{2}}$$

Toto je gamma rozdeleni $\Gamma(\frac{1}{2}, \frac{1}{2}) = \chi^2(1)$.

Věta Necht X_1, X_2, \dots, X_n jsou nezávislé náhodné veličiny s rozděleními $\Gamma(a_i, b)$. Pak $Y = X_1 + X_2 + \dots + X_n$ má rozdeleni $\Gamma(a_1 + a_2 + \dots + a_n, b)$.

Důkaz plyne z výpočtu momentu

$$M_Y(t) = \prod_{i=1}^n M_{X_i}(t) = \frac{b^{a_1}}{(b-t)^{a_1}} \cdot \frac{b^{a_2}}{(b-t)^{a_2}} \dots \frac{b^{a_n}}{(b-t)^{a_n}} = \left(\frac{b}{(b-t)}\right)^{a_1 + a_2 + \dots + a_n}$$

Rozdělení $\chi^2(n)$ - chi kvadrat s n stupni volnosti

Necht $Y = Z_1^2 + Z_2^2 + \dots + Z_n^2$, kde Z_1, \dots, Z_n jsou nezávislé náhodné veličiny

-17-

s rozdělením $N(0, 1)$. Pak Z_i^2 mají
rozdělení $\Gamma(\frac{1}{2}, \frac{1}{2})$ a Y má rozdělení

$$\Gamma\left(\frac{n}{2}, \frac{1}{2}\right).$$

To se nazývá chi-kvadrát s n stupni
volnosti a značí se $\chi^2(n)$. ■