

## Základní pojmy matematické statistiky

### **Motivace:**

Matematická statistika je věda, která analyzuje a interpretuje data především za účelem získání předpovědi a zlepšení rozhodování v různých oborech lidské činnosti. Přitom se řídí principem statistické indukce, tj. na základě znalostí o náhodném výběru z určitého rozložení pravděpodobností se snaží učinit závěry o vlastnostech tohoto rozložení.

Ústředním pojmem matematické statistiky je tedy pojem náhodného výběru.

### Definice náhodného výběru:

- a) Necht'  $X_1, \dots, X_n$  jsou stochasticky nezávislé náhodné veličiny, které mají všechny stejné rozložení  $L(\vartheta)$ . Řekneme, že  $X_1, \dots, X_n$  je **náhodný výběr rozsahu  $n$  z rozložení  $L(\vartheta)$** . (Číselné realizace  $x_1, \dots, x_n$  náhodného výběru  $X_1, \dots, X_n$  uspořádané do sloupcového vektoru odpovídají datovému souboru zavedenému v popisné statistice.)
- b) Necht'  $(X_1, Y_1), \dots, (X_n, Y_n)$  jsou stochasticky nezávislé dvourozměrné náhodné vektory, které mají všechny stejné dvourozměrné rozložení  $L_2(\vartheta)$ . Řekneme, že  $(X_1, Y_1), \dots, (X_n, Y_n)$  je **dvourozměrný náhodný výběr rozsahu  $n$  z dvourozměrného rozložení  $L_2(\vartheta)$** . (Číselné realizace  $(x_1, y_1), \dots, (x_n, y_n)$  náhodného výběru  $(X_1, Y_1), \dots, (X_n, Y_n)$  uspořádané do matice typu  $n \times 2$  odpovídají dvourozměrnému datovému souboru zavedenému v popisné statistice.)
- c) Analogicky lze definovat  $p$ -rozměrný **náhodný výběr rozsahu  $n$  z  $p$ -rozměrného rozložení  $L_p(\vartheta)$** .

### Definice statistiky:

Libovolná funkce  $T = T(X_1, \dots, X_n)$  náhodného výběru  $X_1, \dots, X_n$  (resp.  $T = T(X_1, Y_1, \dots, X_n, Y_n)$  náhodného výběru  $(X_1, Y_1), \dots, (X_n, Y_n)$ ) se nazývá (výběrová) **statistika**.

### Definice důležitých statistik:

a) Necht'  $X_1, \dots, X_n$  je náhodný výběr,  $n \geq 2$ .

Označme  $M = \frac{1}{n} \sum_{i=1}^n X_i$  ... **výběrový průměr**,  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - M)^2$  ... **výběrový rozptyl**,  $S = \sqrt{S^2}$  ... **výběrová směrodatná odchylka**

Pro libovolné, ale pevně dané reálné číslo  $x$  je statistikou též hodnota **výběrové distribuční funkce**  $F_n(x) = \frac{1}{n} \text{card}\{i; X_i \leq x\}$

b) Necht' je dáno  $r \geq 2$  stochasticky nezávislých náhodných výběrů o rozsazích  $n_1 \geq 2, \dots, n_r \geq 2$ .

Celkový rozsah je  $n = \sum_{j=1}^r n_j$ .

Označme  $M_1, \dots, M_r$  výběrové průměry a  $S_1^2, \dots, S_r^2$  výběrové rozptyly jednotlivých výběrů. Necht'  $c_1, \dots, c_r$  jsou reálné konstanty, aspoň jedna nenulová.

$\sum_{j=1}^r c_j M_j$  ... **lineární kombinace výběrových průměrů**,  $S_*^2 = \frac{\sum_{j=1}^r (n_j - 1) S_j^2}{n - r}$  ... **vážený průměr výběrových rozptylů**.

c) Necht'  $(X_1, Y_1), \dots, (X_n, Y_n)$  je náhodný výběr z dvourozměrného rozložení o rozsahu  $n$ .

Označme  $M_1 = \frac{1}{n} \sum_{i=1}^n X_i$ ,  $M_2 = \frac{1}{n} \sum_{i=1}^n Y_i$  výběrové průměry,  $S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_1)^2$ ,  $S_2^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - M_2)^2$  výběrové rozptyly.

$S_{12} = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_1)(Y_i - M_2)$  ... **výběrová kovariance**,  $R_{12} = \begin{cases} \frac{S_{12}}{S_1 S_2} & \text{pro } S_1 S_2 \neq 0 \\ 0 & \text{jinak} \end{cases}$  ... **výběrový koeficient korelace**.

Pro libovolnou, ale pevně zvolenou dvojici reálných čísel  $x, y$  je statistikou též hodnota **výběrové simultánní distribuční funkce**  $F_n(x, y) = \frac{1}{n} \text{card}\{i; X_i \leq x \wedge Y_i \leq y\}$ .

**Upozornění:** Číselné realizace statistik  $M, S^2, S, S_{12}, R_{12}$  odpovídají číselným charakteristikám  $m, s^2, s, s_{12}, r_{12}$  zavedeným v popisné statistice, ale u rozptylu, směrodatné odchylky, kovariance a koeficientu korelace je multiplikační konstanta  $\frac{1}{n-1}$ , nikoliv  $\frac{1}{n}$ , jak tomu bylo v popisné statistice. Jak uvidíme později, uvedené číselné realizace mohou být považovány za **odhady** číselných realizací náhodných veličin zavedených v počtu pravděpodobnosti.

Charakteristika vlastnosti	Počet pravděpodobnosti	Matematická statistika	Popisná statistika
poloha	$E(X) = \mu$	$M$	$m$
variabilita	$D(X) = \sigma^2$	$S^2$	$\frac{n-1}{n} s^2$
variabilita	$\sqrt{D(X)} = \sigma$	$S$	$\sqrt{\frac{n-1}{n}} s$
společná variabilita	$C(X_1, X_2) = \sigma_{12}$	$S_{12}$	$\frac{n-1}{n} s_{12}$
těsnost vztahu	$R(X_1, X_2) = \rho$	$R_{12}$	$r_{12}$
rozložení	$\Phi(x)$	$F_n(x)$	$F(x)$

**Příklad** (výpočet realizací výběrového průměru, výběrového rozptylu a hodnot výběrové distribuční funkce):

Desetkrát nezávisle na sobě byla změřena jistá konstanta  $\mu$ . Výsledky měření byly: 2 1,8 2,1 2,4 1,9 2,1 2 1,8 2,3 2,2. Tyto výsledky považujeme za číselné realizace náhodného výběru  $X_1, \dots, X_{10}$ . Vypočtěte realizaci  $m$  výběrového průměru  $M$ , realizaci  $s^2$  výběrového rozptylu  $S^2$ , realizaci  $s$  výběrové směrodatné odchylky  $S$  a hodnoty výběrové distribuční funkce  $F_{10}(x)$ .

**Řešení:**

$$m = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{10} (2 + 1,8 + \dots + 2,2) = 2,06, s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2 = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - nm^2 \right) = \frac{1}{9} (2^2 + 1,8^2 + \dots + 2,2^2 - 10 \cdot 2,06^2) = 0,0404$$

$$s = \sqrt{s^2} = \sqrt{0,0404} = 0,2011$$

Pro usnadnění výpočtu hodnot výběrové distribuční funkce  $F_{10}(x)$  uspořádáme měření podle velikosti:

1,8 1,8 1,9 2 2 2,1 2,1 2,2 2,3 2,4.

$$x < 1,8 : F_{10}(x) = 0$$

$$1,8 \leq x < 1,9 : F_{10}(x) = \frac{2}{10} = 0,2$$

$$1,9 \leq x < 2 : F_{10}(x) = \frac{3}{10} = 0,3$$

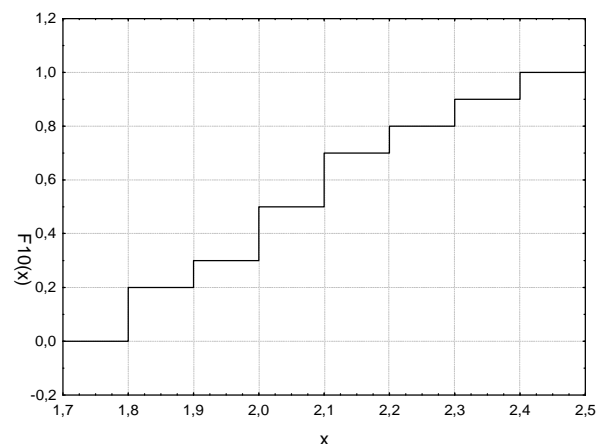
$$2 \leq x < 2,1 : F_{10}(x) = \frac{5}{10} = 0,5$$

$$2,1 \leq x < 2,2 : F_{10}(x) = \frac{7}{10} = 0,7$$

$$2,2 \leq x < 2,3 : F_{10}(x) = \frac{8}{10} = 0,8$$

$$2,3 \leq x < 2,4 : F_{10}(x) = \frac{9}{10} = 0,9$$

$$x \geq 2,4 : F_{10}(x) = 1$$



**Příklad** (výpočet realizace výběrového koeficientu korelace):

U 11 náhodně vybraných aut jisté značky bylo zjišťováno jejich stáří (náhodná veličina  $X$  – v letech) a cena (náhodná veličina  $Y$  – v tisících Kč). Výsledky:

(5, 85), (4, 103), (6, 70), (5, 82), (5, 89), (5, 98), (6, 66), (6, 95), (2, 169), (7, 70), (7, 48).

Vypočítejte a interpretujte číselnou realizaci  $r_{12}$  výběrového koeficientu korelace  $R_{12}$ .

**Řešení:**

$$m_1 = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{11} (5 + 4 + \dots + 7) = 5,28$$

$$m_2 = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{11} (85 + 103 + \dots + 48) = 88,63$$

$$s_1^2 = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - nm_1^2 \right) = \frac{1}{10} (5^2 + 4^2 + \dots + 7^2 - 11 \cdot 5,28^2) = 2,02$$

$$s_2^2 = \frac{1}{n-1} \left( \sum_{i=1}^n y_i^2 - nm_2^2 \right) = \frac{1}{10} (85^2 + 103^2 + \dots + 48^2 - 11 \cdot 88,63^2) = 970,85$$

$$s_{12} = \frac{1}{n-1} \left( \sum_{i=1}^n x_i y_i - nm_1 m_2 \right) = \frac{1}{10} (5 \cdot 85 + 4 \cdot 103 + \dots + 7 \cdot 48 - 11 \cdot 5,28 \cdot 88,63) = -40,89$$

$$r_{12} = \frac{s_{12}}{s_1 \cdot s_2} = \frac{-40,82}{\sqrt{2,02} \cdot \sqrt{970,85}} = -0,92$$

Mezi náhodnými veličinami  $X$  a  $Y$  existuje silná nepřímá lineární závislost. Čím starší auto, tím nižší cena.

## Bodové a intervalové odhady parametrů a parametrických funkcí

Vycházíme z náhodného výběru  $X_1, \dots, X_n$  z rozložení  $L(\vartheta)$ , které závisí na parametru  $\vartheta$ . Množinu všech přípustných hodnot tohoto parametru označíme  $\Xi$ . Tato množina se nazývá **parametrický prostor**.

Např. je-li  $X_1, \dots, X_n$  náhodný výběr z rozložení  $N(\mu, \sigma^2)$ , pak  $\vartheta = (\mu, \sigma^2)$  a v tomto případě parametrický prostor  $\Xi = (-\infty, \infty) \times (0, \infty)$ .

Parametr  $\vartheta$  neznáme a chceme ho odhadnout pomocí daného náhodného výběru (případně chceme odhadnout nějakou **parametrickou funkci**  $h(\vartheta)$ ).

**Bodovým odhadem** parametrické funkce  $h(\vartheta)$  je statistika  $T_n = T(X_1, \dots, X_n)$ , která nabývá hodnot blízkých  $h(\vartheta)$ , ať je hodnota parametru  $\vartheta$  jakákoliv. Existují různé metody, jak konstruovat bodové odhady (např. metoda momentů či metoda maximální věrohodnosti, ale těmi se zde zabývat nebudeme) a také různé typy bodových odhadů. Omezíme se na odhady nestranné, asymptoticky nestranné a konzistentní.

**Intervalovým odhadem** parametrické funkce  $h(\vartheta)$  rozumíme interval  $(D, H)$ , jehož meze jsou statistiky  $D = D(X_1, \dots, X_n)$ ,  $H = H(X_1, \dots, X_n)$  a který s dostatečně velkou pravděpodobností pokrývá  $h(\vartheta)$ , ať je hodnota parametru  $\vartheta$  jakákoliv.

## Typy bodových odhadů

Nechť  $X_1, \dots, X_n$  je náhodný výběr z rozložení  $L(\vartheta)$ ,  $h(\vartheta)$  je parametrická funkce,  $T, T_1, T_2, \dots$  jsou statistiky.

a) Řekneme, že statistika  $T$  je **nestranným odhadem** parametrické funkce  $h(\vartheta)$ , jestliže

$$\forall \vartheta \in \Xi : E(T) = h(\vartheta).$$

(Význam nestrannosti spočívá v tom, že odhad  $T$  nesmí parametrickou funkci  $h(\vartheta)$  systematicky nadhodnocovat ani podhodnocovat. Není-li tato podmínka splněna, jde o vychýlený odhad.)

b) Jsou-li  $T_1, T_2$  nestranné odhady téže parametrické funkce  $h(\vartheta)$ , pak řekneme, že  $T_1$  je **lepší odhad** než  $T_2$ , jestliže

$$\forall \vartheta \in \Xi : D(T_1) < D(T_2).$$

c) Posloupnost  $\{T_n\}_{n=1}^{\infty}$  se nazývá **posloupnost asymptoticky nestranných odhadů** parametrické funkce  $h(\vartheta)$ , jestliže

$$\forall \vartheta \in \Xi : \lim_{n \rightarrow \infty} E(T_n) = h(\vartheta).$$

(Význam asymptotické nestrannosti spočívá v tom, že s rostoucím rozsahem výběru klesá vychýlení odhadu.)

d) Posloupnost  $\{T_n\}_{n=1}^{\infty}$  se nazývá **posloupnost konzistentních odhadů** parametrické funkce  $h(\vartheta)$ , jestliže

$$\forall \vartheta \in \Xi \forall \varepsilon > 0 : \lim_{n \rightarrow \infty} P(|T_n - h(\vartheta)| > \varepsilon) = 0.$$

(Význam konzistence spočívá v tom, že s rostoucím rozsahem výběru klesá pravděpodobnost, že odhad se bude realizovat „daleko“ od parametrické funkce  $h(\vartheta)$ .)

Lze dokázat, že z nestrannosti odhadu vyplývá jeho asymptotická nestrannost a z asymptotické nestrannosti vyplývá konzistence, pokud posloupnost rozptylů odhadu konverguje k nule.



### Vlastnosti důležitých statistik

a) **Případ jednoho náhodného výběru:** Necht'  $X_1, \dots, X_n$  je náhodný výběr z rozložení se střední hodnotou  $\mu$ , rozptylem  $\sigma^2$  a distribuční funkcí  $\Phi(x)$ . Necht'  $n \geq 2$ . Označme  $M_n$  výběrový průměr,  $S_n^2$  výběrový rozptyl a pro libovolné, ale pevně dané  $x \in \mathbf{R}$  označme  $F_n(x)$  hodnotu výběrové distribuční funkce. Pak pro libovolné hodnoty parametrů  $\mu$ ,  $\sigma^2$  a libovolné, ale pevně dané reálné číslo  $x$  platí:

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

$$D(M_n) = \frac{\sigma^2}{n},$$

$$E(S_n^2) = \sigma^2,$$

$$D(S_n^2) = \frac{\gamma_4}{n} - \frac{\sigma^4(n-3)}{n(n-1)}, \text{ kde } \gamma_4 \text{ je 4. centrální moment,}$$

$$E(F_n(x)) = \Phi(x),$$

$$D(F_n(x)) = \frac{\Phi(x)[1 - \Phi(x)]}{n}$$

Znamená to, že  $M_n$  je nestranným odhadem  $\mu$ ,  $S_n^2$  je nestranným odhadem  $\sigma^2$ , pro libovolné, ale pevně dané  $x \in \mathbf{R}$  je výběrová distribuční funkce  $F_n(x)$  nestranným odhadem  $\Phi(x)$ .

Posloupnost  $\{M_n\}_{n=1}^{\infty}$  je posloupnost konzistentních odhadů  $\mu$ ,

$\{S_n^2\}_{n=1}^{\infty}$  je posloupnost konzistentních odhadů  $\sigma^2$ ,

pro libovolné, ale pevně dané  $x \in \mathbf{R}$  je  $\{F_n(x)\}_{n=1}^{\infty}$  posloupnost konzistentních odhadů  $\Phi(x)$ .

b) **Případ  $r \geq 2$  stochasticky nezávislých náhodných výběrů:** Necht'  $X_{11}, \dots, X_{1n_1}, \dots, X_{r1}, \dots, X_{rn_r}$  je  $r$  stochasticky nezávislých náhodných výběrů o rozsazích  $n_1 \geq 2, \dots, n_r \geq 2$  z rozložení se středními hodnotami  $\mu_1, \dots, \mu_r$  a rozptylem  $\sigma^2$ . Celkový rozsah je  $n = \sum_{j=1}^r n_j$ . Necht'  $c_1, \dots, c_r$  jsou reálné konstanty, aspoň jedna nenulová. Pak pro libovolné hodnoty parametrů  $\mu_1, \dots, \mu_r$  a  $\sigma^2$  platí:

$$E\left(\sum_{j=1}^r c_j M_j\right) = \sum_{j=1}^r c_j \mu_j,$$

$$E(S_*^2) = \sigma^2.$$

Znamená to, že lineární kombinace výběrových průměrů  $\sum_{j=1}^r c_j M_j$  je nestranným odhadem lineární kombinace středních hod-

not  $\sum_{j=1}^r c_j \mu_j$  a vážený průměr výběrových rozptylů  $S_*^2 = \frac{\sum_{j=1}^r (n_j - 1) S_j^2}{n - r}$  je nestranným odhadem rozptylu  $\sigma^2$ .

c) **Případ jednoho náhodného výběru z dvourozměrného rozložení:** Necht'  $(X_1, Y_1), \dots, (X_n, Y_n)$  je náhodný výběr z dvourozměrného rozložení s kovariancí  $\sigma_{12}$  a koeficientem korelace  $\rho$ . Pak pro libovolné hodnoty parametrů  $\sigma_{12}$  a  $\rho$  platí:

$$E(S_{12}) = \sigma_{12},$$

$$E(R_{12}) \approx \rho \quad (\text{shoda je vyhovující pro } n \geq 30).$$

Znamená to, že výběrová kovariance  $S_{12}$  je nestranným odhadem kovariance  $\sigma_{12}$ , avšak výběrový koeficient korelace  $R_{12}$  je vychýleným odhadem koeficientu korelace  $\rho$ .

## Pojem intervalu spolehlivosti

Nechť  $X_1, \dots, X_n$  je náhodný výběr z rozložení  $L(\vartheta)$ ,

$h(\vartheta)$  je parametrická funkce,

$\alpha \in (0,1)$ ,

$D = D(X_1, \dots, X_n)$ ,  $H = H(X_1, \dots, X_n)$  jsou statistiky.

a) Interval  $(D, H)$  se nazývá **100(1- $\alpha$ )% (oboustranný) interval spolehlivosti** pro parametrickou funkci  $h(\vartheta)$ ,

jestliže:  $\forall \vartheta \in \Xi : P(D < h(\vartheta) < H) \geq 1-\alpha$ .

b) Interval  $(D, \infty)$  se nazývá **100(1- $\alpha$ )% levostranný interval spolehlivosti** pro parametrickou funkci  $h(\vartheta)$ ,

jestliže:  $\forall \vartheta \in \Xi : P(D < h(\vartheta)) \geq 1-\alpha$ .

c) Interval  $(-\infty, H)$  se nazývá **100(1- $\alpha$ )% pravostranný interval spolehlivosti** pro parametrickou funkci  $h(\vartheta)$ ,

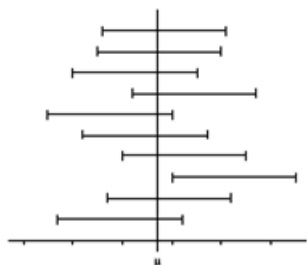
jestliže:  $\forall \vartheta \in \Xi : P(h(\vartheta) < H) \geq 1-\alpha$ .

Číslo  $\alpha$  se nazývá **riziko** (zpravidla  $\alpha = 0,05$ , méně často 0,1 či 0,01), číslo  $1 - \alpha$  se nazývá **spolehlivost**.

### Postup při konstrukci intervalu spolehlivosti

- a) Vyjdeme ze statistiky  $V$ , která je nestranným bodovým odhadem parametrické funkce  $h(\vartheta)$ .
- b) Najdeme tzv. pivotovou statistiku  $W$ , která vznikne transformací statistiky  $V$ , je monotónní funkcí  $h(\vartheta)$  a přitom její rozložení je známé a na  $h(\vartheta)$  nezávisí. Pomocí známého rozložení pivotové statistiky  $W$  najdeme kvantily  $w_{\alpha/2}$ ,  $w_{1-\alpha/2}$ , takže platí:  $\forall \vartheta \in \Xi: P(w_{\alpha/2} < W < w_{1-\alpha/2}) \geq 1 - \alpha$ .
- c) Nerovnost  $w_{\alpha/2} < W < w_{1-\alpha/2}$  převedeme ekvivalentními úpravami na nerovnost  $D < h(\vartheta) < H$ .
- d) Statistiky  $D$ ,  $H$  nahradíme jejich číselnými realizacemi  $d$ ,  $h$  a získáme tak  $100(1-\alpha)\%$  empirický interval spolehlivosti, o němž prohlásíme, že pokrývá  $h(\vartheta)$  s pravděpodobností aspoň  $1 - \alpha$ . (Tvrzení, že  $(d, h)$  pokrývá  $h(\vartheta)$  s pravděpodobností aspoň  $1 - \alpha$  je třeba chápat takto: jestliže mnohonásobně nezávisle získáme realizace  $x_1, \dots, x_n$  náhodného výběru  $X_1, \dots, X_n$  z rozložení  $L(\vartheta)$  a pomocí každé této realizace sestojíme  $100(1-\alpha)\%$  empirický interval spolehlivosti pro  $h(\vartheta)$ , pak podíl počtu těch intervalů, které pokrývají  $h(\vartheta)$  k počtu všech sestojených intervalů bude přibližně  $1 - \alpha$ .)

**Ilustrace:** Jestliže 100x nezávisle na sobě uskutečníme náhodný výběr z rozložení se střední hodnotou  $\mu$  a pokaždé sestrojíme 95% empirický interval spolehlivosti pro  $\mu$ , pak přibližně v 95-ti případech bude ležet parametr  $\mu$  v intervalech spolehlivosti a asi v 5-ti případech interval spolehlivosti  $\mu$  nepokryje.



**Volba oboustranného, jednostranného, nebo jednostranného intervalu:** závisí na konkrétní situaci.

Např. **oboustranný** interval spolehlivosti použije konstruktér, kterého zajímá dolní i horní hranice pro skutečnou délku  $\mu$  nějaké součástky.

**Jednostranný** interval spolehlivosti použije výkupčí drahých kovů, který potřebuje znát dolní mez pro skutečný obsah zlata  $\mu$  v kupovaném slitku.

**Jednostranný** interval spolehlivosti použije chemik, který potřebuje znát horní mez pro obsah nečistot  $\mu$  v analyzovaném vzorku.

**Příklad:** Necht'  $X_1, \dots, X_n$  je náhodný výběr z rozložení  $N(\mu, \sigma^2)$ , kde  $n \geq 2$  a rozptyl  $\sigma^2$  známe. Sestrojte  $100(1-\alpha)\%$  interval spolehlivosti pro neznámou střední hodnotu  $\mu$ .

**Řešení:** V tomto případě parametrická funkce  $h(\vartheta) = \mu$ . Nestranným odhadem střední hodnoty je výběrový průměr  $M = \frac{1}{n} \sum_{i=1}^n X_i$ . Protože  $M$  je lineární kombinací normálně rozložených náhodných veličin, bude mít také normální rozložení se střední hodnotou  $E(M) = \mu$  a rozptylem  $D(M) = \frac{\sigma^2}{n}$ . Pivotovou statistikou  $W$  bude standardizovaná náhodná veličina

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

Kvantil  $w_{\alpha/2} = u_{\alpha/2} = -u_{1-\alpha/2}$ ,  $w_{1-\alpha/2} = u_{1-\alpha/2}$ .

$$\forall \vartheta \in \Xi: 1 - \alpha \leq P(-u_{1-\alpha/2} < U < u_{1-\alpha/2}) = P\left(-u_{1-\alpha/2} < \frac{M - \mu}{\frac{\sigma}{\sqrt{n}}} < u_{1-\alpha/2}\right) = P\left(M - \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2} < \mu < M + \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2}\right).$$

Meze  $100(1-\alpha)\%$  intervalu spolehlivosti pro střední hodnotu  $\mu$  při známém rozptylu  $\sigma^2$  tedy jsou:

$$D = M - \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2}, \quad H = M + \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2}.$$

Při konstrukci jednostranných intervalů spolehlivosti se riziko nepůlí, tedy  $100(1-\alpha)\%$  jednostranný interval spolehlivosti pro  $\mu$  je  $\left(M - \frac{\sigma}{\sqrt{n}} u_{1-\alpha}, \infty\right)$  a pravostranný je  $\left(-\infty, M + \frac{\sigma}{\sqrt{n}} u_{1-\alpha}\right)$ .

Dosadíme-li do vzorců pro dolní a horní mez číselnou realizaci  $m$  výběrového průměru  $M$ , dostaneme  $100(1-\alpha)\%$  empirický interval spolehlivosti. Postup si ukážeme na následujícím numerickém příkladu.

**Příklad:** 10 krát nezávisle na sobě byla změřena jistá konstanta  $\mu$ . Výsledky měření byly:

2 1,8 2,1 2,4 1,9 2,1 2 1,8 2,3 2,2.

Výsledky považujeme za číselné realizace náhodného výběru  $X_1, \dots, X_{10}$  z rozložení  $N(\mu, \sigma^2)$ , kde  $\mu$  neznáme a  $\sigma^2 = 0,04$ .

Najděte 95% empirický interval spolehlivosti pro  $\mu$ , a to

- a) oboustranný,
- b) levostranný,
- c) pravostranný.

**Řešení:**

Vypočteme realizaci výběrového průměru:  $m = 2,06$ . Riziko  $\alpha$  je 0,05. V tabulkách najdeme kvantil  $u_{0,975} = 1,96$  pro oboustranný interval spolehlivosti a kvantil  $u_{0,95} = 1,64$  pro jednostranné intervaly spolehlivosti.

$$\text{ad a) } d = m - \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2} = 2,06 - \frac{0,2}{\sqrt{10}} 1,96 = 1,94$$

$$h = m + \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2} = 2,06 + \frac{0,2}{\sqrt{10}} 1,96 = 2,18$$

$1,94 < \mu < 2,18$  s pravděpodobností aspoň 0,95.

$$\text{ad b) } d = m - \frac{\sigma}{\sqrt{n}} u_{1-\alpha} = 2,06 - \frac{0,2}{\sqrt{10}} 1,64 = 1,96$$

$1,96 < \mu$  s pravděpodobností aspoň 0,95.

$$\text{ad c) } h = m + \frac{\sigma}{\sqrt{n}} u_{1-\alpha} = 2,06 + \frac{0,2}{\sqrt{10}} 1,64 = 2,16$$

$\mu < 2,16$  s pravděpodobností aspoň 0,95.

## Šířka intervalu spolehlivosti

Nechť  $(d, h)$  je  $100(1-\alpha)\%$  empirický interval spolehlivosti pro  $h(\vartheta)$  zkonstruovaný pomocí číselných realizací  $x_1, \dots, x_n$  náhodného výběru  $X_1, \dots, X_n$  z rozložení  $L(\vartheta)$ .

- Při konstantním riziku klesá šířka  $h-d$  s rostoucím rozsahem náhodného výběru.
- Při konstantním rozsahu náhodného výběru klesá šířka  $h-d$  s rostoucím rizikem.

**Příklad:** (stanovení minimálního rozsahu výběru z normálního rozložení)

Nechť  $X_1, \dots, X_n$  je náhodný výběr z  $N(\mu, \sigma^2)$ , kde  $\sigma^2$  známe. Jaký musí být minimální rozsah výběru  $n$ , aby šířka  $100(1-\alpha)\%$  empirického intervalu spolehlivosti pro střední hodnotu  $\mu$  nepřesáhla číslo  $\Delta$ ?

**Řešení:** Požadujeme, aby  $\Delta \geq h - d = m + \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2} - (m - \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2}) = \frac{2\sigma}{\sqrt{n}} u_{1-\alpha/2}$ . Z této podmínky dostaneme, že

$n \geq \frac{4\sigma^2 u_{1-\alpha/2}^2}{\Delta^2}$ . Za rozsah výběru zvolíme nejmenší přirozené číslo vyhovující této podmínce.

**Příklad:** Hloubka moře se měří přístrojem, jehož systematická chyba je nulová a náhodné chyby měření mají normální rozložení se směrodatnou odchylkou  $\sigma = 1$  m. Kolik měření je nutno provést, aby se hloubka stanovila s chybou nejvýše  $\pm 0,25$  m při spolehlivosti 0,95?

**Řešení:** Hledáme rozsah výběru tak, aby šířka 95% intervalu spolehlivosti pro střední hodnotu  $\mu$  nepřesáhla 0,5 m. Přitom  $\sigma$  známe. Z předešlého příkladu vyplývá, že  $n \geq \frac{4\sigma^2 u_{1-\alpha/2}^2}{\Delta^2} = \frac{4 \cdot 1,96^2}{0,5^2} = 61,4656$ . Nejmenší počet měření je tedy 62.



## Základní typy uspořádání pokusů

Metody matematické statistiky často slouží k vyhodnocování výsledků pokusů. Aby mohl být pokus správně vyhodnocen, musí být dobře naplánován. Uvedeme zde nejjednodušší typy uspořádání pokusů.

Předpokládejme například, že sledujeme hmotnostní přírůstky selat téhož plemene při různých výkrmných dietách.

a) **Jednoduché pozorování:** Náhodná veličina  $X$  je pozorována za týchž podmínek. Situace je charakterizována jedním náhodným výběrem  $X_1, \dots, X_n$ .

Náhodně vylosujeme  $n$  selat téhož plemene, podrobíme je jediné výkrmné dietě a zjistíme u každého selete hmotnostní přírůstek. Tím dostaneme realizaci jednoho náhodného výběru.

b) **Dvojné pozorování:** Náhodná veličina  $X$  je pozorována za dvojitých různých podmínek. Existují dvě odlišná uspořádání tohoto pokusu.

**Dvouvýběrové porovnávání:** situace je charakterizována dvěma nezávislými náhodnými výběry  $X_{11}, \dots, X_{1n_1}$  a  $X_{21}, \dots, X_{2n_2}$ .

Náhodně vylosujeme  $n_1$  a  $n_2$  selat téhož plemene, náhodně je rozdělíme na dva soubory o  $n_1$  a  $n_2$  jedincích, první podrobíme výkrmné dietě č. 1 a druhý výkrmné dietě číslo 2. Tak dostaneme realizace dvou nezávislých náhodných výběrů.

**Párové porovnávání:** situace je charakterizována jedním náhodným výběrem  $(X_{11}, X_{12}), \dots, (X_{n1}, X_{n2})$  z dvourozměrného rozložení. Přejdeme k rozdílovému náhodnému výběru  $Z_i = X_{i1} - X_{i2}$ ,  $i = 1, \dots, n$  a tím dostaneme jednoduché pozorování.

Náhodně vylosujeme  $n$  vrhů stejně starých selat téhož plemene, z každého odebereme dva sourozence a náhodně jim přiřadíme první a druhou výkrmnou dietu. Tak dostaneme realizaci jednoho dvourozměrného náhodného výběru, kde první složka odpovídá první dietě a druhá složka druhé dietě.

(Párové porovnávání je efektivnější, protože skutečný rozdíl v účinnosti obou diet je překrýván pouze náhodnými vlivy při samotném krmení a trvání, kdežto vliv různých dědičných vloh, který byl losováním znáhodněn, je u sourozeneckého páru selat částečně vyloučen.)

c) **Mnohonásobné pozorování:** Náhodná veličina  $X$  je pozorována za  $r \geq 3$  různých podmínek. Existují dvě odlišná uspořádání tohoto pokusu.

**Mnohovýběrové porovnávání:** situace je charakterizována  $r$  nezávislými náhodnými výběry  $X_{11}, \dots, X_{1n_1}$  až  $X_{r1}, \dots, X_{rn_r}$ . Náhodně vylosujeme  $n_1, n_2, \dots, n_r$  selat téhož plemene, náhodně je rozdělíme na  $r$  souborů o  $n_1, n_2, \dots, n_r$  jedincích, první podrobíme výkrmné dietě č. 1, druhý výkrmné dietě číslo 2 atd. až  $r$ -tý podrobíme výkrmné dietě číslo  $r$ . Tak dostaneme realizace  $r$  nezávislých náhodných výběrů.

**Blokové porovnávání:** situace je charakterizována jedním náhodným výběrem  $(X_{11}, \dots, X_{1r}), \dots, (X_{n1}, \dots, X_{nr})$  z  $r$ -rozměrného rozložení.

Náhodně vylosujeme  $n$  vrhů stejně starých selat téhož plemene, z každého odebereme  $r$  sourozenců a náhodně jim přiřadíme první až  $r$ -tou výkrmnou dietu. Tak dostaneme realizaci jednoho  $r$ -rozměrného náhodného výběru, kde první složka odpovídá první dietě, druhá složka druhé dietě atd. až  $r$ -tá složka odpovídá  $r$ -té dietě.

## Úvod do testování hypotéz

**Motivace:** Častým úkolem statistika je na základě dat ověřit předpoklady o parametrech nebo typu rozložení, z něhož pochází náhodný výběr. Takovému předpokladu se říká nulová hypotéza. Nulová hypotéza vyjadřuje nějaký teoretický předpoklad, často skeptického rázu a uživatel ji musí stanovit předem, bez přihlídnutí k datovému souboru. Proti nulové hypotéze stavíme alternativní hypotézu, která říká, co platí, když neplatí nulová hypotéza. Alternativní hypotéza je formulována tak, aby mohla platit jenom jedna z těchto dvou hypotéz. Pravdivost alternativní hypotézy by znamenala objevení nějakých nových skutečností, nebo zásadnější změnu v dosavadních představách.

Např. výzkumník by chtěl na základě dat prověřit tezi (nový objev), že pasivní kouření škodí zdraví. Jako nulovou hypotézu tedy položí tvrzení, že pasivní kouření neškodí zdraví a proti nulové hypotéze postaví alternativní, že pasivní kouření škodí zdraví.

Testováním hypotéz se myslí rozhodovací postup, který je založen na daném náhodném výběru a s jehož pomocí rozhodneme o zamítnutí či nezamítnutí nulové hypotézy.

## Nulová a alternativní hypotéza

Nechť  $X_1, \dots, X_n$  je náhodný výběr z rozložení  $L(\vartheta)$ , kde parametr  $\vartheta \in \Xi$  neznáme. Nechť  $h(\vartheta)$  je parametrická funkce a  $c$  daná reálná konstanta.

a) **Oboustranná alternativa:** Tvrzení  $H_0: h(\vartheta) = c$  se nazývá **jednoduchá nulová hypotéza**. Proti nulové hypotéze postavíme **složenou oboustrannou alternativní hypotézu**  $H_1: h(\vartheta) \neq c$ .

b) **Levostranná alternativa:** Tvrzení  $H_0: h(\vartheta) \geq c$  se nazývá **složená pravostranná nulová hypotéza**. Proti jednoduché nebo složené pravostranné nulové hypotéze postavíme **složenou levostrannou alternativní hypotézu**  $H_1: h(\vartheta) < c$ .

c) **Pravostranná alternativa:** Tvrzení  $H_0: h(\vartheta) \leq c$  se nazývá **složená levostranná nulová hypotéza**. Proti jednoduché nebo složené levostranné nulové hypotéze postavíme **složenou pravostrannou alternativní hypotézu**  $H_1: h(\vartheta) > c$ .

**Testováním  $H_0$  proti  $H_1$**  rozumíme rozhodovací postup založený na náhodném výběru  $X_1, \dots, X_n$ , s jehož pomocí zamítneme či nezamítneme platnost nulové hypotézy.

## Chyba 1. a 2. druhu

Při testování  $H_0$  proti  $H_1$  se můžeme dopustit jedné ze dvou chyb: **chyba 1. druhu** spočívá v tom, že  $H_0$  zamítneme, ač ve skutečnosti platí a **chyba 2. druhu** spočívá v tom, že  $H_0$  nezamítneme, ač ve skutečnosti neplatí. Situaci přehledně znázorňuje tabulka:

skutečnost	rozhodnutí	
	$H_0$ nezamítáme	$H_0$ zamítáme
$H_0$ platí	správné rozhodnutí	chyba 1. druhu
$H_0$ neplatí	chyba 2. druhu	správné rozhodnutí

Pravděpodobnost chyby 1. druhu se značí  $\alpha$  a nazývá se **hladina významnosti testu** (většinou bývá  $\alpha = 0,05$ , méně často 0,1 či 0,01). Pravděpodobnost chyby 2. druhu se značí  $\beta$ . Číslo  $1-\beta$  se nazývá **síla testu** a vyjadřuje pravděpodobnost, že bude  $H_0$  zamítnuta za předpokladu, že neplatí. Obvykle se snažíme, aby síla testu byla aspoň 0,8. Obě hodnoty,  $\alpha$  i  $1-\beta$ , závisí na velikosti efektu, který se snažíme detekovat. Čím drobnější efekt, tím musí být větší rozsah náhodného výběru.

skutečnost	rozhodnutí	
	zdravý	nemocný
jsem zdravý	zdravý a neléčený	zdravý a léčený
jsem nemocný	nemocný a neléčený	nemocný a léčený

## Testování pomocí kritického oboru

Najdeme statistiku  $T_0 = T_0(X_1, \dots, X_n)$ , kterou nazveme **testovým kritériem**. Množina všech hodnot, jichž může testové kritérium nabýt, se rozpadá na **obor nezamítnutí nulové hypotézy** (značí se V) a **obor zamítnutí nulové hypotézy** (značí se W a nazývá se též **kritický obor**). Tyto dva obory jsou odděleny kritickými hodnotami (pro danou hladinu významnosti  $\alpha$  je lze najít ve statistických tabulkách).

Jestliže číselná realizace  $t_0$  testového kritéria  $T_0$  padne do kritického oboru W, pak nulovou hypotézu zamítáme na hladině významnosti  $\alpha$  a znamená to skutečné vyvrácení testované hypotézy. Jestliže  $t_0$  padne do oboru nezamítnutí V, pak jde o pouhé mlčení, které platnost nulové hypotézy jenom připouští.

Pravděpodobnosti chyb 1. a 2. druhu nyní zapíšeme takto:

$$P(T_0 \in W/H_0 \text{ platí}) = \alpha, P(T_0 \in V/H_1 \text{ platí}) = \beta.$$

Stanovení kritického oboru pro danou hladinu významnosti  $\alpha$ :

Označme  $t_{\min}$  (resp.  $t_{\max}$ ) nejmenší (resp. největší) hodnotu testového kritéria.

Kritický obor v případě oboustranné alternativy má tvar

$W = (t_{\min}, K_{\alpha/2}(T)) \cup (K_{1-\alpha/2}(T), t_{\max})$ , kde  $K_{\alpha/2}(T)$  a  $K_{1-\alpha/2}(T)$  jsou kvantily rozložení, jímž se řídí testové kritérium  $T_0$ , je-li nulová hypotéza pravdivá.

Kritický obor v případě jednostranné alternativy má tvar:

$$W = (t_{\min}, K_{\alpha}(T)).$$

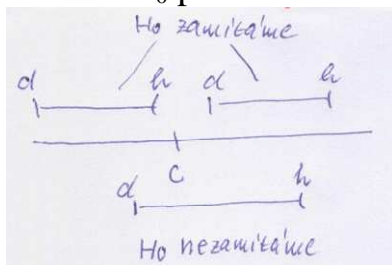
Kritický obor v případě jednostranné alternativy má tvar:

$$W = (K_{1-\alpha}(T), t_{\max}).$$

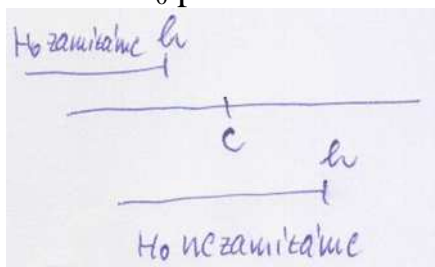
## Testování pomocí intervalu spolehlivosti

Sestrojíme  $100(1-\alpha)\%$  empirický interval spolehlivosti pro parametrickou funkci  $h(\vartheta)$ . Pokryje-li tento interval hodnotu  $c$ , pak  $H_0$  nezamítáme na hladině významnosti  $\alpha$ , v opačném případě  $H_0$  zamítáme na hladině významnosti  $\alpha$ .

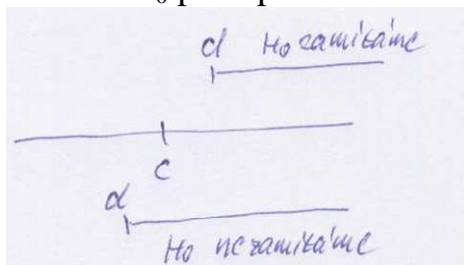
Pro test  $H_0$  proti oboustranné alternativě sestrojíme oboustranný interval spolehlivosti.



Pro test  $H_0$  proti levostranné alternativě sestrojíme pravostranný interval spolehlivosti.



Pro test  $H_0$  proti pravostranné alternativě sestrojíme levostranný interval spolehlivosti.





## Testování pomocí p-hodnoty

**p-hodnota** udává nejnižší možnou hladinu významnosti pro zamítnutí nulové hypotézy. Je to riziko, že bude zamítnuta  $H_0$  za předpokladu, že platí (riziko planého poplachu). Jestliže  $p\text{-hodnota} \leq \alpha$ , pak  $H_0$  zamítáme na hladině významnosti  $\alpha$ , je-li  $p\text{-hodnota} > \alpha$ , pak  $H_0$  nezamítáme na hladině významnosti  $\alpha$ .

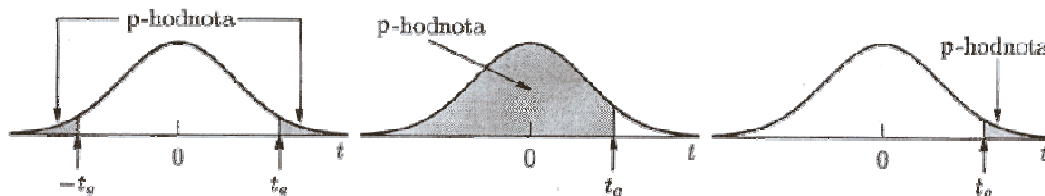
**Způsob výpočtu p-hodnoty:**

Pro oboustrannou alternativu  $p = 2 \min\{P(T_0 \leq t_0), P(T_0 \geq t_0)\}$ .

Pro levostrannou alternativu  $p = P(T_0 \leq t_0)$ .

Pro pravostrannou alternativu  $p = P(T_0 \geq t_0)$ .

Ilustrace významu p-hodnoty pro test nulové hypotézy proti oboustranné, levostranné a pravostranné alternativě:



(Zvonovitá křivka reprezentuje hustotu rozložení, kterým se řídí testové kritérium, je-li nulová hypotéza pravdivá.)

p-hodnota vyjadřuje pravděpodobnost, s jakou číselné realizace  $x_1, \dots, x_n$  náhodného výběru  $X_1, \dots, X_n$  podporují  $H_0$ , je-li pravdivá. Statistické programové systémy poskytují ve svých výstupech p-hodnotu. Její výpočet vyžaduje znalost distribuční funkce rozložení, kterým se řídí testové kritérium  $T_0$ , je-li  $H_0$  pravdivá.

## Doporučený postup při testování hypotéz

1. Stanovíme nulovou hypotézu a alternativní hypotézu. Přitom je vhodné zvolit jako alternativní hypotézu ten předpoklad, jehož přijetí znamená závažné opatření a mělo by k němu dojít jen s malým rizikem omylu.
2. Zvolíme hladinu významnosti  $\alpha$ . Zpravidla volíme  $\alpha = 0,05$ , méně často 0,1 nebo 0,01.
3. Najdeme vhodné testové kritérium a na základě zjištěných dat vypočítáme jeho realizaci.
4.
  - a) Testujeme-li pomocí kritického oboru, pak ho stanovíme. Jestliže realizace testového kritéria padla do kritického oboru, nulovou hypotézu zamítáme na hladině významnosti  $\alpha$  a přijímáme alternativní hypotézu. V opačném případě nulovou hypotézu nezamítáme na hladině významnosti  $\alpha$ .
  - b) Testujeme-li pomocí intervalu spolehlivosti, vypočteme empirický  $100(1-\alpha)\%$  interval spolehlivosti pro parametrickou funkci  $h(\vartheta)$ . Pokud číslo  $c$  padne do tohoto intervalu, nulovou hypotézu nezamítáme na hladině významnosti  $\alpha$ . V opačném případě nulovou hypotézu zamítáme na hladině významnosti  $\alpha$  a přijímáme alternativní hypotézu.
  - c) Testujeme-li pomocí p-hodnoty, vypočteme ji a porovnáme ji s hladinou významnosti  $\alpha$ . Jestliže  $p \leq \alpha$ , pak nulovou hypotézu zamítáme na hladině významnosti  $\alpha$  a přijímáme alternativní hypotézu. Je-li  $p > \alpha$ , pak nulovou hypotézu nezamítáme na hladině významnosti  $\alpha$ .
5. Na základě rozhodnutí, které jsme učinili o nulové hypotéze, provedeme nějaké konkrétní opatření, např. seřídíme obráběcí stroj.

(Při testování hypotéz musíme mít k dispozici odpovídající nástroje, nejlépe vhodný statistický software. Nemáme-li ho k dispozici, musíme znát příslušné vzorce. Dále potřebujeme statistické tabulky a kalkulačku.)

**Příklad:** 10 x nezávisle na sobě byla změřena jistá konstanta  $\mu$ . Výsledky měření byly: 2,1, 1,8, 2,1, 2,4, 1,9, 2,1, 2,1, 1,8, 2,3, 2,2. Tyto výsledky považujeme za číselné realizace náhodného výběru  $X_1, \dots, X_{10}$  z rozložení  $N(\mu, 0,04)$ . Nějaká teorie tvrdí, že  $\mu = 1,95$ .

### 1. Oboustranná alternativa

Proti nulové hypotéze  $H_0: \mu = 1,95$  postavíme oboustrannou alternativu

$H_1: \mu \neq 1,95$ . Na hladině významnosti 0,05 testujte  $H_0$  proti  $H_1$  všemi třemi popsánymi způsoby.

#### Řešení:

$$m = \frac{1}{10}(2 + \dots + 2,2) = 2,06, \sigma^2 = 0,04, n = 10, \alpha = 0,05, c = 1,95$$

a) **Test provedeme pomocí kritického oboru.**

Pro úlohy o střední hodnotě normálního rozložení při známém rozptylu používáme pivotovou statistiku  $U = \frac{M - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$ .

Testové kritérium tedy bude

$T_0 = \frac{M - c}{\frac{\sigma}{\sqrt{n}}}$  a bude mít rozložení  $N(0, 1)$ , pokud je nulová hypotéza pravdivá. Vypočítáme realizaci testového kritéria:

$$t_0 = \frac{2,06 - 1,95}{\frac{0,2}{\sqrt{10}}} = 1,74. \text{ Stanovíme kritický obor:}$$

$$W = (t_{\min}, K_{\alpha/2}(T)) \cup (K_{1-\alpha/2}(T), t_{\max}) = (-\infty, u_{\alpha/2}) \cup (u_{1-\alpha/2}, \infty) = (-\infty, -u_{1-\alpha/2}) \cup (u_{1-\alpha/2}, \infty) = (-\infty, -u_{0,975}) \cup (u_{0,975}, \infty) = (-\infty, -1,96) \cup (1,96, \infty).$$

Protože  $1,74 \notin W$ ,  $H_0$  nezamítáme na hladině významnosti 0,05.

b) Test provedeme pomocí intervalu spolehlivosti.

Meze  $100(1-\alpha)\%$  empirického intervalu spolehlivosti pro střední hodnotu  $\mu$  při známém rozptylu  $\sigma^2$  jsou:

$$(d, h) = \left( m - \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2}, m + \frac{\sigma}{\sqrt{n}} u_{1-\alpha/2} \right).$$

V našem případě dostáváme:

$$d = 2,06 - \frac{0,2}{\sqrt{10}} u_{0,975} = 2,06 - \frac{0,2}{\sqrt{10}} \cdot 1,96 = 1,936,$$

$$h = 2,06 + \frac{0,2}{\sqrt{10}} u_{0,975} = 2,06 + \frac{0,2}{\sqrt{10}} \cdot 1,96 = 2,184.$$

Protože  $1,95 \in (1,936; 2,184)$ ,  $H_0$  nezamítáme na hladině významnosti 0,05.

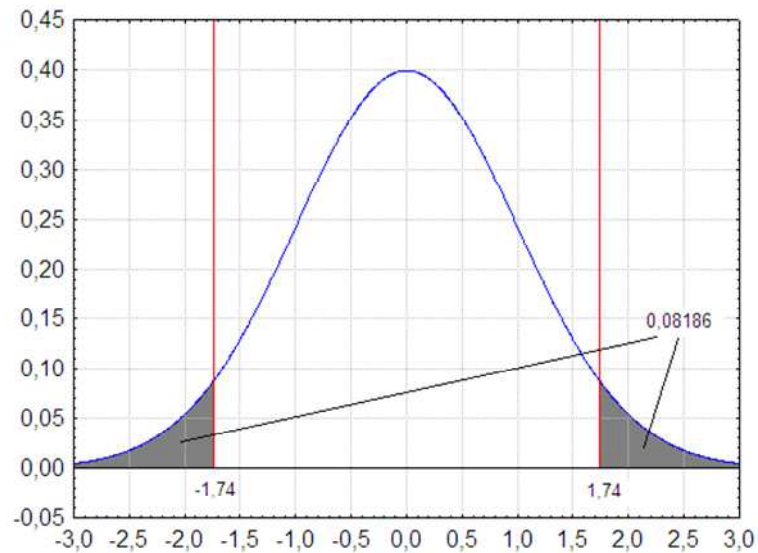
c) Test provedeme pomocí p-hodnoty.

Protože proti nulové hypotéze stavíme oboustrannou alternativu, použijeme vzorec

$$p = 2 \min\{P(T_0 \leq t_0), P(T_0 \geq t_0)\} = 2 \min\{P(T_0 \leq 1,74), P(T_0 \geq 1,74)\} = \\ = 2 \min\{\Phi(1,74), 1 - \Phi(1,74)\} = 2 \min\{0,95907, 1 - 0,95907\} = 0,08186.$$

Jelikož  $0,08186 > 0,05$ , nulovou hypotézu nezamítáme na hladině významnosti  $0,05$ .

Ilustrace významu p-hodnoty pro oboustranný test



## 2. Levostranná alternativa

Proti nulové hypotéze  $H_0: \mu = 1,95$  postavíme levostrannou alternativu

$H_1: \mu < 1,95$ . Na hladině významnosti 0,05 testujte  $H_0$  proti  $H_1$  všemi třemi popsánymi způsoby.

### Řešení:

a) Test provedeme pomocí kritického oboru.

Na rozdíl od oboustranné alternativy bude mít kritický obor tvar

$$W = \langle -\infty, u_\alpha \rangle = \langle -\infty, u_{0,05} \rangle = \langle -\infty, -1,645 \rangle.$$

Protože  $1,74 \notin W$ ,  $H_0$  nezamítáme na hladině významnosti 0,05.

b) Test provedeme pomocí intervalu spolehlivosti.

Meze  $100(1-\alpha)\%$  empirického pravostranného intervalu spolehlivosti pro střední hodnotu  $\mu$  při známém rozptylu  $\sigma^2$  jsou:

$$(-\infty, h) = \left(-\infty, m + \frac{\sigma}{\sqrt{n}} u_{1-\alpha}\right).$$

$$\text{V našem případě dostáváme: } h = 2,06 + \frac{0,2}{\sqrt{10}} u_{0,95} = 2,06 + \frac{0,2}{\sqrt{10}} \cdot 1,645 = 2,164.$$

Protože  $1,95 \in (-\infty; 2,164)$ ,  $H_0$  nezamítáme na hladině významnosti 0,05.

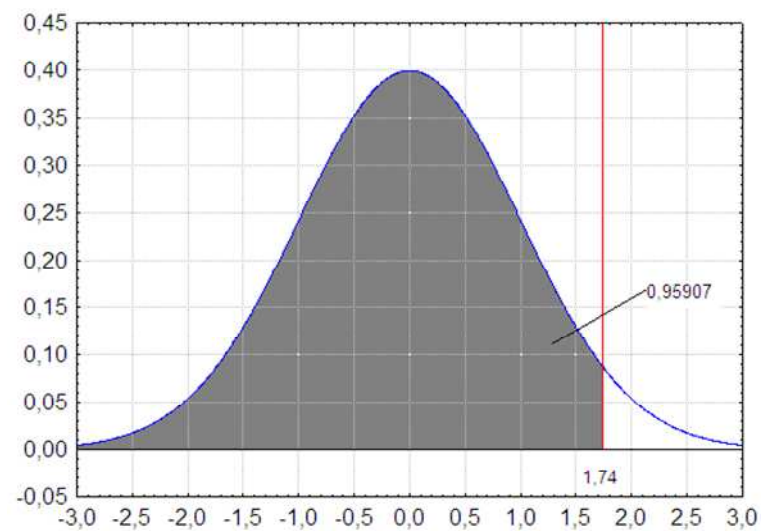
c) Test provedeme pomocí p-hodnoty.

Protože proti nulové hypotéze stavíme levostrannou alternativu, použijeme vzorec

$$p = P(T_0 \leq t_0) = \Phi(1,74) = 0,95907.$$

Jelikož  $0,95907 > 0,05$ , nulovou hypotézu nezamítáme na hladině významnosti 0,05.

Ilustrace významu p-hodnoty pro levostranný test



### 3. Pravostranná alternativa

Proti nulové hypotéze  $H_0: \mu = 1,95$  postavíme pravostrannou alternativu

$H_1: \mu > 1,95$ . Na hladině významnosti 0,05 testujte  $H_0$  proti  $H_1$  všemi třemi popsánymi způsoby.

#### Řešení:

a) Test provedeme pomocí kritického oboru.

Na rozdíl od oboustranné alternativy bude mít kritický obor tvar

$$W = \langle u_{1-\alpha}, \infty \rangle = \langle u_{0,95}, \infty \rangle = \langle 1,645, \infty \rangle.$$

Protože  $1,74 \in W$ ,  $H_0$  zamítáme na hladině významnosti 0,05 ve prospěch pravostranné alternativy.

b) Test provedeme pomocí intervalu spolehlivosti.

Meze 100(1- $\alpha$ )% empirického jednostranného intervalu spolehlivosti pro střední hodnotu  $\mu$  při známém rozptylu  $\sigma^2$  jsou:

$$(d, \infty) = \left( m - \frac{\sigma}{\sqrt{n}} u_{1-\alpha}, \infty \right).$$

$$\text{V našem případě dostáváme: } d = 2,06 - \frac{0,2}{\sqrt{10}} u_{0,95} = 2,06 - \frac{0,2}{\sqrt{10}} \cdot 1,645 = 1,956.$$

Protože  $1,95 \notin (1,956, \infty)$ ,  $H_0$  zamítáme na hladině významnosti 0,05 ve prospěch pravostranné alternativy.



c) Test provedeme pomocí p-hodnoty.

Protože proti nulové hypotéze stavíme pravostrannou alternativu, použijeme vzorec

$$p = P(T_0 \geq t_0) = 1 - \Phi(1,74) = 1 - 0,95907 = 0,04093.$$

Jelikož  $0,04093 \leq 0,05$ , nulovou hypotézu zamítáme na hladině významnosti 0,05 ve prospěch pravostranné alternativy.

Ilustrace významu p-hodnoty pro pravostranný test

