

## 2 Intervalové rozložení četností

### 2.1 Jednorozměrné intervalové rozložení četností

Jestliže v jednorozměrném datovém souboru je počet variant znaku  $X$  blízký rozsahu souboru  $n$ , pak četnosti přiřazujeme nikoliv jednotlivým variantám, ale třídicím intervalům  $(u_1; u_{j+1}), \dots, (u_r; u_{r+1})$ . Hovoříme pak o intervalovém rozložení četností.

Označení:

- $(u_j; u_{j+1})$  –  $j$ -tý třídicí interval znaku  $X$ ,  $j = 1, \dots, r$ ,
- $d_j = u_{j+1} - u_j$  – délka  $j$ -tého třídicího intervalu znaku  $X$ ,
- $x_{[j]} = \frac{u_j + u_{j+1}}{2}$  – střed  $j$ -tého třídicího intervalu znaku  $X$ .

Třídicí intervaly volíme nejčastěji stejně dlouhé. Jejich počet určíme např. pomocí Sturgesova pravidla:  $r \approx 1 + 3,3 \log_{10} n$ , kde  $n$  je rozsah souboru.

Názvy četností jsou podobné jako v 1.2. Navíc se používá četnostní hustota  $j$ -tého třídicího intervalu:  $f_j = \frac{p_j}{d_j}$ . Četnosti a četnostní hustotu zapisujeme do tabulky rozložení četností (viz tabulka 1).

Tabulka 1: Tabulka rozložení četností znaku  $X$

$(u_j; u_{j+1})$	$d_j$	$x_{[j]}$	$n_j$	$p_j$	$N_j$	$F_j$	$f_j$
$(u_1; u_2)$	$d_1$	$x_{[1]}$	$n_1$	$p_1$	$N_1$	$F_1$	$f_1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$(u_r; u_{r+1})$	$d_r$	$x_{[r]}$	$n_r$	$p_r$	$N_r$	$F_r$	$f_r$

Jednorozměrné intervalové rozložení četností graficky znázorníme pomocí histogramu. Nad třídicími intervaly sestrojíme soustavu obdélníků, jejichž plocha je rovna relativním četnostem (a výšky tedy četnostním hustotám) jednotlivých intervalů. Schodovitá čára shora omezující histogram je grafem hustoty četnosti:

$$f(x) = \begin{cases} f_j & \text{pro } u_j < x \leq u_{j+1}, j = 1, \dots, r, \\ 0 & \text{jinak.} \end{cases}$$

Pomocí hustoty četnosti zavedeme intervalovou empirickou distribuční funkci:  $F(x) = \int_{-\infty}^x f(t) dt$ .

#### Příklad 2.1. Řešený příklad

Načtete datový soubor `16-anova-head.txt`. Za předpokladu, že znak  $X$  popisuje délku hlavy v mm (`head.L`) u mužů (a) zjistíte hranice třídicích intervalů; (b) vytvoříte tabulku rozložení četností; (c) nakreslete histogram, resp. graf intervalové empirické distribuční funkce  $F(x)$ .

#### Řešení příkladu 2.1

Datový soubor načteme příkazem `read.delim()`. Z načtených dat vybereme pomocí operátoru `[]` pouze řádky týkající se mužů a sloupec `head.L`. Z výběru odstraníme chybějící údaje příkazem `na.omit()`. Ke stanovení hranic třídicích intervalů potřebujeme znát rozsah datového souboru a minimální a maximální naměřenou hodnotu délky hlavy mužů. Rozsah datového souboru zjistíme pomocí příkazu `length()`, minimální a maximální naměřenou hodnotu pomocí příkazu `range()`. Optimální počet třídicích intervalů stanovíme na základě Sturgesova pravidla.

```
1 data <- read.delim('16-anova-head.txt', sep = '\t', dec = '.')
2 head.LM <- na.omit(data[data$sex == 'm', 'head.L'])
3 n <- length(head.LM) # 75
4 range(head.LM) # 180-214
5 r <- round(3.3 * log10(n) + 1) # 7
6 # 214 - 180 = 34 -> 35; 35 / 7 = 5 -> seq(179, 214, by = 5)
7 b.head.LM <- seq(179, 214, by = 5) # 179, 184, ..., 214
```

Na základě Sturgesova pravidla stanovíme optimální počet třídících intervalů  $r = 7$ . Naměřené hodnoty délky hlavy se pohybují v rozmezí 180 až 214 mm. Vzdálenost mezi minimální naměřenou hodnotou sníženou o 1 a maximální naměřenou hodnotou je  $214 - 179 = 35$  mm. Tato vzdálenost je beze zbytku dělitelná počtem třídících intervalů, tj. sedmi. Optimální šířka jednoho třídícího intervalu  $d = \frac{35}{7} = 5$  mm. Hranice třídících intervalů tedy stanovíme jako posloupnost 179, 184, ..., 214 mm. Celkem získáme osm hranic definujících sedm třídících intervalů o optimální šířce 5 mm.

Tabulku rozložení četností dopočítáme pomocí funkce `trc()` implementované v R-skriptu `AS-sbirka-funkce.R`, který je součástí této publikace. R-skript načteme příkazem `source()`. Vstupními argumenty funkce `trc()` jsou vektor naměřených hodnot  $x$  a hranice třídících intervalů `breaks`. Výstupem funkce `trc()` je tabulka rozložení četností pro znak  $X$ .

```
8 source('AS-sbirka-funkce.R')
9 tab.rel.cet <- trc(head.LM, b.head.LM)
```

	dj	xj	nj	pj	Nj	Fj	fj
1	5	181,5	5	0,0667	5	0,0667	0,0133
2	5	186,5	6	0,0800	11	0,1467	0,0160
3	5	191,5	20	0,2667	31	0,4133	0,0533
4	5	196,5	24	0,3200	55	0,7333	0,0640
5	5	201,5	12	0,1600	67	0,8933	0,0320
6	5	206,5	6	0,0800	73	0,9733	0,0160
7	5	211,5	2	0,0267	75	1,0000	0,0053

10  
11  
12  
13  
14  
15  
16  
17

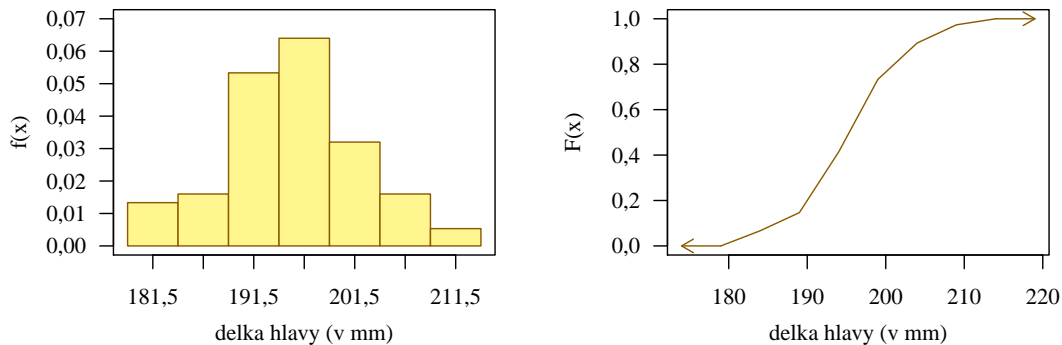
Pro ukázkou si zde uvedeme interpretace vybraných hodnot ze třetího a pátého řádku tabulky rozložení četností: Celkem 20 mužů (26,67%) v datovém souboru mělo délku hlavy v rozmezí 189 až 194 mm, celkem 31 mužů (41,33%) v datovém souboru mělo délku hlavy nejvýše 194 mm. Celkem 12 mužů (16,00%) v datovém souboru mělo délku hlavy v rozmezí 199 až 204 mm, celkem 67 mužů (89,33%) v datovém souboru mělo délku hlavy nejvýše 204 mm.

Histogram vykreslíme pomocí příkazu `hist()` s argumentem `prob = T`. Tento argument zajistí, že obsahem sloupců histogramu budou relativní četnosti a výškou sloupců četnostní hustoty. Hranice třídících intervalů nastavíme pomocí argumentu `breaks`. Okolo grafu dokreslíme rámeček příkazem `box()`. Osu  $x$ , resp. osu  $y$  vykreslíme samostatně pomocí příkazu `axis()` s argumentem `side = 1`, resp. `side = 2`. V souladu se zavedenou konvencí vykreslíme měřítko osy  $x$  ve středech třídících intervalů. Středy třídících intervalů můžeme získat buď ručním výpočtem, nebo jako výstup `mids` funkce `hist()` se specifikací argumentu `plot = F`. Výsledný histogram je zobrazen na obrázku 1 (vlevo).

```
18 h <- hist(head.LM, breaks = b.head.LM, plot = F)
19 hist(head.LM, breaks = b.head.LM, prob = T, include.lowest = F, ylim = c(0, 0.07),
20       col = 'khaki1', border = 'orange4', xlab = 'delka hlavy (v mm)',
21       ylab = 'f(x)', main = '', axes = F)
22 box(bty = 'o')
23 axis(side = 1, at = h$mids)
24 axis(side = 2, las = 1)
```

Graf intervalové empirické distribuční funkce  $F(x)$  vykreslíme pomocí příkazu `plot()` s argumentem `type = 'l'`. Prvním vstupním argumentem příkazu bude vektor hodnot empirické distribuční funkce  $F(x)$  vložený v tabulce rozložení četností `tab.rel.cet` ve sloupci `Fj`. Do grafu dále dokreslíme šipku značící nulovou hodnotu empirické distribuční funkce  $F(x)$  ve všech hodnotách menších než dolní hranice prvního třídícího intervalu ( $u_1 = 179$ ) a dále šipku značící jednotkovou hodnotu funkce  $F(x)$  ve všech hodnotách větších než horní hranice posledního třídícího intervalu ( $u_8 = 214$ ). Šipky vykreslíme pomocí funkce `arrows()`. Výsledný graf intervalové empirické distribuční funkce je zobrazen na obrázku 1 (vpravo).

```
25 Fj <- tab.rel.cet$Fj
26 plot(b.head.LM, c(0, Fj), type = 'l', las = 1, col = 'orange4', xlim = c(174, 219),
27       xlab = 'delka hlavy (v mm)', ylab = 'F(x)')
28 arrows(179, 0, 174, 0, col = 'orange4', length = 0.1)
29 arrows(214, 1, 219, 1, col = 'orange4', length = 0.1)
```



Obrázek 1: Histogram (vlevo) a graf intervalové empirické distribuční funkce  $F(x)$  (vpravo) pro znak  $X$  popisující délku hlavy mužů (v mm)

★

### Příklad 2.2. Neřešený příklad

Načtěte datový soubor `11-two-samples-means-skull.txt`. Za předpokladu, že znak  $X$  popisuje výšku lebky v mm (`skull.H`) u žen (a) zjistěte hranice třídících intervalů; (b) vytvořte tabulku rozložení četností; (c) nakreslete histogram, resp. graf intervalové empirické distribuční funkce  $F(x)$ .

**Výsledky:** (a) optimální počet třídících intervalů  $r = 8$ , optimální šířka jednoho třídícího intervalu  $d = 3$ , hranice třídících intervalů:  $u_1 = 114, u_2 = 117, \dots, u_9 = 138$ ; (b) tabulka rozložení četností viz tabulka 2; (c) histogram, resp. graf intervalové empirické distribuční funkce  $F(x)$  viz obrázek 2.

Tabulka 2: Tabulka rozložení četností

	$d_j$	$x_{[j]}$	$n_j$	$p_j$	$N_j$	$F_j$	$f_j$
1	3,00	115,50	7	0,07	7	0,07	0,02
2	3,00	118,50	12	0,11	19	0,18	0,04
3	3,00	121,50	16	0,15	35	0,33	0,05
4	3,00	124,50	22	0,21	57	0,53	0,07
5	3,00	127,50	27	0,25	84	0,79	0,08
6	3,00	130,50	17	0,16	101	0,94	0,05
7	3,00	133,50	5	0,05	106	0,99	0,02
8	3,00	136,50	1	0,01	107	1,00	0,00

★

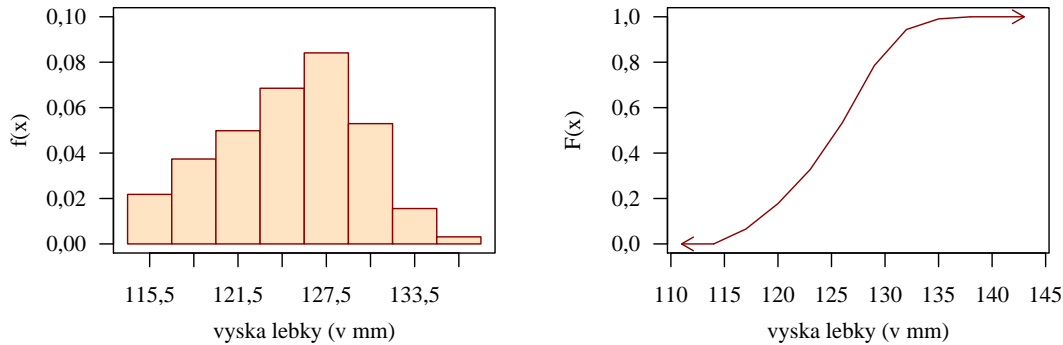
## 2.2 Dvourozměrné intervalové rozložení četností

Ve dvourozměrném datovém souboru o rozsahu  $n$  roztrídíme hodnoty znaku  $X$  do  $r$  třídících intervalů  $(u_j; u_{j+1})$ ,  $j = 1, \dots, r$  s délkami  $d_1, \dots, d_r$  a hodnoty znaku  $Y$  taktéž roztrídíme do  $r$  třídících intervalů  $(v_k; v_{k+1})$ ,  $k = 1, \dots, r$  s délkami  $h_1, \dots, h_r$ . Obdélník  $(u_j; u_{j+1}) \times (v_k; v_{k+1})$  se nazývá  $(j, k)$ -tý dvourozměrný třídící interval.

Názvy četností jsou podobné jako v 1.3. Navíc zavádíme simultánní a marginální četnostní hustoty:

- $f_{jk} = \frac{p_{jk}}{d_j h_k}$  – simultánní četnostní hustota v  $(j, k)$ -tém třídícím intervalu,
- $f_{.j} = \frac{p_{.j}}{d_j}$  – marginální četnostní hustota v  $j$ -tém třídícím intervalu pro znak  $X$ ,
- $f_{.k} = \frac{p_{.k}}{h_k}$  – marginální četnostní hustota v  $k$ -tém třídícím intervalu pro znak  $Y$ .

Kteroukoliv ze simultánních četností zapisujeme do kontingenční tabulky.



Obrázek 2: Histogram (vlevo) a graf intervalové empirické distribuční funkce  $F(x)$  (vpravo) pro znak  $X$  popisující výšku lebky žen (v mm)

Dvourozměrné intervalové rozložení četností graficky znázorníme pomocí stereogramu. Nad dvourozměrnými třídícími intervaly sestrojíme soustavu kvádrů, jejichž objem je roven relativním četnostem (a výšky tedy simultánním četnostním hustotám) jednotlivých intervalů. Schodovitá plocha shora omezující stereogram je grafem simultánní četnostní hustoty:

$$f(x, y) = \begin{cases} f_{jk} \text{ pro } u_j < x \leq u_{j+1}, v_k < y \leq v_{k+1}, j, k = 1, \dots, r, \\ 0 \text{ jinak.} \end{cases}$$

Marginální hustoty četnosti pro znaky  $X$  a  $Y$  odlišíme indexem takto:

$$f_1(x) = \begin{cases} f_{.j} \text{ pro } u_j < x \leq u_{j+1}, j = 1, \dots, r, \\ 0 \text{ jinak,} \end{cases}$$

$$f_2(y) = \begin{cases} f_{.k} \text{ pro } v_k < y \leq v_{k+1}, k = 1, \dots, r, \\ 0 \text{ jinak.} \end{cases}$$

Mezi simultánní hustotou četnosti a marginálními hustotami četnosti platí vztahy:  $f_1(x) = \int_{-\infty}^{\infty} f(x, y) dy$ ,  $f_2(y) = \int_{-\infty}^{\infty} f(x, y) dx$ .

Znaky  $X, Y$  jsou v daném výběrovém souboru četnostně nezávislé při daném dvourozměrném intervalovém rozložení četností, jestliže pro  $\forall j, k = 1, \dots, r$  platí:  $f_{jk} = f_{.j} \cdot f_{.k}$  neboli pro  $\forall (x, y) \in \mathbb{R}^2$ :  $f(x, y) = f_1(x) f_2(y)$ .

### Příklad 2.3. Řešený příklad

Načtete datový soubor 16-anova-head.txt. Za předpokladu, že znak  $X$  popisuje délku hlavy v mm (head.L) a znak  $Y$  popisuje šířku tváře v mm (bizyg.W) u mužů (a) roztříd'te data do dvourozměrných třídících intervalů; (b) vytvořte kontingenční tabulku simultánních absolutních, resp. relativních četností; (c) nakreslete histogram pro znak  $X$ , resp. pro znak  $Y$ ; (d) nakreslete stereogram a dvourozměrný tečkový diagram; (e) zjistěte, zda jsou znaky  $X$  a  $Y$  četnostně nezávislé při daném dvourozměrném intervalovém rozložení četností.

### Řešení příkladu 2.3

Datový soubor načteme příkazem `read.delim()`. Z načtených dat vybereme pomocí operátoru `[]` pouze řádky týkající se mužů a sloupce `head.L` a `bizyg.W`. Z výběru odstraníme chybějící údaje příkazem `na.omit()`. Ke stanovení hranic třídících intervalů potřebujeme znát rozsah dvourozměrného datového souboru a minimální a maximální naměřenou hodnotu každého znaku. Rozsah datového souboru stanovíme příkazem `length()`. Minimální a maximální naměřenou hodnotu znaku  $X$ , resp. znaku  $Y$  zjistíme příkazem `range()`. Optimální počet třídících intervalů každého znaku stanovíme na základě Sturgesova pravidla.

```
30 data <- read.delim('16-anova-head.txt', sep = '\\t', dec = '.')
31 data.M <- data[data$sex == 'm', c('head.L', 'bizyg.W')]
```

```

32 data.M <- na.omit(data.M)
33 head.LM <- data.M$head.L
34 bizyg.WM <- data.M$bizyg.W
35 n <- length(head.LM) # 75
36 range(head.LM) # 180-214
37 range(bizyg.WM) # 113-155
38 r <- round(3.3 * log10(n) + 1) # 7
39 # head.LM: viz priklad 2.1
40 # bizyg.WM: 155 - 112 = 43 / 7 => 49 / 7 = 7 -> seq(110, 159, by = 7)
41 b.head.LM <- seq(179, 214, by = 5)
42 b.bizyg.WM <- seq(110, 159, by = 7)

```

Na základě Sturgesova pravidla stanovíme pro oba znaky  $X$  i  $Y$  optimální počet třídících intervalů  $r = 7$ . Naměřené hodnoty znaku  $X$ , tj. délky hlavy, se pohybují v rozmezí 180 až 214 mm. Analogicky jako v příkladu 2.1 rozdělíme naměřené hodnoty tohoto znaku do 7 třídících intervalů o optimální ekvidistantní šířce  $d = 5$  mm, a to stanovením hranic 179, 184, ..., 214 mm. Naměřené hodnoty znaku  $Y$ , tj. šířky tváře, se pohybují v rozmezí 113 až 155 mm. Vzdálenost mezi minimální naměřenou hodnotou sníženou o 1 a maximální naměřenou hodnotou je  $155 - 112 = 43$  mm. Nejbližší vyšší celé číslo dělitelné beze zbytku počtem třídících intervalů, tj. sedmi, je 49. Optimální šířka jednoho třídícího intervalu  $h = \frac{49}{7} = 7$  mm. Hranice třídících intervalů pro šířku tváře tedy stanovíme jako posloupnost 110, 117, ..., 159 mm. Celkem získáme osm hranic definujících sedm třídících intervalů o optimální šířce 7 mm.

Po stanovení hranic třídících intervalů roztřídíme naměřené hodnoty každého znaku do příslušných třídících intervalů, a to pomocí příkazu `cut()`. Příkaz přiřadí každému měření příslušný třídící interval. Informace o třídícím intervalu přiřazeném každému měření znaku  $X$ , resp. znaku  $Y$  vložíme do proměnné `head.LM.c`, resp. `bizyg.WM.c`. Kontingenční tabulku simultánních absolutních četností vypočítáme pomocí funkce `table()`, jejímiž vstupními argumenty budou vektory `head.LM.c` a `bizyg.WM.c`. Kontingenční tabulku simultánních relativních četností získáme příkazem `prop.table()`.

```

43 head.LM.c <- cut(head.LM, breaks = b.head.LM)
44 bizyg.WM.c <- cut(bizyg.WM, breaks = b.bizyg.WM)
45 KT.abs <- table(head.LM.c, bizyg.WM.c)

```

	(110;117]	(117;124]	(124;131]	(131;138]	(138;145]	(145;152]	(152;159]	
(179;184]	0	0	0	3	1	1	0	46
(184;189]	1	0	0	0	4	0	1	47
(189;194]	0	2	2	7	7	2	0	48
(194;199]	0	0	0	5	12	6	1	49
(199;204]	0	0	1	2	6	2	1	50
(204;209]	0	0	1	0	2	3	0	51
(209;214]	0	0	1	1	0	0	0	52
								53

```

54 KT.rel <- prop.table(KT.abs)

```

	(110;117]	(117;124]	(124;131]	(131;138]	(138;145]	(145;152]	(152;159]	
(179;184]	0,0000	0,0000	0,0000	0,0400	0,0133	0,0133	0,0000	55
(184;189]	0,0133	0,0000	0,0000	0,0000	0,0533	0,0000	0,0133	56
(189;194]	0,0000	0,0267	0,0267	0,0933	0,0933	0,0267	0,0000	57
(194;199]	0,0000	0,0000	0,0000	0,0667	0,1600	0,0800	0,0133	58
(199;204]	0,0000	0,0000	0,0133	0,0267	0,0800	0,0267	0,0133	59
(204;209]	0,0000	0,0000	0,0133	0,0000	0,0267	0,0400	0,0000	60
(209;214]	0,0000	0,0000	0,0133	0,0133	0,0000	0,0000	0,0000	61
								62

Pro ukádku si uvedeme interpretace vybraných hodnot z prvního, druhého a třetího sloupce kontingenční tabulky simultánních absolutních, resp. relativních četností. V datovém souboru se vyskytoval jeden muž (1,33%), jehož délka lebky nabývala hodnoty z intervalu (184; 189) mm a jehož šířka tváře nabývala hodnoty z intervalu (110; 117) mm, dva muži (2,67%), jejichž délka lebky nabývala hodnoty z intervalu (189; 194) mm a jejichž šířka tváře nabývala hodnoty z intervalu (117; 124) mm, dva muži (2,67%), jejichž délka lebky nabývala hodnoty z intervalu (189; 194) mm

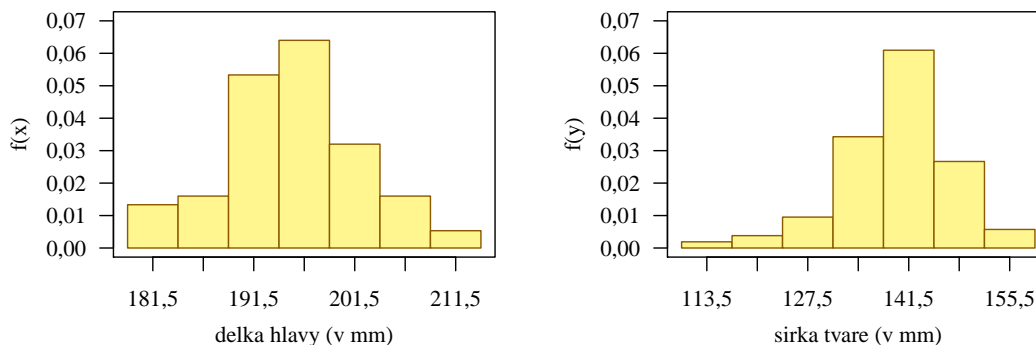
a jejichž šířka tváře nabývala hodnoty z intervalu (124; 131) mm, apod.

Histogram znaku  $X$ , resp. znaku  $Y$  vykreslíme pomocí příkazu `hist()` s argumentem `prob = T`. Hranice třídicích intervalů definujeme explicitně pomocí argumentu `breaks`. Okolo grafu vykreslíme rámeček příkazem `box()`. Měřítko osy  $x$ , resp. osy  $y$  dokreslíme samostatně pomocí příkazu `axis()` s argumentem `side = 1`, resp. `side = 2`. Měřítko osy  $x$  vykreslíme v hodnotách středů třídicích intervalů. Histogramy pro oba znaky jsou zobrazeny na obrázku 3.

```

63 h <- hist(head.LM, breaks = b.head.LM, plot = F)
64 hist(head.LM, breaks = b.head.LM, prob = T, ylim = c(0, 0.07), col = 'khaki1',
65     border = 'orange4', xlab = 'delka hlavy (v mm)', ylab = 'f(x)',
66     main = '', axes = F, include.lowest = F)
67 box(bty = 'o')
68 axis(side = 1, at = h$mids)
69 axis(side = 2, las = 1)
70
71 h <- hist(bizyg.WM, breaks = b.bizyg.WM, plot = F)
72 hist(bizyg.WM, breaks = b.bizyg.WM, prob = T, ylim = c(0, 0.07), col = 'khaki1',
73     border = 'orange4', xlab = 'sirka tvare (v mm)', ylab = 'f(y)',
74     main = '', axes = F, include.lowest = F)
75 box(bty = 'o')
76 axis(side = 1, at = h$mids)
77 axis(2, las = 1)

```



Obrázek 3: Histogram (a) pro znak  $X$  popisující délku hlavy mužů (vlevo); (b) pro znak  $Y$  popisující šířku tváře mužů (vpravo)

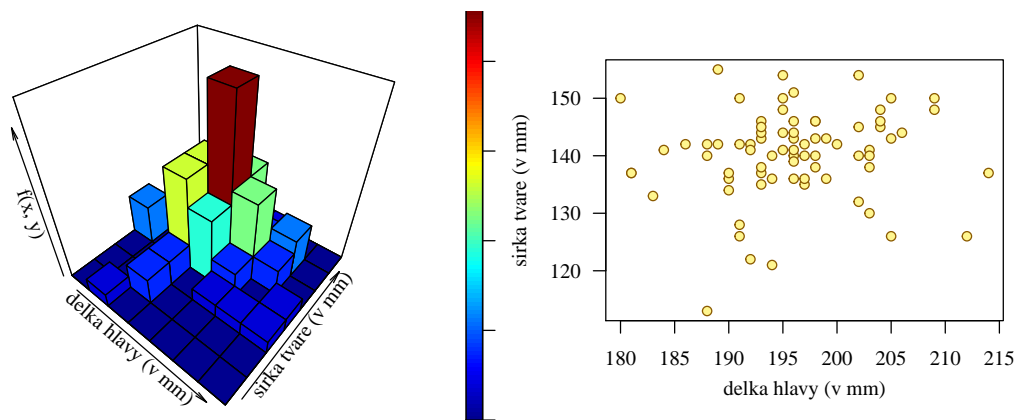
Stereogram vykreslíme pomocí funkce `hist3D()` z knihovny `plot3D`. Jediným povinným vstupním argumentem funkce je matice simultánních četnostních hustot  $f_{jk}$ . Simultánní četnostní hustoty dopočítáme pomocí vzorce  $f_{jk} = \frac{p_{jk}}{d_j h_k}$ , kde  $p_{jk}$  jsou simultánní relativní četnosti (viz proměnná `KT.rel`),  $d_j = 5$  je délka  $j$ -tého třídicího intervalu znaku  $X$  a  $h_k = 7$  je délka  $k$ -tého třídicího intervalu znaku  $Y$ ,  $j, k = 1, \dots, 7$ . Dvourozměrný tečkový diagram vykreslíme příkazem `plot()`. Stereogram i dvourozměrný tečkový diagram jsou zobrazeny na obrázku 4.

```

78 pjk <- KT.rel; dj <- 5; hk <- 7
79 fjk <- pjk / (dj * hk)
80 plot3D::hist3D(z = fjk, border = "black", xlab = 'delka hlavy (v mm)',
81     ylab = 'sirka tvare (v mm)', zlab = 'f(x, y)')
82 plot(head.LM, bizyg.WM, pch = 21, col = 'orange4', bg = 'khaki1', las = 1,
83     xlab = 'delka hlavy (v mm)', ylab = 'sirka tvare (v mm)')

```

Znaky  $X$  a  $Y$  jsou četnostně nezávislé, pokud  $f_{jk} = f_{j.} \cdot f_{.k}$ , kde  $f_{jk}$  jsou simultánní četnostní hustoty,  $f_{j.}$  je marginální četnostní hustota  $j$ -tého třídicího intervalu znaku  $X$  a  $f_{.k}$  je marginální četnostní hustota  $k$ -tého třídicího intervalu



Obrázek 4: (a) Stereogram (vlevo); (b) dvourozměrný tečkový diagram (vpravo) pro znaky  $X$  a  $Y$

znaku  $Y$ ,  $j, k = 1, \dots, 7$ . Simultánní četnostní hustoty  $f_{jk}$  máme vloženy v proměnné `fjk`.

```
84 round(fjk, 4)
```

```
85 format(round(fjk, 4), scientific = F)
```

Marginální četnostní hustoty znaku  $X$  dopočítáme pomocí vzorce  $f_{.j} = \frac{p_{.j}}{d_j}$ , kde  $p_{.j}$  jsou marginální relativní četnosti, které vypočítáme pomocí funkce `apply()` s argumenty `MARGIN = 1` a `FUN = sum` aplikované na matici simultánních relativních četností `pjk`, a  $d_j = 5$  je šířka  $j$ -tého třídicího intervalu znaku  $X$ ,  $j = 1, \dots, 7$ . Marginální četnostní hustoty znaku  $Y$  dopočítáme pomocí vzorce  $f_{.k} = \frac{p_{.k}}{h_k}$ , kde  $p_{.k}$  jsou marginální relativní četnosti, které vypočítáme pomocí funkce `apply()` s argumenty `MARGIN = 2` a `FUN = sum` aplikované na matici simultánních relativních četností `pjk` a  $h_k = 7$  je šířka  $k$ -tého třídicího intervalu znaku  $Y$ ,  $k = 1, \dots, 7$ .

```
86 p.j. <- apply(pjk, 1, sum)
87 p.k <- apply(pjk, 2, sum)
88 f.j. <- p.j. / d.j
89 f.k <- p.k / h.k
90 f.j.f.k <- f.j.%*% t(f.k)
```

	(110;117]	(117;124]	(124;131]	(131;138]	(138;145]	(145;152]	(152;159]	91
(179;184]	0,0000	0,0001	0,0001	0,0005	0,0008	0,0004	0,0001	92
(184;189]	0,0000	0,0001	0,0002	0,0005	0,0010	0,0004	0,0001	93
(189;194]	0,0001	0,0002	0,0005	0,0018	0,0033	0,0014	0,0003	94
(194;199]	0,0001	0,0002	0,0006	0,0022	0,0039	0,0017	0,0004	95
(199;204]	0,0001	0,0001	0,0003	0,0011	0,0020	0,0009	0,0002	96
(204;209]	0,0000	0,0001	0,0002	0,0005	0,0010	0,0004	0,0001	97
(209;214]	0,0000	0,0000	0,0001	0,0002	0,0003	0,0001	0,0000	98

Porovnáním matice simultánních četnostních hustot  $f_{jk}$  s příslušnými součiny marginálních četnostních hustot  $f_{j.f.k}$  vidíme, že znaky  $X$  a  $Y$  nejsou při daném intervalovém rozložení četností četnostně nezávislé. Například  $f_{14} \neq f_{1.f.4}$  ( $0,0011 \neq 0,0005$ ),  $f_{36} \neq f_{3.f.6}$  ( $0,0008 \neq 0,0014$ ),  $f_{65} \neq f_{6.f.5}$  ( $0,0008 \neq 0,0010$ ), apod. ★

### Příklad 2.4. Neřešený příklad

Načtete datový soubor 05-one-sample-correlation-skull-mf.txt. Za předpokladu, že znak  $X$  popisuje největší výšku mozkovny v mm (skull.pH) a znak  $Y$  popisuje morfologickou výšku tváře v mm (face.H) u žen (a) rozřídíte data do dvourozměrných třídicích intervalů; (b) vytvoříte kontingenční tabulku simultánních absolutních, resp. relativních četností; (c) nakreslete histogram pro znak  $X$ , resp. pro znak  $Y$ ; (d) nakreslete stereogram a dvourozměrný tečkový diagram; (e) zjistíte, zda jsou znaky  $X$  a  $Y$  četnostně nezávislé při daném dvourozměrném intervalovém rozložení četností.

**Výsledky:** (a) optimální počet třídicích intervalů  $r = 7$ ; znak  $X$ : optimální šířka jednoho třídicího intervalu  $d_j = 4$ ,  $j = 1, \dots, 7$ , hranice třídicích intervalů:  $u_1 = 117, u_2 = 121, \dots, u_8 = 145$ ; znak  $Y$ : optimální šířka jednoho třídicího intervalu  $h_k = 4$ ,  $k = 1, \dots, 7$ , hranice třídicích intervalů:  $u_1 = 93, u_2 = 97, \dots, u_8 = 121$  mm; (b) kontingenční tabulka simultánních absolutních, resp. relativních četností viz tabulka 3, resp. tabulka 4; (c) histogram znaku  $X$ , resp. znaku  $Y$  viz obrázek 5; (d) stereogram a dvourozměrný tečkový diagram viz obrázek 6; (e) mezi znaky  $X$  a  $Y$  neexistuje při daném dvourozměrném intervalovém rozložení četností četnostní nezávislost (viz porovnání simultánních četnostních hustot  $f_{jk}$  (tabulka 5) se součiny marginálních četnostních hustot  $f_{j \cdot} \cdot f_{\cdot k}$  (tabulka 6)).

Tabulka 3: Tabulka simultánních absolutních četností  $n_{jk}$

	(93; 97)	(97; 101)	(101; 105)	(105; 109)	(109; 113)	(113; 117)	(117; 121)
(117; 121)	0	1	0	0	1	0	0
(121; 125)	1	0	3	2	1	1	1
(125; 129)	0	1	8	4	6	2	1
(129; 133)	2	2	4	7	6	3	0
(133; 137)	0	2	1	3	6	0	3
(137; 141)	0	1	2	2	0	0	0
(141; 145)	0	0	0	0	0	0	1

Tabulka 4: Tabulka simultánních relativních četností  $p_{jk}$

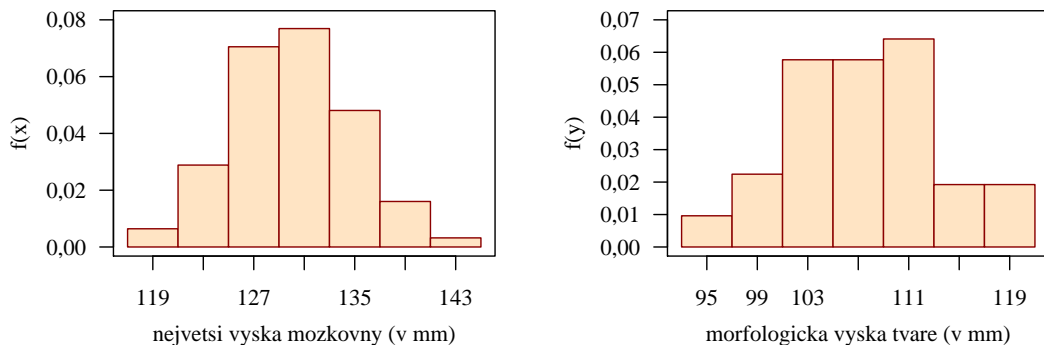
	(93; 97)	(97; 101)	(101; 105)	(105; 109)	(109; 113)	(113; 117)	(117; 121)
(117; 121)	0,0000	0,0128	0,0000	0,0000	0,0128	0,0000	0,0000
(121; 125)	0,0128	0,0000	0,0385	0,0256	0,0128	0,0128	0,0128
(125; 129)	0,0000	0,0128	0,1026	0,0513	0,0769	0,0256	0,0128
(129; 133)	0,0256	0,0256	0,0513	0,0897	0,0769	0,0385	0,0000
(133; 137)	0,0000	0,0256	0,0128	0,0385	0,0769	0,0000	0,0385
(137; 141)	0,0000	0,0128	0,0256	0,0256	0,0000	0,0000	0,0000
(141; 145)	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0128

Tabulka 5: Tabulka simultánních četnostních hustot  $f_{jk}$

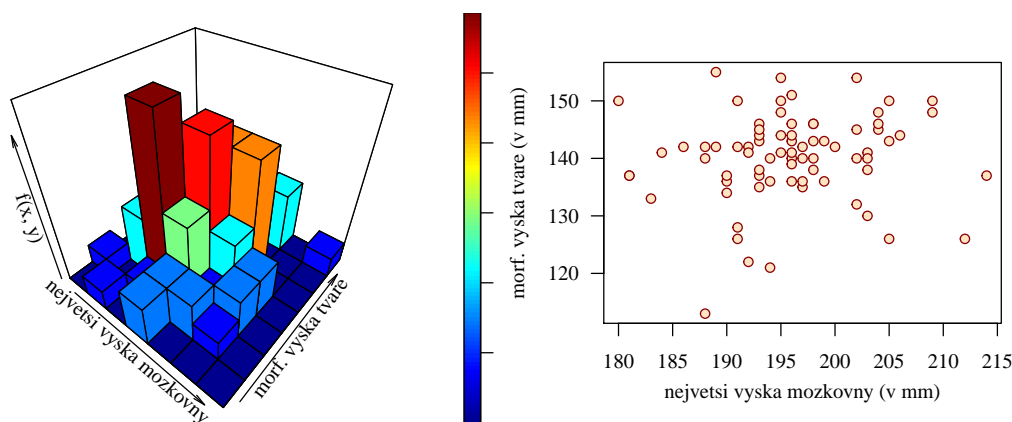
	(93; 97)	(97; 101)	(101; 105)	(105; 109)	(109; 113)	(113; 117)	(117; 121)
(117; 121)	0,0000	0,0004	0,0000	0,0000	0,0004	0,0000	0,0000
(121; 125)	0,0004	0,0000	0,0011	0,0007	0,0004	0,0004	0,0004
(125; 129)	0,0000	0,0004	0,0029	0,0015	0,0022	0,0007	0,0004
(129; 133)	0,0007	0,0007	0,0015	0,0026	0,0022	0,0011	0,0000
(133; 137)	0,0000	0,0007	0,0004	0,0011	0,0022	0,0000	0,0011
(137; 141)	0,0000	0,0004	0,0007	0,0007	0,0000	0,0000	0,0000
(141; 145)	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0004







Obrázek 5: Histogram (a) pro znak  $X$  popisující největší výšku mozkovny žen (vlevo); (b) pro znak  $Y$  popisující morfologickou výšku tváře žen (vpravo)



Obrázek 6: (a) Stereogram (vlevo); (b) dvourozměrný tečkový diagram (vpravo) pro znaky  $X$  a  $Y$

Tabulka 6: Tabulka součinů marginálních četnostních hustot  $f_{j.}f_{.k}$

	(93; 97)	(97; 101)	(101; 105)	(105; 109)	(109; 113)	(113; 117)	(117; 121)
(117; 121)	0,0000	0,0001	0,0002	0,0002	0,0002	0,0001	0,0001
(121; 125)	0,0001	0,0003	0,0008	0,0008	0,0008	0,0003	0,0003
(125; 129)	0,0003	0,0007	0,0019	0,0019	0,0021	0,0006	0,0006
(129; 133)	0,0003	0,0008	0,0020	0,0020	0,0023	0,0007	0,0007
(133; 137)	0,0002	0,0005	0,0013	0,0013	0,0014	0,0004	0,0004
(137; 141)	0,0001	0,0002	0,0004	0,0004	0,0005	0,0001	0,0001
(141; 145)	0,0000	0,0000	0,0001	0,0001	0,0001	0,0000	0,0000