

Matematika III – 12. týden

Centrální limitní věta, příklady důležitých rozdělání, frekvenční a Bayesovská statistika, výběry z populací

Jan Slovák

Masarykova univerzita
Fakulta informatiky

7.12.-11. 12. 2015

Obsah přednášky

- 1 Literatura
- 2 Centrální limitní věta
- 3 Co potkáme
- 4 Matematická statistika
- 5 Výběry z populací
- 6 Intervaly spolehlivosti

Plán přednášky

- 1 Literatura
- 2 Centrální limitní věta
- 3 Co potkáme
- 4 Matematická statistika
- 5 Výběry z populací
- 6 Intervaly spolehlivosti

Kde je dobré číst?

- Karel Zvára, Josef Štěpán, Pravděpodobnost a matematická pravděpodobnost statistika, Matfyzpress, 2006, 230pp.
- J. Slovák, M. Panák, M. Bulant, Matematika drsně a svižně, Muni Press, Brno 2013, v+773 s., elektronická edice www.math.muni.cz/Matematika_drsne_svizne
- Marie Budíková, Štěpán Mikoláš, Pavel Osecký, Teorie pravděpodobnosti a matematická statistika (sbírka příkladů), Masarykova univerzita, 3. vydání, 2004, 117 stran, ISBN 80-210-3313-4.
- Marie Budíková, Tomáš Lerch, Štěpán Mikoláš, Základní statistické metody, Masarykova univerzita, 2005, 170 stran, ISBN 80-210-3886-1.
- Riley, K.F., Hobson, M.P., Bence, S.J. Mathematical Methods for Physics and Engineering, second edition, Cambridge University Press, Cambridge 2004, ISBN 0 521 89067 5, xxiii + 1232 pp.

Plán přednášky

- 1 Literatura
- 2 Centrální limitní věta**
- 3 Co potkáme
- 4 Matematická statistika
- 5 Výběry z populací
- 6 Intervaly spolehlivosti

Uvažme nezávislé náhodné veličiny Y_1, Y_2, \dots , které mají všechny stejné rozdělení se střední hodnotou 0 a rozptylem 1.

Předpokládejme, že třetí absolutní moment $E|Y_i|^3$ je konečný.

Pro náhodnou veličinu $S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$ spočtěme momentovou funkci (koeficient $n^{-1/2}$ je volen tak, aby rozptyl S_n byl stále 1)

$$M_{S_n} = \prod_{i=1}^n E e^{(t/\sqrt{n})Y_i} = (M_Y(t/\sqrt{n}))^n,$$

kde M_Y je společná momentová funkce všech veličin Y_i .

Uvažme nezávislé náhodné veličiny Y_1, Y_2, \dots , které mají všechny stejné rozdělení se střední hodnotou 0 a rozptylem 1.

Předpokládejme, že třetí absolutní moment $E|Y_i|^3$ je konečný.

Pro náhodnou veličinu $S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$ spočtěme momentovou funkci (koeficient $n^{-1/2}$ je volen tak, aby rozptyl S_n byl stále 1)

$$M_{S_n} = \prod_{i=1}^n E e^{(t/\sqrt{n})Y_i} = (M_Y(t/\sqrt{n}))^n,$$

kde M_Y je společná momentová funkce všech veličin Y_i . Nyní

$$M_Y(t/\sqrt{n}) = 1 + 0 \frac{t}{\sqrt{n}} + 1 \frac{t^2}{2n} + o(t^2/n)$$

a v limitě proto dostáváme

$$\lim_{n \rightarrow \infty} M_{S_n}(t) = \lim_{n \rightarrow \infty} \left(1 + \frac{t^2}{2n} + o(1/n) \right)^n = e^{t^2/2}.$$

To je právě momentová funkce pro rozdělení $N(0, 1)$!

Tím jsme skoro dokázali:

Theorem (Centrální limitní věta)

Nechť Y_1, Y_2, \dots jsou nezávislé náhodné veličiny se společnou střední hodnotou $E Y_i = \mu$, rozptylem $\text{var } Y_i = \sigma^2 > 0$ a konečným třetím absolutním momentem $E|Y_i|^3$. Pro distribuční funkce náhodných veličin

$$S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{1}{\sigma} (Y_i - \mu)$$

platí

$$\lim_{n \rightarrow \infty} P(S_n < x) = \Phi(x),$$

kde $\Phi(x)$ je distribuční funkce normálního rozdělení $N(0, 1)$.

Všimněme si: součty $X_n = \sum_{i=1}^n Y_i$ mají střední hodnotu $n\mu$ a rozptyl $n\sigma^2$. Veličiny S_n jsou tedy právě normované veličiny X_n .

Pokud jsou $Y_i \sim A(p)$ nezávislé, pak $E(Y_i)^3 = p < \infty$ a všechny podmínky centrální limitní věty jsou splněny, $\mu = p$, $\sigma^2 = p(1 - p)$.

Pokud jsou $Y_i \sim A(p)$ nezávislé, pak $E(Y_i)^3 = p < \infty$ a všechny podmínky centrální limitní věty jsou splněny, $\mu = p$, $\sigma^2 = p(1 - p)$. Součtové veličiny $X_n = \sum_{i=1}^n Y_i$ pak představují právě binomická rozdělení $\text{Bi}(n, p)$ a příslušné normované veličiny jsou

$$S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\frac{Y_i - p}{\sqrt{p(1-p)}} \right) = \frac{X_n - np}{\sqrt{np(1-p)}}.$$

Podle centrální limitní věty má tato veličina pro velká n rozdělení velmi podobné rozdělení $N(0, 1)$.

Pokud jsou $Y_i \sim A(p)$ nezávislé, pak $E(Y_i)^3 = p < \infty$ a všechny podmínky centrální limitní věty jsou splněny, $\mu = p$, $\sigma^2 = p(1 - p)$. Součtové veličiny $X_n = \sum_{i=1}^n Y_i$ pak představují právě binomická rozdělení $\text{Bi}(n, p)$ a příslušné normované veličiny jsou

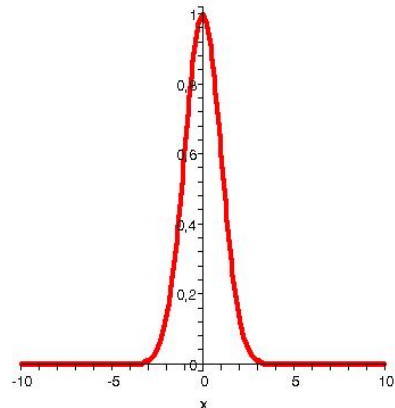
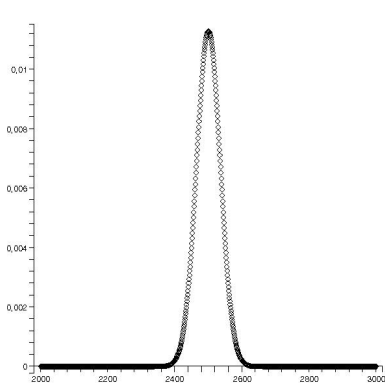
$$S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\frac{Y_i - p}{\sqrt{p(1 - p)}} \right) = \frac{X_n - np}{\sqrt{np(1 - p)}}.$$

Podle centrální limitní věty má tato veličina pro velká n rozdělení velmi podobné rozdělení $N(0, 1)$.

Jinými slovy, rozdělení $\text{Bi}(n, p)$ je velice blízké rozdělení $N(np, np(1 - p))$ pro velká n . To je obsahem tzv.

Laplaceovy–Moivreovy věty. To jsme už viděli minule na obrázcích:

Pro hodnoty $Bi(5000, 0.5)$ je výsledek vidět na obrázku níže. Druhá křivka na obrázku je grafem funkce $f(x) = e^{-x^2/2}$.



Aproximace binomického rozdělení normálním se často považuje v praxi za dostatečnou, jestliže $np(1 - p) > 9$

Při praktických průzkumech zpravidla věříme „zákonu velkých čísel“. Potřebujeme přitom rozhodnout, jak velký vzorek už postačuje.

Typickým příkladem je např. tato úloha: Chceme zjistit poměr p osob s danou krevní skupinou A v populaci. U kolika osob je třeba krevní skupinu skutečně zjistit, abychom měli 90% pravděpodobnost, že naše zjištění se nebude lišit o více než 5%.

Při praktických průzkumech zpravidla věříme „zákonu velkých čísel“. Potřebujeme přitom rozhodnout, jak velký vzorek už postačuje.

Typickým příkladem je např. tato úloha: Chceme zjistit poměr p osob s danou krevní skupinou A v populaci. U kolika osob je třeba krevní skupinu skutečně zjistit, abychom měli 90% pravděpodobnost, že naše zjištění se nebude lišit o více než 5%. Propočítáním zjistíme, že (nezávisle na p) vždy stačí odhadnout $p = X/n$, kde X je náhodná veličina udávající počet osob majících požadovanou skupinu, pro vzorek 270 lidí.

Plán přednášky

- 1 Literatura
- 2 Centrální limitní věta
- 3 Co potkáme**
- 4 Matematická statistika
- 5 Výběry z populací
- 6 Intervaly spolehlivosti

Rozdělení χ^2

Ve statistice budeme pracovat s charakteristikami náhodných vektorů, které budou obdobné výběrovému průměru a rozptylu, ale také s relativními poměry takových charakteristik atd. Podíváme se teď na několik takových případů.

Uvažme $Z \sim N(0, 1)$ a spočtěme hustotu $f_Y(x)$ pro $Y = Z^2$. Evidentě je $f_Y(x) = 0$ pro $x \leq 0$, pro kladná x

$$\begin{aligned} F_Y(x) &= P(Y < x) = P(-\sqrt{x} < Z < \sqrt{x}) \\ &= \int_{-\sqrt{x}}^{\sqrt{x}} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = \int_0^x \frac{1}{\sqrt{2\pi}} t^{-1/2} e^{-t/2} dt. \end{aligned}$$

Hustotu dostaneme derivací

$$f_Y(x) = \frac{d}{dx} F_Y(x) = \frac{1}{\sqrt{2\pi}} x^{-1/2} e^{-x/2}.$$

Tomuto rozdělení se říká χ^2 s **jedním stupněm volnosti**, píšeme $Y \sim \chi^2$.

Gama rozdělení $Y \sim \Gamma(a, b)$

Výběrový rozptyl bude odpovídat součtům takovýchto nezávislých veličin.

Uvažme hustotu (trochu obecnějšího tvaru než u χ^2)

$$f_X(x) = cx^{a-1} e^{-bx}$$

pro $x > 0$, zatímco $f_X(x) = 0$ pro nekladná x (χ^2 odpovídá volbě $a = b = 1/2$).

Je třeba volit $c = \frac{b^a}{\Gamma(a)}$ a jde o rozdělení $\Gamma(a, b)$.

k -tý moment takové veličiny X je

$$\begin{aligned} E X^k &= \int_0^{\infty} x^k \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx} dx \\ &= \frac{\Gamma(a+k)}{\Gamma(a)b^r} \int_0^{\infty} \frac{b^{a+k}}{\Gamma(a+k)} x^{a-1+k} e^{-bx} dx \\ &= \frac{\Gamma(a+k)}{\Gamma(a)b^k} \end{aligned}$$

(protože integrál z hustoty rozdělení $\Gamma(a+k, b)$ v posledním upravovaném výrazu je nutně roven jedné)

Zejména tedy vidíme, že $E X = \frac{\Gamma(a+1)}{b\Gamma(a)} = \frac{a}{b}$, zatímco

$$\text{var } X = E X^2 - (E X)^2 = \frac{\Gamma(a+2)}{b^2\Gamma(a)} - \frac{a^2}{b^2} = \frac{(a+1)a - a^2}{b^2} = \frac{a}{b^2}.$$

Momentová vytvořující funkci pro všechny hodnoty $-b < t < b$ je

$$\begin{aligned} M_X(t) &= \int_0^{\infty} e^{tx} \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx} dx \\ &= \frac{b^a}{(b-t)^a} \int_0^{\infty} \frac{(b-t)^a}{\Gamma(a)} x^{a-1} e^{-(b-t)x} dx \\ &= \frac{b^a}{(b-t)^a}. \end{aligned}$$

Pro součet nezávislých rozdělení $Y = X_1 + \dots + X_n$ s rozděleními $X_i \sim \Gamma(a_i, b)$ tedy okamžitě dostáváme momentovou vytvořující funkci (pro hodnoty $|t| < b$)

$$M_Y(t) = \left(\frac{b}{b-t} \right)^{a_1 + \dots + a_n},$$

tj. $Y \sim \Gamma(a_1 + \dots + a_n, b)$. (Velmi podstatný je přitom předpoklad, že všechna gamma rozdělení sdílí stejnou hodnotu b).

rozdělení χ^2

Jako okamžitý důsledek nyní dostáváme hustotu rozdělení veličiny $Y = Z_1^2 + \dots + Z_n^2$, kde všechna $Z_i \sim N(0, 1)$. Jde totiž o gamma rozdělení $Y \sim \Gamma(n/2, 1/2)$ a má hustotu

$$f_Y(x) = \frac{1}{2^{n/2}\Gamma(n/2)} x^{n/2-1} e^{-x/2}.$$

Tomuto speciálnímu případu gamma rozdělení říkáme rozdělení χ^2 s n stupni volnosti. Značíme jej zpravidla $Y \sim \chi_n^2$.

F-rozdělení

Při prorovnání výběrových rozptylů potkáme veličiny, které jsou dány podílem

$$U = \frac{X/k}{Y/m}$$

$X \sim \chi_k^2$ a $Y \sim \chi_m^2$.

Náhodná veličina $U = \frac{X/k}{Y/m}$ má hustotu $f_U(u)$

$$f_U(u) = \frac{\Gamma((k+m)/2)}{\Gamma(k/2)\Gamma(m/2)} \left(\frac{k}{m}\right)^{k/2} u^{k/2-1} \left(1 + \frac{k}{m}u\right)^{-(k+m)/2}.$$

Takovému rozdělení se říká **Fisherovo-Snedecorovo rozdělení s k a m stupni volnosti**, zkráceně také **F-rozdělení**.

t-rozdělení

Další potřebné rozdělení se objevuje při zkoumání podílu veličin $Z \sim N(0, 1)$ a $\sqrt{X/n}$, kde $X \sim \chi_n^2$ (tj. zajímá nás poměr Z a směrodatné odchytky nějakého výběru).

Dostaneme náhodnou veličinu

$$T = \frac{Z}{\sqrt{X/n}}$$

a hustotou $f_T(t)$

$$f_T(t) = \frac{\Gamma((n+1)/2)}{\Gamma(n/2)\sqrt{n\pi}} \left(1 + \frac{t^2}{n}\right)^{-(n+1)/2}.$$

Tomuto rozdělení říkáme **Studentovo t-rozdělení s n stupni volnosti**.

Plán přednášky

- 1 Literatura
- 2 Centrální limitní věta
- 3 Co potkáme
- 4 Matematická statistika**
- 5 Výběry z populací
- 6 Intervaly spolehlivosti

matematická statistika

Zkoumáme statistiky u nějakého výběru z daného základního souboru (populace).

Matematická statistika se snaží postihnout, do jaké míry jsou zjištěné výsledky relevantní pro celou populaci, případně se ze zjištěných dat pokouší zjistit nebo upřesnit vhodný teoretický model pro chování celého souboru (a z něj pak třeba odhadovat pravděpodobnost nějakého budoucího jevu).

matematická statistika

Zkoumáme statistiky u nějakého výběru z daného základního souboru (populace).

Matematická statistika se snaží postihnout, do jaké míry jsou zjištěné výsledky relevantní pro celou populaci, případně se ze zjištěných dat pokouší zjistit nebo upřesnit vhodný teoretický model pro chování celého souboru (a z něj pak třeba odhadovat pravděpodobnost nějakého budoucího jevu).

Dva základní přístupy:

- **frekvenční statistika** (nebo také klasická statistika)
- **bayesovská statistika**.

frekvenční přístup

- Vychází z matematické abstrakce, že skutečné pravděpodobnosti jsou dány četnostmi výskytů jevů v tak velkých vzorcích dat, že je můžeme dobře aproximovat nekonečnými modely a využít pro odhady spolehlivosti centrální limitní věty.
- Statistik zde na pravděpodobnost pohlíží jako na idealizaci relativní četnosti případů, v nichž se vyskytne určitý výsledek při opakovaných pokusech.

frekvenční přístup

- Vychází z matematické abstrakce, že skutečné pravděpodobnosti jsou dány četnostmi výskytů jevů v tak velkých vzorcích dat, že je můžeme dobře aproximovat nekonečnými modely a využít pro odhady spolehlivosti centrální limitní věty.
- Statistik zde na pravděpodobnost pohlíží jako na idealizaci relativní četnosti případů, v nichž se vyskytne určitý výsledek při opakovaných pokusech.
- Tato zdánlivá výhoda/rigoróznost se může ale rychle stát nevýhodou, jakmile se začneme zabývat spolehlivostí samotných dat a vhodností zvoleného experimentu.
- Stejně tak je obtížné frekvenční statistiku dobře použít pro odhad pravděpodobnosti výskytu jednorázového děje.

Bayesovský přístup

Tento přístup můžeme brát jako příklad matematizace „selského rozumu“. Vstupujeme do procesu s jistým původním přesvědčením, které jsme připraveni postupně pozměňovat ve světle nových dat.

- Jako vstupní předpoklad máme nějaké rozdělení pravděpodobnosti pro odhadovaných parametr,
- samotná data považujeme za konstanty, které hrají roli hypotézy v podmíněné pravděpodobnosti
- výsledkem je upřesnění rozdělení pravděpodobnosti zkoumaného parametru.

Bayesovský přístup

Tento přístup můžeme brát jako příklad matematizace „selského rozumu“. Vstupujeme do procesu s jistým původním přesvědčením, které jsme připraveni postupně pozměňovat ve světle nových dat.

- Jako vstupní předpoklad máme nějaké rozdělení pravděpodobnosti pro odhadovaných parametr,
- samotná data považujeme za konstanty, které hrají roli hypotézy v podmíněné pravděpodobnosti
- výsledkem je upřesnění rozdělení pravděpodobnosti zkoumaného parametru.

Je zajímavé, že historicky byl zjevně první bayesovský přístup (např. Laplace a další již v 18. století), který byl prakticky zcela vystřídán frekvenční statistikou ve 20. století. V posledních desetiletích se však ale bayesovská statistika vrátila, společně s dalšími novými přístupy, do popředí zájmu.

My se jí ale v tomtokurzu nebudeme zabývat.

Plán přednášky

- 1 Literatura
- 2 Centrální limitní věta
- 3 Co potkáme
- 4 Matematická statistika
- 5 Výběry z populací**
- 6 Intervaly spolehlivosti

Máme k dispozici (velký) základní statistický soubor s N jednotkami, který nazýváme **populace**, a zároveň nějaký číselný znak pro každou z jednotek, tj. soubor hodnot (x_1, \dots, x_N) . Z něj ovšem máme k dispozici pouze **výběrový soubor** s hodnotami (X_1, \dots, X_n) .

Máme k dispozici (velký) základní statistický soubor s N jednotkami, který nazýváme **populace**, a zároveň nějaký číselný znak pro každou z jednotek, tj. soubor hodnot (x_1, \dots, x_N) . Z něj ovšem máme k dispozici pouze **výběrový soubor** s hodnotami (X_1, \dots, X_n) .

Abychom se vyhnuli diskusi skutečné velikosti základního statistického souboru s N jednotkami, budeme předpokládat, že vybíráme položky výběrového souboru jednu po druhé a každou vybranou jednotku poté do populace vrátíme. Zároveň předpokládáme, že každá položka má stejnou pravděpodobnost výběru $1/N$. Hovoříme pak o **náhodném výběru**.

Pracujeme tedy s vektorem (X_1, \dots, X_n) nezávislých náhodných veličin a všechny tyto veličiny mají stejné rozdělení pravděpodobnosti. Zejména tedy budou sdílet distribuční funkci $F_X(x)$ a momenty

$$E X_i = \mu, \quad \text{var } X_i = \sigma^2.$$

Dalším naším krokem musí být odvození charakteristik výběrového průměru \bar{X} a **výběrového rozptylu**

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

přičemž následující věta dává hned zdůvodnění, proč volíme koeficient $\frac{1}{n-1}$ místo $\frac{1}{n}$.

Theorem

Pro výběrový průměr \bar{X} spočítaný z náhodného výběru rozsahu n z rozdělení s konečnou střední hodnotou μ a konečným rozptylem σ^2 platí

$$E \bar{X} = \mu, \quad \text{var } \bar{X} = \frac{1}{n} \sigma^2.$$

Pro výběrový rozptyl S^2 platí

$$E S^2 = \sigma^2.$$

Theorem

Pro výběrový průměr \bar{X} spočítaný z náhodného výběru rozsahu n z rozdělení s konečnou střední hodnotou μ a konečným rozptylem σ^2 platí

$$E \bar{X} = \mu, \quad \text{var } \bar{X} = \frac{1}{n} \sigma^2.$$

Pro výběrový rozptyl S^2 platí

$$E S^2 = \sigma^2.$$

Naším úkolem je odhadovat charakteristiky, jako jsou průměr μ hodnot znaku \bar{x} nebo jejich rozptyl σ^2 pro celou populaci pomocí obdobných charakteristik pro náš daleko menší výběr, které budeme značit pomocí velkých písmen, např. \bar{X} , S^2 .

Zde vstupuje do hry pravděpodobnost – budeme chtít znát pravděpodobnost přiblížení hodnot pro náš výběr těm pro celou populaci.

Zde vstupuje do hry pravděpodobnost – budeme chtít znát pravděpodobnost přiblížení hodnot pro náš výběr těm pro celou populaci.

Říkáme, že \bar{X} je nestranným odhadem střední hodnoty znaku pro populaci, zatímco výběrový rozptyl je nestranným odhadem rozptylu.

Zde vstupuje do hry pravděpodobnost – budeme chtít znát pravděpodobnost přiblížení hodnot pro náš výběr těm pro celou populaci.

Říkáme, že \bar{X} je nestranným odhadem střední hodnoty znaku pro populaci, zatímco výběrový rozptyl je nestranným odhadem rozptylu.

V případě, že bychom realizovali výběr z populace bez vracení, bude výběrový průměr stále nestranným odhadem střední hodnoty, výběrový rozptyl ale již ne (vyskočí tam faktor $(N - 1)/N$).

V praktických úlohách je třeba znát nejen číselné charakteristiky výběrového průměru a rozptylu, ale jejich úplné rozdělení pravděpodobnosti. To můžeme samozřejmě odvodit, pouze známe-li konkrétní rozdělení pravděpodobnosti X_i . Jako užitečnou ilustraci se podíváme na náhodný výběr z normálního rozdělení.

Výběrový průměr bude mít normální rozdělení a protože již známe jeho střední hodnotu a rozptyl, bude $\bar{X} \sim N(\mu, \frac{1}{n}\sigma^2)$.

O něco složitější je to s odvozením rozdělení pravděpodobnosti výběrového rozptylu. Uvažme vektor Z normovaných normálních veličin

$$Z_i = \frac{X_i - \mu}{\sigma}.$$

Theorem

Je-li (X_1, \dots, X_n) náhodný výběr z rozdělení $N(\mu, \sigma^2)$, pak jsou \bar{X} a S^2 nezávislé veličiny a platí

$$\bar{X} \sim N\left(\mu, \frac{1}{n}\sigma^2\right), \quad \frac{n-1}{\sigma^2}S^2 \sim \chi_{n-1}^2.$$

Okamžitým důsledkem je, že normalizovaný výběrový průměr

$$T = \sqrt{n} \frac{\bar{X} - \mu}{S}$$

má studentovo t-rozdělení pravděpodobnosti s $n - 1$ stupni volnosti.

Plán přednášky

- 1 Literatura
- 2 Centrální limitní věta
- 3 Co potkáme
- 4 Matematická statistika
- 5 Výběry z populací
- 6 Intervaly spolehlivosti**

Velmi častou úlohou je pro spočtenou hodnotu \bar{X} výběrového průměru určit interval, ve kterém se skutečná hodnota průměru veličiny pro celou populaci nachází s předem danou (vysokou) pravděpodobností.

Velmi častou úlohou je pro spočtenou hodnotu \bar{X} výběrového průměru určit interval, ve kterém se skutečná hodnota průměru veličiny pro celou populaci nachází s předem danou (vysokou) pravděpodobností.

Pro náhodnou veličinu X s normálním rozdělením máme její normovanou veličinu $Z = \frac{X - EX}{\sqrt{\text{var } X}}$. Normovaný výběrový průměr n veličin $X \sim N(0, 1)$ je $\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}}$ a chceme najít takovýto interval pro pravděpodobnost $1 - \alpha$, $\alpha \in (0, 1)$.

Velmi častou úlohou je pro spočtenou hodnotu \bar{X} výběrového průměru určit interval, ve kterém se skutečná hodnota průměru veličiny pro celou populaci nachází s předem danou (vysokou) pravděpodobností.

Pro náhodnou veličinu X s normálním rozdělením máme její normovanou veličinu $Z = \frac{X - EX}{\sqrt{\text{var } X}}$. Normovaný výběrový průměr n veličin $X \sim N(0, 1)$ je $\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}}$ a chceme najít takovýto interval pro pravděpodobnost $1 - \alpha$, $\alpha \in (0, 1)$.

Potřebujeme tedy znát hodnotu $z(\alpha)$ takovou, že $P(Z > z(\alpha)) = \alpha$.

Velmi častou úlohou je pro spočtenou hodnotu \bar{X} výběrového průměru určit interval, ve kterém se skutečná hodnota průměru veličiny pro celou populaci nachází s předem danou (vysokou) pravděpodobností.

Pro náhodnou veličinu X s normálním rozdělením máme její normovanou veličinu $Z = \frac{X - E X}{\sqrt{\text{var } X}}$. Normovaný výběrový průměr n veličin $X \sim N(0, 1)$ je $\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}}$ a chceme najít takovýto interval pro pravděpodobnost $1 - \alpha$, $\alpha \in (0, 1)$.

Potřebujeme tedy znát hodnotu $z(\alpha)$ takovou, že $P(Z > z(\alpha)) = \alpha$.

Je-li $F(x)$ spojitá rostoucí distribuční funkce naší veličiny, pak zjevně $z(\alpha) = F^{-1}(1 - \alpha)$. Pro normální rozdělení splňuje distribuční funkce Φ tento požadavek. Takto definovaným hodnotám $z(\alpha)$ se říká **kritické hodnoty**.

Protože je hustota pro normální rozdělení symetrická kolem jeho střední hodnoty, dostáváme $1 - \alpha = P(|Z| < z(\alpha/2))$.

$$\begin{aligned} 1 - \alpha &= P \left(\left| \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \right| < z(\alpha/2) \right) \\ &= P \left(\bar{X} - \frac{\sigma}{\sqrt{n}} z(\alpha/2) < \mu < \bar{X} + \frac{\sigma}{\sqrt{n}} z(\alpha/2) \right) \end{aligned}$$

což je interval s náhodnými konci, který s námi určenou pravděpodobností pokrývá neznámý parametr μ . V kontextu takových úloh hovoříme o **intervalu spolehlivosti s koeficientem spolehlivosti** $1 - \alpha$.

$$\begin{aligned} 1 - \alpha &= P \left(\left| \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \right| < z(\alpha/2) \right) \\ &= P \left(\bar{X} - \frac{\sigma}{\sqrt{n}} z(\alpha/2) < \mu < \bar{X} + \frac{\sigma}{\sqrt{n}} z(\alpha/2) \right) \end{aligned}$$

což je interval s náhodnými konci, který s námi určenou pravděpodobností pokrývá neznámý parametr μ . V kontextu takových úloh hovoříme o **intervalu spolehlivosti s koeficientem spolehlivosti** $1 - \alpha$.

Pro normální rozdělení je velice populární kritická hodnota $z(0,025) = 1,96$, která odpovídá naší úloze se zvolenou pravděpodobností 95%.

$$\begin{aligned}
 1 - \alpha &= P \left(\left| \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \right| < z(\alpha/2) \right) \\
 &= P \left(\bar{X} - \frac{\sigma}{\sqrt{n}} z(\alpha/2) < \mu < \bar{X} + \frac{\sigma}{\sqrt{n}} z(\alpha/2) \right)
 \end{aligned}$$

což je interval s náhodnými konci, který s námi určenou pravděpodobností pokrývá neznámý parametr μ . V kontextu takových úloh hovoříme o **intervalu spolehlivosti s koeficientem spolehlivosti** $1 - \alpha$.

Pro normální rozdělení je velice populární kritická hodnota $z(0,025) = 1,96$, která odpovídá naší úloze se zvolenou pravděpodobností 95%.

Kritické hodnoty jsou dány pomocí tzv. **kvantilové funkce**

$$F^{-1}(u) = \inf\{x \in \mathbb{R}; F(x) \geq u\}, \quad 0 < u < 1.$$

Kvantilová funkce skutečně dává přímo příslušné kvantily, např. $F^{-1}(0,5)$ je medián, atd.

Example

Před deseti lety byl uskutečněn rozsáhlý výzkum výšky desetiletých chlapců a zjistilo se, že střední výška byla $\mu_0 = 136,1\text{cm}$ se směrodatnou odchylkou $\sigma = 6,4\text{cm}$. Nyní byly na náhodném výběru 15 desetiletých chlapců zjištěny následující výšky: 130, 140, 136, 141, 139, 133, 149, 151, 139, 136, 138, 142, 127, 139, 147. Je známo, že variabilita výšek v populaci se mění velice pomalu, zatímco výšky se mohou měnit rychle. **Otázka: došlo ke změně střední výšky populace desetiletých chlapců?**

Example

Před deseti lety byl uskutečněn rozsáhlý výzkum výšky desetiletých chlapců a zjistilo se, že střední výška byla $\mu_0 = 136,1$ cm se směrodatnou odchylkou $\sigma = 6,4$ cm. Nyní byly na náhodném výběru 15 desetiletých chlapců zjištěny následující výšky: 130, 140, 136, 141, 139, 133, 149, 151, 139, 136, 138, 142, 127, 139, 147. Je známo, že variabilita výšek v populaci se mění velice pomalu, zatímco výšky se mohou měnit rychle. **Otázka: došlo ke změně střední výšky populace desetiletých chlapců?**

Ze zadání předpokládáme, že výběr 15 hodnot je z normálního rozdělení se známým rozptylem σ^2 a otázku si upřesníme tak, že hledáme v jakém intervalu je nyní střední hodnota výšky populace se spolehlivostí 95% :

Example

Před deseti lety byl uskutečněn rozsáhlý výzkum výšky desetiletých chlapců a zjistilo se, že střední výška byla $\mu_0 = 136,1$ cm se směrodatnou odchylkou $\sigma = 6,4$ cm. Nyní byly na náhodném výběru 15 desetiletých chlapců zjištěny následující výšky: 130, 140, 136, 141, 139, 133, 149, 151, 139, 136, 138, 142, 127, 139, 147. Je známo, že variabilita výšek v populaci se mění velice pomalu, zatímco výšky se mohou měnit rychle. **Otázka: došlo ke změně střední výšky populace desetiletých chlapců?**

Ze zadání předpokládáme, že výběr 15 hodnot je z normálního rozdělení se známým rozptylem σ^2 a otázku si upřesníme tak, že hledáme v jakém intervalu je nyní střední hodnota výšky populace se spolehlivostí 95% : $\bar{x} = 139,133$ a tedy interval spolehlivosti je $(139,133 - (6,4/\sqrt{15}), 139,133 + (6,4/\sqrt{15})) = (135,9, 142,4)$.

Example

Před deseti lety byl uskutečněn rozsáhlý výzkum výšky desetiletých chlapců a zjistilo se, že střední výška byla $\mu_0 = 136,1$ cm se směrodatnou odchylkou $\sigma = 6,4$ cm. Nyní byly na náhodném výběru 15 desetiletých chlapců zjištěny následující výšky: 130, 140, 136, 141, 139, 133, 149, 151, 139, 136, 138, 142, 127, 139, 147. Je známo, že variabilita výšek v populaci se mění velice pomalu, zatímco výšky se mohou měnit rychle. **Otázka: došlo ke změně střední výšky populace desetiletých chlapců?**

Ze zadání předpokládáme, že výběr 15 hodnot je z normálního rozdělení se známým rozptylem σ^2 a otázku si upřesníme tak, že hledáme v jakém intervalu je nyní střední hodnota výšky populace se spolehlivostí 95% : $\bar{x} = 139,133$ a tedy interval spolehlivosti je $(139,133 - (6,4/\sqrt{15}), 139,133 + (6,4/\sqrt{15})) = (135,9, 142,4)$.

Protože tento interval pokrývá i populační průměr před deseti lety, nemůžeme na této hladině spolehlivosti tvrdit, že se populační výška změnila.