

3 Charakteristiky polohy a variability a štatistická grafika

Každé rozdelenie pravdepodobnosti býva charakterizované parametrami, ktoré sa odhadujú z realizácií (dát). Tieto parametre označujeme ako **číselné charakteristiky rozdelenia (štatistiky)**. Patria medzi ne **charakteristiky polohy a variability** ako napr. stredná hodnota, medián a ostatné kvantily, rozptyl, rozpätie a pod.

Stredná (očakávaná) hodnota (prvý začiatočný moment) rozdelenia náhodnej veličiny X , ozn. $E[X]$, vypočítame ako (Casella a Berger, 2002)

$$E[X] = \begin{cases} \sum_{i=1}^n x_i \Pr(X = x_i), & \text{ak } X \text{ je diskretná náhodná premenná s realizáciami } x_1, x_2, \dots, x_n, \\ \int_{-\infty}^{\infty} x f_X(x) dx, & \text{ak } X \text{ je spojitá náhodná premenná s hustotou } f_X(x). \end{cases}$$

Rozptyl (druhý centrálny moment) $Var[X] = E[(X - E[X])^2]$, ktorého kladnú odmocninu $SD[X]$, vyjadrujúcu „rozptýlenie“ X okolo $E[X]$, nazývame **smerodajná odchýlka**.

Príklad 110 (charakteristiky normálneho rozdelenia) *Majme náhodný výber X_1, X_2, \dots, X_n z normálneho rozdelenia $N(\mu, \sigma^2)$. Potom $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ nazývame **výberová stredná hodnota** a $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ **výberový rozptyl***

Kedže X je náhodná premenná, \bar{X} nie je aritmetický priemer, ale tzv. výberový priemer; aritmetický priemer \bar{x} je realizácia výberového priemeru. Podobne S^2 je tzv. výberový rozptyl; $\hat{\sigma}^2 = s^2$ je jeho realizácia. Tiež $SD[X] = S$ je výberová smerodajná odchýlka a nie smerodajná odchýlka; smerodajná odchýlka $\hat{\sigma} = s$. Podiel $V_k = S^2/\bar{X}$ nazývame **výberový koeficient variácie**.

Príklad 111 (charakteristiky binomického rozdelenia) *Ak X pochádza z binomického rozdelenia, ozn. $Bin(N, p)$, potom $E[X] = Np$ je **stredná hodnota** a $Var[X] = Np(1-p)$ **rozptyl** náhodnej veličiny X .*

Nech X_1, X_2, \dots, X_n je náhodný výber z nejakého rozdelenia a $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ **usporiadaný náhodný výber** (vo vzrastajúcom poradí) s rozsahom n , kde $X_{(i)}$ nazývame **poriadkové štatistiky**. Sú to náhodné premenné, pre ktoré platí $X_{(1)} < X_{(2)} < \dots < X_{(n)}$. Potom (Casella a Berger, 2002)

$$\begin{aligned} X_{(1)} &= \min_{1 \leq i \leq n} X_i, \\ X_{(2)} &= \text{druhé najmenšie } X_i \text{ (} X_i \text{ v poradí druhé)} \\ &\vdots \\ X_{(n-1)} &= X_i \text{ v poradí } (n-1)\text{-vé,} \\ X_{(n)} &= \max_{1 \leq i \leq n} X_i. \end{aligned}$$

Výberový medián je charakterizovaný číslom $Q_2 = \tilde{X}$, ktoré rozdeľuje náhodný výber tak, že približne polovica X_i je menšia ako táto hodnota a polovica je väčšia ako táto hodnota. Definujeme ho ako²⁴

$$\tilde{X} = \begin{cases} X_{(\frac{n+1}{2})} & \text{ak } n \text{ je nepárne,} \\ \frac{1}{2} \left(X_{(\frac{n}{2})} + X_{(\frac{n}{2}+1)} \right) & \text{ak } n \text{ je párne.} \end{cases}$$

²⁴Ak n je nepárne, potom medián je prostredná hodnota, t.j. hodnota nachádzajúca sa v poradí ako prostredná alebo na mieste X , ktoré zodpovedá poradiu $(n+1)/2$. Ak n je párne, potom medián je uprostred zoradených hodnôt, t.j. hodnota na mieste X , ktoré zodpovedá priemeru X s poradiami $n/2$ a $n/2 + 1$.

Majme číslo $p \in (0, 1)$. Potom **výberovým** $100p$ -tym percentilom náhodného výberu $X_{(\{np\})}$ bude také X_i , kde približne np hodnôt X_i bude menších ako $X_{(\{np\})}$ a $n(1-p)$ väčších ako $X_{(\{np\})}$. Ak $p = 0.5$, ide o 50. percentil alebo medián. Označenie $\{x\}$ v indexe znamená *najbližšie celé číslo*, kde $i - 0.5 \leq x < i + 0.5$, kde $\{x\} = i$. Potom

$$\tilde{X}_p = X_{(\{np\})} = \begin{cases} X_{(\{np\})}, & \text{ak } \frac{1}{2n} < p \leq 0.5, \\ X_{(n+1-\{n(1-p)\})}, & \text{ak } 0.5 < p < 1 - \frac{1}{2n}. \end{cases}$$

Príklad 112 (výberový 100p-ty percentil) Ak $n = 12$, potom výberovým 65. percentilom je $X_{(9)}$, pretože $n(1-p) = 12 \times (1 - 0.65) = 4.2$ a $n + 1 - \{n(1-p)\} = 12 + 1 - 4 = 9$.

Často používanými percentilmi²⁵ sú **výberový dolný kvartil** (25. percentil) $Q_1 = \tilde{X}_{0.25}$ a **výberový horný kvartil** (75. percentil) $Q_3 = \tilde{X}_{0.75}$.

Na výpočet rozptylu nejakého percentilu musíme poznať pravdepodobnostnú funkciu poriadkovej štatistiky, jej hustotu a distribučnú funkciu. Výpočet môžeme výhodne zjednodušiť, keď vieme, aké je asymptotické rozdelenie poriadkovej štatistiky.

Definícia 27 (pravdepodobnostná funkcia poriadkovej štatistiky) Majme náhodný výber X_1, X_2, \dots, X_n z nejakého diskrétného rozdelenia s pravdepodobnostnou funkciou $f_X(x_i) = p_i$, kde $x_1 < x_2 < \dots < x_n$. Nech $P_0 = 0$, $P_1 = p_1$, $P_2 = p_1 + p_2$, \dots , $P_i = p_1 + p_2 + \dots + p_i$. Nech $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ sú poriadkové štatistiky. Potom (Bickel a Doksum, 2006)

$$\Pr(X_{(j)} \leq x_i) = \Pr(Y \geq j) = \sum_{k=j}^n \binom{n}{k} P_i^k (1 - P_i)^{n-k}$$

a

$$\Pr(X_{(j)} = x_i) = \Pr(X_{(j)} \leq x_i) - \Pr(X_{(j)} \leq x_{i-1}) = \sum_{k=j}^n \binom{n}{k} [P_i^k (1 - P_i)^{n-k} - P_{i-1}^k (1 - P_{i-1})^{n-k}],$$

kde Y je náhodná premenná počtu hodnôt z X_1, X_2, \dots, X_n , ktoré sú menšie alebo rovné ako x_i . Nech $\{X_j \leq x_i\}$ je priaznivá udalosť a $\{X_j > x_i\}$ je nepriaznivá udalosť. Ak $i = 1$, potom $\Pr(X_j = x_i) = \Pr(X_j \leq x_i)$, pretože $\Pr_0 = 0$. Teda Y je počet priaznivých udalostí v n pokusoch, t.j. $Y = \text{card}\{X_j \leq x_i\}$. Pravdepodobnosť úspechu je potom $\Pr(X_j \leq x_i)$ pre každý pokus, pretože pokusy sú rovnako rozdelené. Priaznivá a nepriaznivá udalosť j -teho pokusu je nezávislá od výsledku iného pokusu, pretože X_j sú nezávislé od ostatných X_i . Teda $Y \sim \text{Bin}(n, P_i)$.

Ozn. $\text{card}\{\cdot\}$ znamená veľkosť (kardinalitu) množiny.

Ak X_1, X_2, \dots, X_n je náhodný výber zo spojitého rozdelenia, zhody neprichádzajú do úvahy a pravdepodobnosť, že nejaké dve alebo viaceré X_j sú rovnaké, je nula. Teda $\Pr(X_{(1)} < X_{(2)} < \dots < X_{(n)}) = 1$ a výberový priestor pre $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ je $\mathcal{X} = \{(x_1, x_2, \dots, x_n) : x_1 < x_2 < \dots < x_n\}$.

Definícia 28 (hustota a distribučná funkcia poriadkovej štatistiky) Nech $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ sú poriadkové štatistiky náhodného výberu X_1, X_2, \dots, X_n s distribučnou funkciou $F_X(x)$ a hustotou $f_X(x)$. Potom hustota poriadkovej štatistiky $X_{(j)}$ je rovná (Bickel a Doksum, 2006)

$$f_{X_{(j)}}(x) = \frac{n!}{(j-1)!(n-j)!} f_X(x) [F_X(x)]^{j-1} [1 - F_X(x)]^{n-j}$$

²⁵Všeobecne definujeme **kvantil** ako hodnotu skúmanej veličiny, ktorá delí náhodný výber na dve časti pod a nad kvantilom; číslo v dolnom indexe hovorí o časti náhodného výberu pod kvantilom. Podľa toho rozlišujeme aj jeho typ, napr. kvartil, decil, percentil a pod.

a distribučná funkcia

$$F_{X_{(j)}}(x) = \Pr(X_{(j)} \leq x) = \Pr(Y \geq j) = \sum_{k=j}^n \binom{n}{k} [F_X(x)]^k [1 - F_X(x)]^{n-k}.$$

Definícia 29 (asymptotické rozdelenie poriadkovej štatistiky) *Nech $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ sú poriadkové štatistiky náhodného výberu X_1, X_2, \dots, X_n . Majme pravdepodobnosť α , kde $F(t_\alpha) = \alpha$. Asymptoticky platí, že $\sqrt{n}(\frac{j}{n} - \alpha)$ konverguje k 0. Potom je poriadková štatistika $X_{(j)}$ normálne rozdelená so strednou hodnotou $E[X_{(j)}] = t_\alpha$ a rozptylom $\sigma_{X_{(j)}}^2 = \frac{\alpha(1-\alpha)}{f^2(t_\alpha)n}$. Ak $X \sim N(\mu, \sigma^2)$, potom $\sigma_{X_{(j)}}^2 = \sigma^2 \frac{\pi^2}{24 \ln n}$ (Casella a Berger, 2002).*

Príklad 113 (rozptyl poriadkovej štatistiky) *Pomocou delta metódy odvodte rozptyl poriadkovej štatistiky v definícii 29.*

Príklad 114 (rozptyl poriadkovej štatistiky, $X \sim N(\mu, \sigma^2)$) *Pomocou definície 29 odvodte rozptyl poriadkovej štatistiky, ak $X \sim N(\mu, \sigma^2)$.*

Definícia 30 (stredná hodnota a rozptyl mediánu) *Stredná hodnota mediánu $X_{(\frac{n+1}{2})}$ je rovná $E[X_{(\frac{n+1}{2})}] = \tilde{\mu}$ a rozptyl mediánu $\sigma_{X_{(\frac{n+1}{2})}}^2 = \frac{1}{4f^2(\tilde{\mu})n}$, kde n je nepárne. Ak $X \sim N(\mu, \sigma^2)$, potom $\sigma_{X_{(\frac{n+1}{2})}}^2 = \sigma^2 \frac{\pi}{2n}$ (Casella a Berger, 2002).*

Príklad 115 (rozptyl mediánu) *Pomocou delta metódy odvodte rozptyl poriadkovej štatistiky v definícii 30.*

Príklad 116 (rozptyl mediánu, $X \sim N(\mu, \sigma^2)$) *Pomocou definície 30 odvodte rozptyl poriadkovej štatistiky, ak $X \sim N(\mu, \sigma^2)$.*

Ak má náhodná premenná X normálne rozdelenie, výpočet rozptylu mediánu sa zjednoduší, t.j. stačí poznať σ a n . Znalosť rozptylu mediánu je potrebná na výpočet $100 \times (1 - \alpha)\%$ intervalu spoľahlivosti pre medián (pozri kapitolu 4 Testovanie hypotéz). Popis rozdelenia mediánu, ak n je párne, je nad rámec tejto knihy.

Rozpätie náhodného výberu je vzdialenosť medzi najmenšou a najväčšou poriadkovou štatistikou, t.j. $R = X_{(n)} - X_{(1)}$.

Definícia 31 (hodnoty distribučnej funkcie v kvantiloch) *Empirická distribučná funkcia $F_n(x)$ je definovaná nasledovne*

$$F_n(x) = \begin{cases} 0, & \text{ak } x < X_{(1)}, \\ \frac{i}{n}, & \text{ak } X_{(i)} \leq x < X_{(i+1)}, \\ 1, & \text{ak } x \geq X_{(n)}. \end{cases}$$

*Majme transformáciu $T_{(1)} = F_n(X_{(1)})$, $T_{(2)} = F_n(X_{(2)})$, \dots , $T_{(n)} = F_n(X_{(n)})$. Potom $T_{(1)}, T_{(2)}, \dots, T_{(n)}$ sú **poriadkové štatistiky**. Potom platí*

$$\lim_{n \rightarrow \infty} \Pr(\sup_{\forall x \in \mathcal{Y}} [F_n(x) - F(x)] n^{1/2} \leq \lambda) = \Phi(\lambda),$$

kde $F(X)$ je teoretická distribučná funkcia a $\Phi(\lambda) = \sum_{k=-\infty}^{\infty} (-1)^k e^{-2k^2 \lambda^2}$. Potom $100 \times (1 - \alpha)\%$ pás spoľahlivosti pre $F_n(x)$ definujeme ako $F_n(x) \pm \lambda_\alpha 1/n^{1/2}$, kde $\Phi(\lambda_\alpha) = 1 - \alpha$ a $\text{Var}[F_n(x)] = 1/n$ (Kolmogorov, 1933; Smirnov, 1933; Wilks, 1948). Potom môžeme tvrdiť, že $F(X)$ patrí do $100 \times (1 - \alpha)\%$ pásu spoľahlivosti a zároveň je medzi nulou a jednotkou s pravdepodobnosťou $1 - \alpha$.

3.1 Charakteristiky polohy

Realizácie budeme označovať ako x_1, x_2, \dots, x_n , **usporiadané realizácie** budú $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$. Potom môžeme definovať nasledovné odhady charakteristík polohy (výberové charakteristiky polohy) spolu s ich anglickými ekvivalentami:

- **výberové minimum** X_{\min} , ktorého realizácia $x_{\min} = x_{(1)}$;
- **výberové maximum** X_{\max} , ktorého realizácia $x_{\max} = x_{(n)}$;
- **výberový aritmetický priemer** \bar{X} , ktorého realizácia $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{j=1}^{n_j} x_j f_j$, $n_j \leq n$, kde f_j sú frekvencie (počty) prislúchajúcich x_j a $n = \sum_j f_j$;
- **výberový modus** X_{mod} , ktorého realizácia x_{mod} je najčastejšie sa vyskytujúca hodnota (pri diskretnej premennej ide o hodnotu x , v ktorej má pravdepodobnostná funkcia svoje maximum; pri spojitej premennej ide o hodnotu x , v ktorej má hustota svoje maximum);
- **výberový medián** \tilde{X} (robustný odhad polohy), ktorého realizácia

$$\tilde{x} = \begin{cases} x_{(\frac{n+1}{2})} & \text{ak } n \text{ je nepárne,} \\ \frac{1}{2} (x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}) & \text{ak } n \text{ je párne;} \end{cases}$$

rozdelenie je *symetrické*, ak $\bar{x} = \tilde{x} = x_{\text{mod}}$, rozdelenie je *pozitívne zošikmené* (pravostranne), ak $\bar{x} > \tilde{x} > x_{\text{mod}}$ a rozdelenie je *negatívne zošikmené* (ľavostranne), ak $\bar{x} < \tilde{x} < x_{\text{mod}}$;

- **výberové kvartily** poznáme tri
 - **prvý (dolný) kvartil** Q_1 , ktorého realizácia $\tilde{x}_{0.25}$ predstavuje hodnotu, od ktorej je 1/4 dát menšia a 3/4 dát sú väčšie,

$$\Pr [x_{\min}, \tilde{x}_{0.25}] = \Pr [X \leq \tilde{x}_{0.25}] = \frac{1}{4}, \Pr [\tilde{x}_{0.25}, x_{\max}] = \Pr [X \geq \tilde{x}_{0.25}] = \frac{3}{4};$$

- **druhý kvartil** (medián) Q_2 , ktorého realizácia $\tilde{x}_{0.5} = \tilde{x}$ je hodnota, od ktorej je 1/2 dát menšia a 1/2 dát je väčšia,

$$\Pr [x_{\min}, \tilde{x}_{0.5}] = \Pr [X \leq \tilde{x}_{0.5}] = \frac{1}{2}, \Pr [\tilde{x}_{0.5}, x_{\max}] = \Pr [X \geq \tilde{x}_{0.5}] = \frac{1}{2};$$

- **tretí (horný) kvartil** Q_3 , ktorého realizácia $\tilde{x}_{0.75}$ predstavuje hodnotu, od ktorej je 1/4 dát väčšia a 3/4 dát sú menšie,

$$\Pr [x_{\min}, \tilde{x}_{0.75}] = \Pr [X \leq \tilde{x}_{0.75}] = \frac{3}{4}, \Pr [\tilde{x}_{0.75}, x_{\max}] = \Pr [X \geq \tilde{x}_{0.75}] = \frac{1}{4};$$

- **výberové decily** \tilde{X}_k , ktorých realizácie \tilde{x}_k delia súbor na desatiny, t.j. $k/10$ dát je pod decilom a $(10 - k)/10$ nad decilom, kde $k \in \{0, 1, \dots, 10\}$;
- **výberové percentily** \tilde{X}_p (čítame ako 100p-percentil²⁶), ktorých realizácie \tilde{x}_p definujeme ako

$$\tilde{x}_p = \begin{cases} x_{(k+1)} & \text{pre } k \neq np, \\ \frac{1}{2} (x_{(k)} + x_{(k+1)}) & \text{pre } k = np, \end{cases}$$

kde $k = \lfloor np \rfloor$, čo je celá časť čísla np (niekedy sa používa definícia cez $\{x\}$ v indexe, čo znamená najbližšie celé číslo);

²⁶Výraz „100p-percentil“ čítame ako „stokrát p-ty percentil“. Ak $p = 0.75$, potom $100 \times 0.75 = 75$, čo čítame ako „75. percentil“.

- **výberový päťčíselný súhrn** $(X_{\min}, Q_1, Q_2, Q_3, X_{\max})$, ktorého realizáciu označujeme ako $(x_{\min}, \tilde{x}_{0.25}, \tilde{x}_{0.50}, \tilde{x}_{0.75}, x_{\max})$.

Príklad 117 (výšky 10-ročných dievčat) *Majme výšky $n = 12$ náhodne vybraných 10-ročných dievčat v cm usporiadaných podľa veľkosti (poradia ozn. ako r_i pre $x_{(i)}$; pri rovnakých pozorovaniach hovoríme o **strednoporadiach**; strednoporadie sa vypočíta ako priemer poradií realizácií s rovnakou hodnotou).*

Riešenie (pozri tabuľku 16)

Tabuľka 16: Zoradné realizácie x_i a ich poradia r_i pre výšky 10-ročných dievčat

i	1	2	3	4	5	6	7	8	9	10	11	12
$x_{(i)}$	131	132	135	141	141	141	141	142	143	146	146	151
r_i	1	2	3	5.5	5.5	5.5	5.5	8	9	10.5	10.5	12

$\bar{x} \doteq 140.83$, $\tilde{x} = \frac{1}{2}(x_{(6)} + x_{(7)}) = 141$, $Q_1 = \tilde{x}_{0.25} = \frac{1}{2}(x_{(3)} + x_{(4)}) = 138$, kde $k = \lfloor 12 \times 0.25 \rfloor = 3$, $Q_3 = \tilde{x}_{0.75} = \frac{1}{2}(x_{(9)} + x_{(10)}) = 144.5$, kde $k = \lfloor 12 \times 0.75 \rfloor = 9$.

Čo sa stane so spomínanými charakteristikami polohy, keď zmeníme mierku (škálu), napr. gramy na kilogramy alebo namiesto hmotnosti použijeme logaritmus hmotnosti?

Nech $a, b \in \mathbb{R}$ sú nejaké dané konštanty, a je posunutie a b škála. Potom $y_i = a + bx_i$ a pre priemer²⁷

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n (a + bx_i) = \frac{1}{n} \left(na + b \sum_{i=1}^n x_i \right) = a + b\bar{x} = \overline{a + bx}.$$

Pokiaľ $b \in \mathbb{R}^+$, usporiadanie hodnôt x_i sa pri transformácii na $y_i = a + bx_i$ nezmení, teda

$$a + bx_{(1)} \leq a + bx_{(2)} \leq \dots \leq a + bx_{(n)}.$$

Pre medián v tomto prípade platí²⁸

$$\tilde{y} = a + b\tilde{x} = \widetilde{a + bx}.$$

Ľahko sa dá nahliadnuť, že usporiadanie zachová každá rastúca funkcia $g(x)$, teda platí

$$g(x_{(1)}) \leq g(x_{(2)}) \leq \dots \leq g(x_{(n)})$$

a pre medián bude platiť $g(\tilde{x}) \doteq \tilde{g}(x)$. Pre nepárne n platí predchádzajúci vzťah presne, označenie „približnosti“ potrebujeme pre párne n , kde $x_{(\frac{n}{2})} < x_{(\frac{n}{2}+1)}$. V tomto prípade je však $1/2$ hodnôt $g(x_i)$ menšia ako $\tilde{g}(x)$. Teda špeciálne môžeme medián logaritmu (napr. hmotnosti) spočítať ako logaritmus mediánu (napr. hmotnosti). Pokiaľ dôjde v pozorovaniach k posunutiu, dôjde k rovnakému posunutiu aj v charakteristike polohy. Ak zmeníme mierku, potom stačí urobiť rovnakú úpravu aj u charakteristiky polohy.

Robustnou charakteristikou strednej hodnoty (odolnejšou na odlahlé pozorovania) je (Tukey, 1962)

²⁷Rovnosť znamená, že priemer posunutej a preškálovanej veličiny y je rovný posunutému a preškálovanému priemeru pôvodnej veličiny x .

²⁸Rovnosť znamená, že medián posunutej a preškálovanej veličiny y je rovný posunutému a preškálovanému mediánu pôvodnej veličiny x .

- **výberový γ -urezaný aritmetický priemer** \bar{X}_g , ktorého realizáciou je \bar{x}_g a vypočíta sa ako

$$\bar{x}_g = \frac{1}{n-2g} (x_{(g+1)} + x_{(g+2)} + \dots + x_{(n-g)}),$$

kde $g = \{\gamma n\}$, $g = \lfloor \gamma n \rfloor$, $\gamma = 0.1, 0.2$. Viac ako $\gamma 100$ % pozorovaní²⁹ musí byť nahradených, aby sa tento priemer zmenil na malý alebo veľký v porovnaní s pôvodným [³⁰*breakdown point* angl \bar{x}_t je teda γ],

- **výberový γ -winsorizovaný priemer** \bar{X}_w , ktorého realizácia \bar{x}_w je definovaná ako

$$\bar{x}_w = \frac{1}{n} ((g+1)x_{(g+1)} + x_{(g+2)} + \dots + (g+1)x_{(n-g)}).$$

Viac ako $\gamma 100$ % pozorovaní musí byť nahradených, aby sa tento priemer zmenil na malý alebo veľký v porovnaní s pôvodným [*breakdown point* \bar{x}_w je teda γ]. angl

3.2 Charakteristiky variability

Definujeme nasledovné základné odhady charakteristík variability (výberové charakteristiky variability) spolu s ich anglickými ekvivalentami:

- **výberový rozptyl** S^2 , ktorého realizáciou je

$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2;$$

pri lineárnej transformácii sa rozptyl mení nasledovne³¹

$$s_y^2 = s_{a+bx}^2 = b^2 s_x^2,$$

t.j.

$$\begin{aligned} s_y^2 &= s_{a+bx}^2 = \frac{1}{n-1} \sum_{i=1}^n (a + bx_i - \overline{a+bx})^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n (a + bx_i - (a + b\bar{x}))^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n (b(x_i - \bar{x}))^2 = b^2 s_x^2; \end{aligned}$$

výpočtová podoba rozptylu

$$s_x^2 = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) = \frac{1}{n-1} \left(\sum_{j=1}^{n_j} x_j^2 f_j - n\bar{x}^2 \right), n_j \leq n,$$

kde f_j sú frekvencie (počty) prislúchajúcich x_j a $n = \sum_j f_j$;

²⁹Výraz „ $\gamma 100$ % pozorovaní“ čítame ako „gama krát stopercent pozorovaní“.

³⁰*Breakdown point* hovorí o počte pozorovaní, ktoré potrebujeme na to, aby sme výrazne zmenili hodnotu charakteristiky polohy. Pre γ -urezaný a γ -winsorizovaný aritmetický priemer ide o γn pozorovaní, pre medián ide o $n/2$ pozorovaní a pre aritmetický priemer stačí iba jedno pozorovanie (preto hovoríme, že aritmetický priemer je veľmi citlivý na odľahlé pozorovania).

³¹Rovnosť znamená, že rozptyl posunutej a preškáľovanej veličiny y je rovný násobku druhej mocniny škály a rozptylu pôvodnej veličiny x .

- **výberová smerodajná odchýlka** S , ktorej realizáciou je

$$s_x = \sqrt{s_x^2};$$

pri lineárnej transformácii sa smerodajná odchýlka mení nasledovne³²

$$s_y = s_{a+bx} = |b| s_x,$$

teda, ak pripočítame ku všetkým pozorovaniam rovnakú konštantu, miera variability sa nemení; zmena mierky (u pomerovej mierky zmena jednotiek) má za následok rovnakú úpravu jednotlivých pozorovaní i miery variability v podobe smerodajnej odchýlky;

- **výberový koeficient variácie** V_k , ktorého realizácia v_k predstavuje normalizovanú podobu výberového rozptylu (inverzia *signal-to-noise ratio*; podiel variability na priemere) angl

$$v_k = \frac{s_x}{\bar{x}};$$

bezrozmerná veličina, zvyčajne vyjadrovaná v percentách, t.j. $100 \times (s_x/\bar{x}) \%$ a môže sa používať len pre realizácie, ktorých rozsah nadobúda kladné hodnoty; používa sa pri porovnávaní variability súborov s nerovnakými priemermi (napr. pri porovnaní variability výšky detí určitého veku s výškou dospelých určitého veku alebo pri porovnaní variability premenných meraných v rôznych jednotkách);

- **výberový rozptyl aritmetického priemeru** $S_{\bar{X}}^2$, ktorého realizáciou je

$$s_{\bar{x}}^2 = \frac{s_x^2}{n};$$

- **výberová stredná chyba priemeru (štandardná chyba)** $S_{\bar{X}}$, ktorej realizáciou je

$$s_{\bar{x}} = \frac{s_x}{\sqrt{n}};$$

- **výberový koeficient šikmosti** B_1 , ktorého realizáciou je

$$b_1 = \frac{n^{-1} \sum_{i=1}^n (x_i - \bar{x})^3}{[n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2]^{3/2}} = \frac{\sqrt{n} \sum_{i=1}^n (x_i - \bar{x})^3}{[\sum_{i=1}^n (x_i - \bar{x})^2]^{3/2}},$$

kde rozdelenie je *symetrické*, ak $b_1 = 0$, *pozitívne zošikmené* (hustota na ľavej strane stúpa strmšie ako na pravej), ak $b_1 > 0$ a *negatívne zošikmené* (hustota na pravej strane stúpa strmšie ako na ľavej), ak $b_1 < 0$;

- **výberový koeficient špicatosti** B_2 , ktorého realizáciou je

$$b_2 = \frac{n^{-1} \sum_{i=1}^n (x_i - \bar{x})^4}{[n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2]^2} - 3 = \frac{n \sum_{i=1}^n (x_i - \bar{x})^4}{[\sum_{i=1}^n (x_i - \bar{x})^2]^2} - 3,$$

kde rozdelenie je normálne (*mezokurtické*), ak $b_2 = 0$, *špicaté (leptokurtické)*, ak $b_2 > 0$ a *ploché (platykurtické)*, ak $b_2 < 0$;

³²Rovnosť znamená, že smerodajná odchýlka posunutej a preškálovanej veličiny y je rovná násobku absolútnej hodnoty škály a smerodajnej odchýlky pôvodnej veličiny x .

- **výberová suma štvorcov** $\sum_{i=1}^n (X_i - \bar{X})^2$, ktorej realizácia je

$$SS = \sum_{i=1}^n (x_i - \bar{x})^2,$$

kde sa tento čitateľ rozptylu používa napr. v lineárnom regresnom modeli, v modeli ANOVA a pod.;

- **výberová suma absolútnych odchýlok** $\sum_{i=1}^n |X_i - \tilde{X}_{0.5}|$, ktorej realizácia je

$$SAD = \sum_{i=1}^n |x_i - \tilde{x}_{0.5}|;$$

- **výberový priemer absolútnych odchýlok** $\frac{1}{n} \sum_{i=1}^n |X_i - \tilde{X}_{0.5}|$, ktorého realizácia je

$$MAD = SAD/n;$$

- **výberové rozpätie** $X_{\max} - X_{\min}$, ktorého realizáciou je

$$D = x_{\max} - x_{\min};$$

- **výberové medzikvartilové rozpätie** $Q_3 - Q_1$, ktorého realizáciou je

$$D_Q = \tilde{x}_{0.75} - \tilde{x}_{0.25};$$

kde rozdelenie je (medzi kvartilmi) *symetrické*, ak $\tilde{x}_{0.75} - \tilde{x}_{0.50} = \tilde{x}_{0.50} - \tilde{x}_{0.25}$, *pozitívne zošikmené*, ak $\tilde{x}_{0.75} - \tilde{x}_{0.50} > \tilde{x}_{0.50} - \tilde{x}_{0.25}$ a *negatívne zošikmené*, ak $\tilde{x}_{0.75} - \tilde{x}_{0.50} < \tilde{x}_{0.50} - \tilde{x}_{0.25}$;

- **výberové decilové rozpätie** $\tilde{X}_{0.9} - \tilde{X}_{0.1}$, ktorého realizáciou je

$$D_D = \tilde{x}_{0.9} - \tilde{x}_{0.1};$$

- **výberové percentilové rozpätie** $\tilde{X}_{0.99} - \tilde{X}_{0.01}$, ktorého realizáciou je

$$D_P = \tilde{x}_{0.99} - \tilde{x}_{0.01}.$$

Robustnými charakteristikami variability sú (Tukey, 1962)

- **výberový γ -urezaný rozptyl** S_g^2 , ktorého realizácia s_g^2 sa vypočíta ako

$$s_g^2 = \frac{1}{n - 2g - 1} \sum_{i=g+1}^{n-g} x_{(i)};$$

viac ako $\gamma 100$ % pozorovaní musí byť nahradených, aby sa tento rozptyl zmenil na veľký v porovnaní s pôvodným s^2 [*breakdown point* s_g^2 je γ]; platí $s_g^2 < s^2$ pretože urezanie odstráni angl odľahlé hodnoty;

- **výberový γ -winsorizovaný rozptyl** S_w^2 , ktorého realizáciu označujeme ako s_w^2 ; viac ako $\gamma 100$ % pozorovaní musí byť nahradených, aby sa tento rozptyl zmenil na veľký v porovnaní s pôvodným s^2 [*breakdown point* s_w^2 je γ]; platí $s_w^2 < s^2$ pretože winsorizovanie priťahuje angl extrémne hodnoty bližšie k priemeru;

- **výberový kvartilový koeficient variácie** $V_{k,Q}$, ktorého realizáciu $v_{k,Q}$ vypočítame ako

$$v_{k,Q} = \frac{Q_3 - Q_1}{Q_2}.$$

Ďalšie robustné charakteristiky variability (výberové rozpätie) charakterizované pomocou upravených hraníc sú

- **výberové robustné minimum a maximum** („vnútorné hradby“) X_{\min}^* a X_{\max}^* , ktorých realizácie sú definované ako

$$x_{\min}^* = B_D = \tilde{x}_{0.25} - 1.5(\tilde{x}_{0.75} - \tilde{x}_{0.25}) = Q_1 - 1.5D_Q,$$

$$x_{\max}^* = B_H = \tilde{x}_{0.75} + 1.5(\tilde{x}_{0.75} - \tilde{x}_{0.25}) = Q_3 - 1.5D_Q,$$

kde prvky vybočujúce z hradieb sa považujú za *podozrivé, potencionálne odľahlé pozorovania*;

- **výberové robustné minimum a maximum** („vonkajšie hradby“) definované ako $Q_1 - 3(Q_3 - Q_1)$, $Q_3 + 3(Q_3 - Q_1)$, ktorých realizácie sú $B_D^* = \tilde{x}_{0.25} - 3D_Q$, $B_H^* = \tilde{x}_{0.75} + 3D_Q$

– pokiaľ sú nejaké $x_i < B_D^* \vee x_i > B_H^*$, hovoríme, že ide o *vzdialené body*³³,

– ak $x_i \in \langle B_D^*, B_D \rangle \vee \langle B_H, B_H^* \rangle$, ide o *body vonkajšie*,

– ak $x_i \in \langle B_D, B_H \rangle$, ide o *body vnútorné* alebo *body prilahlé mediánu*;

– pre normálne rozdelenie platí $B_H - B_D = Q_3 + 1.5D_Q - Q_1 + 1.5D_Q = 4D_Q \doteq 4.2$; pravdepodobnosť, že $x_i \notin \langle B_D, B_H \rangle$ je potom 0.04;

- **výberové robustné miery šikmosti** B_{1Q} a B_{1O} a ich rozptyly za podmienky asymptotickej normality $B_{1\cdot}$, kde $\cdot = Q$ alebo O , ktorých realizácie sú definované nasledovne

– kvartilový koeficient šikmosti

$$b_{1Q} = \frac{(\tilde{x}_{0.75} - \tilde{x}_{0.50}) - (\tilde{x}_{0.50} - \tilde{x}_{0.25})}{\tilde{x}_{0.75} - \tilde{x}_{0.25}}, \text{Var}_{as}(b_{1Q}) = 1.84,$$

– oktilový koeficient šikmosti

$$b_{1O} = \frac{(\tilde{x}_{0.875} - \tilde{x}_{0.50}) - (\tilde{x}_{0.50} - \tilde{x}_{0.125})}{\tilde{x}_{0.875} - \tilde{x}_{0.125}}, \text{Var}_{as}(b_{1O}) = 1.15.$$

3.3 Detekcia odľahlých pozorovaní

Homogénny náhodný výber je taký výber, v ktorom všetky $x_i, i = 1, 2, \dots, n$, sú realizácie rovnakého rozdelenia pravdepodobnosti s konštantným rozptylom σ^2 . K **nehomogenitám výberu** dochádza všade tam, kde sa vyskytujú výrazné nerovnomernosti v realizáciách, náhle sa menia podmienky experimentu a pod. Nehomogenita môže byť spôsobená aj nevhodne zvoleným výberom subjektov.

Špeciálnym prípadom ovplyvňujúcim homogenitu výberu sú **odľahlé pozorovania** (outliers). angl Takéto pozorovania skresľujú odhady polohy (špeciálne aritmetického priemeru) a variability (hlavne rozptylu), takže môžu znehodnotiť ďalšiu štatistickú analýzu. Pri ich overovaní sa používa mnoho idealizovaných predpokladov. Musíme poznať ich predpokladaný počet, ich rozdelenie a tiež rozdelenie ostatných prvkov náhodného výberu. Navyše je nutné zostrojiť štatistický alebo pravdepodobnostný model, podľa ktorého sa odľahlé pozorovania chovajú. Testovanie odľahlých pozorovaní bez doplnkových informácií je teda málo spoľahlivé.

Jednoduchou technikou, kedy sa predpokladá, že dáta majú normálne rozdelenie, je **modifikácia vnútorných hradieb** B_D a B_H na

³³Ozn. \vee znamená „alebo“ a ozn. \wedge znamená „a súčasne“.

$$B_D^{mod} = \tilde{x}_{0.25} - kD_Q, B_H^{mod} = \tilde{x}_{0.75} + kD_Q,$$

kde sa parameter k volí tak, aby pravdepodobnosť $\Pr(n, k)$ bola dostatočne vysoká, napr. 0.95. $\Pr(n, k)$ je pravdepodobnosť, že žiaden prvok z náhodného výberu z normálneho rozdelenia s rozsahom n nebude mimo intervalu $I = \langle B_D^{mod}, B_H^{mod} \rangle$. Ak $\Pr(n, k) = 0.95$ a $n \in \langle 8, 100 \rangle$ použijeme aproximáciu $k \approx 2.25 - 3.6/n$. Teda prvky mimo I sa považujú za odľahlé. Postup spomenutý vyššie je dostatočne robustný.

Definujme mieru rozptylu ako **kvantilovú odchýlku** $D_{Q^*} = 2D_Q$. Ak urobíme štandardizáciu, dostaneme $D_{Q_{st}^*} = 1$ a **štandardizovaný medián** bude $\tilde{x}_{st} = 0$ a **štandardizovaná kvantilová funkcia** indikujúca tvar bude (Meloun a Militký, 2004)

$$Q_{st}(p) = \frac{\tilde{x}_p - \tilde{x}_{0.5}}{D_{Q^*}}.$$

Hodnoty kvantilov, pre ktoré platí $|Q_{st}(p)| \geq 1$, sú považované za vybočujúce (pre normálne rozdelenie) a hovoríme potom o **identifikátoroch dlhých „chvostov“** (koncov). Hodnoty $Q_{st}(p)$ môžeme použiť na

– identifikáciu miery šikmosti $SQ = Q_{st}(0.25) + Q_{st}(0.75)$, kedy je rozdelenie pravdepodobnosti symetrické medzi kvartilmi, ak SQ je rovné nule,

– identifikáciu dĺžky koncov, kedy

- $Q_{st}(0.95) < 0.5$ hovorí o krátkych koncoch,
- $Q_{st}(0.95) > 1$ hovorí o dlhých koncoch a
- pre stredne dlhé konce bude platiť $Q_{st}(0.95) \in \langle 0.5, 1.0 \rangle$.

Okrem vonkajších a vnútorných hradieb, môžeme odľahlé pozorovania jednoducho detegovať aj nasledovne (Rousseeuw a van Zomeren, 1990)

- $|X - \bar{X}| > 2s$,
- $\hat{\sigma} = MAD/0.6745$, $|X - \tilde{X}_{0.5}| > k \frac{MAD}{0.6745}$ (ako k sa najčastejšie používa 2 alebo 2.24).

Príklad 118 (odľahlé pozorovania; dĺžka kľúčnej kosti) Za predpokladu normality náhodnej premennej X najväčšia dĺžka kosti kľúčnej z pravej strany (**cla.L**; dáta: *more-samples-variances-clavicle.txt*), t.j. $X \sim N(\mu, \sigma^2)$, identifikujte odľahlé pozorovania pomocou (a) vnútorných hradieb B_D a B_H , (b) modifikovaných vnútorných hradieb B_D^{mod} a B_H^{mod} , (c) identifikujte dlhé „chvosty“ rozdelenia tejto premennej na základe štandardizovanej kvantilovej funkcie.

3.4 Z-skóre

V antropológii nás veľmi často zaujímajú normované realizácie, nazývané tiež normované veličiny, a to **z-skóre**, ktoré definujeme ako

$$z_i = \frac{x_i - \bar{x}}{s_x}.$$

Dostaneme ich ako špeciálny prípad lineárnej transformácie $y = a + bx$, kde voľbou $b = 1/s_x$ a $a = -\bar{x}/s_x$. Potom aritmetický priemer z -skóre

$$\bar{z} = -\frac{\bar{x}}{s_x} + \frac{1}{s_x}\bar{x} = 0$$

a rozptyl z -skóre

$$s_z^2 = \left(\frac{1}{s_x}\right)^2 s_x^2 = 1,$$

Pomocou z -skóre môžeme vyjadriť aj koeficient šikmosti a špicatosti, kde

$$b_1 = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x}\right)^3 = \frac{1}{n} \sum_{i=1}^n z_i^3,$$

a za podmienky asymptotickej normality

$$E[b_1] = 0, \text{Var}[b_1] = \frac{n-2}{(n+1)(n+3)};$$

$$b_2 = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x}\right)^4 = \frac{1}{n} \sum_{i=1}^n z_i^4$$

a za podmienky asymptotickej normality

$$E[b_2] = 3 - \frac{6}{n+1}, \text{Var}[b_2] = \frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}.$$

Pokiaľ dáta pochádzajú z normálneho rozdelenia, budú mať oba koeficienty hodnoty približne nulové (pri b_2 po odčítaní konštanty 3).

Musíme si uvedomiť, že prítomnosť odľahlých pozorovaní v realizáciách ovplyvňuje výpočet priemeru a rozptylu, ktoré sú potrebné na výpočet z -skóre, a následne tak ovplyvňuje tiež hodnotu vlastného z -skóre. Ak je rozdelenie dát zošikmenené alebo je normalita porušená inak, nebude z -skóre odrážať situáciu vierohodne a jeho ďalšie použitie je problematické. Ak predpokladáme, že je rozdelenie znaku v populácii normálne, obmedzuje sa použitie z -skóre napr. na zistenie, či nejaké pozorovanie (pacient) patrí svojimi charakteristikami do zdravej populácie.

Príklad 119 (z -skóre; šírka lebky) *Majme náhodnú premennú X šírka lebky (`skull.B`; mm; dáta: `one-sample-mean-skull-mf.txt`) u mužov. Za predpokladu asymptotickej normality X , t.j. $X \sim N(\mu, \sigma^2)$, vypočítajte z -skóre tejto premennej pomocou funkcie `mean()` a `sd()`. Výsledok skontrolujte pomocou funkcie `scale()`.*

3.5 Príklady na charakteristiky polohy a variability

Príklad 120 (argument minima) *Vygenerujte pseudonáhodné čísla $X \sim N(\mu, \sigma^2)$, kde $\mu = 0, \sigma^2 = 1$ a $n = 1000$. Vygenerované čísla ozn. $x_i, i = 1, 2, \dots, 1000$. Nájdite numericky také c , ktoré minimalizuje (a) sumu štvorcov odchýlok $\sum_{i=1}^{1000} (x_i - c)^2$, t.j. $c_1 = \arg \min_c \sum_{i=1}^{1000} (x_i - c)^2$ a (b) sumu absolútnych odchýlok $\sum_{i=1}^{1000} |x_i - c|$, t.j. $c_2 = \arg \min_c \sum_{i=1}^{1000} |x_i - c|$. Za c dosadzujte postupne (1) všetky $x_{(j)}$ ($x_{(j)}$ sú usporiadané x_i podľa veľkosti od najmenšieho po najväčšie) a vybrané charakteristiky polohy ako (2) aritmetický priemer, (3) nejaké kvantily \tilde{x}_p , kde $p \in \langle 0, 1 \rangle$ a pod. Nakreslite obrázok závislosti (a) sumy štvorcov odchýlok na $x_{(j)}$, t.j. body $[x_j, y_j]$, kde $y_j = \sum_{i=1}^{1000} (x_i - x_{(j)})^2$ a (b) sumy absolútnych odchýlok na $x_{(j)}$, t.j. body $[x_{(j)}, y_j]$, kde $y_j = \sum_{i=1}^{1000} |x_i - x_{(j)}|$. Podobné obrázky nakreslite aj pre \tilde{x}_p namiesto $x_{(j)}$.*

Príklad 121 (výšky 10-ročných dievčat, pokrač.) *Vypočítajte základné charakteristiky polohy a variability.*

Riešenie v R (pozri tabuľku 17 a 18)

Dáta:

```
94 | x <- c(131, 132, 135, 141, 141, 141, 141, 142, 143, 146, 146, 151)
```

Minimum, maximum, medián, aritmetický priemer, prvý kvartil, druhý kvartil, tretí kvartil, kvartily (pomocou jednej funkcie), päťčíselný súhrn, rozptyl a smerodajná odchýlka:

```
95 | min(x) # 131
96 | max(x) # 151
97 | median(x) # 141
98 | mean(x) # 140.8333
99 | priemer <- sum(x)/length(x) # 140.8333
100 | q1 <- quantile(x,0.25,type=2) # 138
101 | q2 <- quantile(x,0.50,type=2) # 141
102 | q3 <- quantile(x,0.75,type=2) # 144.5
103 | quantile(x,c(0.25,0.5,0.75),type=2) # 138.0 141.0 144.5
104 | quantile(x,c(0,0.25,0.5,0.75,1),type=2) # 131.0 138.0 141.0 144.5 151.0
105 | var(x) # 33.78788
106 | sd(x) # 5.812734
```

Funkcie na výpočet rozptylu, smerodajnej odchýlky, štandardnej chyby, šikmosti a špicatosti:

```
107 | "rozptyl" <- function(x) sum((x-mean(x))^2)/(length(x)-1)
108 | "smerodch" <- function(x) sqrt(var(x))
109 | "SE" <- function(x) sqrt(var(x) / length(x))
110 | SE(x) # 1.677992
111 | "sikmost" <- function(x) {(1/length(x))*sum(((x-mean(x))/(sqrt(var(x))))^3)}
112 | sikmost(x) # -0.2121993
113 | "spicatost" <- function(x) {(1/length(x))*sum(((x-mean(x))/(
114 | (sqrt(var(x))))^4) - 3}
115 | spicatost(x) # -0.9029347
```

Suma štvorcov, priemerná absolútna odchýlka (priemer absolútnych odchýlok), suma absolútnych odchýlok, rozpätie, medzikvartilové rozpätie, zisťovanie symetrie, robustný výpočet minima a maxima („vnútorné hradby“), robustné rozpätie:

```
116 | var(x)*(length(x)-1) # 371.6667
117 | mad(x) # 5.1891
118 | mad(x)*length(x) # 62.2692
119 | c(min(x),max(x)) # 131 151
120 | range(x) # 131 151
121 | Dq <- quantile(x,0.75,type=2)-quantile(x,0.25,type=2) # 6.5
122 | c(q3 - q2, q2 - q1) # 3.5 3.0
123 | Bd <- q1-1.5*Dq; Bh <- q3+1.5*Dq
124 | c(Bd,Bh) # 128.25 154.25
```

Funkcia na výpočet kvartilovej šikmosti:

```
125 | "kvart.sikmost" <- function(x) {
126 |     q1 <- quantile(x,0.25,type=2)
127 |     q2 <- quantile(x,0.5,type=2)
128 |     q3 <- quantile(x,0.75,type=2)
129 |     ((q3 - q2) - (q2 - q1))/(q3 - q1)
130 | }
131 | kvart.sikmost(x) # 0.07692308
```

Funkcia na výpočet oktilovej šikmosti:

```

132 "oktil.sikmost" <- function(x) {
133     q125 <- quantile(x,0.125,type=2)
134     q2 <- quantile(x,0.5,type=2)
135     q875 <- quantile(x,0.875,type=2)
136     ((q875 - q2) - (q2 - q125))/(q875 - q125)
137 }
138 oktil.sikmost(x) # -0.2857143

```

Funkcia na výpočet výberového γ -urezaného aritmetického priemeru a rozptylu (vytvorenie dátového vektora na ich výpočet):

```

139 "urezanie" <- function(x, gama = 0.1){
140 x <- na.omit(x) # odstranenie NA, ak sa v datach nachadzaju
141 n <- length(x)
142 g.min <- floor(gama*n) # najvacsie cele cislo mensie ako gama*n
143 g.max <- floor((1-gama)*n) # najvacsie cele cislo mensie ako (1-gama)*n
144 x.sort <- sort(x) # zoradenie podľa veľkosti
145 x.min <- x.sort[g.min] # vybratie dolnej hranice
146 x.max <- x.sort[g.max] # vybratie hornej hranice
147 xg <- x[x > x.min & x < x.max]
148 return(xg)
149 }

```

Funkcia na výpočet výberového aritmetického priemeru a rozptylu winsorizovaného pomocou „vnútorných hradieb“ (vytvorenie dátového vektora na ich výpočet):

```

150 "winsorizacia" <- function(x){
151 x <- na.omit(x) # odstranenie NA
152 q1 <- quantile(x,0.25)
153 q3 <- quantile(x,0.75)
154 Dq <- q3 - q1
155 min.x <- q1 - 1.5*Dq
156 max.x <- q3 + 1.5*Dq
157 xw <- x
158 for (i in 1: length(x)) if (x[i] >= max.x) xw[i] <- max.x
159 for (i in 1: length(x)) if (x[i] <= min.x) xw[i] <- min.x
160 return(xw)
161 }

```

Porovnanie troch typov aritmetických priemerov a rozptylov:

```

162 xg <- urezanie(x)
163 xw <- winsorizacia(x)
164 TAB <- rbind(c(length(x), mean(x), sd(x)),
165             c(length(xg), mean(xg), sd(xg)),
166             c(length(xw), mean(xw), sd(xw)))
167 dimnames(TAB)[[1]] <- c("surove_data", "urezane_data", "winsorizovane_data")
168 dimnames(TAB)[[2]] <- c("rozsah", "priemer", "sd")
169 TAB <- round(TAB,2)

```

Tabuľka 17: Rozsah, aritmetický priemer a smerodajná odchýlka pre surové, urezané a winsorizované dáta (výšky 10-ročných dievčat)

	rozsah	aritmetický priemer	smerodajná odchýlka
surové dáta	$n = 12$	$\bar{x} = 140.83$	$s^2 = 5.81$
urezané dáta	$n_g = 8$	$\bar{x}_g = 139.50$	$s_g^2 = 3.85$
winsorizované dáta	$n_w = 12$	$\bar{x}_w = 141.03$	$s_w^2 = 5.21$

Funkcia na výpočet niektorých základných charakteristík:

```

170 "zakl.char" <- function(x, type = 7){
171 # odstranenie chybajucich pozorovani

```

Tabuľka 18: Vybrané charakteristiky polohy a variability pre surové dáta (výšky 10-ročných dievčat)


charakteristika	n	\bar{x}	\tilde{x}_{\min}	$\tilde{x}_{0.25}$	$\tilde{x}_{0.50}$	$\tilde{x}_{0.75}$	\tilde{x}_{\max}	b_1	b_2	s_x	$s_{\bar{x}}$
hodnota	12	140.8	131.0	138.0	141.0	144.5	151.0	-0.21	-0.90	5.81	1.68

```

172 x <- x[!is.na(x)]
173 n <- length(x) # rozsah
174 # kvantily p = 0,0.25,0.5,0.75 a 1
175 kvantily <- quantile(x,c(0,0.25,0.5,0.75,1),type=type)
176 priemer <- mean(x) # priemer
177 SD <- sd(x) # smerodajna odchylka
178 StEr <- SE(x) # standardna chyba
179 sikm <- sikmost(x) # sikmost
180 spic <- spicatost(x) # spicatost
181 # vsetky vysledky spolu
182 vysledky <- c(n,priemer, kvantily, sikm, spic, SD, StEr)
183 # priradenie nazvov vysledkom
184 names(vysledky) <- c("n","priem",names(kvantily),"sik","spic","sd","se")
185 # zaokruhlenie na dve desatinne miesta
186 vysledky <- round(vysledky,2)
187 return(vysledky)
188 }
189 zakl.char(x, type = 2)

```

Príklad 122 (základné charakteristiky polohy a variability) *Vypočítajte základné charakteristiky polohy a variability pre premennú najväčšia dĺžka lebky (`skull.L`) a najväčšia šírka lebky (`skull.B`) u mužov; dáta: `one-sample-mean-skull-mf.txt`. Výsledok uložte pomocou funkcie `write.table()`.*

Riešenie v  (pozri tabuľku 19)

```

190 DATA <- read.table("one-sample-mean-skull-mf.txt",header=TRUE)
191 names(DATA) ## "id" "pop" "sex" "skull.L" "skull.B"
192 attach(DATA)
193 ZCH1 <- zakl.char(skull.L[sex=="m"])
194 ZCH2 <- zakl.char(skull.B[sex=="m"])
195 ZCH <- rbind(ZCH1,ZCH2)
196 dimnames(ZCH)[[1]] <- c("skull.L","skull.B")
197 ZCH
198 write.table(ZCH,"skull-tab-01.txt")

```

Tabuľka 19: Vybrané charakteristiky polohy a variability pre najväčšiu dĺžku lebky

	n	\bar{x}	\tilde{x}_{\min}	$\tilde{x}_{0.25}$	$\tilde{x}_{0.50}$	$\tilde{x}_{0.75}$	\tilde{x}_{\max}	b_1	b_2	s_x	$s_{\bar{x}}$
skull.L	217	182.04	164.00	177.00	182.00	187.00	199.00	-0.06	-0.48	6.37	0.43
skull.B	216	137.19	124.00	134.00	137.00	140.00	149.00	0.08	-0.30	4.82	0.33

Príklad 123 (šikmost' a špicatost') *Naprogramujte funkcie na výpočet rozptylu šikmosti a špicatosti.*

Príklad 124 (základné charakteristiky polohy a variability) *Vypočítajte základné charakteristiky polohy a variability pre nasledujúce premenné:*

(a) *šírko-dĺžkový index lebky vypočítaný ako podiel premenných šírka lebky (`skull.B`; v mm) a dĺžka lebky (`skull.L`; v mm; dáta: `one-sample-mean-skull-mf.txt`) u mužov;*

(b) *stranový rozdiel vertikálneho priemeru diafýzy kľúčnej kosti (`simd.R` a `simd.L`; v mm) na pravej aj ľavej strane tela (dáta: `paired-means-clavicle2.txt`);*

(c) *najväčšia výška mozgovne (`skull.pH`; mm) a morfológická výška tváre (`face.H`; mm; dáta: `one-sample-correlation-skull-mf.txt`).*

3.6 Štatistická grafika

Pokiaľ chceme zobrazíť základné a relevantné grafy (spolu s výpočtom základných charakteristík polohy a variability), hovoríme o **exploratórnej analýze dát (EDA)**; pozri Murrell (2011) a Kabacoff (2011). Grafická interpretácia výberového súboru je možná pomocou stĺpcového diagramu, spojnicového grafu (polygónu početnosti, frekvenčnej krivky a polygónu kumulatívnych početností), bodového grafu, kruhového diagramu, histogramu, empirickej distribučnej funkcie, krabicového diagramu a kvantilového diagramu.

3.6.1 Stĺpcový diagram

Stĺpcový diagram – vyjadruje číselné hodnoty pomocou obdĺžnikových stĺpcov, obyčajne v zvislej, no niekedy aj vo vodorovnej polohe; môže byť neškálovaný, t.j. v *absolútnej škále*, alebo škálovaný, t.j. v *relatívnej škále* (lepšie porovnanie v prípade sledovania viacerých súborov); špeciálnymi variantmi sú **veková pyramída** (strom života, znázorňuje vekové zloženie obyvateľstva) a **histogram**. Stĺpcový diagram nakreslíme pomocou funkcie `barplot(x)`.

Kombinácii stĺpcových diagramov usporiadaných nad sebou pre škálované kategoriálne dáta, t.j. dáta v podobe pravdepodobností, ktorých suma je rovná jednej, sa hovorí aj **spinogram**. Spinogram nakreslíme pomocou funkcií

```
199 | library(vcd)
200 | spine(x)
```

Príklad 125 (stĺpcový diagram; oči vs. vlasy) *Kontingenčná tabuľka predstavuje pravdepodobnosti výskytu rôznych farieb vlasov a očí v populácii. Ide o model multinomického rozdelenia, ozn. $Mult(\mathbf{p}, N)$, pretože všetky pravdepodobnosti dávajú v sume jednotku. Vypočítajte početnosti všetkých buniek tabuľky, ak máme v populácii 1000 jedincov. Prepočítajte pravdepodobnosti na model súčinného multinomického rozdelenia (pravdepodobnosti výskytu pre dve farby očí a tri farby vlasov pozri v tabuľke 20). Použite základné funkcie (pozri Spector, 2008; Venables a kol., 2013; Verzani, 2005) a skontrolujte pomocou funkcií `margin.table(oci)` a `prop.table(oci)`. Nakreslite stĺpcové diagramy pre oba typy súčinných multinomických rozdelení – pre početnosti, ako aj pre pravdepodobnosti. [Marginálne súčty možno pridať pomocou funkcie `addmargins(oci)`; keby sme mali surové dáta (hodnoty 0 a 1 pre každý subjekt), kontingenčnú tabuľku vytvoríme pomocou funkcie `fTable(data)`.]*

Riešenie v  (pozri tabuľky 21 až 24 a obrázok 20)

```
201 | modre.oci <- c(0.12,0.22,0.06) # vektor pravdepodobnosti
202 | hnde.oci <- c(0.15,0.34,0.11) # vektor pravdepodobnosti
203 | oci <- rbind(modre.oci,hnde.oci) # dva vektory spojené do matice
204 | sum(oci) # celkova (totalna) suma = 1
```

Tabuľka 20: Kontingenčná tabuľka 2×3 pravdepodobností výskytu pre dve farby očí a tri farby vlasov

oči/vlasy	blond	hnede	ryšavé
modré	0.12	0.22	0.06
hnede	0.15	0.34	0.11

```

205 # nazvy riadkov a stlpcov
206 dimnames(oci)[[1]] <- c("modre","hnede")
207 dimnames(oci)[[2]] <- c("blond","hnede","rysave")
208 round(addmargins(oci),2) # marginalne pravdepodobnosti

209 oci.pocty <- oci*1000 # fiktivne pocetnosti
210 addmargins(oci.pocty) # marginalne pocetnosti

211 sumy <- apply(oci.pocty,2,sum) # sumy po stlpcoch
212 oci.prav <- oci.pocty
213 oci.prav[1,] <- oci.pocty[1,]/sumy
214 oci.prav[2,] <- oci.pocty[2,]/sumy
215 apply(oci.prav,2,sum) # suma po stlpcoch = 1
216 round(addmargins(oci.prav,1),2) # marginalne pocetnosti po riadkoch

217 sumy1 <- apply(oci.pocty,1,sum) # sumy po riadkoch
218 oci.prav1 <- oci.pocty
219 oci.prav1[,1] <- oci.pocty[,1]/sumy1
220 oci.prav1[,2] <- oci.pocty[,2]/sumy1
221 oci.prav1[,3] <- oci.pocty[,3]/sumy1
222 apply(oci.prav1,1,sum) # suma po riadkoch = 1
223 round(addmargins(oci.prav1,2),4) # marginalne pocetnosti po stlpcoch

224 par(mfcol=c(2,2),mar=c(5.5,2,1,1))
225 barplot(oci.pocty,space=0)
226 barplot(oci.prav,space=0)
227 barplot(t(oci.pocty),space=0)
228 barplot(t(oci.prav1),space=0)
229 # mozne pridanie legendy (nezobrazena)
230 # legend("topright",c("modre oci","hnede oci"),fill=c("black","grey"))

```

Tabuľka 21: Kontingenčná tabuľka 2×3 pravdepodobností výskytu pre dve farby očí a tri farby vlasov spolu s marginálnymi pravdepodobnosťami (multinomické rozdelenie)

oči/vlasy	blond	hnede	ryšavé	suma
modré	0.12	0.22	0.06	0.40
hnede	0.15	0.34	0.11	0.60
suma	0.27	0.56	0.17	1.00

Tabuľka 22: Kontingenčná tabuľka 2×3 početností výskytu pre dve farby očí a tri farby vlasov spolu s marginálnymi početnosťami (multinomické rozdelenie)

oči/vlasy	blond	hnede	ryšavé	suma
modré	120	220	60	400
hnede	150	340	110	600
suma	270	560	170	1000

Tabuľka 23: Kontingenčná tabuľka 2×3 pravdepodobností výskytu pre dve farby očí a tri farby vlasov spolu s marginálnymi stĺpcovými početnosťami (súčinové multinomické rozdelenie; po stĺpcoch)

oči/vlasy	blond	hnede	ryšavé
modré	0.4444	0.3929	0.3529
hnede	0.5556	0.6071	0.6471
suma	1.0000	1.0000	1.0000

Tabuľka 24: Kontingenčná tabuľka 2×3 pravdepodobností výskytu pre dve farby očí a tri farby vlasov spolu s marginálnymi riadkovými početnosťami (súčinové multinomické rozdelenie; po riadkoch)

oči/vlasy	blond	hnede	ryšavé	suma
modré	0.3000	0.5500	0.1500	1.0000
hnede	0.2500	0.5667	0.1833	1.0000

3.6.2 Spojnicový graf, polygón početnosti a frekvenčná krivka

Spojnicový graf – znázorňuje priebeh časového radu a jeho špeciálnymi prípadmi sú **polygón početnosti**, **frekvenčná krivka**, **polygón kumulatívnych početností**.

Polygón početnosti – spojnicový diagram, v ktorom nad stredmi triednych intervalov I_i vztýčime kolmice, s výškou úmernou príslušným triednym početnostiam a koncové body kolmíc pospájame. Ich súradnice sú $[x_i^*, n_i]$. Ide o zobrazenie priebehu početností vnútri každého intervalu, ale plocha uzavretá spojnicou polygónu nie je presne (len približne) úmerná počtu pozorovaní v intervale.

Frekvenčná krivka – vznikne, keď koncové body polygónu početnosti pospájame hladkou krivkou; vystihuje celkom presne priebeh rozdelenia početnosti a plocha v každom mieste ohraničená krivkou je priamo úmerná počtu pozorovaní.

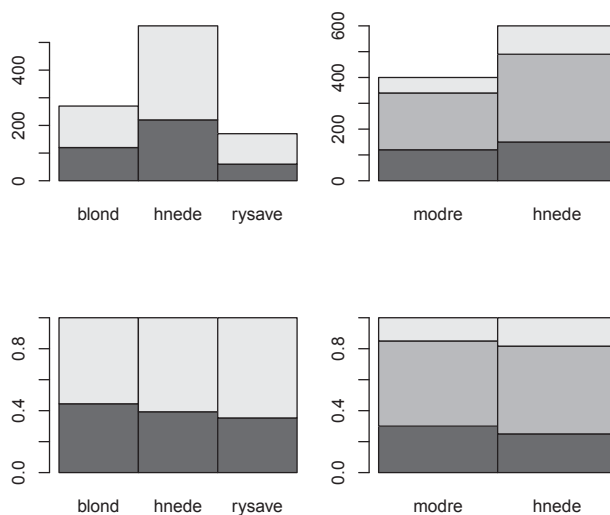
3.6.3 Bodový graf

Bodový (rozptylový) graf – zobrazuje namerané hodnoty v pravouhlej súradnicovej sústave (dvojdimenziálnej, 2D; trojdimenziálnej, 3D), pričom jednotlivé kategórie sa odlišujú pomocou rôznych značiek, farieb a pod.; často sa používa na zobrazenie závislosti dvoch znakov. Dvojdimenziálny a trojdimenziálny rozptylový graf nakreslíme pomocou funkcií

```
231 | plot(x,y) # 2D
232 | library(scatterplot3d)
233 | scatterplot3d(x,y,z) # 3D
```

Argumenty funkcie `plot(x)`:

- `type=` argument kontrolujúci typ grafu
 - `type="p"` – body (prednastavená hodnota; default),
 - `type="l"` – čiary,
 - `type="b"` – body pospájané čiarami,
 - `type="s"` – schodovitá funkcia,
 - `type="n"` – prázdny obrázok;
- argumenty popisu osí a grafu
 - `xlab="retazec"` – popis x -ovej osi,
 - `ylab="retazec"` – popis y -ovej osi,



Obr. 20: Stĺpcové diagramy – početnosti (prvý riadok) pre vlasy (vľavo), pre oči (vpravo); pravdepodobnosti (druhý riadok) pre vlasy (vľavo), pre oči (vpravo)


- `main="retazec"` – hlavný nadpis,
- `sub="retazec"` – podnadpis pod x -ovou osou;
- farba `col="anglicky.nazov"` alebo kódovanie v RGB-škále (funkcia `rgb(cislo1, cislo2, cislo3)` vytvorí RGB-vektor z hodnôt intenzity); transformácia `col2rgb(anglicky.nazov)` vytvorí RGB-vektor z anglického názvu farby, kde RGB-vektor je textový vektor so 7 alebo 9 elementmi, kde za `"#"` nasleduje *red*, *blue* a *green* farebný kanál a voliteľne aj koeficient transparentie α v hexadecimálnej sústave (po preškálovaní na hodnoty $0, \dots, 255$; default je `"black"`);
- veľkosť `cex=k`, $k \in \mathbb{R}$, prednastavená hodnota je 1
- typ bodov `pch=k`, číslo $k = 1, 2, \dots, 20$, prednastavená hodnota je 1 (prázdny krúžok);
- typ čiar `lty=k`, číslo $k = 1, 2, \dots, 20$, prednastavená hodnota je 1 (plná čiara).

Ďalšie funkcie súvisiace s bodovým grafom:

- `points(x,y)` – pridávanie bodov do obrázka;
- `lines(x,y)` – pridávanie čiar do obrázka;
- `text(x, y, labels)` – pridanie textu v bodoch špecifikovaných súradnicami x a y , kde popis `labels[i]` sa zobrazuje v bodoch `(x[i],y[i])`, prednastavené hodnoty sú `1:length(x)`;
`plot(x,y, type="n"); text(x, y, names);`
- `title(main,sub,xlab,ylab)` – dodatočné pridanie nadpisov a popisov osí;
- `legend(x,y,legend)` – dodatočné pridanie legendy, špecificky umiestnenej v súradniciach x a y , kde sú prednastavené nasledovné polohy `"bottomright"`, `"bottom"`, `"bottomleft"`, `"left"`, `"topleft"`, `"top"`, `"topright"`, `"right"` a `"center"` argumenty funkcie
 - `fill="retazec"` – farba výplne (ne)orámovanej legendy,

- `col="retazec"` – farba nakreslených bodov alebo čiar,
 - `lty=k` – typ čiar, $k = 1, 2, \dots, 20$,
 - `lwd=k` – šírka čiar, $k \in \mathbb{N}$, prednastavená hodnota je 1 (plná čiara),
 - `pch=k` – typ bodov, $k = 1, 2, \dots, 20$, prednastavená hodnota je 1 (prázdny krúžok);
- `locator(n,type)` – určenie polohy konkrétneho bodu v grafe (napr. odľahlé pozorovanie) pomocou jedného bodového kliknutia myšou v jeho blízkosti, pričom funkcia bod nielen označí, ale aj vypočíta jeho súradnice `text(locator(1),"retazec")`; použitie v legende `legend(locator(1),...)`;
 - `identify(x,y,labels)` identifikácia bodov, ak poznáme ich súradnice.

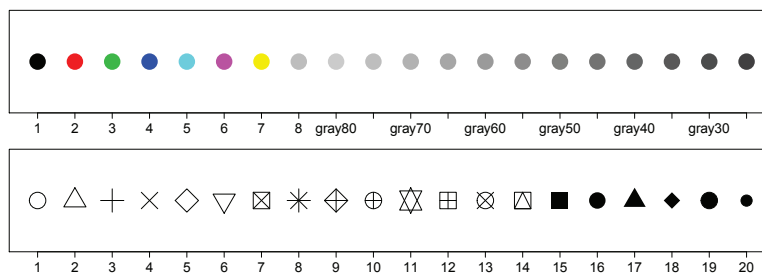
Príklad 126 (typy bodov a základné farby) *Nakreslite obrázok (a) základných dvadsiatich typov bodov a (b) ôsmich typov farieb a dvanástich odtieňov sivej.*

Riešenie v  (pozri obrázok 21)

```


234 windows(14,2.5)
235 par(mar=c(3,0.1,0.1,0.1))
236 plot(1:20,rep(0:1,10),type="n",sub="",xlab="",ylab="",bty="n",axes=FALSE)
237 points(1:20,rep(0.5,20),pch=1:20,cex=4)
238 axis(1,at=1:20,labels=1:20,cex.axis=1.5)
239 box()
240 windows(14,2.5)
241 par(mar=c(3,0.1,0.1,0.1))
242 plot(1:20,rep(0:1,10),type="n",sub="",xlab="",ylab="",bty="n",axes=FALSE)
243 points(1:8,rep(0.5,8),pch=16,col=1:8,cex=4)
244 siva <- paste("gray",rev(seq(25,80,by=5)),sep="")
245 points(9:20,rep(0.5,12),pch=16,col=siva,cex=4)
246 axis(1,at=1:20,labels=c(1:8,siva),cex.axis=1.5)
247 box()

```



Obr. 21: Základné typy bodov (dolný riadok) a farieb (horný riadok)

Príklad 127 (typy bodov a základné farby) *Nakreslite obrázok zpiatich základných typov čiar, ktoré (a) smerujú zvislo, (b) smerujú vodorovne, (c) zvierajú s osou x uhol 45° .*

Riešenie v  (pozri obrázok 22)

```

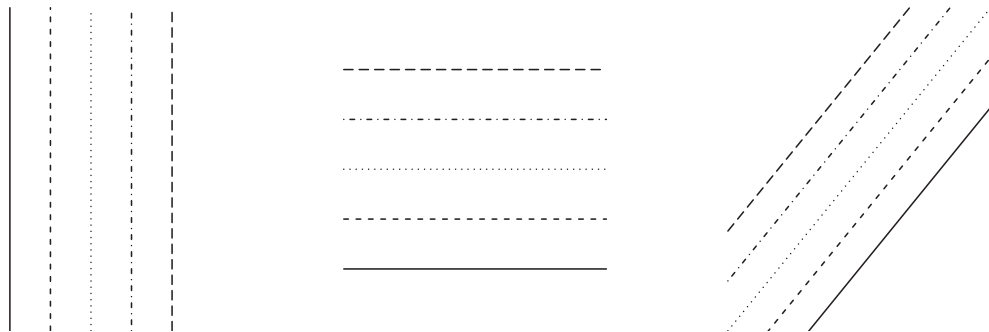
248 windows(12,4)
249 par(mfcol=c(1,3))
250 plot(1,1,type="n",xlab="",ylab="",bty="n",axes=FALSE,xlim=c(-3,3),ylim=c(-3,3))

```

```


251 abline(v=-2:2,lty=1:5)
252 plot(1,1,type="n",xlab="",ylab="",bty="n",axes=FALSE,xlim=c(-3,3),ylim=c(-3,3))
253 abline(h=-2:2,lty=1:5)
254 plot(1,1,type="n",xlab="",ylab="",bty="n",axes=FALSE,xlim=c(-3,3),ylim=c(-3,3))
255 abline(a=-2,b=tan(pi/4),lty=1)
256 abline(a=-1,b=tan(pi/4),lty=2)
257 abline(a=0,b=tan(pi/4),lty=3)
258 abline(a=1,b=tan(pi/4),lty=4)
259 abline(a=2,b=tan(pi/4),lty=5)

```



Obr. 22: Základné typy čiar – zvislo, vodorovne a v uhle 45° (zľava doprava)

Príklad 128 (základy grafiky; dáta iris) Nakreslite rozptylový graf dĺžky a šírky kališných lístkov pre všetky tri taxóny kosatcov pomocou (a) rôznych typov bodov a (b) rôznych farieb. Dokreslite do obrázku (c) regresné priamky pre každý taxón použitím rôzneho typu čiar. Pozri `help(iris)` ohľadom popisu premenných a ďalších detailov o dátach (Fisher, 1936/1971).

Riešenie v  (pozri obrázok 23)

```

260 # zmena datoveho ramca na maticu
261 irisDATA <- as.matrix(iris[,1:4])
262 dimnames(irisDATA)[[2]]
263 # [1] "Sepal.Length" "Sepal.Width"
264 # [3] "Petal.Length" "Petal.Width"
265 irisLABELS <- iris[,5]
266 levels(irisLABELS)
267 # "setosa" "versicolor" "virginica"
268 # rozptylový graf [zle popisujú osi]
269 plot(irisDATA[, "Sepal.Length"], irisDATA[, "Sepal.Width"])
270 # rozsahy oboch premenných
271 x.rozs <- range(irisDATA[, "Sepal.Length"])
272 y.rozs <- range(irisDATA[, "Sepal.Width"])
273 # typy bodov podľa skupín
274 par(mfcol=c(1,3))
275 plot(irisDATA[, "Sepal.Length"], irisDATA[, "Sepal.Width"], type="n",
276      xlab="", ylab="", asp=1)
277 points(irisDATA[irisLABELS=="setosa", "Sepal.Length"],
278        irisDATA[irisLABELS=="setosa", "Sepal.Width"], pch=16)
279 points(irisDATA[irisLABELS=="versicolor", "Sepal.Length"],

```

```

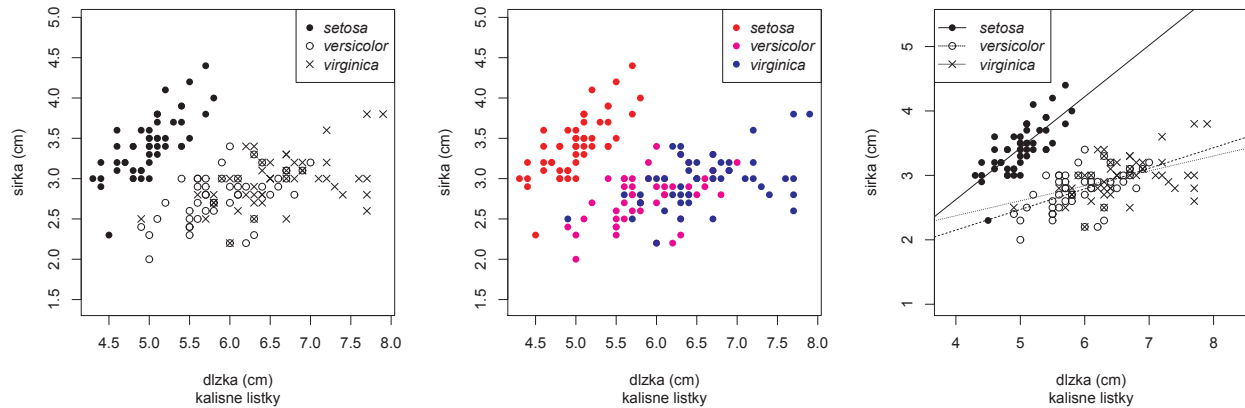
280     irisDATA[irisLABELS=="versicolor","Sepal.Width"],pch=1)
281 points(irisDATA[irisLABELS=="virginica","Sepal.Length"],
282     irisDATA[irisLABELS=="virginica","Sepal.Width"],pch=4)
283 title(xlab="dlzka(cm)",ylab="sirka(cm)",
284     sub="kalisne_listky")
285 legend("topright",c("setosa","versicolor","virginica"),text.font=3,pch=c(16,1,4))
286 # typy farieb podľa skupin
287 plot(irisDATA[, "Sepal.Length"], irisDATA[, "Sepal.Width"], type="n",
288     xlab="", ylab="", asp=1)
289 points(irisDATA[irisLABELS=="setosa","Sepal.Length"],
290     irisDATA[irisLABELS=="setosa","Sepal.Width"], pch=16, col="red")
291 points(irisDATA[irisLABELS=="versicolor","Sepal.Length"],
292     irisDATA[irisLABELS=="versicolor","Sepal.Width"], pch=16,
293     col="magenta")
294 points(irisDATA[irisLABELS=="virginica","Sepal.Length"],
295     irisDATA[irisLABELS=="virginica","Sepal.Width"], pch=16, col="blue")
296 title(xlab="dlzka(cm)",ylab="sirka(cm)",
297     sub="kalisne_listky")
298 legend("topright",c("setosa","versicolor","virginica"),text.font=3,
299     pch=c(16,16,16),col=c("red","magenta","blue"))
300 # typy ciar [presahuju mimo oblakov dat, co nie je statisticky spravne]
301 plot(irisDATA[, "Sepal.Length"], irisDATA[, "Sepal.Width"], type="n",
302     xlab="", ylab="", asp=1, ylim=c(2.0-1, 4.4+1))
303 points(irisDATA[irisLABELS=="setosa","Sepal.Length"],
304     irisDATA[irisLABELS=="setosa","Sepal.Width"], pch=16)
305 points(irisDATA[irisLABELS=="versicolor","Sepal.Length"],
306     irisDATA[irisLABELS=="versicolor","Sepal.Width"], pch=1)
307 points(irisDATA[irisLABELS=="virginica","Sepal.Length"],
308     irisDATA[irisLABELS=="virginica","Sepal.Width"], pch=4)
309 # linearne regresne modely pre vsetky tri taxony zvlast
310 LM1 <- lm(irisDATA[irisLABELS=="setosa","Sepal.Width"] ~
311     irisDATA[irisLABELS=="setosa","Sepal.Length"])
312 LM2 <- lm(irisDATA[irisLABELS=="versicolor","Sepal.Width"] ~
313     irisDATA[irisLABELS=="versicolor","Sepal.Length"])
314 LM3 <- lm(irisDATA[irisLABELS=="virginica","Sepal.Width"] ~
315     irisDATA[irisLABELS=="virginica","Sepal.Length"])
316 # ciary a ich typy pre linearny regresny model
317 abline(LM1,lty=1,lwd=2)
318 abline(LM2,lty=2,lwd=2)
319 abline(LM3,lty=3,lwd=2)
320 title(xlab="dlzka(cm)",ylab="sirka(cm)", sub="kalisne_listky")
321 legend("topleft",c("setosa","versicolor","virginica"),text.font=3,pch=c(16,1,4),
322     lty=c(1,2,3))

```

3.6.4 Kruhový diagram


Kruhový (výsekový, koláčový) diagram – zachytáva štruktúru dát takým spôsobom, že celá plocha kruhu predstavuje celý súbor a kruhové výseky jej jednotlivé časti, pričom polomery zvierajúce uhol 3.6° vymedzujú plochu odpovedajúcu 1 % celého obsahu. Kruhový diagram nakreslíme pomocou funkcie `pie(x)`. Argumenty funkcie `pie(x)`:

- `x` vektor relatívnych hodnôt (pravdepodobností), ktoré v súčte dávajú 1, teda i -ta položka bude $x[i]/\text{sum}(x)$ kruhu (ale aj početností); graf začína horizontálnou čiarou doprava a pokračuje proti smeru hodinových ručičiek;
- `names` vektor mien prislúchajúcich jednotlivým položkám grafu;
- `col` vektor farieb, ktorými sú jednotlivé položky vyfarbené;
- `labels="retazec"` – pomenovania kruhových výsekov



Obr. 23: Rozptylové grafy

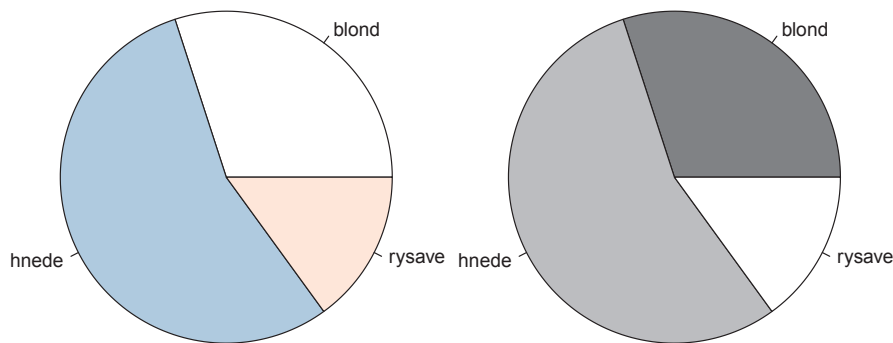
Príklad 129 (kruhový diagram) Vytvorte kruhový diagram znázorňujúci pravdepodobnosti výskytu rôznych farieb vlasov so súčasným výskytom modrých očí (dáta z príkladu `oci vs. vlasy`, tabuľka `oci.prav1`). Vytvorte rovnaký diagram v odtieňoch sivej farby (`gray()`). Použite tabuľku `oci.prav1` z príkladu 125.

Riešenie v  (pozri obrázok 24)

```

323 windows(8,4)
324 par(mfcol=c(1,2),mar=c(0,0,0,0))
325 pie(oci.prav1[1,])
326 pie(oci.prav1[1,],col=gray(seq(0.4,1.0,length=3)))

```

Obr. 24: Kruhový diagram (dáta `oci vs. vlasy`)

Príklad 130 (kruhový diagram) Nakreslite tiež kruhový diagram 24

- odtieňov sivej (`gray(sekvencia)`, kde `sekvencia` sú čísla z intervalu $\langle 0, 1 \rangle$),
- odtieňov farieb dúhy (`rainbow(k)`),

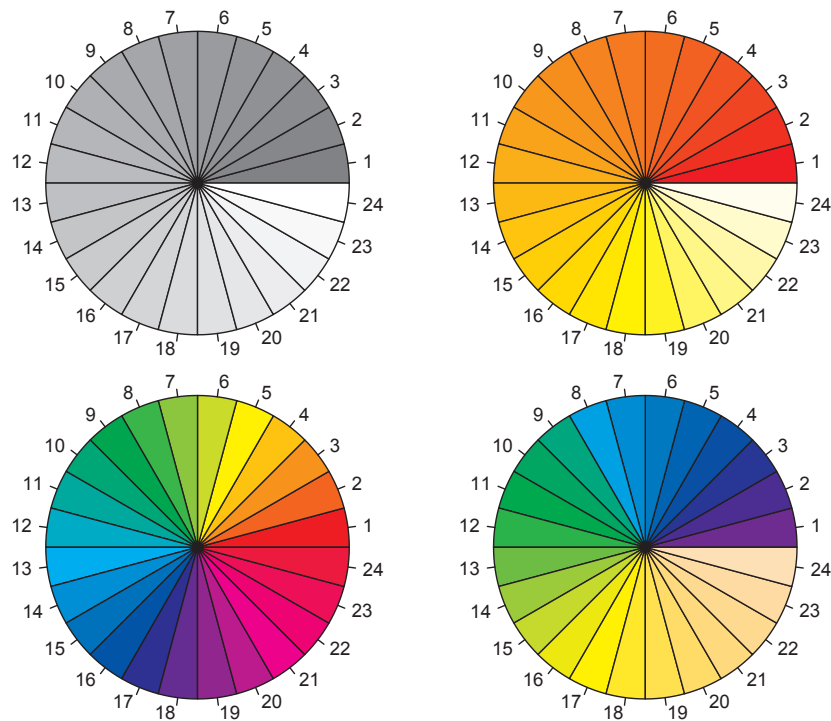
- *teplých farieb* (`heat.colors(k)`),
- *topografických farieb* (`topo.colors(k)` alebo `terrain.colors(k)`).

Riešenie v \mathbb{R} (pozri obrázok 25)

```

327 windows(8,8)
328 par(mfcol=c(2,2),mar=c(0,0,0,0))
329 pie(rep(1,24),col=gray(seq(0.4,1.0,length=24)),radius=0.9)
330 pie(rep(1,24),col=rainbow(24),radius=0.9)
331 pie(rep(1,24),col=heat.colors(24),radius=0.9)
332 pie(rep(1,24),col=topo.colors(24),radius=0.9)

```



Obr. 25: Kruhový diagram (farebné škály)


3.6.5 Histogram

Histogram – predstavuje stĺpcový diagram s k stĺpcami, ktorých základňa sa rovná šírke intervalu $I_i = (x_i, x_{i+1})$ a výška i -teho stĺpca jeho početnosti ($i = 1, 2, \dots, k$). Zobrazuje početnosti pozorovaní v jednotlivých intervaloch v *absolútnej škále* (na osi y sú zobrazené početnosti) a v *relatívnej škále* (obsah histogramu je rovný jednej). Histogram možno opísať pomocou *frekvenčnej tabuľky*, ktorá obsahuje početnosti a relatívne početnosti. Množstvo intervalov volí príslušný štatistický softvér alebo aj sám užívateľ. Potrebných je aspoň 12 triednych intervalov (ich počet nesmie klesnúť pod 6). Šírka jedného je minimálne $h_{\min} = 0.08(x_{\max} - x_{\min})$. Musí obsahovať minimálne 5 meraní. *Počet tried* histogramu je rovný $k = \log_2 n + 1 \doteq 2 + 3.3 \log_{10} n$ (**Sturgesova formula**; Becker a kol. (1988)), intervaly sú definované ako $\langle x_0, x_1 \rangle, \langle x_1, x_2 \rangle, \dots, \langle x_{n-1}, x_n \rangle$. Šírka intervalov je teda $h = D / (\log_2 n + 1)$ pre realizácie z *normálneho rozdelenia*. Teraz už vlastne nepracujeme s realizáciami x_i , ale so stredmi intervalov $x_i^* = (x_i + x_{i+1}) / 2$. Počty hodnôt n_i , ktoré sa v intervale I_i nachádzajú, sa nazývajú *triedne početnosti*. Pokiaľ realizácie nemajú normálne rozdelenie, treba

použiť robustné algoritmy. Taktiež odľahlé pozorovania môžu dramaticky nafúknuť rozpätie, čo môže spôsobiť nárast šírky intervalov. Preto sa využívajú dva algoritmy ako kompromis medzi výchytkou (bias) a rozptylom realizácií pochádzajúcich z normálneho rozdelenia. Potom šírky triednych intervalov budú $h_1 = 3.49\hat{\sigma}n^{-1/3}$, $\hat{\sigma} = s$ (**Scottova formula**; Becker a kol. (1988)), $h_2 = 2D_Qn^{-1/3}$ (robustnejšia, **Freedman-Diaconisova formula**, ktorá je nezávislá od odľahlých pozorovaní a vyberá menšie intervaly ako Scottova formula; Venables a Ripley (2002)). Pre symetrické rozdelenia platí $h_3 = \lfloor 2\sqrt{n} \rfloor$ alebo $h_4 = \lfloor 2.46 \times (n-1)^{0.4} \rfloor$. Pokiaľ sa neočakáva príliš zošíkmené rozdelenie, je šírka triednych intervalov h konštantná. V prípade komplikovanejších tvarov výberových rozdelení treba zväčšiť počet triednych intervalov alebo použiť špeciálne postupy na hľadanie nekonštantne dlhých triednych intervalov (Meloun a Militký, 2004). Histogram nakreslíme pomocou funkcie `hist(x)`. Argumenty funkcie `hist(x)`:

- `prob=FALSE` (prednastavená hodnota) – v **absolútnej škále**, kde na osi y sú početnosti;
- `prob=TRUE` – v **relatívnej škále** – suma obsahov stĺpcov (obdĺžnikov) je rovná jednej;
- `breaks="Sturges"` (**Sturgesova formula**, prednastavená hodnota), ďalšie možnosti sú "Scott" (**Scottova formula**) a "FD" alebo "Freedman-Diaconis" (**Freedman-Diaconisova formula**);
- `nclass=k`, $k \in \mathbb{N}$ – počet triednych intervalov;
- `plot=FALSE` – ak chceme vypísať číselné detaily, ale nechceme obrázok.

Príklad 131 (histogram a hustota normálneho rozdelenia) Nakreslite histogram najväčšej dĺžky lebky (v mm) `skull.L` u mužov (dáta `one-sample-mean-skull-mf.txt`) a superponujte ho s (a) krivkou hustoty normálneho rozdelenia (červená farba) a (b) krivkou hustoty vypočítanej z dát (čierna farba). Pod bázu histogramu nakreslite tzv. „koberec“, ktorý charakterizuje polohu realizácií.


Riešenie v  (pozri obrázok 26)

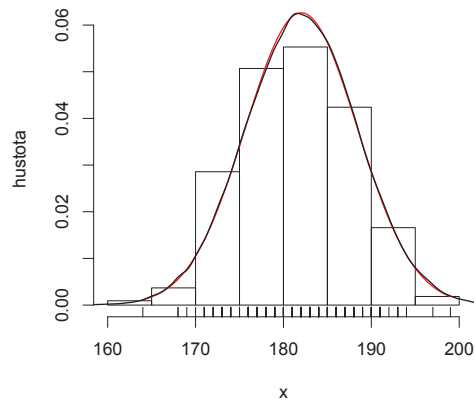
```

333 DATA <- read.table("one-sample-mean-skull-mf.txt",header=TRUE)
334 attach(DATA)
335 x <- na.omit(skull.L[sex=="m"]) # odstranenie chybajúceho pozorovania NA
336 priemer <- mean(x)
337 SD <- sd(x)
338 # 100 kvantilov od minima po maximum x a hodnoty hustoty v nich
339 kvant <- seq(min(x),max(x),length=100)
340 # doplnenie hustoty normalneho rozdelenia
341 hust <- dnorm(kvant,mean = priemer, sd=SD)
342 hist(x,ylim=c(0,range(hust)[2]),prob=TRUE,main="",xlab="x",ylab="hustota")
343 lines(kvant,hust,col="red",lwd=2)
344 # alternativne cez MC simuláciu
345 x.seq <- rnorm(100000,mean=mean(x),sd=sd(x))
346 lines(density(x.seq),lwd=2)
347 # koberec
348 rug(x)

```

Príklad 132 (dva histogramy) Nakreslite dva histogramy tak, aby sa svojimi bázami dotýkali. Aplikujte na dáta `two-samples-means-birth.txt`.

Riešenie v  (pozri obrázok 27)

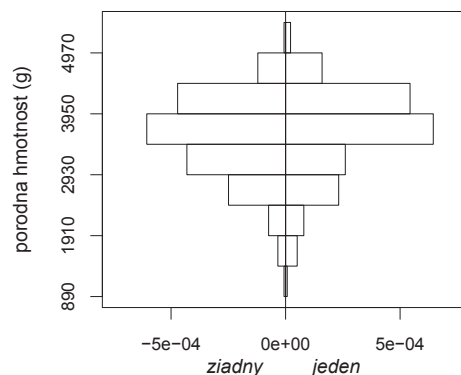


Obr. 26: Histogram so superponovanou krivkou hustoty normálneho rozdelenia (červená farba) a hustoty vypočítanej z dát (čierna farba); pod histogramom je tzv. „koberec“

```


349 DATA <- read.table("two-samples-means-birth.txt",header=TRUE)
350 names(DATA)
351 # [1] "o.sib.N" "birth.W"
352 attach(DATA) # zmena typu objektu na faktor
353 o.sib.N.faktor <- as.factor(o.sib.N) # hladiny faktora
354 # oznacenie hladin a ich zmena
355 levels(o.sib.N.faktor)
356 o.sib.N.faktor1 <- factor(o.sib.N.faktor,labels=c("ziadny","jeden"))
357 rozsah <- range(birth.W)
358 library(Hmisc) # nacistanie kniznice
359 xx <- histbackback(split(birth.W,o.sib.N.faktor1),
360                    probability=TRUE,xlab="",ylab="",axes=FALSE)
361 title(ylab="porodna_hmotnost(g)")
362 axis(1,cex.axis=0.9)
363 axis(2,at=seq(0,8,length=5),labels=seq(rozsah[1],rozsah[2],length=5),
364      cex.axis=0.9,las=1)
365 mtext("ziadny",side=1,line=2,at=-mean(xx$left),font=3)
366 mtext("jeden",side=1,line=2,at=mean(xx$right),font=3)

```



Obr. 27: Dva histogramy s priloženými bázami (základňami)

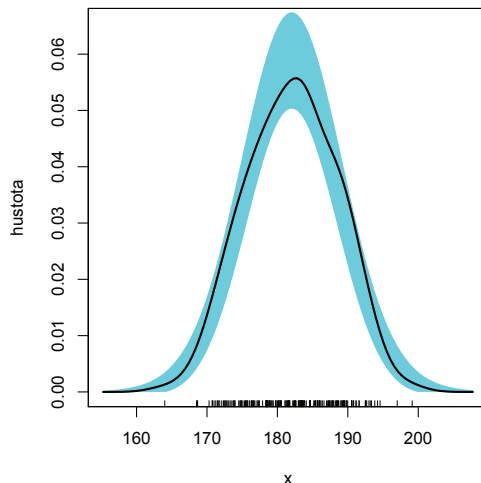
Príklad 133 (MC simulácia hustoty normálneho rozdelenia) Superponujte hustotu najväčšej dĺžky lebky (v mm) `skull.L` (dáta `one-sample-mean-skull-mf.txt`) s hustotou normálneho rozdelenia vypočítanou pomocou MC simulácie s 95% pásom spoľahlivosti normálneho rozdelenia so strednou hodnotou rovnou aritmetickému priemeru a rozptylom rovným výberovému rozptylu (Bowman a Azzalini, 1997).

Riešenie v  (pozri obrázok 28)

```

367 DATA <- read.table("one-sample-mean-skull-mf.txt",header=TRUE)
368 attach(DATA)
369 x <- na.omit(skull.L[sex=="m"])
370 library(sm)
371 yy <- sm.density(x, model="Normal")
372 windows(5,5)
373 par(mar=c(4.5,4.5,1,1))
374 sm.density(x,model="Normal",ylim=c(0,0.0655),ylab="hustota",cex.lab=1.2)
375 lines(yy$eval.points,yy$estimate,lwd=2)

```



Obr. 28: Hustota so superponovanou hustotou normálneho rozdelenia v podobe 95% pásom spoľahlivosti

3.6.6 Empirická distribučná funkcia

Histogram kumulatívnych početností (súčtový histogram) – namiesto početností budeme nad jednotlivými intervalmi I_i zakreslovať obdĺžniky s výškou rovnajúcou sa príslušným kumulatívnym početnostiam $N_i = \sum_{j=1}^i n_j$. Kumulatívne relatívne početnosti definujeme ako N_i/n , čomu zodpovedá **empirická distribučná funkcia**, definovaná pre zvolené číslo x ako relatívna početnosť v intervale $(-\infty, x)$, teda ako N_i -tina hodnôt x_i menších alebo rovných ako x , t.j. (Wasserman, 2006)


$$\hat{F}_n(x) = \frac{\#x_i < x}{n} = \sum_{i=1}^n I(x_i < x)/n,$$

kde $I(\cdot)$ je indikačná funkcia. Obrázok empirickej distribučnej funkcie nakreslíme pomocou funkcie `plot(ecdf(x),verticals = TRUE,do.points=FALSE)`. Argumenty funkcie `plot(ecdf(x))`:

- `verticals = FALSE` je prednastavená hodnota; ak `verticals = TRUE`, je nakreslená schodovitá funkcia;

- `do.points = TRUE` je prednastavená hodnota; vkreslí do obrázka aj body, v ktorých je funkcia počítaná;
- `col.01line="gray70"` (prednastavená hodnota) – farba horizontálnych priamok v bodoch 0 a 1.

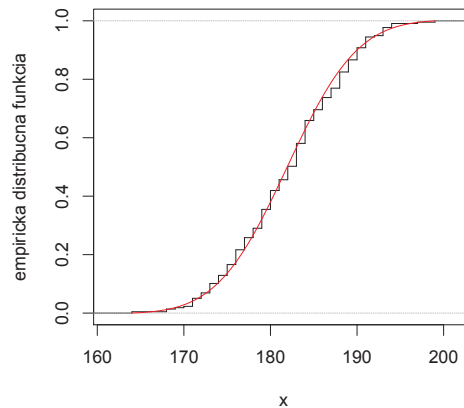
Príklad 134 (distribučná funkcia a hustota normálneho rozdelenia) *Nakreslite empirickú distribučnú funkciu najväčšej dĺžky lebky (v mm) `skull.L` u mužov (dáta `one-sample-mean-skull-mf.txt`) a superponujte ju s krivkou distribučnej funkcie normálneho rozdelenia červenou farbou.*

Riešenie v  (pozri obrázok 29)

```

376 DATA <- read.table("one-sample-mean-skull-mf.txt",header=TRUE)
377 attach(DATA)
378 x <- na.omit(skull.L[sex=="m"])
379 priemer <- mean(x)
380 SD <- sd(x)
381 ROZP <- range(x)
382 # zoradených 200 hodnot (ekvidistantne) medzi min a max x
383 x1 <- seq(ROZP[1],ROZP[2],length=200)
384 # teoretická CDF (normalne rozdelenie)
385 y <- dnorm(x1,mean=priemer,sd=SD)
386 y <- cumsum(y)/sum(y)
387 # zobrazenie krivky distribucnej funkcie
388 # teoretickeho normalneho rozdelenia
389 plot(ecdf(x),verticals=TRUE,do.points=FALSE,xlab="x",
390      ylab="empiricka_distribucna_funkcia",main="")
391 lines(x1,y,col="red",lwd=2)

```



Obr. 29: Emirická distribučná funkcia superponovaná krivkou distribučnej funkcie normálneho rozdelenia (červená farba)

3.6.7 Krabicový diagram

Krabicový diagram – predstavuje grafické znázornenie päťčíselného súhrnu, t.j. zobrazuje v poradí zdola nahor hodnoty x_{\min} , $\tilde{x}_{0.25}$, $\tilde{x}_{0.50}$, $\tilde{x}_{0.75}$, x_{\max} . Umožňuje tiež doplnenie hodnoty aritmetického

priemeru a tak zvýrazniť prípadné odchýlky od normality, identifikovať symetriu rozdelenia medzi kvartilmi, symetriu rozdelenia v koncoch rozdelenia, či odhaliť odľahlé pozorovania. Ak $\tilde{x}_{0.50} = \bar{x}$ ide o **symetrické rozdelenie**, ak $\tilde{x}_{0.50} < \bar{x}$, ide o **pravostranne zošikmené rozdelenie**, ak $\tilde{x}_{0.50} > \bar{x}$, ide o **ľavostranne zošikmené rozdelenie**. Často sa používa na grafické porovnanie dvoch a viacerých skupín. *Šírka krabičiek* je proporčná k odmocnine z rozsahu výberového súboru \sqrt{n} . Ak hovoríme o **krabicových diagramoch so zárezom**, ide o zárez charakterizujúci 95% *empirický interval spoľahlivosti* (pozri kap. Testovanie hypotéz) mediánu $\tilde{\mu}$, ozn. (d, h) , kde

$$d = \tilde{x} - 1.57 \frac{D_Q}{\sqrt{n}}, h = \tilde{x} + 1.57 \frac{D_Q}{\sqrt{n}}.$$

Odhadom rozptylu mediánu je $\hat{\sigma}_{\tilde{x}}^2 = D_Q/1.349$, kde vo všeobecnosti platí (pre akékoľvek rozdelenie pravdepodobnosti), že

$$\sigma_{\tilde{x}}^2 = \frac{1}{4nf^2(\tilde{x})},$$

kde f je hustota rozdelenia pravdepodobnosti (Casella a Berger, 2002). Pre normálne rozdelenie bude platiť $\sigma_{\tilde{x}}^2 = \sigma_x^2 \frac{\pi}{2n}$, kde $\tilde{X} \sim N(\tilde{\mu}, \sigma_x^2)$.

Obrázok krabicového diagramu nakreslíme pomocou funkcie `boxplot(x)`.


Argumenty funkcie `boxplot(x)`:

- `varwidth` je argument relatívnej šírky jednotlivých krabičiek; prednastavená hodnota je `FALSE` s rovnakou šírkou všetkých krabičiek; ak ju zmeníme na `TRUE`, šírka krabičiek bude odpovedať druhej odmocnine z počtu pozorovaní;
- `notch=TRUE` znamená zobrazenie zárezov krabičiek, ktoré odpovedajú 95% intervalom spoľahlivosti pre medián (prednastavená hodnota je `FALSE`);
- `col` určuje farbu vnútri krabičiek;
- `border` určuje farbu hraníc krabičiek;
- `names` je vektor pomenovaní pre jednotlivé zobrazované skupiny, ak ho vynecháme, použijú sa názvy z atribútu `names` z dátového rámca;
- `pch` typ bodu na zobrazenie odľahlých pozorovaní; prednastavená hodnota je 1;
- `horizontal=FALSE` je prednastavená hodnota;
- `plot=FALSE` – ak chceme iba vypísať číselné charakteristiky (prednastavená hodnota je `TRUE`).

Najčastejšie používanou kombináciou argumentov je `boxplot(x, varwidth=TRUE, notch=TRUE, outpch=16)`.

Hodnoty aritmetického priemeru sa do krabicových diagramov dokresľujú pomocou príkazu `points(priemer, pch=16, col="red")`.

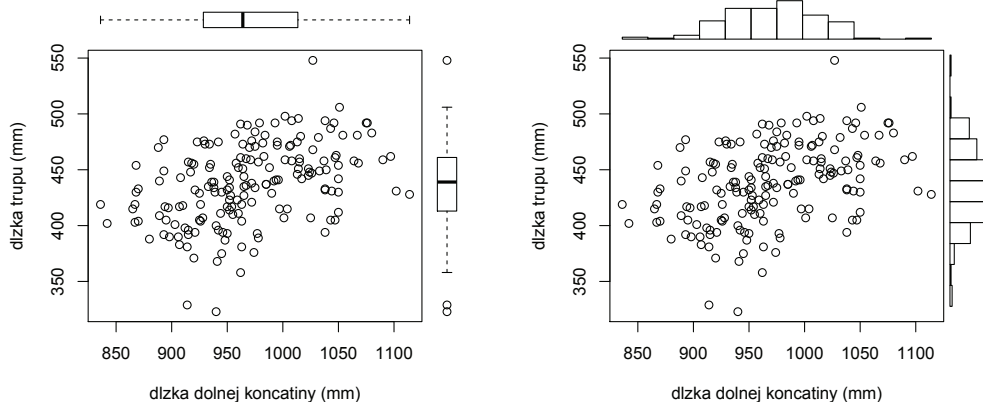
Príklad 135 (bodový graf a marginálne grafy) *Naprogramujte bodový graf dvoch premenných (1) s krabicovými diagramami a (2) histogramami pre marginálne dáta. Aplikujte na dáta `two-samples-correlations-trunk.txt`.*

Riešenie v  (pozri obrázok 30)

```

392 DATA <- read.table("two-samples-correlations-trunk.txt",header=TRUE)
393 names(DATA) # [1] "sex" "lowex.L" "tru.L"
394 attach(DATA)
395 par(fig=c(0,0.9,0,0.9))
396 plot(lowex.L, tru.L, xlab="dlzka_dolnej_koncatiny_(mm)",
397      ylab="dlzka_trupu_(mm)")
398 par(fig=c(0,0.9,0.55,1),new=TRUE)
399 boxplot(lowex.L,horizontal=TRUE,axes=FALSE)
400 par(fig=c(0.65,1,0,0.9),new=TRUE)
401 boxplot(tru.L,axes=FALSE)
402
403 H <- hist(tru.L,plot=FALSE)
404 k1 <- 0
405 k2 <- 5
406 par(fig=c(0,0.9,0,0.9))
407 plot(lowex.L,tru.L,xlab="dlzka_dolnej_koncatiny_(mm)",
408      ylab="dlzka_trupu_(mm)")
409 par(fig=c(0,0.9,0.55,1), new=TRUE)
410 plot(NULL,type="n",ylim=c(0,max(H$counts)+k1),xlim=c(range(H$breaks)),
411      xlab="",ylab="",main="",bty="n",axes=FALSE)
412 rect(H$breaks[1:(length(H$breaks)-1)]+k1,0+k1,
413      H$breaks[2:length(H$breaks)]+k1,H$counts+k1)
414 par(fig=c(0.65,1,0,0.9),new=TRUE)
415 plot(NULL,type="n",xlim=c(0,max(H$counts)+k2),ylim=c(range(H$breaks)),
416      xlab="",ylab="",main="",bty="n",axes=FALSE)
417 rect(0+k2,H$breaks[1:(length(H$breaks)-1)]+k2,H$counts+k2,
418      H$breaks[2:length(H$breaks)]+k2)

```



Obr. 30: Bodový graf s marginálnymi krabicovými diagramami (vľavo) a s histogramami (vpravo)

3.6.8 Kvantilový diagram

Kvantilový diagram (*qq-diagram*) – zobrazuje body so súradnicami $[\Phi^{-1}(i/(n+1)), x_{(i)}]$, kde $\Phi^{-1}(p)$ je kvantilová funkcia normovaného normálneho rozdelenia definovaná nasledovne


$$\Pr(Z \leq \Phi^{-1}(p)) = p.$$

Šikmosť $b_1 > 0$ sa prejaví v podobe **konvexného** usporiadania hodnôt, $b_1 < 0$ sa prejaví v podobe **konkávneho** usporiadania hodnôt. Taktiež zviditeľňuje **dĺžku „chvostov“** (ľavého a pravého konca krivky) rozdelenia, pričom **esovitým usporiadaním** bodov sa prejavujú krátke „chvosty“ a

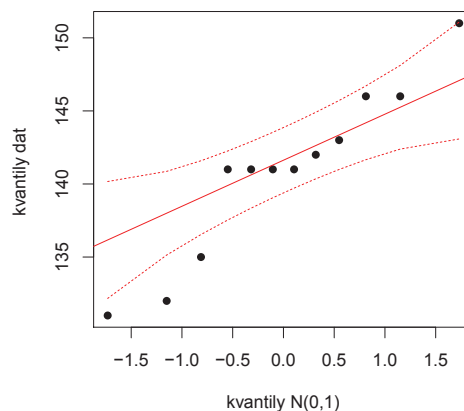
inverzne esovitým usporiadaním bodov dlhé „chvosty“. Zreľná je tiež prípadná bimodalita rozdelenia. Štatisticky možno testovať normalitu rozdelenia pomocou simulácií z $N(0, 1)$ a vytvorením 95% **Atkinsonovej obálky spoľahlivosti**. Pre normálne rozdelenie bude platiť $\sigma_{\tilde{x}_p}^2 = \sigma_x^2 \frac{\pi^2}{24 \ln n}$, kde $\tilde{X}_p \sim N(\tilde{\mu}_p, \sigma_{\tilde{x}_p}^2)$.

Kvantilový diagram vytvoríme pomocou funkcie `qqnorm(x)` a `qqline(x)`. Druhá funkcia dokreslí do grafu priamku prechádzajúcu bodmi charakterizujúcimi prvý a tretí kvartil teoretických a empirických kvantilov.

Príklad 136 (Atkinsonova obálka qq-diagramu) *Nakreslite 95% Atkinsonovu obálku spoľahlivosti (Atkinson, 1981; Flack a Flores, 1989). Aplikujte na dáta výšky 10-ročných dievčat.*


Riešenie v  (pozri obrázok 31)


```
420 library(car)
421 qqPlot(x, distribution="norm", xlab="kvantily_N(0,1)", ylab="kvantily_dat",
422        main="", envelope=.95, col.lines="red", lwd=2, pch=16, cex=1, grid=FALSE)
```



Obr. 31: qq-diagram normovanej spojitej premennej výška 10-ročných dievčat (mm) so superponovanou obálkou normálneho rozdelenia

3.7 Príklady zo štatistickej grafiky

Príklad 137 (všetky základné grafy pre jeden výber ako funkcia) *Naprogramujte v  do jedného obrázka 2×2 štvoricu nasledovných grafov (1) histogram v relatívnej škále so superponovanou krivkou hustoty normálneho rozdelenia, (2) krabicový diagram so zakresleným priemerom, (3) empirickú (kumulatívnu) distribučnú funkciu superponovanú s teoretickou distribučnou funkciou normálneho rozdelenia a (4) qq-diagram so superponovanou qq-priamkou. Aplikujte na dáta výška 10-ročných dievčat.*

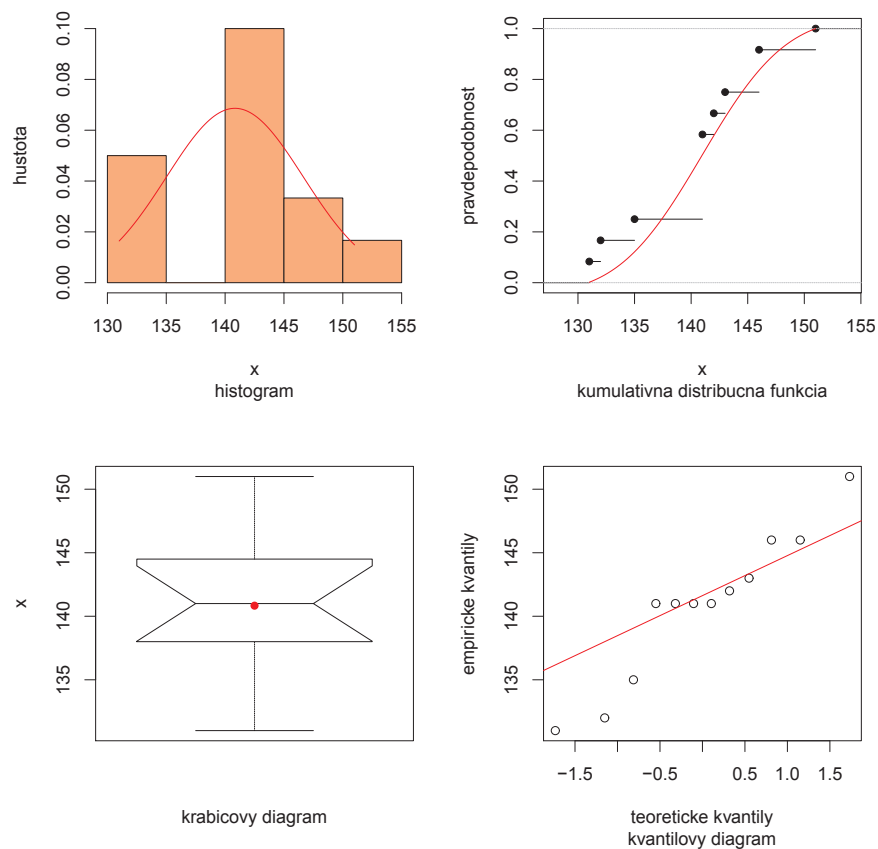
Riešenie v  (pozri obrázok 32)

```
423 "grafy.jeden.vyber" <- function(x){
424 par(mfcol=c(2,2),mar=c(5,5,3,1)) # nastavenie typu okna a jeho okrajov
425 priemer <- mean(x)
```

```

426 SD <- sd(x)
427 kvant <- seq(min(x),max(x),length=100)
428 hust <- dnorm(kvant,priemer,SD)
429 ROZP <- range(x)
430 y1 <- seq(ROZP[1],ROZP[2],length=200)
431 y <- dnorm(y1,mean=mean(y1),sd=sd(y1))
432 y <- cumsum(y)/sum(y)
433 # histogram a teoreticka hustota
434 hist(x,prob=TRUE,main="",xlab="x",ylab="hustota",col="lightsalmon")
435 title(sub="histogram")
436 lines(kvant,hust,col="red",lwd=2)
437 # krabicovy diagram
438 boxplot(x,varwidth=TRUE,notch=TRUE,outpch=16,xlab="krabicovy_diagram",ylab="x")
439 points(priemer,pch=16,col="red")
440 # empiricka a teoreticka CDF
441 plot(ecdf(x),main="",xlab="x",ylab="pravdepodobnost")
442 lines(y1,y,col="red",lwd=2)
443 title(sub="kumulativna_distribucna_funkcia")
444 # kvantilovy diagram
445 qqnorm(x,main="",xlab="",ylab="")
446 qqline(x,col="red",lwd=2)
447 title(sub="kvantilovy_diagram",xlab="teoreticke_kvantily",
448       ylab="empiricke_kvantily")
449 }
450 # graf
451 grafy.jeden.vyber(x)

```



Obr. 32: Základná štvorica grafov pre spojitú premennú výška 10-ročných dievčat (mm)

Príklad 138 (všetky základné grafy pre dva výbery ako funkcia) V \mathbb{R} naprogramujte do jedného obrázka 1×3 trojicu nasledovných grafov (1) superponované krivky hustôt, (2) superponované krivky empirických (kumulatívnych) distribučných funkcií, (3) krabicové diagramy so zakresleným priemerom. Aplikujte na dáta *two-samples-means-birth.txt*.

Riešenie v \mathbb{R} (pozri obrázok 33)

```

452 "grafy.dva.vybery" <- function(x1,x2){
453 hust.1 <- density(x1)$y
454 x1.sekv <- density(x1)$x
455 CDF.1 <- cumsum(hust.1)/sum(hust.1)
456 hust.2 <- density(x2)$y
457 x2.sekv <- density(x2)$x
458 CDF.2 <- cumsum(hust.2)/sum(hust.2)
459 obe.sekv <- c(x1.sekv,x2.sekv)
460 obe.hust <- c(hust.1,hust.2)
461 obe.CDF <- c(CDF.1,CDF.2)
462 ## dve hustoty
463 par(mfcol=c(1,3))
464 plot(obe.sekv,obe.hust,type="n",xlab="x",ylab="hustota",sub="hustoty")
465 lines(x1.sekv,hust.1,col="red",lwd=2)
466 lines(x2.sekv,hust.2,col="blue",lwd=2)
467 legend("topright",c("hust.1","hust.2"),lty=c(1,1),col=c("red","blue"),lwd=2)
468 ## dve CDF
469 plot(obe.sekv,obe.CDF,type="n",xlab="x",ylab="pravdepodobnost",
470      sub="kumulativna_distribucna_funkcia")
471 lines(x1.sekv,CDF.1,col="red",lwd=2)
472 lines(x2.sekv,CDF.1,col="blue",lwd=2)
473 legend("bottomright",c("CDF.1","CDF.2"),lty=c(1,1),col=c("red","blue"),lwd=2)
474 x <- c(x1,x2)
475 kod <- c(rep("x1",length(x1)),rep("x2",length(x2)))
476 ## dva krabicove diagramy
477 boxplot(x ~ kod,varwidth=TRUE,notch=TRUE,outpch=16,sub="krabicove_diagramy")
478 priem <- tapply(x,kod,mean)
479 points(priem,pch=16,col=c("red","blue"))
480 }
481 # grafy
482 DATA <- read.table("two-samples-means-birth.txt",header=TRUE)
483 names(DATA) # [1] "o.sib.N" "birth.W"
484 attach(DATA)
485 o.sib.N.faktor <- as.factor(o.sib.N) # zamena typu objektu na faktor
486 levels(o.sib.N.faktor) # kontrola hladin faktora
487 grafy.dva.vybery(birth.W[o.sib.N.faktor==0],birth.W[o.sib.N.faktor==1])

```

Príklad 139 (stĺpcový diagram) (a) Vypočítajte podmienené pravdepodobnosti (pre každý riadok kontingenčnej tabuľky); t.j. za predpokladu súčinového multinomického rozdelenia. Premennú v riadkoch budeme označovať X (prediktor) a premennú v stĺpcoch ako Y (závisle premenná).

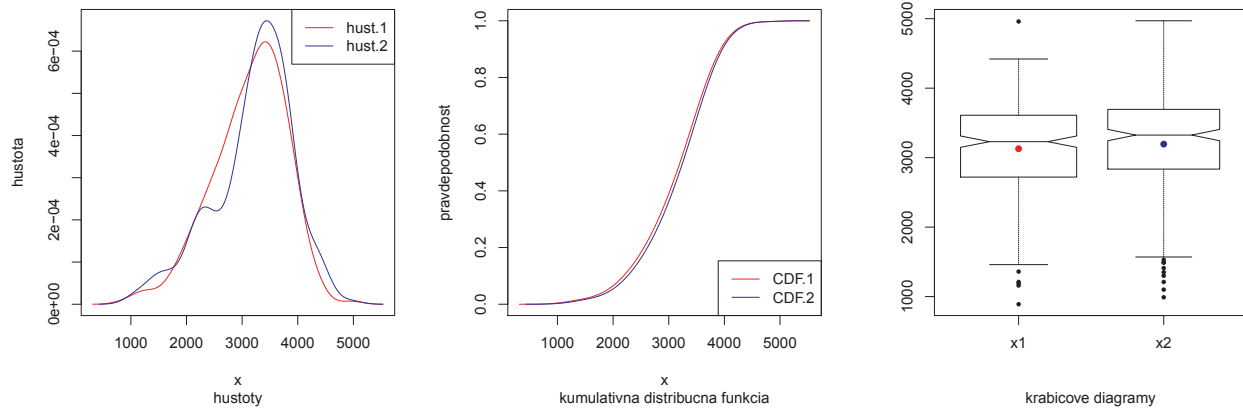
(1) Y : farba dúhovky, X : radiálne útvary v štruktúre dúhovky (dáta: *multinom-iris-color.txt*),

(2) Y : zakončenie troch hlavných dlaňových línií, X : farba vlasov (dáta: *multinom-palmar-lines.txt*),

(3) Y : príľahlosť ušného laloka, X : pohlavie (dáta: *multinom-earlobe.txt*),

(4) Y : krvná skupina, X : mesto (dáta: *multinom-blood-groups.txt*).

(b) Nakreslite stĺpcové diagramy podmienených pravdepodobností pre (1) až (4).

Obr. 33: Základná trojica grafov pre spojitú premennú pre dáta `two-samples-means-birth.txt`

Príklad 140 (štatistická grafika) Nakreslite (1) histogram v relatívnej škále so superponovanou krivkou hustoty normálneho rozdelenia, (2) krabicový diagram so zakresleným priemerom, (3) empirickú (kumulatívnu) distribučnú funkciu superponovanú s teoretickou distribučnou funkciou normálneho rozdelenia a (4) qq-diagram so superponovanou qq-priamkou pre nasledujúce premenné:

(a) stranový rozdiel vertikálneho priemeru diafýzy kľúčnej kosti (`simd.R` a `simd.L`; v mm) na pravej a ľavej strane tela (dáta: `paired-means-clavicle2.txt`),

(b) najväčšia výška mozgovne (`skull.pH`; v mm; dáta: `one-sample-correlation-skull-mf.txt`) a

(c) morfológická výška tváre (`face.H`; v mm; dáta: `one-sample-correlation-skull-mf.txt`).

Príklad 141 (bodový graf a krabicové diagramy) Nakreslite bodový graf spolu s krabicovými diagramami pre marginálne dáta pre premenné vertikálny priemer diafýzy kľúčnej kosti na pravej a ľavej strane tela (`simd.R` a `simd.L`; v mm; dáta: `paired-means-clavicle2.txt`).

Príklad 142 (bodový graf a krabicové diagramy) Nakreslite bodový graf spolu s krabicovými diagramami pre marginálne dáta pre

(a) premenné najväčšia výška mozgovne (`skull.pH`; v mm) a morfológická výška tváre (`face.H`; v mm) a

(b) ich z-skóre (dáta: `one-sample-correlation-skull-mf.txt`).

Príklad 143 (krabicové diagramy) Nakreslite krabicové diagramy pre obe pohlavia pre nasledovné premenné:

(a) dĺžky lebky (`skull.H`; v mm; `sex`; dáta: `two-samples-means-skull.txt`),

(b) dĺžka dolnej končatiny (`lowex.L`; v mm; `sex`; dáta: `two-samples-correlations-trunk.txt`),

(c) dĺžka trupu (`tru.L`; v mm; `sex`; dáta: `two-samples-correlations-trunk.txt`).

Príklad 144 (krabicové diagramy) *Nakreslite krabicové diagramy pre nasledovné premenné:*

(a) *výška hornej časti tváre pre všetky populácie (upface.H; v mm; pop; dáta: anova-means-skull.txt);*

(b) *najväčšia dĺžka klúčnej kosti pravej strany (cla.L; v mm; population; dáta: more-samples-variances-clavicle.txt).*