

## 6 Bodové a intervalové odhady parametrů

### 6.1 Úvod do matematické statistiky

V rámci kapitol 2 a 3 jsme se seznámili se základními metodami popisné statistiky. Připomeňme si, že tyto metody slouží výhradně k seznámení se s datovým souborem, k pochopení podstaty předložených dat a zjištění jejich základních vlastností. Je důležité si uvědomit, že výsledky a závěry metod popisné statistiky se vztahují pouze a jedině k předloženému datovému souboru a jeho hranice nikdy nepřekročí.

Snahou každého výzkumníka je však naopak poznat a používat metody, které jsou schopné hranice datového souboru překročit a umožnit mu rozšíření informací získaných na základě datového souboru na celou zkoumanou populaci. V praxi totiž častokrát nemáme možnost zkoumat výskyt nějaké vlastnosti v celé populaci, neboť zkoumaná populace může být velmi rozsáhlá a nasbírání hodnot od každého subjektu z této populace by bylo časově i finančně velmi náročné. Z tohoto důvodu je pro nás mnohem jednodušší sestavit pouze reprezentativní vzorek subjektů ze zkoumané populace, který svým složením jednak dostatečně pokrývá celou populaci a jednak dostatečně reprezentuje její stěžejní rysy. Tento reprezentativní vzorek potom vyhodnotíme pomocí vhodných statistických metod a závěry platné pro reprezentativní vzorek následně rozšíříme na celou populaci (tento krok si v případě, že vybraný vzorek je skutečně reprezentativním vzorkem celé populace, můžeme dovolit).

Reprezentativní vzorek, ve statistické terminologii nazývaný jako *náhodný výběr*, je soubor  $n$  stochasticky nezávislých náhodných veličin  $X_1, \dots, X_n$ , které se řídí stejným modelem  $L$  s parametry  $\theta$ , tj.  $X_1 \sim L(\theta), \dots, X_n \sim L(\theta)$ . Protože každá díleč náhodná veličina se řídí stejným modelem  $L(\theta)$ , můžeme předpokládat, že celý náhodný výběr  $X_1, \dots, X_n$  se také řídí modelem  $L(\theta)$ . V praxi může být modelem  $L(\theta)$  například alternativní model  $\text{Alt}(p)$ , kde  $\theta = p$ , binomický model  $\text{Bin}(N, p)$ , kde  $\theta = (N, p)^T$ , normální model  $N(\mu, \sigma^2)$ , kde  $\theta = (\mu, \sigma^2)^T$ , apod. Konkrétní číselné realizace náhodného výběru  $X_1, \dots, X_n$  (značíme je malými písmeny  $x_1, \dots, x_n$ ), tvoří *datový soubor*.

V souvislosti s náhodným výběrem definujeme také pojem *statistika*, jako libovolnou funkci  $T = T(X_1, \dots, X_n)$  náhodného výběru, která žádným způsobem nezávisí na parametru  $\theta$ . *Realizací statistiky*  $t$  potom označujeme statistiku  $T$  vyhodnocenou v realizaci náhodného výběru, tj.  $t = T(x_1, \dots, x_n)$ .

#### Příklad 6.1. Repezentativní vzorek

Předpokládejme, že chceme provést studii zkoumající výšku žen ve věku 25–35 let v Jihomoravském kraji. V ideálním případě bychom oslovili všechny ženy v požadovaném věku s trvalým pobytem v Jihomoravském kraji, změřili jejich výšku, zaznamenali ji do tabulky a nasbíraná data statisticky vyhodnotili. Takový výzkum by byl však časově i finančně náročný a navíc není pravděpodobné, že bychom do studie dokázali zahrnout úplně všechny ženy. Proto raději vytvoříme reprezentativní vzorek žen z Jihomoravského kraje o rozsahu například  $n = 1000$ . Aby byl vzorek reprezentativní, měl by rovnoměrně pokrývat ženy z celého Jihomoravského kraje. S využitím multihypergeometrického modelu modelu popsáno v kapitole 4 můžeme vypočítat, že pro zachování rovnoměrného pokrytí celého Jihomoravského kraje bychom měli oslovit přibližně 92 žen okresu Blansko, 320 žen z okresu Brno-město, 187 žen z okresu Brno-venkov, 98 žen z okresu Břeclav, 130 žen z okresu Hodonín, 77 žen z okresu Vyškov a 96 žen z okresu Znojmo. Volba žen v každém okrese by měla být čistě náhodná a měla by pokrývat celou věkovou kategorii 25–35 let. Reprezentativní vzorek, neboli náhodný výběr bude sestávat z  $n = 1000$  náhodných veličin  $X_1, \dots, X_{1000}$ , kde veličina  $X_1$  bude popisovat výšku první ženy,  $\dots$ ,  $X_{1000}$  bude popisovat výšku tisící ženy. O každé náhodné veličině předpokládáme, že se řídí normálním modelem, tj.  $X_1 \sim N(\mu, \sigma^2), \dots, X_{1000} \sim N(\mu, \sigma^2)$ , kde střední hodnota  $\mu$  i rozptyl  $\sigma^2$  jsou shodné pro všechny náhodné veličiny. Potom tedy také o celém náhodném výběru předpokládáme, že se řídí normálním modelem, tj.  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . Nyní se dostáváme do fáze, kdy všechny ženy změříme a zjistíme například, že první změřená žena měří 165 cm, druhá žena měří 168 cm,  $\dots$ , tisící žena měří 163 cm. Zaznamenáním naměřených hodnot do tabulky získáme realizace náhodných veličin  $x_1 = 165, x_2 = 168, \dots, x_{1000} = 163$ , které společně tvoří datový soubor. ★

#### Příklad 6.2. Jednorozměrné statistiky

Mějme jeden náhodný výběr  $X_1, \dots, X_n$  o rozsahu  $n \geq 2$ . Příkladem statistiky pro tento náhodný výběr může být například výběrový průměr

$$M = \frac{1}{n} \sum_{i=1}^n X_i. \quad (6.1)$$

Všimněme si, že ve vzorci ?? výběrového průměru vystupují pouze hodnoty náhodného výběru  $X_1, \dots, X_n$  a rozsah

náhodného výběru  $n$ . Libovolný parametr  $\theta$  (např.  $\mu$ ,  $\sigma^2$ ,  $p$ , apod.) se ve vzorci nevyskytuje.

Dalším příkladem jednorozměrné statistiky je výběrový rozptyl

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - M)^2. \quad (6.2)$$

Opět si všimněme, že ve vzorci ?? výběrového rozptylu se vyskytují pouze hodnoty náhodného výběru  $X_1, \dots, X_n$ , rozsah náhodného výběru  $n$  a výběrový průměr  $M$ , který, jakožto statistika, je funkcí náhodného výběru. Žádný parametr se ve vzorci nevyskytuje. Výběrový rozptyl je tedy opět pouze funkcí náhodného výběru. Z výběrového rozptylu vychází další statistika nazývaná výběrová směrodatná odchylka, která definovaná jako odmocnina z výběrového rozptylu, tj.

$$S = \sqrt{S^2}. \quad (6.3)$$

Výběrová směrodatná odchylka je statistikou, neboť jde pouze o odmocninou statistiky nazývané výběrový rozptyl.

Posledním příkladem statistiky, který si uvedeme, je výběrový variační koeficient. Tento koeficient vyjadřuje míru směrodatné odchylky vysvětlené aritmetickým průměrem (viz kapitola 3).

$$V = \frac{S}{M}. \quad (6.4)$$

Protože výběrový variační koeficient není nic jiného, než podíl dvou statistik, a sice výběrové směrodatné odchylky a výběrového průměru, je sám také statistikou. ★

### Příklad 6.3. Dvourozměrné statistiky

Nechť  $(X_1, Y_1), \dots, (X_n, Y_n)$  je náhodný výběr z dvourozměrného rozdělení,  $M_1$  a  $M_2$  jsou výběrové průměry a  $S_1^2$  a  $S_2^2$  jsou výběrové rozptyly. Příkladem dvourozměrné statistiky je výběrová kovariance

$$S_{12} = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_1)(Y_i - M_2). \quad (6.5)$$

Všimněme si, že ve vzorci ?? výběrové kovariance vystupují kromě hodnot dvourozměrného náhodného výběru pouze výběrové průměry  $M_1$  a  $M_2$ , které jsou statistikami a tedy funkcemi náhodného výběru, a rozsah  $n$ . Proto je výběrová kovariance rovněž statistika. Druhým příkladem dvourozměrné statistiky je výběrový korelační koeficient

$$R_{12} = \frac{S_{12}}{\sqrt{S_1^2 S_2^2}} = \frac{S_{12}}{S_1 S_2}. \quad (6.6)$$

Výběrový korelační koeficient je definován jako podíl výběrové kovariance, která je statistikou, a odmocniny ze součinu výběrových rozptylů  $S_1$  a  $S_2$ , které jsou rovněž statistikami. Žádný parametr ve vzorci ?? nefiguruje, proto je výběrový korelační koeficient též statistikou. ★

### Příklad 6.4. Dvouvýběrové a vícevýběrové statistiky

Nechť  $X_{11}, \dots, X_{1n_1}$  a  $X_{21}, \dots, X_{2n_2}$ ,  $n_1 \geq 2$  a  $n_2 \geq 2$ , jsou dva náhodné výběry z jednozměrných rozdělení,  $M_1$  a  $M_2$  jsou jejich výběrové průměry a  $S_1^2$  a  $S_2^2$  jejich výběrové rozptyly. Příkladem dvouvýběrové statistiky je vážený průměr dvou výběrových rozptylů

$$S_*^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} \quad (6.7)$$

Všimněme si, že vzorec váženého průměru dvou výběrových rozptylů ?? sestává z rozsahů náhodných výběrů  $n_1$  a  $n_2$  a výběrových rozptylů  $S_1^2$  a  $S_2^2$ , které jsou statistikami. Proto i vážený průměr dvou výběrových rozptylů  $S_*^2$  je statistikou. Vzorec váženého průměru dvou výběrových rozptylů ?? můžeme zobecnit na vážený průměr  $m$  náhodných výběrů.

Nechť  $X_{11}, \dots, X_{1n_1}$ ,  $X_{21}, \dots, X_{2n_2}$ ,  $\dots$ ,  $X_{m1}, \dots, X_{mn_m}$ ,  $n_1 \geq 2$ ,  $n_2 \geq 2$ ,  $\dots$ ,  $n_m \geq 2$  je  $m$  náhodných výběrů z jednozměrných rozdělení,  $M_1, M_2, \dots, M_m$  jsou jejich výběrové průměry a  $S_1^2, S_2^2, \dots, S_m^2$  jejich výběrové rozptyly. Příkladem vícevýběrové statistiky je vážený průměr  $m$  výběrových rozptylů

$$S_*^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2 + \dots + (n_m - 1)S_m^2}{n_1 + n_2 + \dots + n_m - m} \quad (6.8)$$

Analogicky jako ve vzorci ?? vystupují ve vzorci váženého průměru  $m$  výběrových rozptylů ?? pouze rozsahy náhodných výběrů  $n_1, n_2, \dots, n_m$  a výběrové rozptyly  $S_1^2, S_2^2, \dots, S_m^2$ , které jsou statistikami. Proto také vážený průměr  $m$  výběrových rozptylů je statistikou. ★

### Příklad 6.5. Jednorozměrné statistiky

Mějme k dispozici datový soubor 15-anova-means-skull.txt obsahující původní kranio-metrické údaje o výšce horní části tváře (proměnná upface.H) mužů z pěti populací (německé, čínské, bantuské, peruánské a malajské). Pro čínskou populaci vypočítejte (a) výběrový průměr; (b) výběrový rozptyl; (c) výběrovou směrodatnou odchylku; (d) výběrový koeficient variace. Všechny vypočítané hodnoty statistik řádně interpretujte.

### Řešení příkladu ??

Celkem máme k dispozici  $n = 18$  náhodných veličin  $X_1, \dots, X_{18}$ , přičemž veličina  $X_1$  popisuje výšku horní části tváře u prvního mužského skeletu čínské populace,  $\dots$ ,  $X_{18}$  popisuje výšku horní části tváře u osmnáctého mužského skeletu čínské populace. Naměřením hodnoty výšky horní části tváře každého skeletu jsme získali celkem 18 realizací náhodných veličin, konkrétně  $x_1 = 77, \dots, x_{18} = 70$ . Těchto 18 realizací tvoří společně datový soubor.

Výběrový průměr výšky horní části tváře vypočítáme dosazením hodnot do vzorce ??, tj.

$$m = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{18} (77 + 71 + \dots + 75 + 70) = \frac{1296}{18} = 72.$$

Výběrový rozptyl výšky horní části tváře dopočítáme dosazením do vzorce ??, tj.


$$\begin{aligned} s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2 \\ &= \frac{1}{18-1} ((77-72)^2 + (71-72)^2 + \dots + (75-72)^2 + (70-72)^2) \\ &= \frac{1}{17} (5^2 + (-1)^2 + \dots + 3^2 + (-2)^2) \\ &= \frac{354}{17} \\ &= 4.563281 \doteq 4.5633. \end{aligned}$$

Výběrovou směrodatnou odchylku výšky horní části tváře dopočítáme jako odmocninu z výběrového rozptylu (viz vzorec ??), tj.

$$s = \sqrt{s^2} = \sqrt{4.563281} = 2.13619 \doteq 2.1362.$$

Konečně, výběrový koeficient variace výšky horní části tváře dopočítáme jako podíl výběrové směrodatné odchylky a výběrového rozptylu (viz vzorec ??), tj.

$$v = \frac{s}{m} = \frac{2.13619}{72} = 0.029669 \doteq 0.0297.$$

Výpočet provedeme také pomocí softwaru . Nejprve načteme datový soubor příkazem read.delim() a odstraníme NA hodnoty příkazem na.omit(). Pomocí operátoru [] vybereme z tabulky data pouze řádky týkající se čínské populace (data\$pop == 'cin') a sloupec s naměřenými výškami horní části tváře 'upface.H'.

```
1 data <- read.delim('15-anova-means-skull.txt')
2 data <- na.omit(data)
3 upface.HC <- data[data$pop == 'cin', 'upface.H']
4 n <- length(upface.HC) # 18
5 upface.HC
```

```
[1] 77 71 76 75 60 75 72 67 68 75 75 78 70 67 70 75 75 70
```

Vidíme, že datový soubor skutečně obsahuje naměřené výšky horní části tváře 18 mužů čínské populace. Nyní vypočítáme hodnoty všech čtyř požadovaných statistik, a to nejprve přímým opisem vzorců `??`, `??`, `??` a `??`, a následně pomocí `R`-kovských funkcí. Výběrový průměr spočítáme pomocí funkce `mean()`, výběrový rozptyl pomocí funkce `var()` a výběrovou směrodatnou odchylku pomocí funkce `sd()`.

```
7 m.HC <- 1 / n * sum(upface.HC)
8 s2.HC <- 1 / (n - 1) * sum((upface.HC - m.HC) ^ 2)
9 s.HC <- sqrt(s2.HC)
10 v.HC <- s.HC / m.HC
11
12 mm.HC <- mean(upface.HC)
13 ss2.HC <- var(upface.HC)
14 ss.HC <- sd(upface.HC)
15
16 (tab <- data.frame(prumer = m.HC, rozptyl = s2.HC, sm.odch = s.HC, koef.var = v.HC))
```

```
prumer rozptyl sm.odch koef.var
1      72 20.82353 4.563281 0.0633789
```

17  
18

**Interpretace výsledků:** Výběrový průměr výšky horní části tváře u mužů čínské populace  $\bar{x} = 72$  mm, výběrový rozptyl  $s^2 = 20.8235$  mm<sup>2</sup>, výběrová směrodatná odchylka  $s = 4.5633$  mm a výběrový koeficient variace  $v = 0.06338$  (6.34%). ★

### Příklad 6.6. Dvourozměrné statistiky

Máme k dispozici datový soubor `30-goldman-alaska.csv` obsahující antropometrické údaje o délce kosti stehenní v mm (proměnná `femur.R`) a acetabulární výšce v mm (proměnná `acetab.R`) z pravé strany u skeletů jedinců z aljašské populace (muži a ženy z kmene Tigara). Pro tento dvourozměrný náhodný výběr vypočítejte (a) výběrovou kovarianci; (b) výběrový korelační koeficient pro skelety mužského pohlaví. Obě hodnoty řádně interpretujte.

#### Řešení příkladu ??

Celkem máme k dispozici  $n = 24$  dvojic náhodných veličin  $(X_1, Y_1), \dots, (X_{24}, Y_{24})$ , přičemž veličina  $X_1$  popisuje délku stehenní kosti z pravé strany u prvního mužského skeletu,  $\dots$ ,  $X_{24}$  popisuje délku stehenní kosti z pravé strany u čtyřicátého mužského skeletu. Podobně potom náhodná veličina  $Y_1$  popisuje acetabulární výšku z pravé strany u prvního mužského skeletu,  $\dots$ ,  $Y_{24}$  popisuje acetabulární výšku z pravé strany u čtyřicátého mužského skeletu. Naměřením hodnot délky stehenní kosti z pravé strany a acetabulární výšky z pravé strany každého skeletu jsme získali celkem 24 dvojic realizací náhodných veličin, konkrétně  $(x_1, y_1) = (452, 52.24)$ ,  $\dots$ ,  $(x_{24}, y_{24}) = (415, 53.35)$ . Těchto 24 dvojic realizací tvoří společně datový soubor. Výběrovou kovarianci dopočítáme dosazením do vzorce `??`, tj.

$$s_{12} = \frac{1}{n-1} \sum_{i=1}^n (x_i - m_1)(y_i - m_2),$$

kde výběrový průměr délky stehenní kosti  $m_1$  a výběrový průměr acetabulární výšky  $m_2$  vypočítáme dosazením do vzorce `??`.

$$m_1 = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{24} (452 + 415 + \dots + 440 + 494) = \frac{10\,269.5}{24} = 427.8958.$$

$$m_2 = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{24} (52.24 + 53.35 + \dots + 52.05 + 60.49) = \frac{1\,246}{24} = 51.92708 \doteq 51.9271.$$

$$\begin{aligned}
s_{12} &= \frac{1}{24-1}((452-427.8958)(52.24-51.92708) + (415-427.8958)(53.35-51.92708) + \dots \\
&\quad \dots + (440-427.8958)(52.05-51.92708) + (60.49-427.8958)(186-51.92708)) \\
&= \frac{1}{23}(7.54268626 - 18.34969174 + \dots + 1.48784826 + 566.04497626) \\
&= \frac{968.5627}{23} = 42.11142 \doteq 42.1114.
\end{aligned}$$

Výběrový korelační koeficient vypočítáme dosazením do vzorce ??, tj.



$$r_{12} = \frac{s_{12}}{\sqrt{s_1^2 s_2^2}} = \frac{s_{12}}{s_1 s_2}, \quad (6.9)$$

kde  $s_{12}$  je výběrové kovariance (viz výše), a výběrový rozptyl délky stehenní kosti  $s_1^2$  a výběrový rozptyl acetabulární výšky  $s_2^2 = 40.7664$  vypočítáme pomocí vzorce ??.

$$\begin{aligned}
s_1^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2 \\
&= \frac{1}{24-1} ((452-427.8958)^2 + (415-427.8958)^2 + \dots + (440-427.8958)^2 + (494-427.8958)^2) \\
&= \frac{1}{23} (24.1042^2 + (-12.8958)^2 + \dots + 12.1042^2 + 66.1042^2) \\
&= \frac{12400.99}{23} \\
&= 539.1735.
\end{aligned}$$

$$\begin{aligned}
s_2^2 &= \frac{1}{n-1} \sum_{i=1}^n (y_i - m)^2 \\
&= \frac{1}{24-1} ((52.24-51.92708)^2 + (53.35-51.92708)^2 + \dots + (52.05-51.92708)^2 + (60.49-51.92708)^2) \\
&= \frac{1}{17} (0.3129^2 + 1.4229^2 + \dots + 0.1229^2 + 8.5629^2) \\
&= \frac{224.6287}{23} \\
&= 9.766465 \doteq 9.7665.
\end{aligned}$$

$$r_{12} = \frac{42.11142}{\sqrt{539.1735} \sqrt{9.766465}} = \frac{42.11142}{23.22011 \times 3.125134} = \frac{42.11142}{72.56596} = 0.5803192.$$

Hodnotu výběrové kovariance, resp. výběrového korelačního koeficientu dále spočítáme pomocí softwaru  přepisem vzorce ??, resp. ?? s použitím funkcí `sum()` a `sqrt()`. Následně obě statistiky vypočítáme pomocí  funkcí, a sice výběrovou kovarianci pomocí funkce `cov()` a výběrový korelační koeficient pomocí funkce `cor()`.

	kovariance	korelacni.koeficient
1	42.11142	0.5803192

19  
20

**Interpretace výsledků:** Výběrová kovariance délky klíční kosti z pravé strany a acetabulární výšky z pravé strany u skeletů mužského pohlaví z aljašské populace kmene Tigara  $s_{12} = 42.1114$ , výběrový korelační koeficient  $r_{12} = 0.58032$ .



### Příklad 6.7. Dvourozměrné statistiky

Máme k dispozici datový soubor obsahující antropometrické údaje o délce kosti pažní v mm (znak  $X$ ) a délce kosti stehenní v mm (znak  $Y$ ) z levé strany u 18 skeletů mužského pohlaví a 20 skeletů ženského pohlaví z římského pohřebiště v Poundbury. Ze zadaných údajů byly dopočítány následující statistiky pro skelety ženského pohlaví: výběrové průměry:  $m_X = 288.9500$  mm,  $m_Y = 411.4000$  mm; výběrové směrodatné odchylky:  $s_X = 10.3287$  mm,  $s_Y = 16.1323$  mm; výběrová kovariance:  $s_{12} = 104.7579$ . Na základě uvedených údajů vypočítejte výběrový korelační koeficient pro skelety ženského pohlaví a vypočítanou hodnotu řádně interpretujte.

#### Řešení příkladu ??

V tomto příkladu vypočítáme hodnotu výběrového korelačního koeficientu délky pažní kosti z levé strany a délky stehenní kosti z levé strany u skeletů ženského pohlaví, ovšem bez znalosti datového souboru. K dispozici máme pouze rozsah dvourozměrného náhodného výběru ( $n = 20$ ) a hodnoty výběrových průměrů, výběrových směrodatných odchylek a výběrové kovariance. To nám však k výpočtu stačí, ba dokonce některé údaje uvedené v zadání ani nevyužijeme. Dosazením výběrové kovariance a výběrových korelačních koeficientů do vzorce ?? získáme hodnotu výběrového korelačního koeficientu.

$$r_{12} = \frac{s_{12}}{s_1 s_2} = \frac{104.7579}{10.3287 \times 16.1323} = \frac{104.7579}{166.6257} = 0.6287019 \doteq 0.6287.$$

*Poznámka:* K výpočtu jsme tentokrát použili variantu vzorce pro výběrový korelační koeficient se směrodatnými odchylkami ve jmenovateli, a to proto, že v zadání příkladu jsme měli uvedené právě hodnoty výběrových směrodatných odchylek.

```
21 s12 <- 104.7579
22 s1 <- 10.3287
23 s2 <- 16.1323
24 (r12 <- s12 / (s1 * s2))
```

```
[1] 0.628702
```

25

**Interpretace výsledků:** Výběrový korelační koeficient délky pažní a stehenní kosti z levé strany u skeletů ženského pohlaví z pohřebiště v Poundbury  $r_{12} = 0.6287$ .



### Příklad 6.8. Dvouvýběrové statistiky

Máme k dispozici datový soubor 21-goldman-measures obsahující naměřené údaje o délce kyčelní kosti (v mm) z levé strany u mužských skeletů tří japonských populací (9 skeletů z populace Tsugumo Shell Mound, 7 skeletů z populace Yoshigo Shell Mound a 3 skelety z populace Yasaki Shell Mound). Vypočítejte (a) vážený průměr výběrových rozptylů pro každou dvojici uvedených japonských populací; (b) vážený průměr výběrových rozptylů všech tří uvedených populací. Všechny vypočítané hodnoty vážených průměrů řádně interpretujte.

#### Řešení příkladu ??

Celkem máme k dispozici  $n = 9$  náhodných veličin  $X_{11}, \dots, X_{19}$ , přičemž veličina  $X_{11}$  popisuje délku kyčelní kosti z levé strany u prvního mužského skeletu z populace Tsugumo Shell Mound,  $\dots, X_{19}$  popisuje délku kyčelní kosti z levé strany u devátého mužského skeletu z populace Tsugumo Shell Mound. Dále máme k dispozici  $n = 7$  náhodných veličin  $X_{21}, \dots, X_{27}$ , přičemž veličina  $X_{21}$  popisuje délku kyčelní kosti z levé strany u prvního mužského skeletu z populace Yoshigo Shell Mound,  $\dots, X_{27}$  popisuje délku kyčelní kosti z levé strany u sedmého mužského skeletu z populace Yoshigo Shell Mound. Konečně, máme k dispozici  $n = 3$  náhodné veličiny  $X_{31}, \dots, X_{33}$ , přičemž veličina  $X_{31}$  popisuje délku kyčelní kosti z levé strany u prvního mužského skeletu z populace Yasaki Shell Mound,  $\dots, X_{33}$  popisuje délku kyčelní kosti z levé strany u třetího mužského skeletu z populace Yasaki Shell Mound. Naměřením hodnot délky kyčelní kosti z levé strany každého mužského skeletu populace Tsugumo Shell Mound jsme získali celkem 9 realizací náhodných veličin  $x_{11} = 152, \dots, x_{19} = 137$ , naměřením hodnot délky kyčelní kosti z levé strany každého mužského skeletu populace Yoshigo Shell Mound jsme získali celkem 9 realizací náhodných veličin  $x_{21} = 142, \dots, x_{27} = 152$  a naměřením hodnot délky kyčelní kosti z levé strany každého mužského skeletu populace Yasaki Shell Mound jsme získali celkem 3 realizace náhodných veličin  $x_{31} = 156, \dots,$

$x_{33} = 154$ . Těchto 9, 7 a 3 realizace tvoří tři datové soubory.

K výpočtu vážených průměrů výběrových rozptylů potřebujeme znát hodnoty výběrových rozptylů. Nejprve si tedy vypočítáme hodnoty výběrových průměrů a hodnoty výběrových rozptylů pro všechny tři populace a nakonec vypočítáme vážené průměry výběrových rozptylů pro různé dvojice populací podle vzorce ?? a vážený průměr výběrových rozptylů pro všechny tři populace podle vzorce ?? upraveného pro  $m = 3$ .

Výběrové průměry a výběrové rozptyly vypočítáme rovnou pomocí R funkcí. Načteme datový soubor a vytvoříme postupně vektory délek kyčelních kostí z levé strany (iblade.LL) pro skelety mužského pohlaví (sex == 'm') populace Tsugumo Shell Mound (pop == 'Tsugumo Shell Mound'), Yoshigo Shell Mound (pop == 'Yoshigo Shell Mound') a Yasaki Shell Mound (pop == 'Yasaki Shell Mound'). Jednotlivé vektory postupně pojmenujeme iblade.LLT, iblade.LLYo, iblade.LLYa a následně z každého z nich odstraníme chybějící pozorování.

```
26 data <- read.delim('00-Data//21-goldman-measures.csv', sep = ';', dec = '.')
27 iblade.LLT <- data[data$pop == 'Tsugumo Shell Mound' & data$sex == 'm',
28                   'iblade.LL']
29 iblade.LLYo <- data[data$pop == 'Yoshigo Shell Mound' & data$sex == 'm',
30                   'iblade.LL']
31 iblade.LLYa <- data[data$pop == 'Yasaki Shell Mound' & data$sex == 'm',
32                   'iblade.LL']
33 iblade.LLT <- na.omit(as.numeric(iblade.LLT))
34 iblade.LLYo <- na.omit(as.numeric(iblade.LLYo))
35 iblade.LLYa <- na.omit(as.numeric(iblade.LLYa))
```

Nyní příkazem length stanovíme počet pozorování, příkazem mean() vypočítáme výběrové průměry a příkazem var() vypočítáme výběrové rozptyly pro každou populaci.

```
36 n.T <- length(iblade.LLT)
37 n.Yo <- length(iblade.LLYo)
38 n.Ya <- length(iblade.LLYa)
39
40 m.T <- mean(iblade.LLT)
41 m.Yo <- mean(iblade.LLYo)
42 m.Ya <- mean(iblade.LLYa)
43
44 s2.T <- var(iblade.LLT)
45 s2.Yo <- var(iblade.LLYo)
46 s2.Ya <- var(iblade.LLYa)
47
48 tab <- data.frame(n = c(n.T, n.Yo, n.Ya),
49                  prumery = c(m.T, m.Yo, m.Ya),
50                  rozptyly = c(s2.T, s2.Yo, s2.Ya),
51                  row.names = c('Tsugumo', 'Yoshigo', 'Yasaki'))
52 tab
```

	n	prumery	rozptyly
Tsugumo	9	149.2222	32.19444
Yoshigo	7	151.0000	20.66667
Yasaki	3	154.0000	4.00000

53  
54  
55  
56

Výběrový průměr délek kyčelních kostí z levé strany pro skelety mužského pohlaví pro populaci Tsugumo Shell Mound  $m_1 = 149.222$  mm, pro populaci Yoshigo Shell Mound  $m_2 = 151$  mm a pro populaci Yasaki Shell Mound  $m_3 = 154$  mm. Výběrový rozptyl délek kyčelních kostí z levé strany pro skelety mužského pohlaví pro populaci Tsugumo Shell Mound  $s_1^2 = 32.1944$  mm<sup>2</sup>, pro populaci Yoshigo Shell Mound  $s_2^2 = 20.6667$  mm<sup>2</sup> a pro populaci Yasaki Shell Mound  $s_3^2 = 4$  mm<sup>2</sup>.

Vážený průměr výběrových rozptylů délek kyčelních kostí z levé strany pro skelety mužského pohlaví pro populace Tsugumo Shell Mound a Yoshigo Shell Mound vypočítáme dosazením příslušných výběrových rozptylů a rozsahů výběrů do vzorce ??.

$$\begin{aligned}
S_{*12}^2 &= \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} \\
&= \frac{(9 - 1)32.19444 + (7 - 1)20.66667}{9 + 7 - 2} \\
&= \frac{8 \times 32.19444 + 6 \times 20.66667}{14} \\
&= \frac{381.5555}{14} = 27.25396 \doteq 27.2540
\end{aligned}$$


Analogicky vypočítáme vážený proměr výběrových rozptylů pro populace Tsugumo Shell Mound a Yasaki Shell Mound a pro populace Yoshigo Shell Mound a Yasaki Shell Mound.

$$\begin{aligned}
s_{*13}^2 &= \frac{(n_1 - 1)s_1^2 + (n_3 - 1)s_3^2}{n_1 + n_3 - 2} \\
&= \frac{(9 - 1)32.19444 + (3 - 1)4}{9 + 3 - 2} \\
&= \frac{8 \times 32.19444 + 2 \times 4}{10} \\
&= \frac{265.5555}{10} = 26.55555 \doteq 26.55556
\end{aligned}$$

$$\begin{aligned}
s_{*23}^2 &= \frac{(n_2 - 1)s_2^2 + (n_3 - 1)s_3^2}{n_2 + n_3 - 2} \\
&= \frac{(7 - 1)20.66667 + (3 - 1)4}{7 + 3 - 2} \\
&= \frac{6 \times 20.66667 + 2 \times 4}{8} \\
&= \frac{132}{8} = 16.5
\end{aligned}$$

Nakonec dosazením rozsahů náhodných výběrů a výběrových rozptylů všech tří populací do vzorce ?? vypočítáme hodnotu váženého průměru výběrových rozptylů pro všechny tři populace Tsugumo, Yoshigo a Yakasi Shell Mound.

$$\begin{aligned}
s_*^2 &= \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2 + (n_3 - 1)s_3^2}{n_1 + n_2 + n_3 - 3} \\
&= \frac{(9 - 1)32.19444 + (7 - 1)20.66667 + (3 - 1)4}{9 + 7 + 3 - 3} \\
&= \frac{8 \times 32.19444 + 6 \times 20.66667 + 2 \times 4}{16} \\
&= \frac{389.5555}{16} = 24.34722 \doteq 24.3472
\end{aligned}$$

Výpočet provedeme také pomocí softwaru .

```

57 sh.TYo <- ((n.T - 1) * s2.T + (n.Yo - 1) * s2.Yo) / (n.T + n.Yo - 2) # 27.25397
58 sh.TYa <- ((n.T - 1) * s2.T + (n.Ya - 1) * s2.Ya) / (n.T + n.Ya - 2) # 26.55556
59 sh.YoYa <- ((n.Yo - 1) * s2.Yo + (n.Ya - 1) * s2.Ya) / (n.Yo + n.Ya - 2) # 16.5
60
61 sh.TYoYa <- ((n.T - 1) * s2.T + (n.Yo - 1) * s2.Yo + (n.Ya - 1) * s2.Ya) /
62 (n.T + n.Yo + n.Ya - 3) # 24.34722

```



**Interpretace výsledků:** Vážený průměr výběrových rozptylů délek kyčelních kostí z levé strany pro skelety mužského pohlaví pro populace Tsugumo Shell Mound a Yoshigo Shell Mound  $s_{*12}^2 = 27.2540\text{m}^2$ , pro populace Tsugumo Shell Mound a Yasaki Shell Mound  $s_{*13}^2 = 26.5556\text{m}^2$  a pro populace Yoshigo Shell Mound a Yakasi Shell Mound  $s_{*23}^2 = 16.5\text{m}^2$ . Vážený průměr výběrových rozptylů všech tří populací Tsugumo, Yoshigo a Yasaki Shell Mound  $s_*^2 = 24.3472\text{mm}^2$ . ★

### Příklad 6.9. Dvouvýběrové statistiky

Máme k dispozici naměřené údaje o acetabulární výšce (v mm) z pravé strany u mužských skeletů ze tří pohřebišť na území Nového Mexika (19 skeletů s pohřebišť Hawikuh, 4 skelety z pohřebišť Pueblo Bonito a 7 skeletů z pohřebišť Puye). Ze zadaných údajů byly dopočítány následující charakteristiky: (a) Hawikuh: výběrový průměr:  $m_1 = 47.98\text{ mm}$ ; výběrová směrodatná odchylka:  $s_1 = 2.15\text{ mm}$ ; (b) Pueblo Bonito:  $m_2 = 51.08\text{ mm}$ ;  $s_2 = 1.83\text{ mm}$ ; (c) Puye:  $m_3 = 46.20\text{ mm}$ ;  $s_3 = 2.73\text{ mm}$ . Vypočítejte (a) vážený průměr výběrových rozptylů pro každou dvojici uvedených japonských populací; (b) vážený průměr výběrových rozptylů všech tří uvedených populací. Všechny vypočítané hodnoty vážených průměrů řádně interpretujte.

### Řešení příkladu ??

V tomto příkladu vypočítáme hodnotu vážených průměrů výběrových rozptylů acetabulární výšky z pravé strany u skeletů mužského pohlaví, ovšem bez znalosti datového souboru. K dispozici máme pouze rozsahy náhodných výběrů a hodnoty výběrových průměrů a výběrových směrodatných odchylek. To nám však k výpočtu stačí, ba dokonce některé údaje uvedené v zadání ani nevyužijeme. Dosazením rozsahů náhodných výběrů a výběrových směrodatných odchylek do vzorce ?? vypočítáme hodnoty vážených průměrů výběrových rozptylů pro skelety z pohřebišť Hawikuh a Pueblo Bonito, z pohřebišť Hawikuh a Puye a z pohřebišť Pueblo Bonito a Puye.

$$\begin{aligned} S_{*12}^2 &= \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} \\ &= \frac{(19 - 1)2.15^2 + (4 - 1)1.83^2}{19 + 4 - 2} \\ &= \frac{18 \times 4.6225 + 3 \times 3.3489}{21} \\ &= \frac{93.2517}{21} = 4.440557 \doteq 4.44056 \end{aligned}$$

$$\begin{aligned} s_{*13}^2 &= \frac{(n_1 - 1)s_1^2 + (n_3 - 1)s_3^2}{n_1 + n_3 - 2} \\ &= \frac{(19 - 1)2.15^2 + (7 - 1)2.73^2}{19 + 7 - 2} \\ &= \frac{18 \times 4.6225 + 6 \times 7.4529}{24} \\ &= \frac{127.9224}{24} = 5.3301 \end{aligned}$$


$$\begin{aligned} s_{*23}^2 &= \frac{(n_2 - 1)s_2^2 + (n_3 - 1)s_3^2}{n_2 + n_3 - 2} \\ &= \frac{(4 - 1)1.83^2 + (7 - 1)2.73^2}{4 + 7 - 2} \\ &= \frac{3 \times 3.3489 + 6 \times 7.4529}{9} \\ &= \frac{54.7641}{9} = 6.0849 \end{aligned}$$

*Poznámka:* Na rozdíl od předchozího příkladu ?? jsme nyní do vzorce vážených průměrů výběrových rozptylů vkládali hodnoty výběrových směrodatných odchylek. Proto jsme ve všech vzorcích tyto hodnoty umocňovali na

druhou.

Nakonec dosazením rozsahů náhodných výběrů a výběrových rozptylů všech tří populací do vzorce ?? vypočítáme hodnotu váženého průměru výběrových rozptylů pro populace ze všech tří pohřebišť Hawikuh, Pueblo Bonito a Puye.

$$\begin{aligned} s_*^2 &= \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2 + (n_3 - 1)s_3^2}{n_1 + n_2 + n_3 - 3} \\ &= \frac{(19 - 1)2.15^2 + (4 - 1)1.83^2 + (7 - 1)2.73^2}{19 + 4 + 7 - 3} \\ &= \frac{18 \times 4.6225 + 3 \times 3.3489 + 6 \times 7.4529}{27} \\ &= \frac{137.9691}{27} = 5.109967 \doteq 5.1100 \end{aligned}$$

Výpočet vážených průměrů výběrových rozptylů provedeme také pomocí softwaru . Hodnoty rozsahů náhodných výběrů a směrodatných odchylek všech tří populací vhojíme do proměnných n.H, n.PB, n.P, s.H, s.PB a s.P a následně přepíšeme vzorec ?? vypočítáme vážené průměry výběrových rozptylů.

```
63 n.H <- 19
64 n.PB <- 4
65 n.P <- 7
66
67 s.H <- 2.15
68 s.PB <- 1.83
69 s.P <- 2.73
70
71 sh.HPB <- ((n.H - 1) * s.H ^ 2 + (n.PB - 1) * s.PB ^ 2) / (n.H + n.PB - 2) # 4.440557
72 sh.HP <- ((n.H - 1) * s.H ^ 2 + (n.P - 1) * s.P ^ 2) / (n.H + n.P - 2) # 5.3301
73 sh.PBP <- ((n.PB - 1) * s.PB ^ 2 + (n.P - 1) * s.P ^ 2) / (n.PB + n.P - 2) # 6.0849
74 sh.HPBP <- ((n.H - 1) * s.H ^ 2 + (n.PB - 1) * s.PB ^ 2 + (n.P - 1) * s.P ^ 2) /
75 (n.H + n.PB + n.P - 3) # 5.109967
```

**Interpretace výsledků:** Vážený průměr výběrových rozptylů acetabulární výšky z pravé strany pro skelety mužského pohlaví pro populace z pohřebišť Hawikuh a Pueblo Bonito  $s_{*12}^2 = 4.44056\text{m}^2$ , pro populace z pohřebišť Hawikuh a Puye  $s_{*13}^2 = 5.3301\text{m}^2$  a pro populace z pohřebišť Pueblo Bonito a Puye  $s_{*23}^2 = 6.0849\text{m}^2$ . Vážený průměr výběrových rozptylů populací z pohřebišť Hawikuh, Pueblo Bonito a Puye  $s_*^2 = 5.10997\text{mm}^2$



## 6.2 Bodové odhady parametrů

Předpokládejme nyní, že náhodný výběr  $X_1, \dots, X_n$  se řídí nějakým modelem  $L$  s parametrem  $\theta$ , tj.  $X_1, \dots, X_n \sim L(\theta)$ . Skutečnou hodnotu parametru  $\theta$  neznáme a bohužel ji nikdy znát nebudeme. Jde o teoretickou hodnotu, kterou není možné přesně stanovit. Hodnotu parametru  $\theta$  můžeme ale na základě datového souboru alespoň odhadnout, přičemž můžeme stanovit buď bodový nebo intervalový odhad parametru  $\theta$ . Bodovým odhadem parametru  $\theta$  je statistika  $T = T(X_1, \dots, X_n)$ , která nabývá hodnot blízkých hodnotě parametru  $\theta$ , ať je hodnota tohoto parametru jakákoli. Neformálně vzato, bodovým odhadem parametru  $\theta$  je jedno konkrétní číslo, které získáme jako realizaci nějaké statistiky.

V praxi se můžeme setkat s různými typy bodových odhadů. Nejlepším odhadem je tzv. *nestranný* bodový odhad parametru  $\theta$ . Tento odhad skutečnou hodnotu parametru  $\theta$  ani nepodhodnocuje, ani nenadhodnocuje a proto je nejlepším možným typem bodového odhadu. Opakem nestranného odhadu je *vychýlený* odhad. Takový odhad skutečnou hodnotu parametru  $\theta$  buď systematicky podhodnocuje, nebo systematicky nadhodnocuje. Třetím typem odhadu je tzv. *asymptoticky nestranný* odhad. Asymptoticky nestranný odhad parametru  $\theta$ , stanovený na základě náhodného výběru s malým rozsahem  $n$ , je vychýlený, ale s rostoucím rozsahem náhodného výběru  $n$  jeho vychýlení klesá. Čím je tedy rozsah náhodného výběru použitého ke stanovení asymptotického odhadu parametru  $\theta$  větší, tím více se stanovený odhad blíží k nestrannému odhadu.

### Příklad 6.10. Bodový odhad parametru $\mu$ a parametru $\sigma^2$

Předpokládejme nyní, že  $X_1, \dots, X_n$ ,  $n \geq 2$  je náhodný výběr řídicí se modelem  $L$  se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $\mu$  a  $\sigma^2$  jsou parametry rozdělení  $L$ , jejichž přesnou hodnotu nebudeme nikdy znát. Nechť dále  $M$  je výběrový průměr a  $S^2$  je výběrový rozptyl, tj.  $M$  a  $S^2$  jsou statistiky vypočítané na základě náhodného výběru  $X_1, \dots, X_n$ . Potom výběrový průměr  $M$  je nestranným odhadem parametru  $\mu$  a výběrový rozptyl je nestranným odhadem parametru  $\sigma^2$ . ★

### Příklad 6.11. Bodový odhad parametru $\sigma_{12}$ a parametru $\rho$

Nechť  $(X_1, Y_1), \dots, (X_n, Y_n)$  je náhodný výběr řídicí se dvourozměrným modelem  $L_2$  s kovariancí  $\sigma_{12}$  a koeficientem korelace  $\rho$ , tj.  $\sigma_{12}$  a  $\rho$  jsou parametry rozdělení  $L_2$ , jejichž skutečnou hodnotu nebudeme nikdy znát. Nechť dále  $S_{12}$  je výběrové kovariance a  $R$  je výběrový korelační koeficient, tj.  $S_{12}$  a  $R$  jsou statistiky vypočítané na základě dvourozměrného náhodného výběru. Potom výběrové kovariance je nestranným odhadem parametru  $\sigma_{12}$ , zatímco výběrový korelační koeficient je asymptoticky nestranným odhadem parametru  $\rho$ . ★

### Příklad 6.12. Bodové odhady parametrů $\mu$ , $\sigma^2$ a $\sigma$ normálního rozdělení

Načtete datový soubor 01-one-sample-mean-skull-mf.txt a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *největší šířku mozkovny* u skeletů mužského pohlaví. Za předpokladu, že se náhodná veličina  $X$  řídí normálním modelem se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , stanovte nestranný bodový odhad (a) střední hodnoty  $\mu$ ; (b) rozptylu  $\sigma^2$ ; (c) směrodatné odchylky  $\sigma$ .

### Řešení příkladu ??

Celkem máme k dispozici  $n = 216$  náhodných veličin  $X_1, \dots, X_{216}$ , přičemž veličina  $X_1$  popisuje největší šířku mozkovny u prvního skeletu,  $\dots$ ,  $X_{216}$  popisuje největší šířku mozkovny u dvěstěšestnáctého skeletu. Předpokládáme, že všechny náhodné veličiny se řídí normálním modelem se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X_1 \sim N(\mu, \sigma^2), \dots, X_{216} \sim N(\mu, \sigma^2)$ . Protože se všechny náhodné veličiny řídí stejným modelem  $N(\mu, \sigma^2)$ , předpokládáme, že celý datový soubor se řídí tímž modelem  $N(\mu, \sigma^2)$ . Naměřením hodnoty největší šířky mozkovny každého skeletu jsme získali celkem 216 realizací náhodných veličin, konkrétně  $x_1 = 145, \dots, x_{216} = 137$ . Těchto 216 realizací tvoří společně datový soubor. Skutečnou hodnotu parametrů  $\mu$  a  $\sigma^2$  (resp.  $\sigma$ ) nebudeme nikdy znát. Jejich hodnoty ale můžeme odhadnout pomocí nestranných bodových odhadů. Bodový odhad parametru  $\mu$  stanovíme pomocí výběrového průměru, tj.

$$m = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{216} (124 + 127 + \dots + 149 + 149) = \frac{29\,632}{216} = 137.1852.$$

Bodový odhad parametru  $\sigma^2$  stanovíme pomocí výběrového rozptylu, tj.

$$\begin{aligned}
s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2 \\
&= \frac{1}{215} ((124 - 137.1852)^2 + (127 - 137.1852)^2 + \dots + (149 - 137.1852)^2 + (149 - 137.1852)^2) \\
&= \frac{1}{215} ((-13.1852)^2 + (-10.1852)^2 + \dots + 11.8148^2 + 11.8148^2) \\
&= 23.27717 \doteq 23.2772.
\end{aligned}$$

Konečně bodový odhad parametru  $\sigma$  stanovíme pomocí výběrové směrodatné odchylky, neboli jako odmocninu z výběrového rozptylu, tj.

$$s = \sqrt{s^2} = \sqrt{23.27717} = 4.824642 \doteq 4.8246.$$

Datový soubor načteme příkazem `read.delim()` a NA hodnoty odstraníme příkazem `na.omit()`. Pomocí operátoru `[]` vybereme z tabulky `data` pouze ty řádky, které se vztahují k mužským skeletům (`data$sex == 'm'`) a sloupec obsahující údaje o největší šířce mozkovny 'skull.B'. Hodnotu výběrového průměru, resp. výběrového rozptylu můžeme dopočítat pomocí softwaru  $\mathbb{R}$  přepisem vzorce `??`, resp. `??` s použitím funkce `sum()`. Hodnotu výběrové směrodatné odchylky získáme odmocněním výběrového rozptylu s využitím funkce `sqrt()`. Druhou možností je vypočítat výběrový průměr pomocí funkce `mean()`, výběrový rozptyl pomocí funkce `var()` a výběrovou směrodatnou odchylku pomocí funkce `sd()`.

```

76 data <- read.delim('01-one-sample-mean-skull-mf.txt')
77 data <- na.omit(data)
78 # head(data)
79 skull.BM <- data[data$sex == 'm', 'skull.B']
80 n <- length(skull.BM)
81 m.BM <- 1 / n * sum(skull.BM)
82 s2.BM <- 1 / (n - 1) * sum((skull.BM - m.BM) ^ 2)
83 s.BM <- sqrt(s2.BM)
84
85 mm.BM <- mean(skull.BM)
86 ss2.BM <- var(skull.BM)
87 ss.BM <- sd(skull.BM)
88
89 (tab <- data.frame(prumer = m.BM, rozptyl = s2.BM, sm.odch = s.BM))

```

	prumer	rozptyl	sm.odch
1	137.1852	23.27717	4.824642

90  
91

**Interpretace výsledků:** Nestranný odhad střední hodnoty největší šířky mozkovny pro skelety mužského pohlaví je 137.19 mm. Nestranný odhad rozptylu (resp. směrodatné odchylky) největší šířky mozkovny pro skelety mužského pohlaví je 23.28 mm<sup>2</sup> (resp. 4.82 mm). To znamená, že největší šířka mozkovny skeletů mužského pohlaví se pohybuje okolo hodnoty 137.19 mm se směrodatnou odchylkou 4.82 mm.

*Poznámka:* Všimněme si, že hodnota výběrového průměru vypočítaná v příkladu `??` je totožná s hodnotou aritmetického průměru vypočítanou v příkladu `??`. Rozdíl je však v přístupu k výsledné hodnotě. V příkladu `??` jsme aritmetický průměr uvažovali jako hodnotu vztahující se pouze k datovému souboru. V příkladu `??` již pracujeme s informací, že výběrový průměr je nestranným odhadem střední hodnoty  $\mu$  normálního rozdělení a tedy je možné ji brát jako výsledek relevantní pro celou populaci skeletů mužského pohlaví starověké egyptské populace.

Naopak srovnáme-li hodnotu výběrového rozptylu vypočítanou v příkladu `??` s hodnotou rozptylu vypočítanou v příkladu `??`, vidíme, že výsledky se mírně liší. Konkrétně hodnoty vypočítané v příkladu `??` jsou nepatrně vyšší než hodnoty vypočítané v příkladu `??`. Rozdíly v hodnotách jsou způsobeny použitím odlišných vzorců v obou příkladech. V příkladu `??` jsme k výpočtu rozptylu použili vzorec  $\frac{1}{n} \sum (X_i - M)^2$ , který má sice lepší interpretaci (jde o aritmetický průměr kvadrátů odchylek naměřených hodnot  $X_i$  od průměrné hodnoty  $M$ ), ale není nestranným odhadem parametru  $\sigma^2$ . Jde o odhad vychýlený, který systematicky skutečnou hodnotu parametru  $\sigma^2$  podhodnocuje.

Naopak v příkladu ?? jsme k výpočtu rozptylu použili vzorec  $\frac{1}{n-1} \sum (X_i - M)^2$ , který je nestranným odhadem parametru  $\sigma^2$ .

Analogicky vidíme, že hodnota výběrové směrodatné odchylky vypočítané v příkladu ?? je nepatrně vyšší než hodnota směrodatné odchylky vypočítaná v příkladu ?? . Směrodatná odchylka vypočítaná v příkladu ?? je opět vychýleným odhadem parametru  $\sigma$ , který skutečnou hodnotu parametru systematicky podhodnocuje. Naopak výběrová směrodatná odchylka vypočítaná v příkladu ?? je nestranným odhadem parametru  $\sigma$ . ★

### Příklad 6.13. Bodové odhady parametrů $\sigma_{12}$ a $\rho$ dvourozměrného normálního rozdělení

Načtete datový soubor 01-one-sample-mean-skull-mf.txt a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *největší šířku mozkovny* a náhodnou veličinu  $Y$  popisující *největší délku mozkovny* u skeletu mužského pohlaví. Za předpokladu, že se náhodný vektor  $(X, Y)^T$  řídí dvourozměrným normálním modelem, tj.  $(X, Y)^T \sim N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , kde  $\boldsymbol{\mu}$  je vektor středních hodnot a  $\boldsymbol{\Sigma}$  je varianční matice, stanovte (a) nestranný bodový odhad kovariance  $\sigma_{12}$ ; (b) asymptoticky nestranný bodový odhad korelačního koeficientu  $\rho$ .

#### Řešení příkladu ??

Celkem máme k dispozici  $n = 216$  dvojic náhodných veličin  $(X_1, Y_1), \dots, (X_{216}, Y_{216})$ , přičemž veličina  $X_1$  popisuje největší šířku mozkovny u prvního skeletu,  $\dots$ ,  $X_{216}$  popisuje největší šířku mozkovny u dvěstěšestnáctého skeletu a veličina  $Y_1$  popisuje největší délku mozkovny u prvního skeletu,  $\dots$ ,  $Y_{216}$  popisuje největší délku mozkovny u dvěstěšestnáctého skeletu. Předpokládáme, že všechny dvojice náhodných veličin se řídí dvourozměrným normálním modelem s vektorem středních hodnot  $\boldsymbol{\mu}$  a varianční maticí  $\boldsymbol{\Sigma}$ , tj.  $(X_1, Y_1) \sim N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \dots, (X_{216}, Y_{216}) \sim N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Protože se všechny dvojice náhodných veličin řídí stejným modelem  $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , předpokládáme, že celý datový soubor se řídí tímž modelem  $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Naměřením hodnot největší šířky a největší délky mozkovny každého skeletu jsme získali celkem 216 dvojic realizací náhodných veličin, konkrétně  $(x_1, y_1) = (145, 188), \dots, (x_{216}, y_{216}) = (137, 186)$ . Těchto 216 dvojic realizací tvoří společně datový soubor. Skutečnou hodnotu parametrů  $\sigma_{12}$  a  $\rho$  nebudeme nikdy znát. Jejich hodnoty ale můžeme odhadnout pomocí bodových odhadů. Nestranný bodový odhad parametru  $\sigma_{12}$  stanovíme pomocí výběrové kovariance, tj.

$$s_{12} = \frac{1}{n-1} \sum_{i=1}^n (x_i - m_1)(y_i - m_2),$$

kde  $m_1 = 137.1851$  je výběrový průměr největší šířky mozkovny (viz příklad ??) a  $m_2 = 182.0324$  je výběrový průměr největší délky mozkovny. Hodnotu výběrového průměru  $m_2$  získáme analogickým postupem uvedeným v příkladu ?? . Výběrovou kovarianci potom dopočítáme jako


$$\begin{aligned} s_{12} &= \frac{1}{215} ((145 - 137.1851)(188 - 182.0324) + (139 - 137.1851)(172 - 182.0324) + \dots \\ &\quad \dots + (142 - 137.1851)(183 - 182.0324) + (137 - 137.1851)(186 - 182.0324)) \\ &= \frac{1}{215} (46.6356 - 18.2068 + \dots + 4.6588 - 0.7348) \\ &= \frac{1113.7037}{215} = 5.1800172265 \doteq 5.1800. \end{aligned}$$

Asymptoticky nestranný bodový odhad parametru  $\rho$  stanovíme pomocí výběrového korelačního koeficientu, tj.

$$r_{12} = \frac{s_{12}}{\sqrt{s_1^2 s_2^2}} = \frac{s_{12}}{s_1 s_2}, \quad (6.10)$$

kde  $s_{12}$  je výběrové kovariance (viz výše),  $s_1^2 = 23.2772$  je výběrový rozptyl největší šířky mozkovny (viz příklad ??) a  $s_2^2 = 40.7664$  je výběrový rozptyl největší délky mozkovny. Hodnotu výběrového rozptylu  $s_2^2$  získáme analogickým postupem uvedeným v příkladu ?? . Výběrový korelační koeficient potom dopočítáme jako

$$r_{12} = \frac{5.1800}{\sqrt{23.2772} \sqrt{40.7664}} = \frac{5.1800}{4.8246 \times 6.3849} = \frac{5.1800}{30.8046} = 0.1682.$$

Hodnotu výběrové kovariance, resp. výběrového korelačního koeficientu můžeme dopočítat pomocí softwaru  přepisem vzorce ??, resp. ?? s použitím funkcí `sum()` a `sqrt()`. Druhou možností je vypočítat výběrovou kovarianci pomocí funkce `cov()` a výběrový korelační koeficient pomocí funkce `cor()`.

	kovariance	korelacni_koeficient
1	5.180017	0.168157

92  
93

**Interpretace výsledků:** Nestranný odhad kovariance největší šířky a délky mozkovny pro skelety mužského pohlaví je  $5.18 \text{ mm}^2$ . Asymptoticky nestranný odhad korelačního koeficientu největší šířky a délky mozkovny pro skelety mužského pohlaví je  $0.1682$ . To znamená, že mezi největší šířkou a délkou mozkovny u skeletů mužského pohlaví existuje nízký stupeň přímé lineární závislosti.

*Poznámka:* Všimněme si, že hodnota výběrového korelačního koeficientu vypočítaná v příkladu ?? je totožná s hodnotou korelačního koeficientu vypočítanou v příkladu ?? . Rozdíl je však v přístupu k výsledné hodnotě. V příkladu ?? jsme korelační koeficient uvažovali jako hodnotu vztahující se pouze k datovému souboru. V příkladu ?? již pracujeme s informací, že výběrový korelační koeficient je asymptoticky nestranným odhadem parametru  $\rho$  dvourozměrného normálního rozdělení a tedy je možné jej brát jako výsledek relevantní pro celou populaci skeletů mužského pohlaví starověké egyptské populace. Neměli bychom však zapomínat na to, že výběrový korelační koeficient je pouze asymptoticky nestranným odhadem parametru  $\rho$  a tedy jeho vychýlení klesá s rozsahem náhodného výběru. Rozsah náhodného výběru mužských skeletů,  $n = 216$ , je však dostatečně vysoký a tedy odhad parametru  $\rho$  můžeme považovat za nestranný. ★

#### Příklad 6.14. Bodový odhad vektoru středních hodnot $\mu$ a varianční matice $\Sigma$ dvourozměrného normálního rozdělení

Načtete datový soubor `01-one-sample-mean-skull-mf.txt` a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *největší šířku mozkovny* a náhodnou veličinu  $Y$  popisující *největší délku mozkovny* u skeletů mužského pohlaví. Za předpokladu, že se náhodný vektor  $(X, Y)^T$  řídí dvourozměrným normálním modelem, tj.  $(X, Y)^T \sim N_2(\mu, \Sigma)$ , kde  $\mu$  je vektor středních hodnot a  $\Sigma$  je varianční matice, stanovte (a) nestranný bodový odhad vektoru středních hodnot  $\mu$ ; (b) asymptoticky nestranný bodový odhad varianční matice  $\Sigma$ .

#### Řešení příkladu ??

Předpokládáme, že náhodný vektor  $(X, Y)^T$  se řídí dvourozměrným normálním modelem, tj.  $(X, Y)^T \sim N_2(\mu, \Sigma)$  s vektorem středních hodnot  $\mu = (\mu_1, \mu_2)^T$  a varianční maticí

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix},$$

kde  $\mu_1$  je střední hodnota náhodné veličiny  $X$ ,  $\mu_2$  je střední hodnota náhodné veličiny  $Y$ ,  $\sigma_1^2$  je rozptyl náhodné veličiny  $X$ ,  $\sigma_2^2$  je rozptyl náhodné veličiny  $Y$  a  $\rho$  je korelační koeficient popisující vztah mezi veličinami  $X$  a  $Y$ .

Z příkladů ?? a ?? známe výběrové průměry  $m_1 = 137.1851$  a  $m_2 = 182.0324$ , které jsou nestrannými odhady středních hodnot  $\mu_1$  a  $\mu_2$ , výběrové rozptyly  $s_1^2 = 23.2772$  a  $s_2^2 = 40.7664$ , které jsou nestrannými odhady rozptylů  $\sigma_1^2$  a  $\sigma_2^2$ , výběrové směrodatné odchylky  $s_1 = 4.8246$  a  $s_2 = 6.3842$ , které jsou nestrannými odhady odchylek  $\sigma_1$  a  $\sigma_2$  a výběrový korelační koeficient  $r_{12} = 0.1681$ , který je nestranným odhadem korelačního koeficientu  $\rho$ . Nestranným odhadem vektoru středních hodnot  $\mu$  je potom vektor  $(137.1851, 182.0324)^T$ . Asymptoticky nestranným odhadem varianční matice  $\Sigma$  je matice

$$\begin{pmatrix} 23.2772 & 0.1681 \times 4.8246 \times 6.3842 \\ 0.1681 \times 4.8246 \times 6.3842 & 40.7664 \end{pmatrix} = \begin{pmatrix} 23.2772 & 5.1800 \\ 5.1800 & 40.7664 \end{pmatrix}.$$

★

#### Příklad 6.15. Bodový odhad parametru $p$ alternativního modelu

Načtete datový soubor `17-anova-newborns.txt` a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *ženské pohlaví* novorozenců. Za předpokladu, že náhodná veličina  $X$  pochází z alternativního rozdělení s parametrem  $p$ , tj.  $X \sim Alt(p)$ , kde  $p$  je pravděpodobnost narození holčičky, stanovte bodový odhad parametru  $p$ .

## Řešení příkladu ??

Celkem máme k dispozici 1382 náhodných veličin  $X_1, \dots, X_{1382}$ , přičemž veličina  $X_1$  popisuje výskyt události narození holčičky ( $X_1 = 1$ ; úspěch), nebo výskyt události narození chlapečka ( $X_1 = 0$ ; neúspěch) u první matky,  $\dots$ ,  $X_{1382}$  popisuje výskyt události narození holčičky ( $X_{1382} = 1$ ; úspěch) nebo chlapečka ( $X_{1382} = 0$ ; neúspěch) u tisíci třísté osmdesáté druhé matky. Za předpokladu, že všechny náhodné veličiny se řídí alternativním modelem se stejným parametrem  $p$ , tj.  $X_1 \sim \text{Alt}(p)$ ,  $\dots$ ,  $X_{1382} \sim \text{Alt}(p)$ , se také celý náhodný výběr řídí tímž alternativním modelem, tj.  $X \sim \text{Alt}(p)$ . Parametr  $p$  určuje pravděpodobnost narození holčičky u jedné matky. Skutečnou hodnotu parametru  $p$  nebudeme nikdy znát, můžeme ji ale odhadnout pomocí nestranného bodového odhadu.

Vektor  $X$  obsahující údaje o pohlaví novorozenců je soubor 1 (narození holčičky) a 0 (narození chlapečka). Odhad parametru  $p$  získáme opět pomocí výběrového průměru, tj.

$$m = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{1382} (0 + 0 + \dots + 1 + 0) = \frac{663}{1382} = 0.4797. \quad (6.11)$$

Všimněme si, že odhad parametru  $p$  není nic jiného, než celkový počet narozených holčiček (čitatel vzorce ??) ku celkovému počtu všech novorozenců (jmenovatel vzorce ??), což je vlastně relativní četnost výskytu holčiček v datovém souboru. Výběrový průměr sestavený nad vektorem nul a jedniček je tedy roven relativní četnosti.

Nejprve načteme datový soubor příkazem `read.delim()` a odstraníme NA hodnoty příkazem `na.omit()`. Dále do proměnné `sex` vložíme údaje o pohlaví. Bližším prozkoumáním vektoru `sex` zjistíme, že jde o proměnnou typu `factor`, která nabývá dvou úrovní, a sice úrovně 1 (s popiskem 'f' (female)) a úrovně 2 (s popiskem 'm' (male)).

```
94 data <- read.delim('17-anova-newborns.txt')
95 data <- na.omit(data)
96 sex <- data$sex.C
97 head(sex)
```

```
[1] m m f m m m
Levels: f m
```

98  
99

Protože chceme vektor `sex` použít k odhadu parametru  $p$ , upravíme si jej nejprve do vhodné číselné podoby. Pomocí funkce `as.numeric()` převedeme faktor na číselný vektor a vložíme jej do proměnné `pohlavi`. Vidíme, že vektor `pohlavi` si zachoval původní kódování 1 = female, 2 = male. Ve vektoru `pohlavi` tedy změníme všechny hodnoty 2 na hodnoty 0, čímž dostaneme požadované kódování 0 = male, 1 = female.

```
100 pohlavi <- as.numeric(sex)
101 pohlavi[pohlavi == 2] <- 0
102 head(pohlavi)
```

```
[1] 0 0 1 0 0 0
```

103

Odhad parametru  $p$  nyní získáme buď přepisem vzorce ??, nebo funkcí `mean()`. Nakonec si ověříme že hodnota výběrového průměru je shodná s hodnotou relativní četnosti vypočítané pomocí původního faktoru `sex`.

```
104 N <- length(pohlavi)
105 m <- 1 / N * sum(pohlavi)
106 mm <- mean(pohlavi)
107 p <- sum(sex == 'f') / N
108 tab <- data.frame(m, mm, p)
109 round(tab, 4)
```

```
      m      mm      p
1 0.4797 0.4797 0.4797
```

110  
111

**Interpretace výsledků:** Nestranný odhad pravděpodobnosti narození holčičky je 0.4797. To znamená, že k narození holčičky u jedné matky dojde s pravděpodobností 47.97% ★

### 6.3 Intervalové odhady parametru

Hodnotu parametru  $\theta$  modelu  $L$ , ze kterého pochází náhodný výběr  $X_1, \dots, X_n$  zkusíme nyní odhadnout pomocí tzv. intervalového odhadu. Zatímco bodový odhad parametru  $\theta$  je jedno číslo (vypočítané na základě vhodné statistiky), intervalový odhad parametru  $\theta$  je interval  $(D, H)$ , který s dostatečně velkou pravděpodobností pokrývá hodnotu parametru  $\theta$ . Hranice intervalového odhadu tvoří opět vhodné statistiky, neboli funkce náhodného výběru, tj.  $D = D(x_1, \dots, X_n)$  a  $H = H(X_1, \dots, X_n)$ . Intervalový odhad nazýváme ve statistické terminologii jako interval spolehlivosti. Všechny zde prezentované intervaly spolehlivosti jsou intervaly spolehlivosti Waldova typu, nazývané zkráceně Waldovy intervaly spolehlivosti.

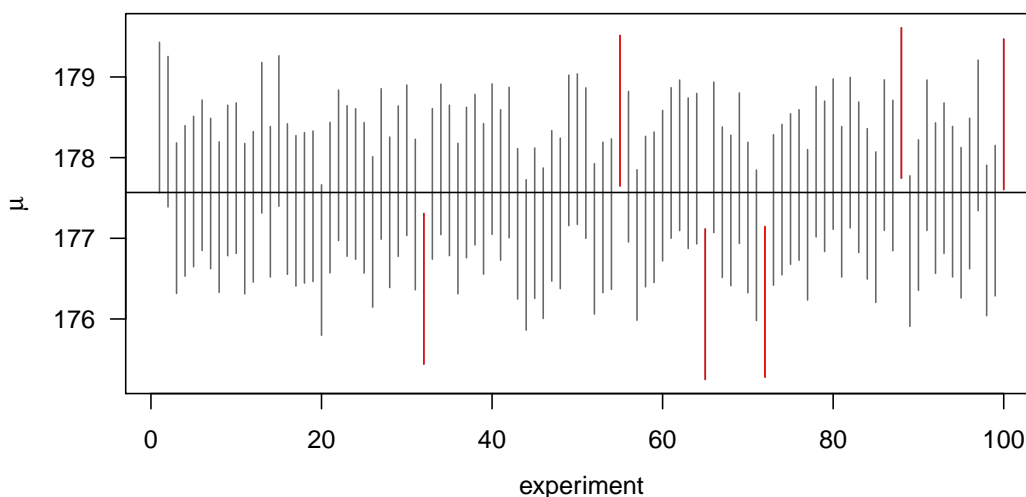
Mějme nyní riziko  $\alpha$ , což je koeficient nabývající hodnoty z intervalu  $(0, 1)$ . Tento koeficient určuje pravděpodobnost, s jakou interval spolehlivosti nepokrývá hodnotu parametru  $\theta$ . Doplnkem k riziku  $\alpha$  je tzv. koeficient spolehlivosti  $(1 - \alpha)$  určující pravděpodobnost, s jakou interval spolehlivosti pokrývá hodnotu parametru  $\theta$ . Podle potřeby volíme nejčastěji hodnotu rizika  $\alpha = 0.1$  (koeficient spolehlivosti  $1 - \alpha = 0.90$ , tj. pravděpodobnost, že interval spolehlivosti pokrývá hodnotu parametru  $\theta$  je 90%),  $\alpha = 0.05$  (koeficient spolehlivosti  $1 - \alpha = 0.95$ , tj. pravděpodobnost, že interval spolehlivosti pokrývá hodnotu parametru  $\theta$  je 95%) nebo  $\alpha = 0.01$  (koeficient spolehlivosti  $1 - \alpha = 0.99$ , tj. pravděpodobnost, že interval spolehlivosti pokrývá hodnotu parametru  $\theta$  je 99%).

$(1 - \alpha)\%$  pravděpodobnost, že interval spolehlivosti pokrývá hodnotu parametru  $\theta$  chápeme ve smyslu, že kdybychom nasbírali  $n$  náhodných výběrů a na základě každého z nich vypočítali intervalový odhad zkoumaného parametru  $\theta$ , potom alespoň v  $(1 - \alpha) \times n$  případech by vypočítaný interval spolehlivosti pokrýval (obsahoval) skutečnou hodnotu parametru  $\theta$  a ve zbylých  $(\alpha \times n$  a méně) případech by interval spolehlivosti skutečnou hodnotu parametru  $\theta$  nepokrýval (neobsahoval).

#### Příklad 6.16. Pravděpodobnost pokrytí $100 \times (1 - \alpha)\%$ intervalu spolehlivosti

Předpokládejme, že náhodná veličina  $X$  popisující největší šířku lebky novověké egyptské mužské populace, se řídí normálním modelem se střední hodnotou  $\mu = 177.568$  a rozptylem  $\sigma^2 = 7.526^2$ , tj.  $X \sim N(177.568, 7.526^2)$ . Zde tedy na chvíli připustíme, že známe skutečnou hodnotu parametru  $\mu$  i skutečnou hodnotu parametru  $\sigma^2$ .

Představme si nyní, že jsme v rámci jednoho experimentu vybrali náhodný vzorek 250 mužů z novodobé egyptské populace a změřili největší šířku jejich lebky. Získali jsme náhodný výběr  $\mathbf{X}_1 = (X_{1,1}, \dots, X_{1,250})$ . Na základě náhodného výběru jsme stanovili 95% Waldův empirický interval spolehlivosti pro parametr  $\mu$  a následně jsme zkontrolovali, zda skutečná hodnota parametru  $\mu = 177.568$  náleží do vypočítaného intervalu spolehlivosti, nebo nikoli. Analogický experiment jsme následně zopakovali stokrát. V alespoň 95 případech ze 100 experimentů vypočítaný Waldův empirický interval spolehlivosti pokrývá skutečnou hodnotu parametru  $\mu = 177.568$ , zatímco v pěti a méně případech ze 100 experimentů vypočítaný interval spolehlivosti skutečnou hodnotu parametru  $\mu$  nepokrývá (viz obrázek ??). ★



Obrázek 1: Pravděpodobnost pokrytí 95% Waldova empirického intervalů spolehlivosti pro parametr  $\mu$  normálního rozdělení při známém rozptylu  $\sigma^2$

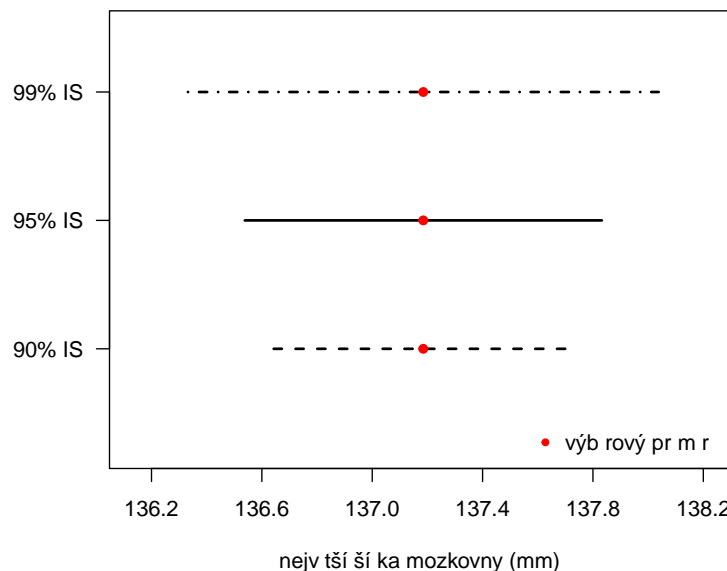


Naším cílem je nalézt takový interval spolehlivosti, který je jednak rozumně široký a který pokrývá skutečnou hodnotu parametru  $\theta$  s co největší pravděpodobností. Bohužel s rostoucí pravděpodobností pokrytí parametru  $\theta$  intervalem spolehlivosti roste také šířka tohoto intervalu. Proto volba výše pokrytí parametru  $\theta$  intervalem spolehlivosti je vždy otázkou kompromisu. Představme si, že bychom chtěli sestrojít 100% interval spolehlivosti, tedy interval, který pokrývá parametr  $\theta$  se 100% pravděpodobností. Takový interval existuje a má tvar  $(-\infty; \infty)$  (se 100% pravděpodobností bude skutečná hodnota parametru  $\theta$  nabývat hodnoty mezi  $-\infty$  a  $\infty$ ). Takový interval spolehlivosti nám však příliš platný není. Kdybychom se naopak rozhodli, že sestrojíme co nejpřesnější interval spolehlivosti, tj. interval, který bude mít co nejmenší šířku, sestrojili bychom 1% interval spolehlivosti. Šířka tohoto intervalu by byla tak malá, že by se intervalový odhad blížil bodovému odhadu. Ovšem pravděpodobnost, že skutečná hodnota parametru  $\theta$  náleží do toho úzkého intervalu by byla pouhé 1% (s 1% pravděpodobností skutečná hodnota parametru  $\theta$  náleží do intervalu spolehlivosti, ale s 99% pravděpodobností ne). Tento interval je tedy také ne příliš užitečný. Proto volíme hodnotu rizika  $\alpha$  jako 0.1, 0.05 nebo 0.01, protože odpovídající koeficient spolehlivosti  $1 - \alpha$  (0.9, 0.95 nebo 0.99) zajišťuje přijatelnou šířku intervalu spolehlivosti při zachování velmi vysoké pravděpodobnosti pokrytí parametru  $\theta$ .

Jak je uvedeno výše, šířka intervalu spolehlivosti roste s rostoucím koeficientem spolehlivosti  $1 - \alpha$ , neboli s rostoucí pravděpodobností pokrytí parametru  $\theta$ . Porovnáme-li tedy navzájem 90%, 95% a 99% interval spolehlivosti, bude šířka 90% intervalu spolehlivosti menší než šířka 95% intervalu spolehlivosti a ta bude menší než šířka 99% intervalu spolehlivosti.

### Příklad 6.17. Porovnání šířky Waldových empirických intervalů spolehlivosti

Načtete datový soubor 01-one-sample-mean-skull-mf.txt a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *největší šířku mozkovny* u skeletů mužského pohlaví. Za předpokladu, že náhodná veličina  $X$  pochází z normálního rozdělení se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , vypočítejte 90%, 95% a 99% Waldův empirický interval spolehlivosti pro parametr  $\mu$ . Následně porovnejte šířky těchto intervalů.



Obrázek 2: Porovnání šířky 90%, 95% a 99% Waldova empirického intervalu spolehlivosti pro parametr  $\mu$  normálního modelu

*Poznámka:* Přesný postup výpočtu intervalů spolehlivosti si názorně ukážeme později, v příkladu ??.

Z obrázku ?? vidíme, že skutečně nejméně široký je 90% Waldův empirický interval spolehlivosti. Naopak nejširší je 99% Waldův empirický interval spolehlivosti. 95% Waldův empirický interval spolehlivosti má šířku větší než 90% interval spolehlivosti, ale menší než 99% interval spolehlivosti.



Rozlišujeme tři základní typy intervalů spolehlivosti, a sice  $100 \times (1 - \alpha)\%$  oboustranný interval spolehlivosti  $(D, H)$

pro parametr  $\theta$ ,  $100 \times (1 - \alpha)\%$  levostranný interval spolehlivosti  $(D, \infty)$  pro parametr  $\theta$  a  $100 \times (1 - \alpha)\%$  pravostranný interval spolehlivosti  $(-\infty, H)$  pro parametr  $\theta$ . Ve všech třech případech je pravděpodobnost, že parametr  $\theta$  náleží do intervalu spolehlivosti, alespoň  $(1 - \alpha) \times 100\%$ , tj.  $\Pr(\theta \in IS) \geq (1 - \alpha) \times 100\%$ . V rámci tohoto textu se zaměříme na konstrukci (oboustranného / levostranného / pravostranného) intervalu spolehlivosti pro parametry  $\mu$  a  $\sigma^2$  normálního modelu a parametru  $p$  alternativního modelu. Pro parametr  $\rho$  dvourozměrného normálního modelu existuje několik typů intervalů spolehlivosti. Jejich konstrukcí se budeme zabývat v sekcích ?? a ??.

Nyní si představíme konkrétní tvary jednotlivých intervalů spolehlivosti. Na níže uvedené vzorce se potom budeme odkazovat v řešených příkladech.

Předpokládejme nejprve, že náhodný výběr  $X_1, \dots, X_n$  je náhodný výběr, který se řídí normálním modelem se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , kde parametr  $\sigma^2$  známe. Necht  $m$  značí realizaci výběrového průměru,  $\sigma$  směrodatnou odchylku vypočítanou jako  $\sigma = \sqrt{\sigma^2}$ ,  $n$  je rozsah náhodného výběru a  $u_\alpha$  (resp.  $u_{\alpha/2}$ ,  $u_{1-\alpha/2}$ ,  $u_{1-\alpha}$ ) je  $\alpha$ -kvantil (resp.  $\frac{\alpha}{2}$ -kvantil,  $(1 - \frac{\alpha}{2})$ -kvantil,  $(1 - \alpha)$ -kvantil) standardizovaného normálního modelu.  $100 \times (1 - \alpha)\%$  oboustranný Waldův empirický interval spolehlivosti pro parametr  $\mu$ , když  $\sigma^2$  známe, má tvar


$$(d, h) = (m - \frac{\sigma}{\sqrt{n}}u_{1-\alpha/2}, m - \frac{\sigma}{\sqrt{n}}u_{\alpha/2}). \quad (6.12)$$

$100 \times (1 - \alpha)\%$  levostranný Waldův empirický interval spolehlivosti pro parametr  $\mu$ , když  $\sigma^2$  známe, má tvar

$$(d, \infty) = (m - \frac{\sigma}{\sqrt{n}}u_{1-\alpha}, \infty). \quad (6.13)$$

$100 \times (1 - \alpha)\%$  pravostranný Waldův empirický interval spolehlivosti pro parametr  $\mu$ , když  $\sigma^2$  známe, má tvar

$$(-\infty, h) = (-\infty, m - \frac{\sigma}{\sqrt{n}}u_\alpha). \quad (6.14)$$

Hodnotu  $\alpha$ -kvantilu standardizovaného normálního rozdělení  $u_\alpha$  vypočítáme pomocí softwaru  příkazem `qnorm(alpha)`. Analogicky můžeme získat hodnoty kvantilů  $u_{\alpha/2}$ ,  $u_{1-\alpha}$  a  $u_{1-\alpha/2}$ .

*Poznámka:* V praxi se se situací, kdy odhadujeme parametr střední hodnoty  $\mu$  a přitom známe skutečnou hodnotu rozptylu  $\sigma^2$ , spíše neseťkáme. Jak jsme si uvedli, skutečná hodnota parametru  $\sigma^2$  nám není známá. Intervalové odhady ??, ?? a ?? se tedy spíše využívají při simulačních studiích, v rámci kterých zkoumáme různé vlastnosti těchto odhadů (např. změnu polohy intervalu spolehlivosti při měnící se hodnotě parametru  $\mu$ , nebo změnu šířky intervalu spolehlivosti při měnící se hodnotě parametru  $\sigma^2$ , rozsahu náhodného výběru  $n$  nebo koeficientu spolehlivosti  $1 - \alpha$ , apod.). Dále budeme tyto intervaly spolehlivosti využívat v kapitole ?? zabývající se testováním hypotéz. Zde budeme interval spolehlivosti pro parametr  $\mu$  když  $\sigma^2$  známe používat v případech, kdy budeme porovnávat střední hodnotu našeho náhodného výběru se střední hodnotou publikovanou, společně se svým rozptylem, v literatuře.

Předpokládejme nyní, že náhodný výběr  $X_1, \dots, X_n$  je náhodný výběr, který se řídí normálním modelem se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , kde parametr  $\sigma^2$  neznáme. Necht  $m$  značí realizaci výběrového průměru,  $s$  značí realizaci výběrové směrodatné odchylky,  $n$  je rozsah náhodného výběru a  $t_{n-1}(\alpha)$  (resp.  $t_{n-1}(\alpha/2)$ ,  $t_{n-1}(1 - \alpha/2)$ ,  $t_{n-1}(1 - \alpha)$ ) je  $\alpha$ -kvantil (resp.  $\frac{\alpha}{2}$ -kvantil,  $(1 - \frac{\alpha}{2})$ -kvantil,  $(1 - \alpha)$ -kvantil) Studentova modelu o  $n - 1$  stupních volnosti.  $100 \times (1 - \alpha)\%$  oboustranný Waldův empirický interval spolehlivosti pro parametr  $\mu$ , když  $\sigma^2$  neznáme, má tvar

$$(d, h) = (m - \frac{s}{\sqrt{n}}t_{n-1}(1 - \alpha/2), m - \frac{s}{\sqrt{n}}t_{n-1}(\alpha/2)). \quad (6.15)$$

$100 \times (1 - \alpha)\%$  levostranný Waldův empirický interval spolehlivosti pro parametr  $\mu$ , když  $\sigma^2$  neznáme, má tvar

$$(d, \infty) = (m - \frac{s}{\sqrt{n}}t_{n-1}(1 - \alpha), \infty). \quad (6.16)$$

$100 \times (1 - \alpha)\%$  pravostranný Waldův empirický interval spolehlivosti pro parametr  $\mu$ , když  $\sigma^2$  neznáme, má tvar

$$(-\infty, h) = (-\infty, m - \frac{s}{\sqrt{n}}t_{n-1}(\alpha)). \quad (6.17)$$

Hodnotu  $\alpha$ -kvantilu Studentova modelu s  $n - 1$  stupni volnosti  $t_{n-1}(\alpha)$  vypočítáme pomocí softwaru  $\mathbb{R}$  příkazem `qt(alpha, n-1)`. Analogicky můžeme získat hodnoty kvantilů  $t_{n-1}(\alpha/2)$ ,  $t_{n-1}(1 - \alpha)$  a  $t_{n-1}(1 - \alpha/2)$ .

Předpokládejme dále, že  $X_1, \dots, X_n$  je náhodný výběr, který se řídí normálním modelem se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , kde parametr  $\mu$  neznáme. Nechť  $s^2$  značí realizaci výběrového rozptylu,  $n$  je rozsah náhodného výběru a  $\chi_{n-1}^2(\alpha)$  (resp.  $\chi_{n-1}^2(\alpha/2)$ ,  $\chi_{n-1}^2(1 - \alpha/2)$ ,  $t\chi_{n-1}^2(1 - \alpha)$ ) je  $\alpha$ -kvantil (resp.  $\frac{\alpha}{2}$ -kvantil,  $(1 - \frac{\alpha}{2})$ -kvantil,  $(1 - \alpha)$ -kvantil)  $\chi^2$  modelu o  $n - 1$  stupních volnosti.  $100 \times (1 - \alpha)\%$  oboustranný Waldův empirický interval spolehlivosti pro parametr  $\sigma^2$ , když  $\mu$  neznáme, má tvar

$$(d, h) = \left( \frac{(n-1)s^2}{\chi_{n-1}^2(1-\alpha/2)}, \frac{(n-1)s^2}{\chi_{n-1}^2(\alpha/2)} \right). \quad (6.18)$$

$100 \times (1 - \alpha)\%$  levostranný Waldův empirický interval spolehlivosti pro parametr  $\sigma^2$ , když  $\mu$  neznáme, má tvar

$$(d, \infty) = \left( \frac{(n-1)s^2}{\chi_{n-1}^2(1-\alpha)}, \infty \right). \quad (6.19)$$

$100 \times (1 - \alpha)\%$  pravostranný Waldův empirický interval spolehlivosti pro parametr  $\sigma^2$ , když  $\mu$  neznáme, má tvar

$$(0, h) = \left( 0, \frac{(n-1)s^2}{\chi_{n-1}^2(\alpha)} \right). \quad (6.20)$$

Hodnotu  $\alpha$ -kvantilu  $\chi^2$  modelu s  $n - 1$  stupni volnosti  $\chi_{n-1}^2(\alpha)$  vypočítáme pomocí softwaru  $\mathbb{R}$  příkazem `qchisq(alpha, n-1)`. Analogicky můžeme získat hodnoty kvantilů  $\chi_{n-1}^2(\alpha/2)$ ,  $\chi_{n-1}^2(1 - \alpha)$  a  $\chi_{n-1}^2(1 - \alpha/2)$ .

Protože parametr rozptylu  $\sigma^2$  je z definice vždy větší než 0, stanovuje se hodnota dolní hranice  $100 \times (1 - \alpha)\%$  pravostranného Waldova empirického intervalu spolehlivosti jako nula namísto nekonečna.

Konečně předpokládejme, že  $X_1, \dots, X_n$  je náhodný výběr, který se řídí alternativním modelem s parametrem  $p$ , tj.  $X \sim \text{Alt}(p)$ . Nechť  $m$  značí realizaci výběrového průměru,  $N$  je rozsah náhodného výběru a  $u_\alpha$  (resp.  $u_{\alpha/2}$ ,  $u_{1-\alpha/2}$ ,  $u_{1-\alpha}$ ) je  $\alpha$ -kvantil (resp.  $\frac{\alpha}{2}$ -kvantil,  $(1 - \frac{\alpha}{2})$ -kvantil,  $(1 - \alpha)$ -kvantil) standardizovaného normálního modelu.  $100 \times (1 - \alpha)\%$  oboustranný Waldův empirický interval spolehlivosti pro parametr  $p$  má tvar

$$(d, h) = \left( m - \sqrt{\frac{m(1-m)}{N}} u_{1-\alpha/2}, m - \sqrt{\frac{m(1-m)}{N}} u_{\alpha/2} \right). \quad (6.21)$$

$100 \times (1 - \alpha)\%$  levostranný Waldův empirický interval spolehlivosti pro parametr  $p$  má tvar

$$(d, 1) = \left( m - \sqrt{\frac{m(1-m)}{N}} u_{1-\alpha}, 1 \right). \quad (6.22)$$

$100 \times (1 - \alpha)\%$  pravostranný Waldův empirický interval spolehlivosti pro parametr  $p$  má tvar

$$(0, h) = \left( 0, m - \sqrt{\frac{m(1-m)}{N}} u_\alpha \right). \quad (6.23)$$

Hodnotu  $\alpha$ -kvantilu standardizovaného normálního rozdělení  $u_\alpha$  vypočítáme pomocí softwaru  $\mathbb{R}$  příkazem `qnorm(alpha)`. Analogicky můžeme získat hodnoty kvantilů  $u_{\alpha/2}$ ,  $u_{1-\alpha}$  a  $u_{1-\alpha/2}$ .

Protože parametr  $p$  značí pravděpodobnost úspěchu v jednom pokusu, platí, že  $p \in (0, 1)$ , a tedy horní hranice levostranného Waldova empirického intervalu spolehlivosti je 1 (viz vzorec ??). Analogicky dolní hranice pravostranného Waldova empirického intervalu spolehlivosti je 0 (viz vzorec ??).

### Příklad 6.18. Intervalový odhad parametru $\mu$ normálního modelu

Načtěte datový soubor `01-one-sample-mean-skull-mf.txt` a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *největší šířku mozkovny* u skeletů mužského pohlaví. Za předpokladu, že náhodná veličina  $X$

pochází z normálního rozdělení se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , vypočítejte 90%, 95% a 99% Waldův empirický interval spolehlivosti pro parametr  $\mu$ .

### Řešení příkladu ??

Není-li v textu specifikován typ požadovaného intervalu spolehlivosti, počítáme vždy oboustranný interval spolehlivosti. Ze zadání víme, že chceme spočítat interval spolehlivosti pro parametr  $\mu$ . Dále si všimneme, že ze zadání příkladu není známá skutečná hodnota rozptylu  $\sigma^2$ . Naším úkolem je tedy vypočítat 90% (resp. 95%, či 99%) Waldův empirický oboustranný interval spolehlivosti pro parametr  $\mu$  když  $\sigma^2$  neznáme.

Při výpočtu intervalů spolehlivosti budeme vycházet ze vzorce ???. Z příkladu ?? víme, že realizace výběrového průměru  $m = 137.1852$  a rozsah náhodného výběru  $n = 216$ . Hodnotu směrodatné odchylky  $\sigma$  odhadneme pomocí výběrové směrodatné odchylky  $s = 4.8246$  (viz příklad ??). Zbývá stanovit hodnotu kvantilu  $t_{n-1}(\alpha/2)$  a kvantilu  $t_{n-1}(1 - \alpha/2)$  Studentova modelu. K tomu je potřeba nejprve dopočítat koeficient  $\alpha$ . Ten vyjádříme, v případě výpočtu 90% Waldova empirického intervalu spolehlivosti, postupnými kroky z rovnice  $100 \times (1 - \alpha) \% = 90 \%$ .

$$100 \times (1 - \alpha) \% = 90 \%$$

$$100 \times (1 - \alpha) = 90$$

$$1 - \alpha = 0.90$$

$$1 - 0.90 = \alpha$$

$$\alpha = 0.10$$

Pomocí softwaru  $\mathbb{R}$  a funkce `qt()` nyní stanovíme hodnotu kvantilu  $t_{n-1}(\alpha/2) = t_{215}(0.10/2) = t_{215}(0.05) = \text{qt}(0.05, 215) = -1.6520$  a kvantilu  $t_{n-1}(1 - \alpha/2) = t_{215}(1 - 0.10/2) = t_{215}(0.95) = \text{qt}(0.95, 215) = 1.6520$ . Nyní již známe všechny potřebné hodnoty a můžeme dosadit do vzorce ???.

$$\begin{aligned} (d, h) &= \left( m - \frac{s}{\sqrt{n}} t_{n-1}(1 - \alpha/2), m - \frac{s}{\sqrt{n}} t_{n-1}(\alpha/2) \right) \\ &= \left( 137.1852 - \frac{4.8246}{\sqrt{215}} 1.6520, 137.1852 - \frac{4.8246}{\sqrt{215}} (-1.6520) \right) \\ &= (137.1852 - 0.5436, 137.1852 - (-0.5436)) \\ &= (136.6416, 137.7288) \end{aligned}$$

V případě výpočtu 95% Waldova empirického intervalu spolehlivosti, vyjádříme koeficient  $\alpha$  z rovnice  $100 \times (1 - \alpha) \% = 95 \%$ .

$$100 \times (1 - \alpha) \% = 95 \%$$

$$100 \times (1 - \alpha) = 95$$

$$1 - \alpha = 0.95$$

$$1 - 0.95 = \alpha$$

$$\alpha = 0.05$$

Pomocí softwaru  $\mathbb{R}$  a funkce `qt()` stanovíme hodnotu kvantilu  $t_{n-1}(\alpha/2) = t_{215}(0.05/2) = t_{215}(0.025) = \text{qt}(0.025, 215) = -1.9711$  a kvantilu  $t_{n-1}(1 - \alpha/2) = t_{215}(1 - 0.05/2) = t_{215}(0.975) = \text{qt}(0.975, 215) = 1.9711$ . Hodnoty  $m$ ,  $s$  a  $n$  jsme stanovili výše, zbývá tedy dosadit do vzorce ???.

$$\begin{aligned} (d, h) &= \left( m - \frac{s}{\sqrt{n}} t_{n-1}(1 - \alpha/2), m - \frac{s}{\sqrt{n}} t_{n-1}(\alpha/2) \right) \\ &= \left( 137.1852 - \frac{4.8246}{\sqrt{215}} 1.9711, 137.1852 - \frac{4.8246}{\sqrt{215}} (-1.9711) \right) \\ &= (137.1852 - 0.6486, 137.1852 - (-0.6486)) \\ &= (136.5366, 137.8338) \end{aligned}$$

Konečně, v případě výpočtu 99% Waldova empirického intervalu spolehlivosti, vyjádříme koeficient  $\alpha$  z rovnice  $100 \times (1 - \alpha) \% = 99 \%$ .

$$\begin{aligned} 100 \times (1 - \alpha) \% &= 99 \% \\ 100 \times (1 - \alpha) &= 99 \\ 1 - \alpha &= 0.99 \\ 1 - 0.99 &= \alpha \\ \alpha &= 0.01 \end{aligned}$$

Pomocí softwaru  $\mathbb{R}$  stanovíme hodnotu kvantilu  $t_{n-1}(\alpha/2) = t_{215}(0.01/2) = t_{215}(0.005) = \text{qt}(0.005, 215) = -2.5989$  a kvantilu  $t_{n-1}(1 - \alpha/2) = t_{215}(1 - 0.01/2) = t_{215}(0.995) = \text{qt}(0.995, 215) = 2.5989$ . Hranice intervalu spolehlivosti dopočítáme analogicky jako v předchozích dvou případech dosazením do vzorce ??.

$$\begin{aligned} (d, h) &= \left( m - \frac{s}{\sqrt{n}} t_{n-1}(1 - \alpha/2), m - \frac{s}{\sqrt{n}} t_{n-1}(\alpha/2) \right) \\ &= \left( 137.1852 - \frac{4.8246}{\sqrt{215}} 2.5989, 137.1852 - \frac{4.8246}{\sqrt{215}} (-2.5989) \right) \\ &= (137.1852 - 0.8551, 137.1852 - (-0.8551)) \\ &= (136.3301, 138.0403) \end{aligned}$$

Datový soubor načteme příkazem `read.delim()` a NA hodnoty odstraníme příkazem `na.omit()`. Pomocí operátoru `[]` vybereme z tabulky `data` pouze ty řádky, které se vztahují k mužským skeletům (`data$sex == 'm'`) a sloupec obsahující údaje o největší šířce mozkovny 'skull.B'. Hodnotu výběrového průměru a výběrové směrodatné odchylky dopočítáme pomocí funkcí `mean()` a `sd()`, rozsah náhodného výběru stanovíme funkcí `length()`. Do proměnné `alpha` si vložíme všechny tři hodnoty koeficientů  $\alpha$ , tj. 0.1, 0.05 i 0.01, najednou. Nyní přepísem vzorce ??, kde funkci `qt()` využijeme na výpočet  $\alpha/2$ , resp.  $1 - \alpha/2$  kvantilů Studentova rozdělení, a to současně pro všechny tři koeficienty  $\alpha$ , získáme dolní, resp. horní hranice všech tří Waldových empirických intervalů spolehlivosti.

```
112 data <- read.delim('01-one-sample-mean-skull-mf.txt')
113 data <- na.omit(data)
114
115 skull.BM <- data[data$sex == 'm', 'skull.B']
116 m <- mean(skull.BM)
117 s <- sd(skull.BM)
118 n <- length(skull.BM)
119 alpha <- c(0.1, 0.05, 0.01)
120
121 dh <- m - s / sqrt(n) * qt(1 - alpha / 2, n - 1)
122 hh <- m - s / sqrt(n) * qt(alpha / 2, n - 1)
123
124 tab <- data.frame(d = dh, h = hh, row.names = c('90% DIS', '95% DIS', '99% DIS'))
125 round(tab, 4)
```

	d	h
90% DIS	136.6429	137.7275
95% DIS	136.5381	137.8322
99% DIS	136.3320	138.0383

126  
127  
128  
129

**Interpretace výsledků:** 90% Waldův empirický interval spolehlivosti pro parametr  $\mu$  má tvar (136.64, 137.73) mm. To znamená, že  $136.64 \text{ mm} < \mu < 137.73 \text{ mm}$  s pravděpodobností 90%. V 90 případech ze sta bude střední hodnota největší šířky mozkovny u skeletů mužského pohlaví nabývat hodnoty z intervalu (136.64, 137.73) mm.

95% Waldův empirický interval spolehlivosti pro parametr  $\mu$  má tvar (136.54, 137.83) mm. To znamená, že  $136.54 \text{ mm} < \mu < 137.83 \text{ mm}$  s pravděpodobností 95%. V 95 případech ze sta bude střední hodnota největší šířky mozkovny u skeletů mužského pohlaví nabývat hodnoty z intervalu (136.54, 137.83) mm.

99% Waldův empirický interval spolehlivosti pro parametr  $\mu$  má tvar (136.33, 138.04) mm. To znamená, že  $136.33 \text{ mm} < \mu < 138.04 \text{ mm}$  s pravděpodobností 99%. V 99 případech ze sta bude střední hodnota největší šířky mozkovny u skeletů mužského pohlaví nabývat hodnoty z intervalu (136.33, 138.04) mm. ★

### Příklad 6.19. Intervalový odhad parametru $\mu$ normálního modelu

Máme datový soubor 21-goldman-measures.csv obsahující údaje o délce vřetenní kosti na pravé straně (radius.LR) a na levé straně (radius.LL) u mužů (sex = 'm' a žen (sex = 'f') pohřbených v oblasti Indian Knoll v Kentucky (viz sekce ??). Mějme náhodnou veličinu  $X$  popisující *délku vřetenní kosti* u skeletů ženského pohlaví z oblasti Indian Knoll v Kentucky. Za předpokladu, že náhodná veličina  $X$  pochází z normálního rozdělení se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , vypočítejte (a) 95% oboustranný intervalový odhad; (b) 99% jednostranný intervalový odhad; (c) 90% jednostranný intervalový odhad pro parametr  $\mu$ .

### Řešení příkladu ??

Ze zadání víme, že chceme spočítat interval spolehlivosti pro parametr  $\mu$ . Dále si všimněme, že ze zadání příkladu není známá skutečná hodnota rozptylu  $\sigma^2$ . Naším úkolem je tedy vypočítat 95% (resp. 99%, či 90%) Waldův empirický oboustranný (resp. jednostranný či jednostranný) interval spolehlivosti pro parametr  $\mu$  když  $\sigma^2$  neznáme.

Při výpočtu intervalů spolehlivosti budeme vycházet ze vzorců ??, ?? a ?. Nejprve si tedy dopočítáme hodnoty výběrového průměru  $m$  a výběrové směrodatné odchylky  $s$  a následně vyjádříme hodnotu koeficientu  $\alpha$ .

Budeme tedy vycházet ze vzorce ?. Hodnotu parametru  $\sigma^2$  odhadneme pomocí výběrového rozptylu. Nejprve tedy spočítáme výběrový průměr a následně výběrový rozptyl, tj.


$$m = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{27} (223 + 231 + \dots + 218.5 + 221) = \frac{6042.5}{27} = 223.7963.$$

$$\begin{aligned} s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2 \\ &= \frac{1}{27-1} ((223 - 223.7963)^2 + (231 - 223.7963)^2 + \dots + (218.5 - 223.7963)^2 + (221 - 223.7963)^2) \\ &= \frac{1}{26} ((-0.7963)^2 + 7.2037^2 + \dots + (-5.2963)^2 + (-2.7963)^2) \\ &= \frac{2534.13}{26} = 97.46652 \doteq 97.4665 \end{aligned}$$

$$s = \sqrt{s^2} = \sqrt{97.46652} = 9.872513 \doteq 9.8725$$

Výběrový průměr  $m = 223.7963 \text{ mm}$ , výběrový rozptyl  $s^2 = 97.4665 \text{ mm}^2$ , výběrová směrodatná odchylka  $s = 9.8725 \text{ mm}$ , rozsah náhodného výběru  $n = 27$ . Zbývá stanovit hodnotu kvantilu  $t_{n-1}(\alpha/2)$  a kvantilu  $t_{n-1}(1 - \alpha/2)$  Studentova rozdělení. K tomu je potřeba nejprve dopočítat koeficient  $\alpha$ , a to vyjádřením z rovnice  $100 \times (1 - \alpha)\% = 95\%$ , analogicky jako v příkladu ??

$$\begin{aligned} 100 \times (1 - \alpha)\% &= 95\% \\ 100 \times (1 - \alpha) &= 95 \\ 1 - \alpha &= 0.95 \\ 1 - 0.95 &= \alpha \\ \alpha &= 0.05 \end{aligned}$$

Pomocí softwaru  a funkce qt() nyní stanovíme hodnotu kvantilu  $t_{n-1}(\alpha/2) = t_{27-1}(0.05/2) = t_{26}(0.025) = \text{qt}(0.025, 26) = -2.055529$  a kvantilu  $t_{n-1}(1 - \alpha/2) = t_{27-1}(1 - 0.05/2) = t_{26}(0.975) = \text{qt}(0.975, 26) = 2.055529$ . Nyní již známe všechny potřebné hodnoty a můžeme dosadit do vzorce ??

$$\begin{aligned}
(d, h) &= \left( m - \frac{s}{\sqrt{n}} t_{n-1}(1 - \alpha/2), m - \frac{s}{\sqrt{n}} t_{n-1}(\alpha/2) \right) \\
&= \left( 223.7963 - \frac{9.872514}{\sqrt{27}} 2.055529, 223.7963 - \frac{9.872514}{\sqrt{27}} (-2.055529) \right) \\
&= (223.7963 - 3.905436, 223.7963 - (-3.905436)) \\
&= (219.8909, 227.7017)
\end{aligned}$$

V případě výpočtu 99% Waldova empirického intervalu spolehlivosti, vyjádříme koeficient  $\alpha$  z rovnice  $100 \times (1 - \alpha) \% = 99 \%$ .

$$\begin{aligned}
100 \times (1 - \alpha) \% &= 99 \% \\
100 \times (1 - \alpha) &= 99 \\
1 - \alpha &= 0.99 \\
1 - 0.99 &= \alpha \\
\alpha &= 0.01
\end{aligned}$$

Pomocí softwaru  $\mathbb{R}$  a funkce `qt()` stanovíme hodnotu kvantilu  $t_{n-1}(1 - \alpha) = t_{27-1}(1 - 0.01) = t_{26}(0.99) = \text{qt}(0.99, 26) = 2.47863$ . Hodnoty  $m$ ,  $s$  a  $n$  jsme stanovili výše, zbývá tedy dosadit do vzorce ??.

$$\begin{aligned}
(d, \infty) &= \left( m - \frac{s}{\sqrt{n}} t_{n-1}(1 - \alpha), \infty \right) \\
&= \left( 223.7963 - \frac{9.872514}{\sqrt{26}} 2.47863, \infty \right) \\
&= (223.7963 - 4.799021, \infty) \\
&= (218.9973, \infty)
\end{aligned}$$

Konečně, v případě výpočtu 90% Waldova empirického intervalu spolehlivosti, vyjádříme koeficient  $\alpha$  z rovnice  $100 \times (1 - \alpha) \% = 90 \%$ .

$$\begin{aligned}
100 \times (1 - \alpha) \% &= 90 \% \\
100 \times (1 - \alpha) &= 90 \\
1 - \alpha &= 0.90 \\
1 - 0.90 &= \alpha \\
\alpha &= 0.10
\end{aligned}$$

Pomocí softwaru  $\mathbb{R}$  stanovíme hodnotu kvantilu  $t_{n-1}(\alpha) = t_{27-1}(0.10) = \text{qt}(0.10, 26) = -1.314972$ . Hranice intervalu spolehlivosti dopočítáme dosazením do vzorce ??.

$$\begin{aligned}
(-\infty, h) &= \left( -\infty, m - \frac{s}{\sqrt{n}} t_{n-1}(\alpha/2) \right) \\
&= \left( -\infty, 223.7963 - \frac{9.872514}{\sqrt{26}} (-1.314972) \right) \\
&= (-\infty, 223.7963 - (-2.545995)) \\
&= (-\infty, 226.3423)
\end{aligned}$$

Datový soubor načteme příkazem `read.delim()`. Z datové tabulky vybereme pouze délky vřetených kostí z pravé strany (`radius.LR`) skeletů ženského pohlaví (`sex == 'f'`) z oblasti Indian Knoll v Kentucky (`pop == 'Indian Knoll'`).

Chybějící hodnoty odstraníme příkazem `na.omit()`. Hodnotu výběrového průměru a výběrové směrodatné odchylky dopočítáme pomocí funkcí `mean()` a `sd()`, rozsah náhodného výběru stanovíme funkcí `length()`. Protože nejprve stanovujeme hranice 95% Waldova empirického oboustranného intervalu spolehlivosti, vložíme do proměnné `alpha` nejprve hodnotu 0.05. Dále přepíšeme vzorec ??, kde pomocí funkce `qt()` vypočítáme  $\alpha/2$ , resp.  $1 - \alpha/2$  kvantil Studentova rozdělení, získáme dolní, resp. horní hranici 95% Waldova empirického oboustranného intervalu spolehlivosti. Následně změníme hodnotu koeficientu `alpha` na hodnotu 0.01 a přepíšeme vzorec ?? dopočítáme dolní hranici 99% Waldova empirického jednostranného intervalu spolehlivosti. Nakonec předefinujeme hodnotu koeficientu `alpha` na hodnotu 0.10 a přepíšeme vzorec ?? spočítáme horní hranici 90% Waldova empirického jednostranného intervalu spolehlivosti.



```

130 data <- read.delim('00-Data//21-goldman-measures.csv', sep = ';', dec = '.')
131 radius.LRF <- data[data$sex == 'f' & data$pop == 'Indian Knoll', 'radius.LR']
132 radius.LRF <- as.numeric(na.omit(radius.LRF))
133
134 m <- mean(radius.LRF) # 223.7963
135 s <- sd(radius.LRF) # 9.872514
136 n <- length(radius.LRF) # 27
137 alpha <- 0.05
138 dh <- m - s / sqrt(n) * qt(1 - alpha / 2, n - 1) # 219.8909
139 hh <- m - s / sqrt(n) * qt(alpha / 2, n - 1) # 227.7017
140
141 alpha <- 0.01
142 DH <- m - s / sqrt(n) * qt(1 - alpha, n - 1) # 219.087
143
144 alpha <- 0.10
145 HH <- m - s / sqrt(n) * qt(alpha, n - 1) # 226.2947
146
147 (tab <- data.frame(d = c(round(c(dh, DH), 4), '-inf'),
148                   h = c(round(hh, 4), 'inf', round(HH, 4)),
149                   row.names = c('95% DIS', '99% LIS', '90% PIS')))

```

	d	h
95% DIS	219.8909	227.7017
99% LIS	219.087	inf
90% PIS	-inf	226.2947

150  
151  
152  
153

*Poznámka:* Horní hranice jednostranného Waldova empirického intervalu spolehlivosti je nekonečno. Abychom mohli tuto hodnotu zanesť do tabulky výsledků, musíme ji zadat jako textový řetězec. Datové tabulky v softwaru  však mají tu vlastnost, že hodnoty ve stejném sloupci musí být stejného typu, tedy buď jsou všechny hodnoty numerické, nebo jsou všechny hodnoty textové. Proto v okamžiku, kdy do sloupce `h` tabulky `tab` vložíme textový řetězec `inf`, software  všechny numerické hodnoty v tomto sloupci automaticky převede na textové řetězce a tak se k nim od tohoto okamžiku také chová. Protože text nejde zaokrouhlit, není potom možné tabulku `tab` zaokrouhlit na čtyři desetinná místa. Proto musíme všechny hodnoty v tabulce zaokrouhlit dříve, než je do tabulky vložíme společně s textovým řetězcem `inf`.

**Interpretace výsledků:** 95% Waldův empirický oboustranný interval spolehlivosti pro parametr  $\mu$  má tvar (219.89, 227.70) mm. To znamená, že  $219.89 \text{ mm} < \mu < 227.70 \text{ mm}$  s pravděpodobností 95 %. V 95 případech ze sta bude střední hodnota délky vřetenní kosti z pravé strany u skeletů ženského pohlaví z oblasti Indian Knoll v Kentucky nabývat hodnoty z intervalu (219.89, 227.70) mm.

99% Waldův empirický jednostranný interval spolehlivosti pro parametr  $\mu$  má tvar (219.087,  $\infty$ ) mm. To znamená, že  $\mu > 219.087 \text{ mm}$  s pravděpodobností 99 %. V 99 případech ze sta bude střední hodnota délky vřetenní kosti z pravé strany u skeletů ženského pohlaví populace Ipituaq nabývat hodnoty z oblasti Indian Knoll v Kentucky intervalu (219.087,  $\infty$ ) mm.

90% Waldův empirický jednostranný interval spolehlivosti pro parametr  $\mu$  má tvar ( $-\infty$ , 226.29) mm. To znamená, že  $226.29 \text{ mm} > \mu$  s pravděpodobností 90 %. V 90 případech ze sta bude střední hodnota délky vřetenní kosti z pravé strany u skeletů ženského pohlaví z oblasti Indian Knoll v Kentucky nabývat hodnoty z intervalu ( $-\infty$ , 226.29) mm.



*Poznámka:* Mírné odchylky ve výsledných hranicích intervalů spolehlivosti spočítaných ručním výpočtem a výpočtem pomocí softwaru  $\mathbb{R}$  jsou způsobeny průměrným zaokrouhlováním výsledků při ručním výpočtu. Přesnější jsou tedy výsledky získané přímým výpočtem pomocí softwaru  $\mathbb{R}$ . Proto právě výsledky získané pomocí softwaru  $\mathbb{R}$  bereme jakofinální a jejich hodnoty interpretujeme.



### Příklad 6.20. Intervalový odhad parametru $\sigma^2$ normálního modelu

Načtete datový soubor 01-one-sample-mean-skull-mf.txt a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *největší délku mozkovny* u skeletů ženského pohlaví. Za předpokladu, že náhodná veličina  $X$  pochází z normálního rozdělení se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , vypočítejte (a) 99%, oboustranný intervalový odhad; (b) 90% levostranný intervalový odhad; (c) 95% pravostranný intervalový odhad rozptylu  $\sigma^2$ .

#### Řešení příkladu ??

Celkem máme k dispozici  $n = 109$  náhodných veličin  $X_1, \dots, X_{109}$ , přičemž veličina  $X_1$  popisuje největší délku mozkovny u prvního ženského skeletu,  $\dots$ , veličina  $X_{109}$  popisuje výšku největší délku mozkovny u sto devátého ženského skeletu. Předpokládáme, že všechny náhodné veličiny pochází z normálního rozdělení, tj.  $X_1 \sim N(\mu, \sigma^2)$ ,  $\dots$ ,  $X_{109} \sim N(\mu, \sigma^2)$ , a tedy i celý náhodný výběr pochází ze stejného normálního rozdělení, tj.  $X_1, \dots, X_{109} \sim N(\mu, \sigma^2)$ . Naměřením hodnoty největší délky mozkovny každého skeletu jsme získali celkem 109 realizací náhodných veličin  $X_1, \dots, X_{109}$ . Těchto 109 realizací tvoří dohromady datový soubor. Skutečné hodnoty parametrů  $\mu$  a  $\sigma^2$  nebudeme nikdy znát, ale můžeme je odhadnout pomocí bodových nebo intervalových odhadů.

V rámci tohoto příkladu se máme zaměřit na parametr rozptylu  $\sigma^2$  a odhadnout jej pomocí všech tří typů intervalových odhadů. Všimněme si, že ze zadání neznáme skutečnou hodnotu parametru  $\mu$ . Budeme tedy sestrojovat 99% Waldovy empirické intervaly spolehlivosti pro rozptyl  $\sigma^2$  když parametr  $\mu$  neznáme, a tedy budeme vycházet ze vzorců ??, ?? a ??.

Hodnotu parametru  $\sigma^2$  odhadneme pomocí výběrového rozptylu. K jeho výpočtu potřebujeme dopočítat výběrový průměr, tj.


$$m = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{109} (168 + 174 + \dots + 162 + 170) = \frac{19\,024}{109} = 174.5321.$$

Dosazením výběrového průměru  $m = 174.53$  do vzorce ?? získáme hodnotu výběrového rozptylu, tj.

$$\begin{aligned} s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2 \\ &= \frac{1}{108} ((168 - 174.5321)^2 + (174 - 174.5321)^2 + \dots + (162 - 174.5321)^2 + (170 - 174.5321)^2) \\ &= \frac{1}{108} ((-6.5321)^2 + (-0.5321)^2 + \dots + (-12.5321)^2 + (-4.5321)^2) \\ &= 38.6772 \end{aligned}$$

Výběrový rozptyl  $s^2 = 38.6772$ , rozsah náhodného výběru  $n = 109$ . Zbývá stanovit hodnotu kvantilu  $\chi_{n-1}^2(\alpha/2)$  a kvantilu  $\chi_{n-1}^2(1 - \alpha/2)$   $\chi^2$  modelu. K tomu je potřeba nejprve dopočítat koeficient  $\alpha$ , a to vyjádřením z rovnice  $100 \times (1 - \alpha)\% = 99\%$ .

$$\begin{aligned} 100 \times (1 - \alpha)\% &= 99\% \\ 100 \times (1 - \alpha) &= 99 \\ 1 - \alpha &= 0.99 \\ 1 - 0.99 &= \alpha \\ \alpha &= 0.01 \end{aligned}$$

Pomocí softwaru  a funkce `qchisq()` nyní stanovíme hodnotu kvantilu  $\chi_{n-1}^2(\alpha/2) = \chi_{108}^2(0.01/2) = \chi_{108}^2(0.005) = \text{qchisq}(0.005, 108) = 73.8989$  a kvantilu  $\chi_{n-1}^2(1 - \alpha/2) = \chi_{108}^2(1 - 0.01/2) = \chi_{108}^2(0.995) = \text{qchisq}(0.995, 108) = 149.5994$ .

Hranice 99% Waldova empirického oboustranného intervalu spolehlivosti vypočítáme dosazením do vzorce ??.

$$\begin{aligned}(d, h) &= \left( \frac{(n-1)s^2}{\chi_{n-1}^2(1-\alpha/2)}, \frac{(n-1)s^2}{\chi_{n-1}^2(\alpha/2)} \right) \\ &= \left( \frac{108 \times 38.6772}{149.5994}, \frac{108 \times 38.6772}{73.8989} \right) \\ &= (27.9222, 56.5250)\end{aligned}$$

Pro výpočet levostranného intervalu spolehlivosti stanovíme nejprve hodnotu kvantilu  $\chi_{n-1}^2(1 - \alpha)$ . Protože chceme znát hranice 90% Waldova empirického levostranného intervalu spolehlivosti bude  $\alpha = 0.10$  a tedy  $\chi_{n-1}^2(1 - \alpha) = \chi_{108}^2(1 - 0.10) = \chi_{108}^2(0.90) = \text{qchisq}(0.90, 108) = 127.2111$ . Hranice 90% Waldova empirického levostranného intervalu spolehlivosti vypočítáme dosazením hodnot do vzorce ??.

$$\begin{aligned}(d, \infty) &= \left( \frac{(n-1)s^2}{\chi_{n-1}^2(1-\alpha)}, \infty \right) \\ &= \left( \frac{108 \times 38.6772}{127.2111}, \infty \right) \\ &= (32.83627, \infty)\end{aligned}$$

Pro výpočet pravostranného intervalu spolehlivosti stanovíme nejprve hodnotu kvantilu  $\chi_{n-1}^2(\alpha)$ . Protože koeficient  $\alpha$  je tentokrát rovný hodnotě 0.05, bude kvantil  $\chi_{n-1}^2(\alpha) = \chi_{108}^2(0.05) = \text{qchisq}(0.05, 108) = 85.0149$ . Hranice 95% Waldova empirického pravostranného intervalu spolehlivosti vypočítáme dosazením hodnot do vzorce ??.

$$\begin{aligned}(0, h) &= \left( 0, \frac{(n-1)s^2}{\chi_{n-1}^2(\alpha)} \right) \\ &= \left( 0, \frac{108 \times 38.6772}{85.0149} \right) \\ &= (0, 49.13418)\end{aligned}$$

Datový soubor načteme příkazem `read.delim()` a NA hodnoty odstraníme příkazem `na.omit()`. Pomocí operátoru `[]` vybereme z tabulky `data` pouze ty řádky, které se vztahují k ženským skeletům (`data$sex == 'f'`) a sloupec obsahující údaje o největší délce mozkovny 'skull.L'. Hodnotu výběrové směrodatné odchylky dopočítáme funkcí `sd()`, rozsah náhodného výběru stanovíme pomocí funkce `length()`. Do proměnné `alpha` vložíme hodnotu koeficientu  $\alpha = 0.01$ . Nyní přepíšeme vzorců ??, ?? a ??, kde funkci `qchisq()` využijeme na výpočet  $\alpha/2$ ,  $1 - \alpha/2$ ,  $\alpha$  a  $1 - \alpha$  kvantilů  $\chi^2$  modelu získáme dolní a horní hranice všech tří Waldových empirických intervalů spolehlivosti.

```
154 data <- read.delim('00-Data//01-one-sample-mean-skull-mf.txt')
155 data <- na.omit(data)
156
157 skull.LF <- data[data$sex == 'f', 'skull.L']
158 n <- length(skull.LF)
159 s.LF <- sd(skull.LF)
160 alpha <- 0.01
161
162 dh <- (n - 1) * s.LF ^ 2 / qchisq(1 - alpha / 2, n - 1) # 27.92216
163 hh <- (n - 1) * s.LF ^ 2 / qchisq(alpha / 2, n - 1) # 56.52505
164
165 alpha <- 0.10
166 DH <- (n - 1) * s.LF ^ 2 / qchisq(1 - alpha, n - 1) # 32.83628
167
```

```

168 alpha <- 0.05
169 HH <- (n - 1) * s.LF ^ 2 / qchisq(alpha, n - 1) # 49.13418
170
171 (tab <- data.frame(d = round(c(dh, DH, 0), 4),
172                   h = c(round(hh, 4), 'inf', round(HH, 4)),
173                   row.names = c('99% DIS', '99% LIS', '99% PIS')))

```

	d	h
99% DIS	27.9222	56.525
99% LIS	32.8363	inf
99% PIS	0.0000	49.1342

174  
175  
176  
177

**Interpretace výsledků:** 99% Waldův empirický oboustranný interval spolehlivosti pro parametr  $\sigma^2$  má tvar  $(27.92, 56.53) \text{ mm}^2$ . To znamená, že  $27.92 \text{ mm}^2 < \sigma^2 < 56.53 \text{ mm}^2$  s pravděpodobností 99%. V 99 případech ze sta bude rozptyl největší délky mozkovny u skeletů ženského pohlaví větší než  $27.92 \text{ mm}^2$  a menší než  $56.53 \text{ mm}^2$ .

90% Waldův empirický levostranný interval spolehlivosti pro parametr  $\sigma^2$  má tvar  $(32.8363, \infty) \text{ mm}^2$ . To znamená, že  $\sigma^2 > 32.84 \text{ mm}^2$  s pravděpodobností 90%. V 90 případech ze sta bude rozptyl největší délky mozkovny u skeletů ženského pohlaví větší než  $32.84 \text{ mm}^2$ .

95% Waldův empirický pravostranný interval spolehlivosti pro parametr  $\sigma^2$  má tvar  $(0, 49.1342) \text{ mm}^2$ . To znamená, že  $\sigma^2 < 49.13 \text{ mm}^2$  s pravděpodobností 95%. V 95 případech ze sta bude rozptyl největší délky mozkovny u skeletů ženského pohlaví menší než  $49.13 \text{ mm}^2$ . ★

### Příklad 6.21. Intervalový odhad parametru $\sigma$ normálního modelu

Načtete datový soubor 01-one-sample-mean-skull-mf.txt a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *největší šířku mozkovny* u skeletů ženského pohlaví. Za předpokladu, že náhodná veličina  $X$  pochází z normálního rozdělení se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , vypočítejte (a) 95%, oboustranný intervalový odhad směrodatné odchylky  $\sigma$ .

#### Řešení příkladu ??

V rámci tohoto příkladu se máme zaměřit na parametr směrodatné odchylky  $\sigma$ . Žádný interval spolehlivosti pro směrodatnou odchylku ale neznáme. Vystačíme si tedy s tím, co máme k dispozici. Vypočítáme hranice 95% oboustranného Waldova empirického intervalu spolehlivosti pro parametr rozptylu  $\sigma^2$ . Jejich odmocněním získáme hranice 95% oboustranného Waldova empirického intervalu spolehlivosti pro parametr  $\sigma$ .

Budeme tedy vycházet ze vzorce ???. Hodnotu parametru  $\sigma^2$  odhadneme pomocí výběrového rozptylu. Nejprve tedy spočítáme výběrový průměr a následně výběrový rozptyl, tj.

$$m = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{109} (130 + 134 + \dots + 138 + 140) = \frac{14622}{109} = 134.1468.$$

$$\begin{aligned}
s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2 \\
&= \frac{1}{108} ((130 - 134.1468)^2 + (134 - 134.1468)^2 + \dots + (138 - 134.1468)^2 + (140 - 134.1468)^2) \\
&= \frac{1}{108} ((-4.1468)^2 + (-0.1468)^2 + \dots + (3.8532)^2 + (5.8532)^2) \\
&\doteq 22.0523
\end{aligned}$$

Výběrový rozptyl  $s^2 = 22.0523$ , rozsah náhodného výběru  $n = 109$ . Zbývá stanovit hodnotu kvantilu  $\chi_{n-1}^2(\alpha/2)$  a kvantilu  $\chi_{n-1}^2(1 - \alpha/2)$   $\chi^2$  modelu. K tomu je potřeba nejprve dopočítat koeficient  $\alpha$ , a to vyjádřením z rovnice  $100 \times (1 - \alpha)\% = 95\%$ .


$$100 \times (1 - \alpha) \% = 95 \%$$

$$100 \times (1 - \alpha) = 95$$

$$1 - \alpha = 0.95$$

$$1 - 0.95 = \alpha$$

$$\alpha = 0.05$$

Pomocí softwaru  a funkce `qchisq()` nyní stanovíme hodnotu kvantilu  $\chi_{n-1}^2(\alpha/2) = \chi_{108}^2(0.05/2) = \chi_{108}^2(0.025) = \text{qchisq}(0.025, 108) = 81.1329$  a kvantilu  $\chi_{n-1}^2(1 - \alpha/2) = \chi_{108}^2(1 - 0.05/2) = \chi_{108}^2(0.975) = \text{qchisq}(0.975, 108) = 138.6506$ .

Hranice 95% Waldova empirického oboustranného intervalu spolehlivosti pro parametr  $\sigma^2$  vypočítáme dosazením do vzorce ??.

$$\begin{aligned} (d, h)_{\sigma^2} &= \left( \frac{(n-1)s^2}{\chi_{n-1}^2(1-\alpha/2)}, \frac{(n-1)s^2}{\chi_{n-1}^2(\alpha/2)} \right) \\ &= \left( \frac{108 \times 22.0523}{138.6506}, \frac{108 \times 22.0523}{81.1329} \right) \\ &= (17.1773, 29.3549) \end{aligned}$$

Konečně, hranice 95% Waldova empirického oboustranného intervalu spolehlivosti pro parametr  $\sigma$  získáme odmocněním hranic intervalu spolehlivosti (17.1773, 29.3549), tj.

$$\begin{aligned} (d, h)_{\sigma} &= \sqrt{(d, h)_{\sigma^2}} \\ &= \sqrt{(17.1773, 29.3549)} \\ &= (4.1446, 5.4180) \end{aligned}$$

Datový soubor načteme příkazem `read.delim()` a NA hodnoty odstraníme příkazem `na.omit()`. Pomocí operátoru `[]` vybereme z tabulky `data` pouze ty řádky, které se vztahují k ženským skeletům (`data$sex == 'f'`) a sloupec obsahující údaje o největší šířce mozkovny 'skull.L'. Hodnotu výběrové směrodatné odchylky dopočítáme funkcí `sd()`, rozsah náhodného výběru stanovíme pomocí funkce `length()`. Do proměnné `alpha` vložíme hodnotu koeficientu  $\alpha = 0.05$ . Nyní přepisem vzorce ??, kde pomocí funkce `qchisq()` vypočítáme  $\alpha/2$ ,  $1 - \alpha/2$  kvantily  $\chi^2$  modelu získáme dolní a horní hranici 95% Waldova empirického oboustranného intervalu spolehlivosti pro parametr  $\sigma^2$ . Odmocněním obou hranic pomocí funkce `sqrt()` získáme hranice 95% Waldova empirického oboustranného intervalu spolehlivosti pro parametr  $\sigma$ .

```
178 data <- read.delim('01-one-sample-mean-skull-mf.txt')
179 data <- na.omit(data)
180
181 skull.BF <- data[data$sex == 'f', 'skull.B']
182 n <- length(skull.BF)
183 s.BF <- sd(skull.BF)
184 alpha <- 0.05
185
186 dh <- (n - 1) * s.BF ^ 2 / qchisq(1 - alpha / 2, n - 1) # 17.17736
187 hh <- (n - 1) * s.BF ^ 2 / qchisq(alpha / 2, n - 1) # 29.35493
188
189 dh.s <- sqrt(dh) # 4.144558
190 hh.s <- sqrt(hh) # 5.418019
191
192 tab <- data.frame(d = c(dh, dh.s), h = c(hh, hh.s),
193                 row.names = c('95% DIS pro rozptyl', '95% DIS pro sm.odchylku'))
194 round(tab, 4)
```

	d	h
95% DIS pro rozptyl	17.1774	29.3549
95% DIS pro sm.odchylku	4.1446	5.4180

195  
196  
197

**Interpretace výsledků:** 95% Waldův empirický oboustranný interval spolehlivosti pro parametr  $\sigma$  má tvar (4.14, 5.42) mm. To znamená, že  $4.14 \text{ mm} < \sigma < 5.42 \text{ mm}$  s pravděpodobností 95%. V 95 případech ze sta bude směrodatná odchylka největší šířky mozkovny u skeletů ženského pohlaví větší než 4.14 mm a menší než 5.42 mm. ★

### Příklad 6.22. Intervalový odhad parametru $\sigma$ normálního modelu

Načtěte datový soubor 01-one-sample-mean-skull-mf.txt a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující *největší délku mozkovny* u skeletů ženského pohlaví. Za předpokladu, že náhodná veličina  $X$  pochází z normálního rozdělení se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj.  $X \sim N(\mu, \sigma^2)$ , vypočítejte (a) 99%, oboustranný intervalový odhad; (b) 90% levostranný intervalový odhad; (c) 95% pravostranný intervalový odhad směrodatné odchylky  $\sigma$ .

### Řešení příkladu ??

Hranice všech tří intervalů spolehlivosti získáme odmocněním hranic Waldových empirických intervalů spolehlivosti pro rozptyl  $\sigma^2$  vypočítaných v rámci příkladu ???. Konkrétně hranice 99% Waldova empirického oboustranného intervalu spolehlivosti pro parametr  $\sigma$  vypočítáme jako

$$(d, h) = \sqrt{(27.9222, 56.5250)} \\ = (5.2842, 7.5183)$$

Hranice 90% Waldova empirického levostranného intervalu spolehlivosti pro parametr  $\sigma$  vypočítáme jako

$$(d, \infty) = \sqrt{(32.8363, \infty)} \\ = (5.730297, \infty)$$

Hranice 95% Waldova empirického pravostranného intervalu spolehlivosti pro parametr  $\sigma$  vypočítáme jako

$$(0, h) = \sqrt{(0, 49.1342)} \\ = (0, 7.009579)$$

Odmocniny dolních resp. horních hranic Waldových empirických intervalů spolehlivosti získáme pomoc funkce `sqrt()`.

```
198 dh <- sqrt(dh) # 5.284142
199 hh <- sqrt(hh) # 7.518314
200
201 DH <- sqrt(DH) # 5.730295
202 HH <- sqrt(HH) # 7.009578
203
204 (tab <- data.frame(d = round(c(dh, DH, 0), 4),
205                    h = c(round(hh, 4), 'inf', round(HH, 4)),
206                    row.names = c('99% DIS', '99% LIS', '99% PIS')))
```

	d	h
99% DIS	5.2841	7.5183
99% LIS	5.7303	inf
99% PIS	0.0000	7.0096

207  
208  
209  
210

**Interpretace výsledků:** 99% Waldův empirický oboustranný interval spolehlivosti pro parametr  $\sigma$  má tvar (5.28, 7.52) mm. To znamená, že  $5.28 \text{ mm} < \sigma < 7.52 \text{ mm}$  s pravděpodobností 99%. V 99 případech ze sta bude rozptyl největší délky mozkovny u skeletů ženského pohlaví větší než 5.28 mm a menší než 7.52 mm.


90% Waldův empirický levostranný interval spolehlivosti pro parametr  $\sigma$  má tvar  $(5.73, \infty)$  mm. To znamená, že  $\sigma > 5.73$  mm s pravděpodobností 90%. V 90 případech ze sta bude rozptýl největší délky mozkovny u skeletů ženského pohlaví větší než 5.73 mm.

95% Waldův empirický pravostranný interval spolehlivosti pro parametr  $\sigma$  má tvar  $(0, 7.01)$  mm. To znamená, že  $\sigma < 7.01$  mm s pravděpodobností 95%. V 95 případech ze sta bude rozptýl největší délky mozkovny u skeletů ženského pohlaví menší než 7.01 mm. ★

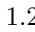
### Příklad 6.23. Intervalový odhad parametru $p$ alternativního modelu

Načtete datový soubor 17-anova-newborns.txt a odstraňte z načtených dat NA hodnoty. Mějme náhodnou veličinu  $X$  popisující ženské pohlaví novorozenců. Za předpokladu, že náhodná veličina  $X$  pochází z alternativního rozdělení s parametrem  $p$ , tj.  $X \sim Alt(p)$ , kde  $p$  je pravděpodobnost narození holčičky, stanovte (a) 95% oboustranný intervalový odhad parametru  $p$ ; (b) 90% levostranný intervalový odhad parametru  $p$ ; (c) 99% pravostranný intervalový odhad parametru  $p$ .

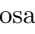
### Řešení příkladu ??

95% oboustranný intervalový odhad získáme pomocí 95% oboustranného Waldova empirického intervalu spolehlivosti. Z příkladu ?? víme, že realizace výběrového průměru  $m = 0.4797$  a rozsah náhodného výběru  $N = 1382$ . Zbývá nám tedy dopočítat hodnotu  $\alpha/2$ -kvantilu a  $1 - \alpha/2$ -kvantilu standardizovaného normálního rozdělení. Z rovnice  $100 \times (1 - \alpha)\% = 95\%$  dopočítáme, že  $\alpha = 0.05$ . Pomocí softwaru  zjistíme hodnoty kvantilů  $u_{\alpha/2} = u_{0.05/2} = u_{0.025} = -1.9600$  a  $u_{1-\alpha/2} = u_{1-0.05/2} = u_{0.975} = 1.9600$ . Dosazením hodnot do vzorce ?? získáme realizaci 95% oboustranného Waldova empirického intervalu spolehlivosti pro parametr  $p$ , tj.

$$\begin{aligned} (d, h) &= \left( m - \sqrt{\frac{m(1-m)}{N}} u_{1-\alpha/2}, m - \sqrt{\frac{m(1-m)}{N}} u_{\alpha/2} \right) \\ &= \left( 0.4797 - \sqrt{\frac{0.4797(1-0.4797)}{1382}} 1.9600, 0.4797 - \sqrt{\frac{0.4797(1-0.4797)}{1382}} (-1.9600) \right) \\ &= (0.4797 - 0.0263, 0.4797 - (-0.0263)) \\ &= (0.4534, 0.5060) \end{aligned}$$

90% levostranný intervalový odhad získáme pomocí 90% levostranného Waldova empirického intervalu spolehlivosti. Realizace výběrového průměru  $m = 0.4797$  a rozsah náhodného výběru  $N = 1382$ . Zbývá nám tedy dopočítat hodnotu  $1 - \alpha$ -kvantilu standardizovaného normálního rozdělení. Z rovnice  $100 \times (1 - \alpha)\% = 90\%$  dopočítáme, že  $\alpha = 0.1$ . Pomocí softwaru  zjistíme hodnotu kvantilu  $u_{1-\alpha} = u_{1-0.1} = u_{0.9} = 1.2816$ . Dosazením hodnot do vzorce ?? získáme realizaci 90% levostranného Waldova empirického intervalu spolehlivosti pro parametr  $p$ , tj.

$$\begin{aligned} (d, 1) &= \left( m - \sqrt{\frac{m(1-m)}{N}} u_{1-\alpha}, 1 \right) \\ &= \left( 0.4797 - \sqrt{\frac{0.4797(1-0.4797)}{1382}} 1.2816, 1 \right) \\ &= (0.4797 - 0.0172, 1) \\ &= (0.4625, 1) \end{aligned}$$

99% pravostranný intervalový odhad získáme pomocí 99% pravostranného Waldova empirického intervalu spolehlivosti. Realizace výběrového průměru  $m = 0.4797$  a rozsah náhodného výběru  $N = 1382$ . Zbývá nám tedy dopočítat hodnotu  $\alpha$ -kvantilu standardizovaného normálního rozdělení. Z rovnice  $100 \times (1 - \alpha)\% = 99\%$  dopočítáme, že  $\alpha = 0.01$ . Pomocí softwaru  zjistíme hodnotu kvantilu  $u_{\alpha} = u_{0.01} = -2.3263$ . Dosazením hodnot do vzorce ?? získáme realizaci 99% pravostranného Waldova empirického intervalu spolehlivosti pro parametr  $p$ , tj.

$$\begin{aligned}
(0, h) &= \left( 0, m - \sqrt{\frac{m(1-m)}{N}} u_\alpha \right) \\
&= \left( 0, 0.4797 - \sqrt{\frac{0.4797(1-0.4797)}{1382}} (-2.3263) \right) \\
&= (0, 0.4797 - (-0.03126)) \\
&= (0, 0.5110)
\end{aligned}$$

Nejprve načteme datový soubor příkazem `read.delim()` a odstraníme NA hodnoty příkazem `na.omit()`. Dále do proměnné `sex` vložíme údaje o pohlaví. Protože proměnná `sex` je typu `factor` s dvěma úrovněmi, úrovní 1 ('f'; (female)) a úrovní 2 ('m', (male)), převedeme jej nejprve pomocí funkce `as.numeric()` na číselný vektor, který vložíme do proměnné `pohlavi`. Následně změňme všechny hodnoty 2 na hodnoty 0, čímž dostaneme požadované kódování 0 = male, 1 = female (viz příklad ??).

```

211 data <- read.delim('17-anova-newborns.txt')
212 data <- na.omit(data)
213 sex <- data$sex.C
214 pohlavi <- as.numeric(sex)
215 pohlavi[pohlavi == 2] <- 0

```

Hranice 95% oboustranného Waldova empirického intervalu spolehlivosti nyní získáme přepisem vzorce ?? s použitím funkcí `mean()`, `sd()`, `length()` a `qnorm()`.

```

216 alpha <- 0.05
217 N <- length(pohlavi)
218 m <- mean(pohlavi)
219 s <- sd(pohlavi)
220
221 dh <- m - sqrt(m * (1 - m) / N) * qnorm(1 - alpha / 2)
222 hh <- m - sqrt(m * (1 - m) / N) * qnorm(alpha / 2)

```

Dolní hranici 90% jednostranného Waldova empirického intervalu spolehlivosti nyní získáme přepisem vzorce ??.

```

223 DH <- m - sqrt(m * (1 - m) / N) * qnorm(1 - alpha)

```

Horní hranici 99% jednostranného Waldova empirického intervalu spolehlivosti nyní získáme přepisem vzorce ??.

Výsledné hranice všech tří intervalů spolehlivosti vložíme do jedné tabulky.

```

224 HH <- m - sqrt(m * (1 - m) / N) * qnorm(alpha)
225 tab <- data.frame(d = c(dh, DH, 0), h = c(hh, 1, HH),
226                   row.names = c('95% OIS', '90% LIS', '99% PIS'))
227 round(tab, 4)

```

	d	h
95% OIS	0.4534	0.5061
90% LIS	0.4576	1.0000
99% PIS	0.0000	0.5018

228  
229  
230  
231

**Interpretace výsledků:** 95% oboustranný Waldův empirický interval spolehlivosti pro parametr  $p$  má tvar  $(0.4534, 0.5061)$ , což znamená, že  $0.4534 < p < 0.5061$  s pravděpodobností 95%. S 95% pravděpodobností se pravděpodobnost narození holčičky pohybuje v rozmezí 45.34% – 50.61%.

90% jednostranný Waldův empirický interval spolehlivosti pro parametr  $p$  má tvar  $(0.4576, 1)$ , což znamená, že  $0.4576 < p < 1$  s pravděpodobností 90%. S pravděpodobností 90% je pravděpodobnost narození holčičky větší než 45.76%.

99% jednostranný Waldův empirický interval spolehlivosti pro parametr  $p$  má tvar  $(0, 0.5018)$ , což znamená, že  $0 < p < 0.5018$  s pravděpodobností 99%. S pravděpodobností 99% je pravděpodobnost narození holčičky menší než 50.61%.

