

Statistické metody

Základní metody:

- **odhad parametrů**
- **testování hypotéz**

Pokročilejší metody:

- shluková (klastrová) analýza
- faktorová analýza
- analýza hlavních komponent (PCA)
- ...

Statistické metody

Základní soubor (populace) je příliš velký a nemůžeme ho celý „proměřit“.

Proto dělám reprezentativní výběr, ten změřím, tedy náhodným procesem získávám konkrétní hodnoty náhodných veličin.

Spočítám výběrové charakteristiky souboru.

Tyto výběrové charakteristiky chci vztáhnout na celý základní soubor. Musím nějak **kvantifikovat jistotu či nejistotu**, že **moje odhady se potkávají s neznámou skutečností**.

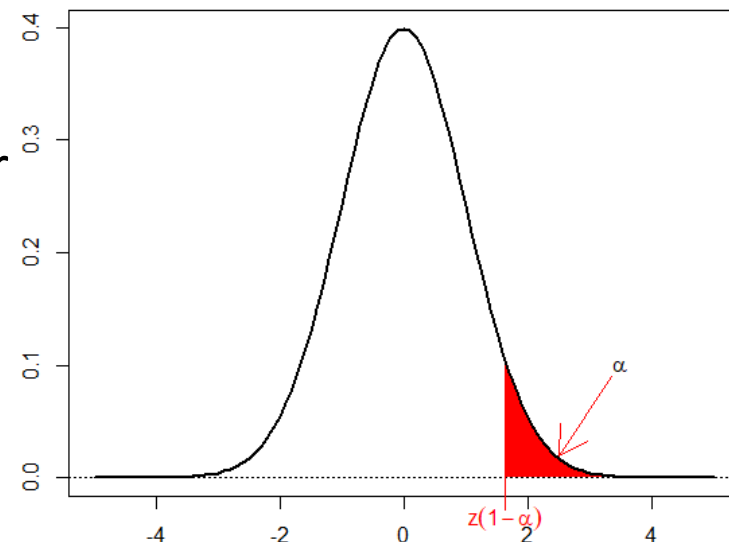
Připomínka značení:

μ vs. \bar{X} ... skutečný neznámý parametr vs. náš odhad

σ^2 vs. S^2

$z(1 - \alpha)$... $(1 - \alpha)\%$ kvantil rozdělení prstí, pro který platí

$t_{df}(1 - \alpha)$

$$P(X > z(1 - \alpha)) = \alpha$$


Bodový odhad parametru [point estimate of the parameter]

Základní předpoklad dalšího odvozování:

mám výběr n hodnot $(X_1, X_2, X_3, \dots, X_n)$, které jsou **iid.**, tedy vzájemně nezávislé a všechny pocházejí ze stejného rozdělení prstí.

K odhadu typické hodnoty (charakteristika polohy) nejčastěji používáme

výběrový průměr $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ [sample mean]

Protože výběrový průměr je náhodná veličina, má smysl se ptát:

- jaká je jeho střední hodnota [expected value of the estimate]
- jaký je jeho rozptyl [variance of the estimate]
- jaká je jeho směrodatná odchylka [standard error of the estimate]

Populační charakteristiky průměru (odvození dále):

$$E\bar{X} = \mu \qquad \text{var } \bar{X} = \frac{\sigma^2}{n} \qquad \text{sd } \bar{X} = \frac{\sigma}{\sqrt{n}}$$

Odvození pro výběrový průměr:

(a) Střední hodnota výběrového průměru:

$$E\bar{X} = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n EX_i = \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} \cdot n \cdot \mu = \mu$$

vlastnost střední hodnoty: $E(X + Y) = EX + EY$

- tento odhad je nestranný, protože $E\bar{X} = \mu$

Odvození pro výběrový průměr:

(b) Rozptyl výběrového průměru:

$$\text{var}\bar{X} = \text{var}\left(\frac{1}{n}\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{var}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \left(\sum_{i=1}^n \text{var}X_i\right) = \frac{1}{n^2} \cdot n \cdot \sigma^2 = \frac{\sigma^2}{n}$$

$$\text{var}(\beta \cdot X) = \beta^2 \cdot \text{var}X$$

- (1) všechna X_i jsou *iid.*, proto $\text{cov}(X_i, X_j) = 0$ pro $\forall i, j$
 (2) $\text{var}(X + Y) = \text{var}X + \text{var}Y + 2\text{cov}(X, Y)$

$$\text{var}\bar{X} = \frac{\sigma^2}{n}$$

- $n = 1 \rightarrow \text{var}\bar{X}_1 = \sigma^2$
- větší $n \rightarrow$ menší rozptyl \bar{X}
- problém: σ^2 většinou neznáme

Odvození pro výběrový průměr:

(c) Směrodatná odchylka výběrového průměru:

$$S. E. (\bar{X}) = \sqrt{\text{var } \bar{X}} = \frac{\sigma}{\sqrt{n}}$$

- říkáme jí **střední chyba průměru** [standard error of mean, SEM]
- často se uvádí ve výsledcích článků
- charakterizuje „přesnost“ odhadu (pozor: přesnost odhadu ve smyslu střední kvadratické chyby (viz dále) zahrnuje i vychýlení odhadu)
- platí: čím větší výběr (n), tím přesnější odhad
- *SEM* závisí na parametru σ , který většinou neznáme a nahrazujeme ho vhodným odhadem, např. výběrovým rozptylem (za chvíli). Slovní označení „střední chyba“ se používá i tehdy, když místo σ použijí odhad.

Bodový odhad variance – výběrový rozptyl

K odhadu variability hodnot v populaci nejčastěji používáme

výběrový rozptyl $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ [sample variance]

- střední hodnota výběrového rozptylu:

$$ES^2 = \sigma^2$$

- rozptyl výběrového rozptylu běžně nepotřebujeme, proto neuvádím

- Jiný bodový odhad variability hodnot v populaci:

výběrový momentový rozptyl $S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ většinou

nepoužíváme, protože o $\frac{1}{n}$ podhodnocuje skutečný parametr σ^2 (dále)

Bodový odhad populační SD – výběrová směrodatná odchylka

$$S = \sqrt{S^2}$$

- tento odhad je vychýlený, skutečnou směr. odchylku v průměru podhodnocuje, protože platí $ES < \sigma$.

Dodatek: jiný tvar výběrového rozptylu:

užitečný tvar pro „ruční“ výpočet, používá se v algoritmech (je rychlejší):

$$\begin{aligned} S^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i^2 - 2\bar{X}X_i + \bar{X}^2) = \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2 \right) = \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - 2\bar{X} \cdot n \frac{\sum X_i}{n} + n\bar{X}^2 \right) = \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - 2\bar{X} \cdot n\bar{X} + n\bar{X}^2 \right) = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right) \end{aligned}$$

Dodatek: odvození výpočtu střední hodnoty výběrového rozptylu

$$\begin{aligned}
 ES^2 &= E\left(\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right) = E\left(\frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2\right)\right) = \\
 &= \frac{1}{n-1} \left(\sum_{i=1}^n E X_i^2 - n \cdot E \bar{X}^2\right) = \\
 &= \frac{1}{n-1} \left(n(\sigma^2 + \mu^2) - n \cdot \left(\frac{\sigma^2}{n} + \mu^2\right)\right) = \frac{1}{n-1} (n\sigma^2 + n\mu^2 - \sigma^2 - n\mu^2) = \\
 &= \frac{1}{n-1} \cdot (n-1)\sigma^2 = \sigma^2
 \end{aligned}$$

$$E(\beta \cdot X) = \beta \cdot EX$$

$$\rightarrow \text{var } X_i = E(X_i - EX_i)^2 = E(X_i^2 - 2X_i EX_i + (EX_i)^2) = E X_i^2 - 2 \cdot EX_i \cdot EX_i + (EX_i)^2 = E X_i^2 - (EX_i)^2$$

$$\text{odtud: } E X_i^2 = \text{var } X_i + (EX_i)^2 = \sigma^2 + \mu^2$$

$$\rightarrow \text{podobně: } \text{var } \bar{X} = E(\bar{X} - E\bar{X})^2 = \dots = E(\bar{X})^2 - (E\bar{X})^2$$

$$\text{odtud: } E\bar{X}^2 = \text{var } \bar{X} + (E\bar{X})^2 = \frac{\sigma^2}{n} + \mu^2$$

Vlastnosti bodového odhadu

Nestranný odhad (nevychýlený, nezkreslený) [unbiased estimation]

- když střední hodnota odhadu = teoretickému parametru
- právě jsme měli: $E\bar{X} = \mu$ a $ES^2 = \sigma^2$
- nestranný odhad systematicky nenadhodnocuje ani nepodhodnocuje odhadovaný parametr
- příklad vychýleného odhadu – výběrový momentový rozptyl:

$$ES_n^2 = E\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right) = \dots = \frac{n-1}{n} \sigma^2 = \sigma^2 - \frac{1}{n} \sigma^2$$

vychýlení značíme $B(\sigma^2, S_n^2) = ES_n^2 - \sigma^2 = \frac{n-1}{n} \sigma^2 - \frac{n}{n} \sigma^2 = -\frac{1}{n} \sigma^2$

S_n^2 podhodnocuje skutečný parametr σ^2 .

Vlastnosti bodového odhadu

Asymptoticky nestranný odhad

- když odhad je sice vychýlený, ale se zvyšujícím se rozsahem výběru n se vychýlení zmenšuje až k nule
- to je případ výběrového momentového rozptylu:

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

$$ES_n^2 = \sigma^2 - \frac{1}{n} \sigma^2$$

$$\text{Vychýlení: } -\frac{1}{n} \sigma^2$$

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \sigma^2 \rightarrow -\frac{1}{\infty} \sigma^2 = 0$$

Vlastnosti bodového odhadu

Konzistentní odhad [consistent estimation]

- pokud se s rostoucím rozsahem výběru n odhad zpřesňuje

- $E(\text{odhadu}) = \text{parametr}$

- a zároveň $\text{var}(\text{odhadu}) \xrightarrow{n \rightarrow \infty} 0$

- platí např. pro výběrový průměr:

$$E\bar{X} = \mu$$

$$\text{var}\bar{X} = \frac{\sigma^2}{n} \xrightarrow{n \rightarrow \infty} \frac{\sigma^2}{\infty} = 0$$

Vlastnosti bodového odhadu

Vydatný, eficientní, nejlepší nestranný odhad [efficient estimation]

- má nejmenší rozptyl mezi všemi nestrannými odhady téhož parametru

Přesnost, kvalita odhadu [quality of the estimation]

- měříme pomocí střední kvadratické chyby odhadu
- výběrová chyba odhadu: $odhad - parametr$
- zkratka $MSE(odhadu)$ [mean squared error] (! \neq SEM, stand. error of mean)
- Kromě variability zahrnuje i vychýlení odhadu. Pro nestranné odhady (vychýlení = 0) je to totéž jako $var(odhadu)$ a potažmo $S.E.(odhadu)$
- $MSE(odhadu) = E(odhad - parametr)^2 = var(odhadu) + B^2(odhadu) = E(odhad - E(odhadu))^2 + (E(odhadu) - parametr)^2$
- příklad: $MSE(S_n^2) = E(S_n^2 - \sigma^2)^2 = \dots$

Ze statistického slovníku:

Robustní = odolný

přibližně řečeno je to schopnost spočítat „spolehlivý“ výsledek, přestože jsou narušeny předpoklady testu, odhadu apod.

Konečnostní násobitel

Většinou zahrnuje náš výběr méně než 5 % jedinců z celé populace, proto můžeme takovou populaci považovat za nekonečnou.

Pokud ovšem vybíráme z menší konečné populace a rozsah výběru je větší než 5 % všech jedinců, potom výběrový průměr \bar{X} zůstává nestranným odhadem populačního průměru, ale odhad rozptylu \bar{X} bude nadhodnocený. Aby byly odhadované vlastnosti \bar{X} správné, je třeba rozptyl vynásobit konečnostním násobitelem $\frac{N-n}{N-1}$.

Tedy:

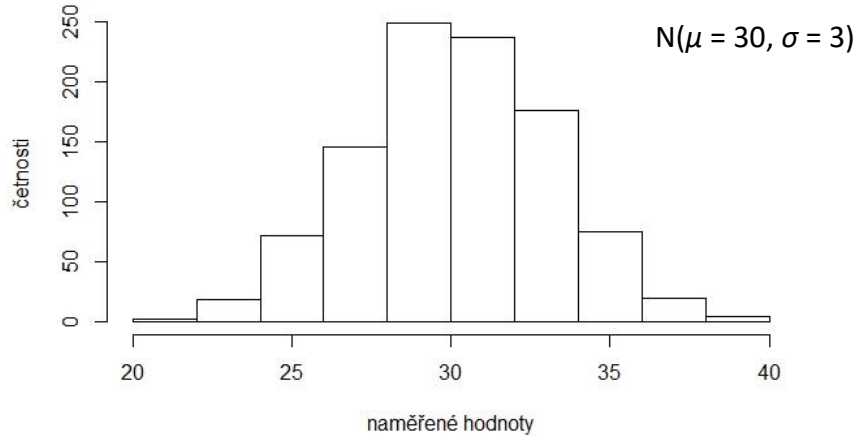
$$E\bar{X} = \mu \quad \dots \text{to je stejné}$$

$$\text{var } \bar{X} = \frac{N-n}{N-1} \cdot \frac{\sigma^2}{n} \quad \Rightarrow \quad S.E.(\bar{X}) = \frac{\sigma}{\sqrt{n}} \cdot \sqrt{\frac{N-n}{N-1}}$$

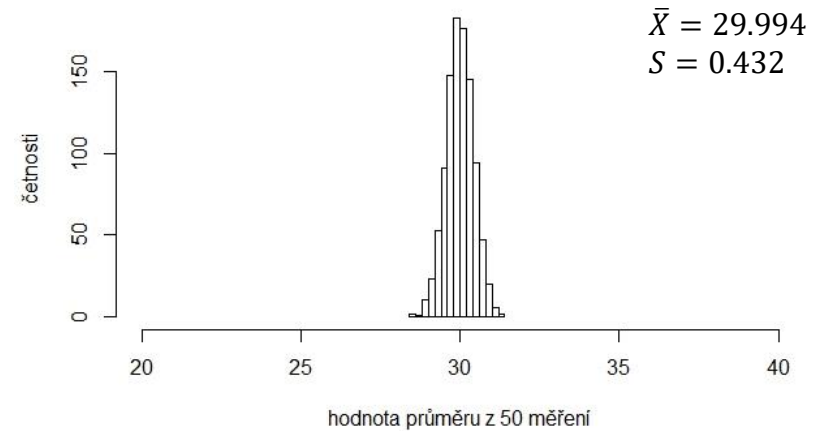
(Citace: Zvára, Karel: Biostatistika. Karolinum, Praha 2008.)

Centrální limitní věta graficky

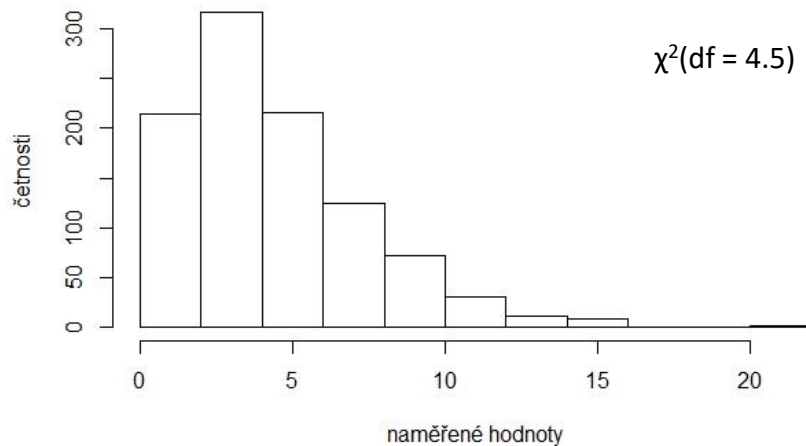
Základní data mají normální rozdělení



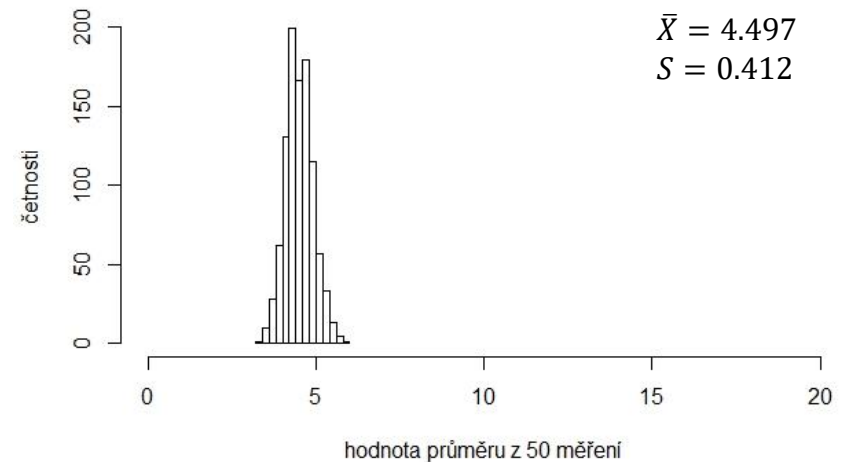
Průměry mají také normální rozdělení



Základní data NEMají normální rozdělení



Průměry PŘESTO mají normální rozdělení



Centrální limitní věta (CLV) [central limit theorem]

Když hodnoty ve výběru mají normální rozdělení $N(\mu, \sigma^2)$
potom také jejich průměr \bar{X} má normální rozdělení $N\left(\mu, \frac{\sigma^2}{n}\right)$.
Toho využíváme pro výpočet intervalu spolehlivosti nebo v testech.

Ale co když hodnoty ve výběru **nemají** normální rozdělení?

**Mám-li „dostatečně velký“ výběr n , potom se rozdělení průměru \bar{X}
blíží normálnímu s parametry odvozenými z výběrových dat $N\left(\mu, \frac{\sigma^2}{n}\right)$.**

Toto tvrzení je matematicky zpracováno v centrální limitní větě.

„Dostatečně velké“ n je v praxi alespoň 30 a více. Mám-li hodnot ve výběru méně, musím věnovat větší pozornost předpokladům parametrických testů.

Neplatí vždycky, ale lze aplikovat na průměr, relativní četnost či součet pořadí, také na testy o střední hodnotě nějakého rozdělení.

Použití CLV na aproximaci binomického rozdělení

$Y \sim \mathbf{Bi}(n, p)$, kde $Y = \sum_{i=1}^n X_i$ a $X_i \sim \mathbf{Alt}(p)$

víme, že $\mathbf{E}X_i = p$ a $\mathbf{var}X_i = p(1 - p)$

tedy $\mathbf{E}Y = n \cdot p$ a $\mathbf{var}Y = n \cdot p \cdot (1 - p)$

Podle CLV má náh. vel. $\mathbf{Z} = \frac{Y - np}{\sqrt{np(1-p)}} \sim \mathbf{N}(0, 1)$ pro velká n .

Proto $Y \sim \mathbf{Bi}(n, p)$ může být pro velká n aproximována $\sim \mathbf{N}(np, np(1 - p))$.

Zkušenosti starších říkají, že aproximace je dobře použitelná pro

$np(1 - p) > 9$ nebo

p	→	n
0.5		≥ 30
0.4 a 0.6		≥ 50
0.3 a 0.7		≥ 80
0.2 a 0.8		≥ 200
0.1 a 0.9		≥ 600

Intervalový odhad parametru [confidence interval of the parameter]

také konfidenční interval či interval spolehlivosti.

Konstrukci intervalu provedeme na příkladu výběrového průměru, teorie však platí pro odhady všech parametrů.

- Výběrový průměr \bar{X} je náhodná veličina, má tedy i své rozdělení pravděpodobností. Tvar rozdělení je dán rozdělením hodnot X_i a rozsahem výběru n .
- Víme, že $E\bar{X} = \mu$ a $var \bar{X} = \frac{\sigma^2}{n}$ (skutečné, ale neznámé parametry).
- Pokud **výběr pochází z normálního rozdělení** $N(\mu, \sigma^2)$, potom také náhodná veličina \bar{X} má normální rozdělení s parametry $N\left(\mu, \frac{\sigma^2}{n}\right)$.
- Když **výběr nepochází** z normálního rozdělení (histogram je šikmý nebo hrbatý), potom záleží na velikosti výběru. Při rozumně velkém výběru n funguje **centrální limitní věta** (dále) a podle té má $\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$ i když původní data nejsou z normálního rozdělení.

Intervalový odhad parametru

Teoreticky: hodnoty, kterých může nabývat průměr \bar{X} jsou popsány normálním rozdělením $N\left(\mu, \frac{\sigma^2}{n}\right)$:



Chceme sestavit interval takový, aby pokrýval „rozumné“ hodnoty \bar{X} a abychom znali pravděpodobnost chybného tvrzení o tomto intervalu.

Zvolíme velikost možné chyby $\alpha = 0,05$, tj. 5 % (například).

Pomůžeme si normovaným tvarem $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$ se známými kvantily:

$$P\left(-z(1 - \alpha/2) < \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} < z(1 - \alpha/2)\right) = 0,95$$

$N(0, 1)$ je souměrné, proto $z(1 - \alpha/2) = -z(\alpha/2)$.

$$\approx P\left(\bar{X} - z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}}\right) = 0,95$$

Intervalový odhad parametru

$$\rightarrow P\left(\bar{X} - z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha = 0,95$$

Tedy jsem zpět $v \sim N\left(\mu, \frac{\sigma^2}{n}\right)$.

$$\text{Jiný tvar: } P\left(\mu \in \left\{\bar{X} - z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}} ; \bar{X} + z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}}\right\}\right) = 1 - \alpha$$

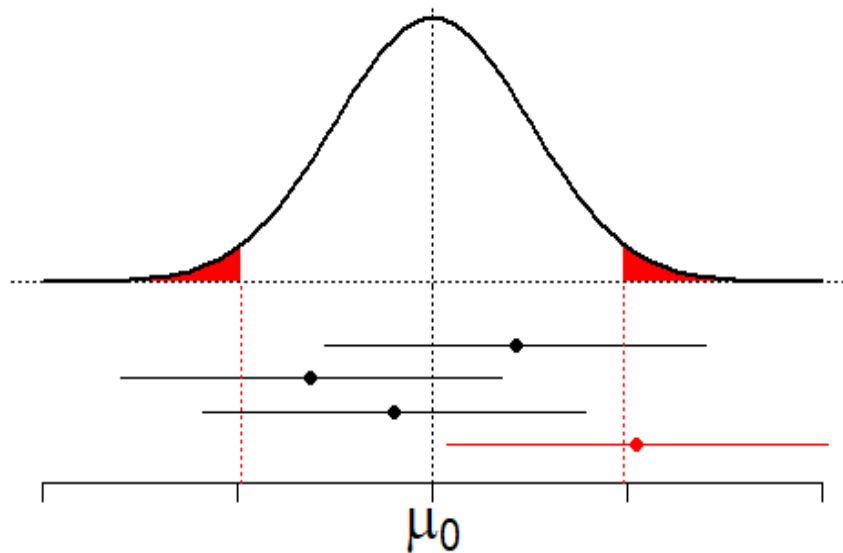
konfidenční interval odhadu parametru μ na hladině $\alpha = 0.05$.

Další způsob zápisu: $\bar{X} \pm z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}}$

Výsledek 15.3 ± 3.65 čteme např. takto: *střední hodnotu odhadujeme hodnotou 15.3, přičemž skutečná hodnota střední hodnoty leží na 95 % v rozmezí 15.3 – 3.65 a 15.3 + 3.65.* Je třeba uvádět také pravděp. nebo α .

Intervalový odhad parametru – graficky

$$P\left(\mu \in \left\{\bar{X} - z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}} ; \bar{X} + z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}}\right\}\right) = 1 - \alpha = 0,95$$



Červený interval je to „chybné tvrzení o intervalu spolehlivosti“. Červený interval nezahrnuje (nepokrývá) skutečnou hodnotu μ_0 .
Pravděpodobnost této chyby je α (5 %).

R: qnorm (pravděpodobnost, mean, sd) ... spočte takovou hodnotu na x-ové ose, pro kterou je $P(X \leq x) =$ zadaná pravděpodobnost.

pnorm (x, mean, sd) ... spočte $P(X \leq x)$ pro zadané x.

dnorm (x, mean, sd) ... spočte hustotu normálního rozdělení pro zadané x.

rnorm (n, mean, sd) ... vygeneruje n náhodných hodnot ze zadaného $N(\text{mean}, \text{sd})$.

Testování hypotéz [hypotheses testing]

Příklady:

- Z histogramu vidím, že data mají zhruba normální rozdělení. Ale tvrzení, že výběr pochází z normálního rozdělení, musím podepřít testem.
- Mám data o hmotnosti samců a samic nějakého druhu a z grafické prezentace je vidět, že samci jsou těžší. Statistický test řekne, zda je rozdíl mezi pohlavími „systematický“ nebo zda bylo věcí náhody, že někteří samci byli těžší a posunuli průměr napravo.

Základní poučka metodologie vědy: shoda dat s hypotézou ještě neznamená, že hypotéza je pravdivá; na druhou stranu data odporující hypotéze ukazují, že hypotéza pravdivá není.

**Proto hypotézu nelze na základě dat dokázat,
ale hypotézu lze na základě dat vyvrátit.**

Ad příklad 2) chci vyvrátit tvrzení, že samci i samice mají stejnou hmotnost.

Formulujeme **nulovou hypotézu H_0** [null hypothesis] a její negaci, tzv. **alternativní hypotézu H_1** , příp. **H_A** [alternative hypothesis].

Příklad.

H_0 : dva datové soubory mají stejnou střední hodnotu, $\mu_1 = \mu_2$;

H_1 : střední hodnoty se liší, $\mu_1 \neq \mu_2$.

H_0 : výběr pochází z normálního rozdělení;

H_1 : výběr nepochází z normálního rozdělení

Máme 2 možná rozhodnutí: H_0 zamítáme nebo H_0 nezamítáme.

Následují 4 možné situace:

	SKUTEČNOST	
NAŠE ROZHODNUTÍ	H_0 platí	H_0 neplatí (platí H_1)
H_0 zamítáme	Chyba 1. druhu: α Prst. chyby $\leq \alpha$	SPRÁVNÉ ROZHODNUTÍ $P = 1 - \beta$ síla testu
H_0 nezamítáme	SPRÁVNÉ ROZHODNUTÍ ($P \geq 1 - \alpha$)	Chyba 2. druhu: β β většinou neznáme

Testování hypotéz

Nulová hypotéza souvisí s nějakým předem daným uspořádáním dat. Toto uspořádání je popsáno nějakým teoretickým rozdělením prstí nějaké náhodné veličiny.

Naše výběrová data tedy porovnáváme s určeným teoretickým rozdělením pomocí odhadu určené náhodné veličiny.

Nulovou hypotézu zamítáme tehdy, když naše uspořádání výběrového souboru je za předpokladu platnosti H_0 velmi nepravděpodobné.

Příklad: **Test hypotézy o střední hodnotě normálního rozdělení.**

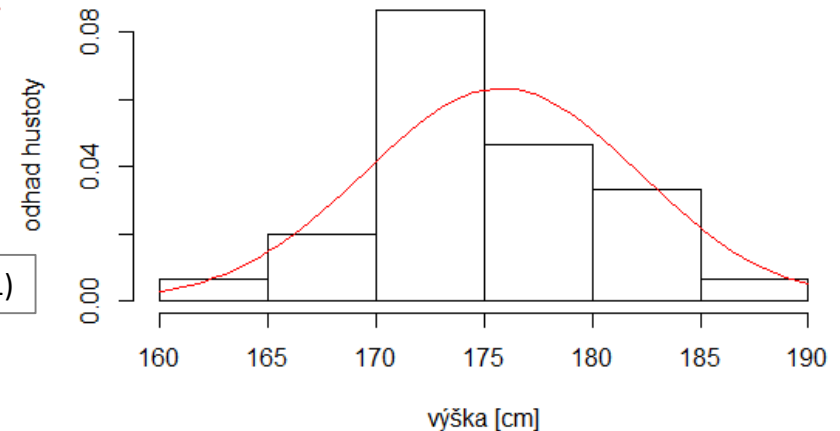
Data „Stulong“: výška mužů

Populační výška = 175.8 cm

Výběrová výška = 178.1 cm Pro výběr set.seed(21)

$H_0: \mu_{\text{vyber}} = 175.8$; $H_1: \mu_{\text{vyber}} \neq 175.8$

Výška mužů ve výběru (rel. četnosti)



Test hypotézy o střední hodnotě normálního rozdělení

Data „Stulong“: „populační“ výška mužů

- Mám výběr X_1, X_2, \dots, X_n *Výška mužů v našem výběru*
- **Je můj výběr reprezentativní?** *Výpočet v R-skriptu*
- Předpokládám, že $X_i \sim N(\mu_X, \sigma_X^2)$ a jsou iid. $N(175.8 \text{ cm}, \sigma = 6.3 \text{ cm})$
- Testuji, zda $\mu_X = \mu_0 \dots \mu_0$ nějaké číslo $\mu_0 = 175.8 \text{ cm}$
- Hypotéza $H_0: \mu_X = \mu_0, H_1: \mu_X \neq \mu_0$ $H_0: \mu_X = 175.8; H_1: \mu_X \neq 175.8$
- μ_X odhadnu pomocí \bar{X} , protože vím, $E\bar{X} = \mu$ $\bar{X} = 178.1 \text{ cm}$
- Rozhodovací pravidlo: $|\bar{X} - \mu_0| \dots$ bude-li velký rozdíl, H_0 zamítnu
- Jak velký musí být rozdíl $|\bar{X} - \mu_0|$, abych H_0 zamítla?
- Podle toho, jakou dovolím pravděpodobnost α chyby 1. druhu $\alpha = 0.05$
- Dopočítám kvantily pro $P \sim 2.5 \%$ a $P \sim 97.5 \%$ $x_{\alpha/2} = 173.6$ a $x_{1-\alpha/2} = 178.0$
- Je $\bar{X} \leq x_{\alpha/2}$ nebo $\bar{X} \geq x_{1-\alpha/2}$? Potom zamítám H_0 .
 $\bar{X} = 178.1 \text{ cm} \geq 178.0 \Rightarrow \text{zamítám } H_0$

Test hypotézy o střední hodnotě normálního rozdělení

Obecně bývá výpočet převeden z $N(\mu, \frac{\sigma}{\sqrt{n}})$ na $N(0,1)$:

$$Z = \frac{\bar{X} - \mu_{\bar{X}}}{\sigma_{\bar{X}}} = \frac{\bar{X} - \mu_X}{\frac{\sigma_X}{\sqrt{n}}} \sim N(0,1)$$

Zpracujeme předpoklad $H_0: \mu_X = \mu_0$

$$\rightarrow Z = \frac{\bar{X} - \mu_0}{\frac{\sigma_X}{\sqrt{n}}} = \frac{\bar{X} - \mu_0}{\sigma_X} \sqrt{n} \sim N(0,1)$$

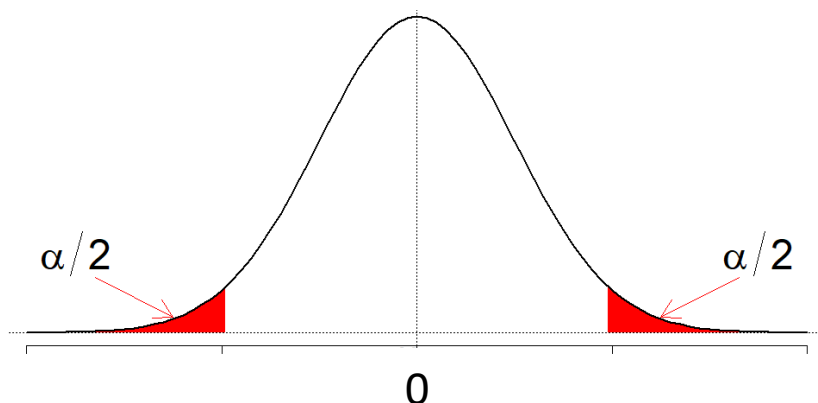
Test hypotézy o střední hodnotě normálního rozdělení – neznámé σ

Odvodili jsme testovou statistiku Z , která má – za platnosti H_0 – rozdělení $N(0, 1)$:

$$Z = \frac{\bar{X} - \mu_0}{\sigma_x} \sqrt{n} \sim N(0, 1) \quad \dots \mu_0 = \text{známé číslo}$$

V tuto chvíli otazník jen u σ_x $\left\{ \begin{array}{l} \text{a) známe} \\ \text{b) neznáme} \end{array} \right.$

a) σ_x známe: rozhod. pravidlo bude $|Z| \geq z(1 - \alpha/2)$, protože $H_1: \mu_X \neq \mu_0$



oboustranná alternativa
[two-tailed test]

b) σ_x neznáme: nahradíme ho odhadem $\sqrt{S_X^2} = S_X$

test. statistika $t = \frac{\bar{X} - \mu_0}{S_X} \sqrt{n} \sim t_{n-1}$ a rozhod. pravidlo $|t| \geq t_{n-1} \left(1 - \frac{\alpha}{2}\right)$.

t-test hypotézy o střední hodnotě normálního rozdělení v číslech:

Příklad: Data „Stulong“: výška mužů v našem výběru

$$H_0: \mu_{\text{vyber}} = 175.8 \text{ cm}; H_1: \mu_{\text{vyber}} \neq 175.8 \text{ cm}$$

$$\bar{X} = 178.1 \text{ cm}, \sigma_x \text{ neznáme} \rightarrow \text{odhad } S = 7.1 \text{ cm}$$

$$\text{Testová statistika: } t = \frac{178.1 - 175.8}{7.1} \sqrt{30} = 1.77$$

$$\text{Kvantil } t_{(29)}(1 - 0,025) = 2.05$$

Rozhodnutí: $|1.77| < 2.05$, proto nezamítám H_0 , že skutečná $\mu_{\text{vyber}} = 175.8 \text{ cm}$.

P-hodnota provedeného testu $p = 0.087$, tj. 8.7 %

Hladina testu α

Jiný název pro zvolenou chybu 1. druhu α .

Dosažená hladina významnosti testu

Také **p-hodnota** [p-value]

Je to pravděpodobnost, které odpovídá testová statistika coby kvantil.

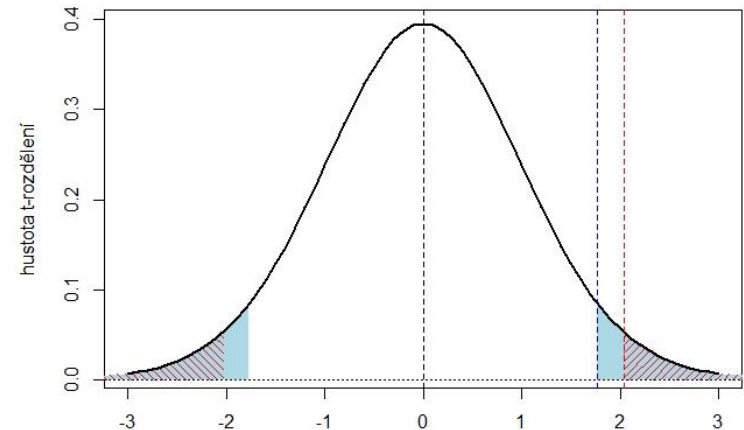
Dnes je toto číslo velmi cennou informací v publikacích, proto je častou součástí výsledků.

Na grafu: 97.5% kvantil t -rozdělení ($df=29$) = 2.04

testová statistika $t = 1.77$... modrá plocha $p = 2 * 0.044 = 0.088$

Co nastává: zvolili jsme $\alpha = 0.05$ (5 %) a ...

- p-hodnota vyjde 0.0023, tj. 0.23 %. Výsledek je tedy hluboko za kritickou hranicí, výsledek (rozdíl) je evidentně průkazný. Hurá!
- $p = 0.049$, tedy zamítám H_0 , ale jen velmi těsně.
- $p = 0.052$, tedy nezamítám H_0 , ale také velmi těsně.
- $p = 0.43$, tedy H_0 nezamítám a je zřejmé, že se výsledek hranici 5 % ani zdaleka neblíží.



Formulace nulové hypotézy

- a) Vidím, že samci a samičky mají skoro stejnou charakteristiku a chci je spojit do jedné skupiny. Potřebuji testem ukázat, že v datech není rozpor se „sjednocením“.
- hledaný výsledek: „*nezamítám H_0* “, „*rozdíl mezi samci a samičkami je neprůkazný*“, apod. Tvrzení podporuje velká p-hodnota, např. 0.3 a větší.
- b) Chci ukázat, že dvě skupiny se v nějaké charakteristice liší. Potom H_0 formuluji tak, abych ji na základě svých dat mohla zamítnout.
- Hledaný výsledek: „*zamítám H_0 o tom, že mezi charakteristikami první a druhé skupiny není rozdíl*“. Tvrzení musí mít p-hodnotu $\leq \alpha$.
- Nezamítnutí H_0 s p-hodnotou kolem 0.1 ~ 10 % znamená spíše nedostatek důkazů pro zamítnutí, než potvrzení platnosti H_0 .
 - **Vědecký důkaz = zamítnutí hypotézy (H_0)**
 - Nezamítnutí hypotézy nic nedokazuje, jen říká, že data nejsou v rozporu
 - Odpověď při neúspěchu: „*Na základě dat nemůžeme zamítnout H_0 .*“
 - Nelze napsat: „*dokázali jsme nulovou hypotézu...*“ **CHYBA!!**

Formulace nulové hypotézy

ad (b) Chci ukázat, že v datech je rozdíl:

Dobrá hypotéza je vyvratitelná.

Příklad: *V parku jsou lišky. // V parku nejsou žádné lišky.*

Nepřítomnost důkazů není důkaz nepřítomnosti.

[Absence of evidence is not evidence of absence.]

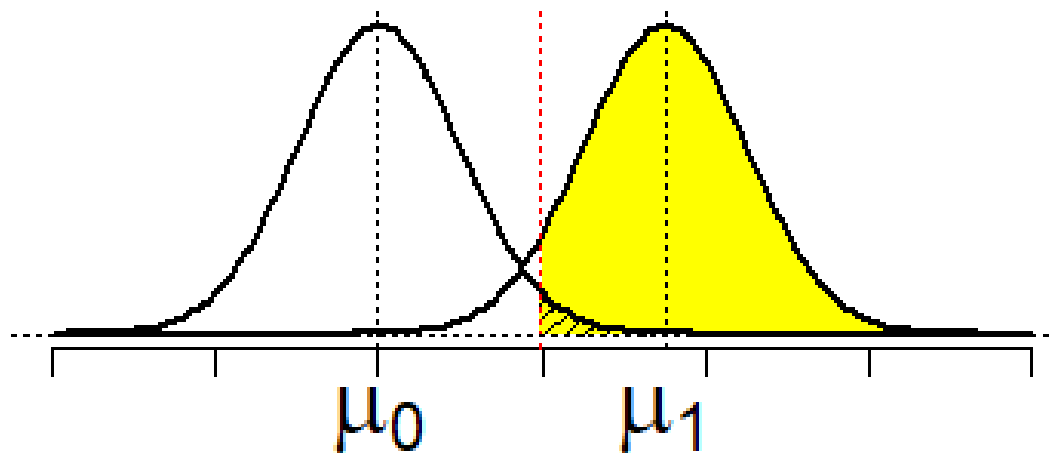
Neprůkaznost rozdílu, který jsme očekávali, je nejčastěji důsledkem toho, že buď rozdíly neexistují, nebo máme málo dat.

Poznámky k postupu

- Statistik má nejdříve formulovat hypotézu, zvolit rozhodovací pravidlo, určit hladinu testu, podle toho spočítat minimální rozsah výběru, a pak teprve sbírat data.
- Biolog nasbívá data, polovinu jich vyřadí a pak se ptá, co z toho lze otestovat 😊
- Přesto máme pokusy, kdy je třeba o rozsahu výběru i o hladině testu uvažovat předem -> plánování experimentů, výpočet potřebného rozsahu výběru tak, aby bylo možné dosáhnout potřebné hladiny testu α .

Síla testu ($1 - \beta$) [power of the test]

= pravděpodobnost, že nulovou hypotézu zamítneme, když ona neplatí
= pravděpodobnost, s jakou odhalíme neplatnost hypotézy \rightarrow ta žlutá prst.

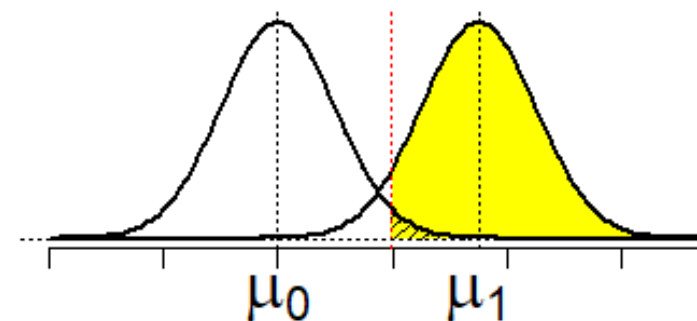


Šrafovaná část: pro takové hodnoty \bar{X} zamítám hypotézu, že $\mu_X = \mu_0$.
Typicky je pravděpodobnost této plochy $\alpha/2 \sim 2.5\%$.
Žlutá plocha: pravděpodobnost $1 - \beta$, tedy síla testu.

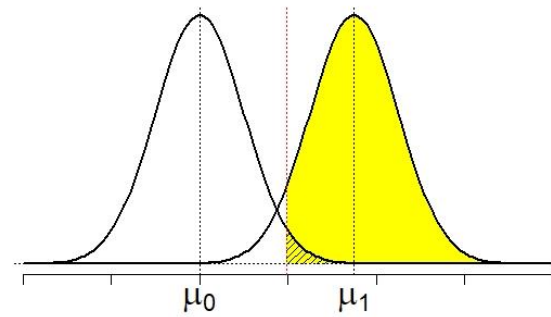
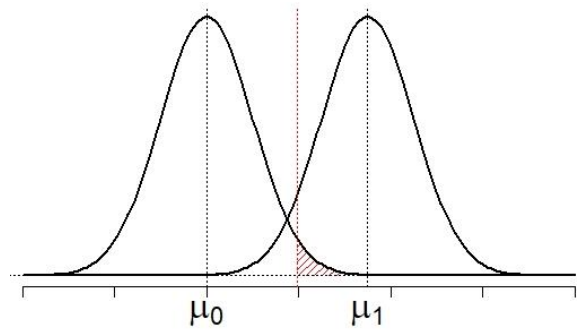
- Sílu testu většinou neznáme. Závisí na skutečném rozdělení výběrového souboru.
- Víme ale, že síla testu roste s odchylkou od nulové hypotézy a také s počtem pozorování (rozsahem výběru).
- Také platí, že čím menší je α , tím větší bude β .

Síla testu ($1 - \beta$) : myšlenkový pochod

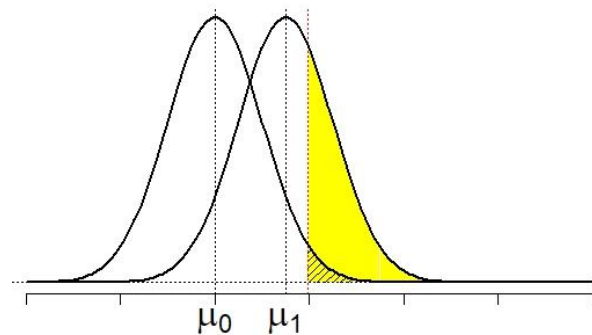
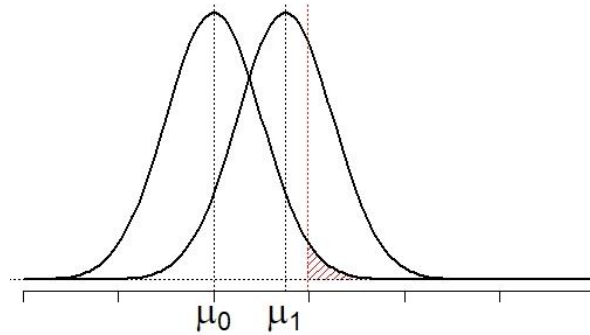
- Mám data, spočítám \bar{X} .
- Ptám se, zda \bar{X} patří do rozdělení s $EX = \mu_0$ nebo spíš do posunutého rozdělení s $EX = \mu_1$.
- Víme, že pro rozumná n má \bar{X} normální rozdělení, porovnááme tedy $N(\mu_0, \frac{\sigma}{\sqrt{n}})$ nebo $N(\mu_1, \frac{\sigma}{\sqrt{n}})$.
- Z dat odhadnu $se(\bar{X}) = \frac{S}{\sqrt{n}}$ a můžu vykreslit (odhad) tvaru hustoty pro \bar{X} .
- Podle hladiny testu α vyznačím příslušné kvantily na vodorovné ose (a mohu vyznačit šrafovanou pravděpodobnost $\alpha/2$).
- Pokud vyjde \bar{X} větší než kvantil $x_{\alpha/2}$, zamítám hypotézu $\mu_X = \mu_0$ ve prospěch alternativy $\mu_X = \mu_1$ (žlutá část vodorovné osy). A pravděpodobnost, že toto nastane, se jmenuje **síla testu** (žlutá plocha).



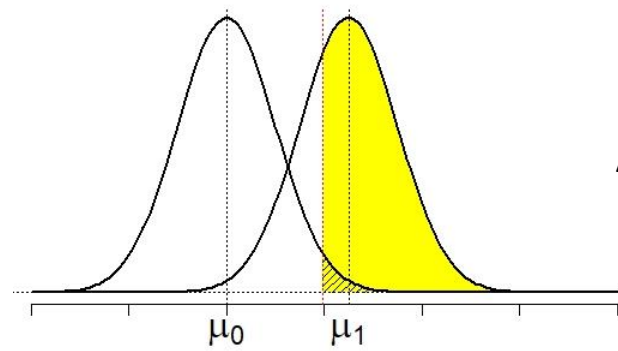
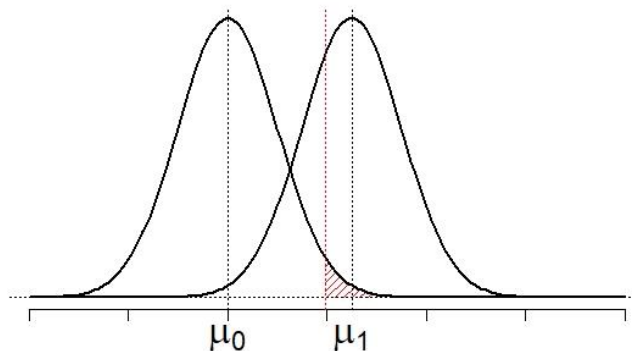
Síla testu podle vzdálenosti μ_0 a μ_1



Daleko od sebe
~ velká síla testu

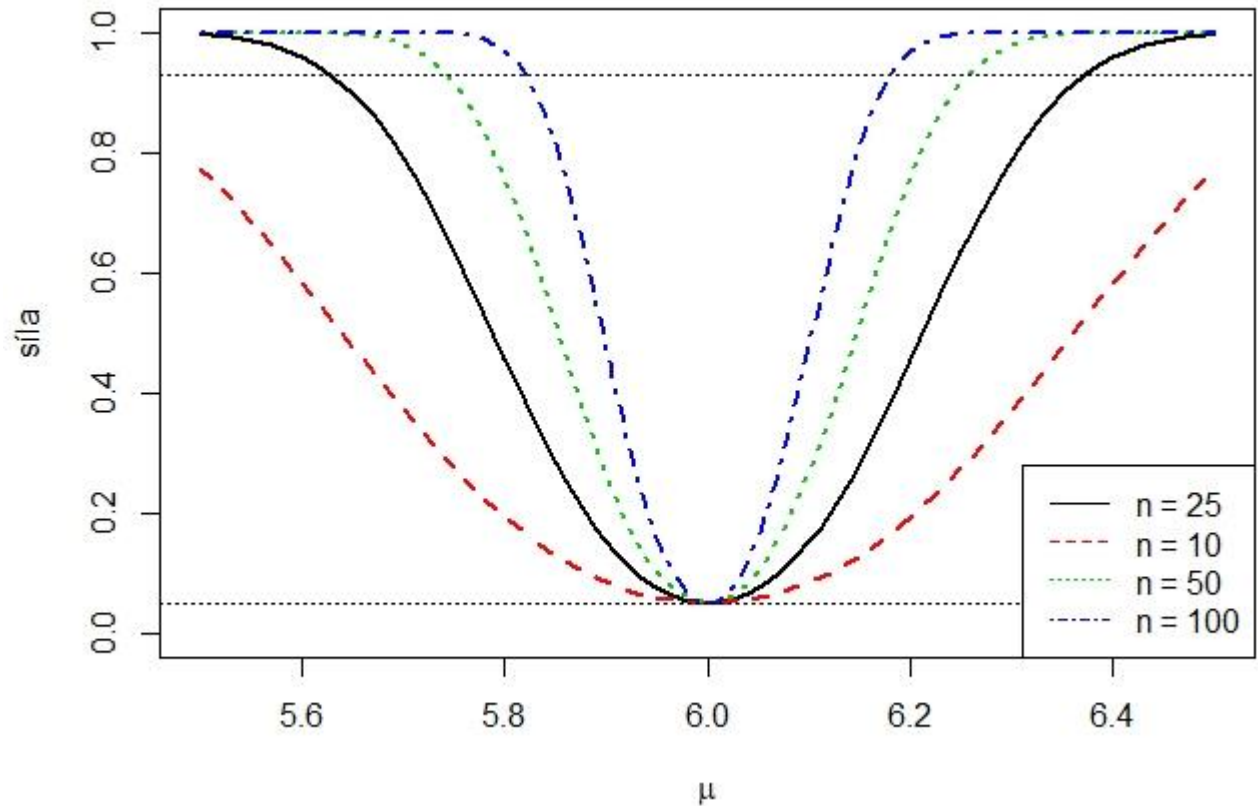
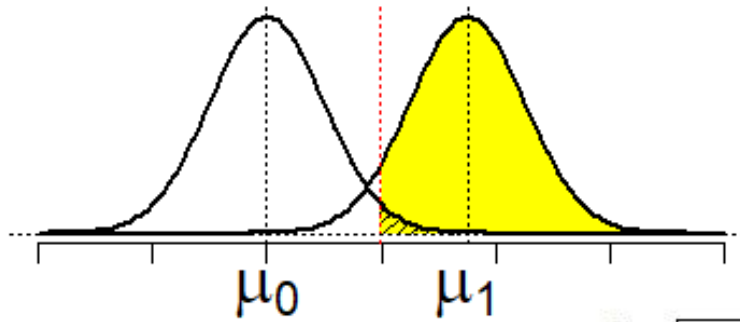


Blízko u sebe
~ malá síla testu



A celá škála mezi tím 😊

Jak se mění síla testu se vzdáleností od μ_0 a s počtem pozorování n



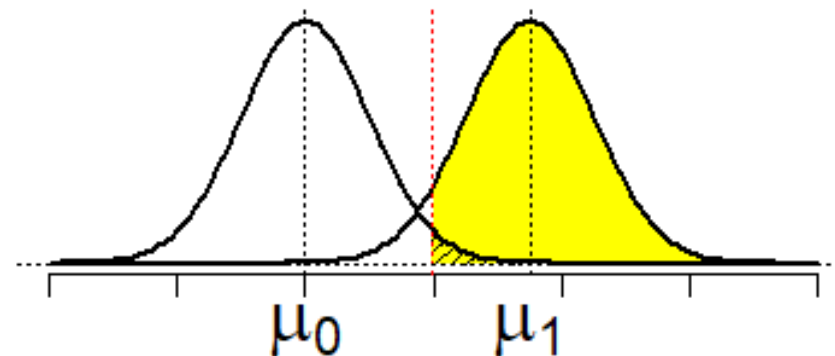
Jak spočítám sílu testu v R

```
power.t.test (n = NULL, delta = NULL, sd = 1, sig.level = 0.05,  
             power = NULL, type = c("two.sample", "one.sample", "paired"),  
             alternative = c("two.sided", "one.sided"), ...)
```

`power.anova`

`power.prop.test`

- Právě jeden z parametrů *n*, *delta*, *power*, *sd* nebo *sig.level* musí být neznámý (= NULL). Tento parametr se pak dopočítává z ostatních, které naopak musí být zadány, specifikovány.
- *type*: musím specifikovat, z kolika výběrů test počítám
- *alternative*: mám oboustranný nebo jednostranný test?



Síla testu ($1 - \beta$)

Různé typy testů mají také různé síly, tím se zabývá teorie.

Nás pak zajímají praktické poznámky typu

- „test B je silnější než běžně používaný test A“
- „test C je silný, ale je citlivý na porušení předpokladů o normalitě dat“ (tzn. mám pěkná data z normálního rozd. => беру test C)
- „test D je spíše slabý, ale je robustní k narušení předpokladů“ (tzn. použiju ho tam, kde data nejsou zrovna příkladně gaussovská).

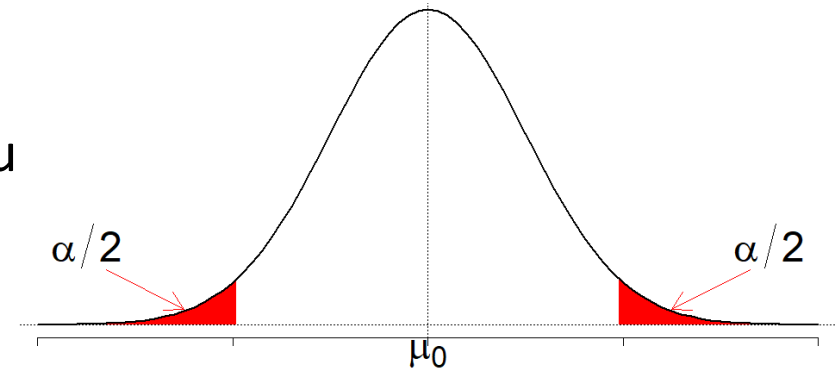
Oboustranná vs. jednostranná alternativa

[two-tailed vs. one-tailed alternative]

Oboustranná alternativa

$$H_0: \mu_X = \mu_0, \quad H_1: \mu_X \neq \mu_0$$

... tedy μ_X může být větší. Teorie případu nenapovídá nic o tom, na kterou stranu se rozdělení dat může posunout (přestože nám to napovídají čísla!)

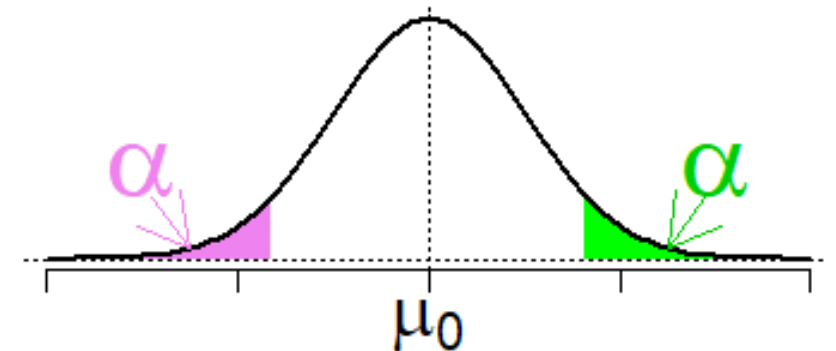


Jednostranná alternativa

Pokud z povahy případu vyplývá, že pokud se střední hodnota změní, může být jedině menší (větší) než testovaná hodnota μ_0 , zapracuju tento fakt do H_1 :

$$H_0: \mu_X \geq \mu_0, \quad H_1: \mu_X < \mu_0 \quad \text{nebo} \quad H_0: \mu_X \leq \mu_0, \quad H_1: \mu_X > \mu_0$$

Rozhodovací pravidlo: $T < t_{n-1}(\alpha)$ nebo $T > t_{n-1}(1 - \alpha)$



Testování hypotéz - slovníček

Chyba 1. druhu – Type I error

Chyba 2. druhu – Type II error

Síla testu – power of the test

Hladina testu – significance level

Zamítnout hypotézu – to reject hypothesis

Oboustranný test – two-tailed test

Jednostranný test – one-tailed test, left/right-tailed test

Kritický obor – takové výsledky testové statistiky, kdy H_0 zamítáme

Obor přijetí – takové výsledky testové statistiky, kdy H_0 nezamítáme