

Statistika

semestrální projekt

**Zuzana Foltýnová
324542**

U 96 náhodně vybraných studentů VŠE v Praze byly zjištěny následující údaje:

Pohlaví (0 – žena, 1 – muž)

Výška (tělesná výška v cm)

Hmotnost (tělesná hmotnost v kg)

Známka (známka z matematiky v 1. semestru)

Úkol 1.

Zjistěte absolutní a relativní četnosti proměnných Pohlaví a Známka, přičemž pro proměnnou známka zjistěte též kumulativní absolutní a relativní četnosti. Pro proměnnou Pohlaví vytvořte sloupkový diagram, pro proměnnou Známka polygon četností.

Řešení:

I. Proměnná Pohlaví - absolutní a relativní četnost

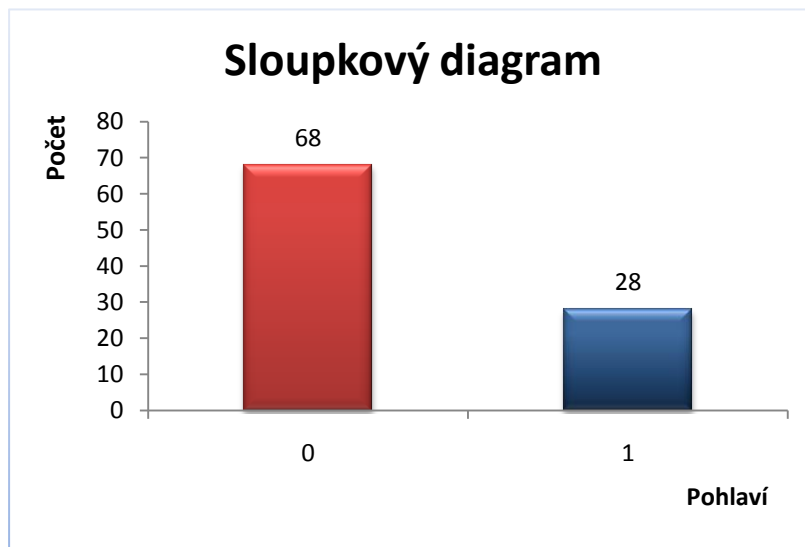
Tabulka četností		
$x_{[j]}$	n_j	p_j
1	28	0,29
0	68	0,71

$x_{[j]}$ - proměnná pohlaví

n_j - absolutní četnost varianty $x_{[j]}$

p_j - relativní četnost varianty $x_{[j]}$

II. Proměnná Pohlaví - sloupkový diagram



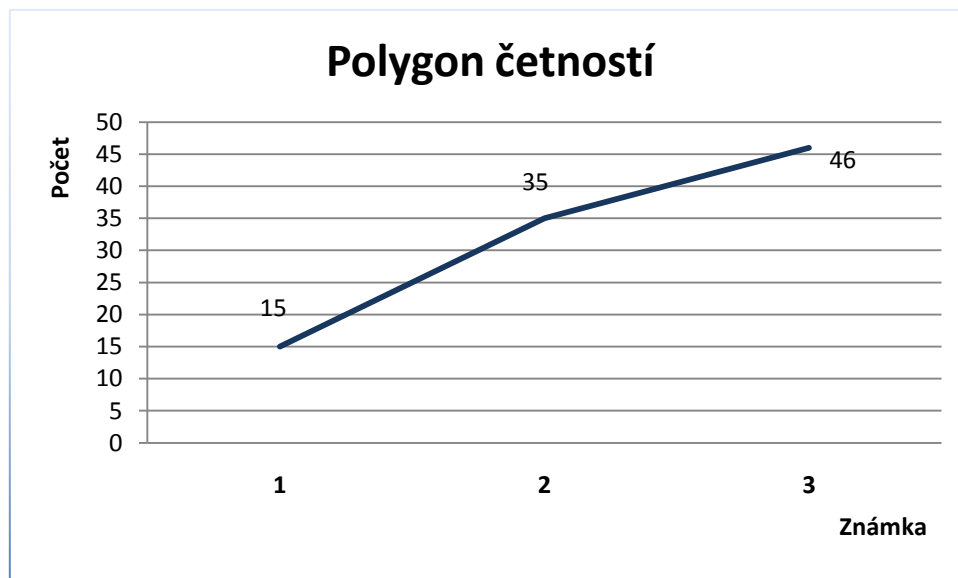
Komentář: Ve zkoumaném datovém souboru jsou ženy (0) zastoupeny v nadpoloviční většině oproti mužům (1).

III. Proměnná Znamka - tabulka četností (absolutní a relativní četnost, kumulativní absolutní a relativní četnost)

Tabulka četností				
$x_{[j]}$	n_j	p_j	N_j	F_j
1	15	0,16	15	0,16
2	35	0,36	50	0,52
3	46	0,48	96	1

$x_{[j]}$ - proměnná pohlaví
 n_j - absolutní četnost varianty $x_{[j]}$
 p_j - relativní četnost varianty $x_{[j]}$
 N_j - kumulativní absolutní četnost prvních j variant
 F_j - kumulativní relativní četnost prvních j variant

IV. Proměnná Znamka - polygon četností



Komentář: Nejčastěji obdržena známka je 3, s 46 výskyty. Nejméně obdržena známka je 1, a to pouze s 15 výskyty.

Úkol 2.

Pro proměnné Pohlaví a Známkou sestavte kontingenční tabulky absolutních a relativních četností, sloupcově a řádkově podmíněných relativních četností. Kolik procent žen má z matematiky jedničku? Kolik procent studentů, kteří mají jedničku, jsou muži?

Řešení:

I. Proměnná Pohlaví a Známkou - kontingenční tabulka absolutních a relativních četností

Kontingenční tabulka absolutních četností					
-	y	1	2	3	$n_{j.}$
x	n_{jk}				
0		11	27	30	68
1		4	8	16	28
$n_{.k}$		15	35	46	n = 96

x - proměnná pohlaví
y - proměnná známka
 n_{jk} - absolutní četnost dvojice $(x_{[j]}, y_{[k]})$
 $n_{j.}$ - absolutní četnost varianty $x_{[j]}$
 $n_{.j}$ - absolutní četnost varianty $y_{[k]}$

Kontingenční tabulka relativních četností					
-	y	1	2	3	$p_{j.}$
x	p_{jk}				
0		0,12	0,28	0,31	0,71
1		0,04	0,08	0,17	0,29
$p_{.k}$		0,16	0,36	0,48	1,00

x - proměnná pohlaví
y - proměnná známka
 p_{jk} - relativní četnost dvojice $(x_{[j]}, y_{[k]})$
 $p_{j.}$ - relativní četnost varianty $x_{[j]}$
 $p_{.j}$ - relativní četnost varianty $y_{[k]}$

II. Proměnná Pohlaví a Známkou - kontingenční tabulka sloupcově a řádkově podmíněných relativních četností

Kontingenční tabulka sloupcově podmíněných relativních četností				
-	y	1	2	3
x	$p_{j(k)}$			
0		0,73	0,77	0,65
1		0,27	0,23	0,35
Σ		1,00	1,00	1,00

x - proměnná pohlaví
y - proměnná známka
 $p_{j(k)}$ - sloupcově podmíněná relativní četnost varianty $x_{[j]}$ za předpokladu $y_{[k]}$

Komentář: 27% studentů, kteří mají jedničku, jsou muži.

Kontingenční tabulka řádkově podmíněných relativních četností					
-	y	1	2	3	Σ
x	n_{jk}				
0		0,16	0,40	0,44	1,00
1		0,14	0,29	0,57	1,00

x - proměnná pohlaví
y - proměnná známka
 $p_{j(k)}$ - řádkově podmíněná relativní četnost varianty $y_{[k]}$ za předpokladu $x_{[j]}$

Komentář: 16 % žen má z matematiky jedničku

Úkol 3.

Podle Sturgersova pravidla stanovte optimální počet třídících intervalů pro proměnné Výška a Hmotnost a nakreslete jejich histogramy, a to

- pro celý soubor
- pro ženy
- pro muže.

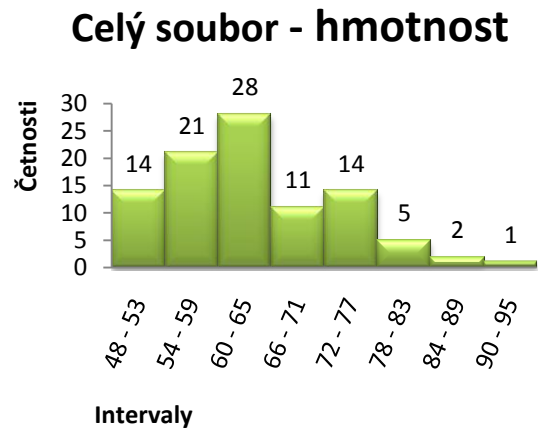
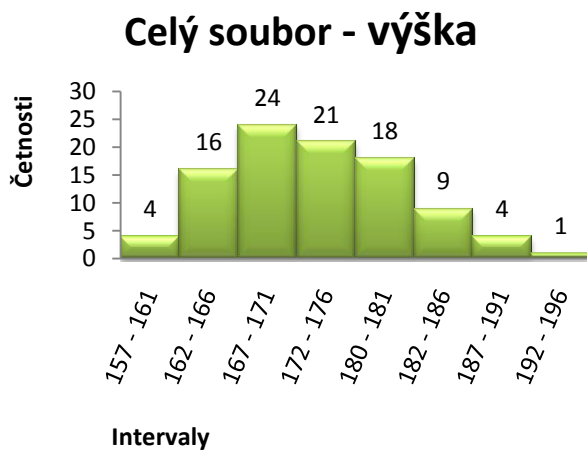
Řešení:

I. Histogramy pro celý soubor

Celkový rozsah souboru je 96, tedy podle Sturgesova pravidla je optimální počet třídících intervalů $r = 1 + 3,3 * \log(96) \approx 8$.

Rozsah intervalů (výška) = $(\max - \min) / 8 = (193 - 160) / 8 \approx 5$

Rozsah intervalů (hmotnost) = $(\max - \min) / 8 = (91 - 48) / 8 \approx 6$



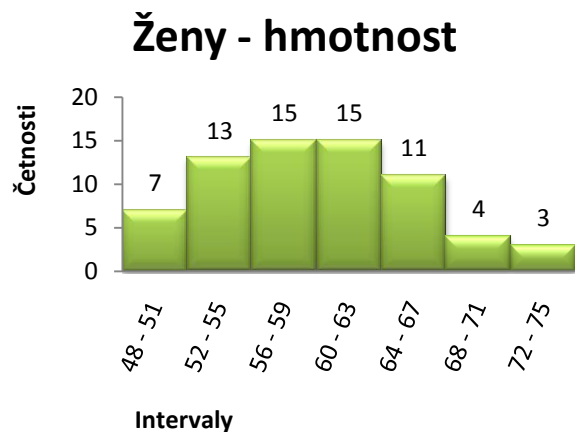
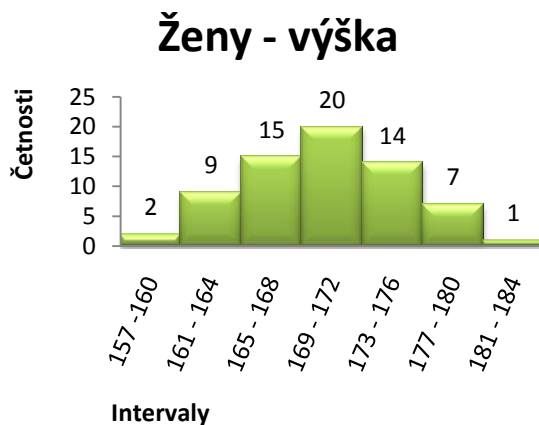
Komentář: Ve zkoumaném souboru převládají osoby s výškou v intervalu 167 - 171 cm či váhou mezi 60 - 65 kg. Naopak nejméně jsou zastoupeny osoby s výškou pod 161 cm a nad 192 cm a s váhou větší než 84 kg.

II. Histogramy pro ženy

Počet žen v souboru je 68. Podle Sturgesova pravidla je optimální počet třídících intervalů $r = 1 + 3,3 * \log(68) \approx 7$

Rozsah intervalů (výška) = $(\max - \min) / 7 = (182 - 160) / 7 \approx 4$

Rozsah intervalů (hmotnost) = $(\max - \min) / 7 = (76 - 48) / 7 \approx 4$



Komentář: Ve zkoumaném souboru převládají ženy s výškou v intervalu 169 - 172 cm či váhou mezi 59 - 63 kg. Naopak nejméně jsou zastoupeny ženy s výškou pod 160 cm a nad 181 cm a s váhou větší než 72 kg.

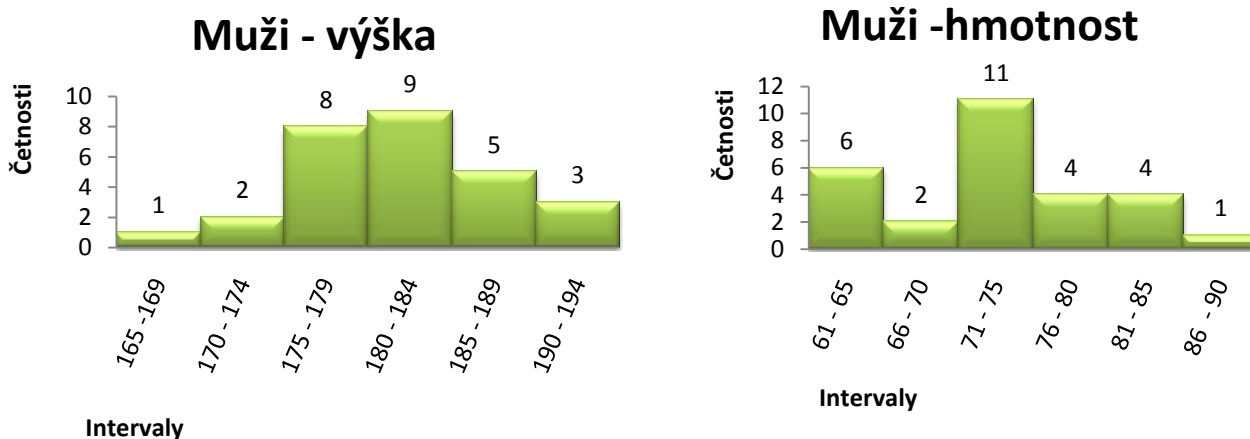
III. Histogramy pro muže

Počet mužů v souboru je 28. Podle Sturgesova pravidla je optimální počet třídících intervalů

$$r = 1 + 3,3 \cdot \log(28) \approx 6$$

$$\text{Rozsah intervalů (výška)} = (\max - \min)/6 = (193 - 168)/6 \approx 5$$

$$\text{Rozsah intervalů (hmotnost)} = (\max - \min)/6 = (91 - 61)/6 \approx 5$$



Komentář: Co se týče mužů, nejvíce jsou zastoupeni ti s výškou v intervalu 180 - 184 cm či s váhou mezi 71 - 75 kg. Naopak nejméně jsou zastoupeni muži s výškou pod 169 cm či s váhou větší než 86 kg.

Úkol 4.

Vypočítejte minimum, maximum, medián, průměr, směrodatnou odchylku, šikmost a špičatost proměnných Výška a Hmotnost

- pro celý soubor
- pro ženy
- pro muže.

Řešení:

I. Pro celý soubor

Hodnoty pro celý soubor								
-	počet	minimum	maximum	medián	průměr	směrodatná odchylka	šikmost	špičatost
výška	96	160	192	172,0	173,19	7,336	0,428	-0,299
hmotnost	96	48	90	62,5	63,38	9,327	0,588	-0,135

II. Pro ženy

Hodnoty pro ženy								
-	počet	minimum	maximum	medián	průměr	směrodatná odchylka	šikmost	špičatost
výška	68	160	181	170	169,94	5,063	0,017	-0,638
hmotnost	68	48	75	58	59,26	6,392	0,322	-0,446

III. Pro muže

Hodnoty pro muže								
-	počet	minimum	maximum	medián	průměr	směrodatná odchylka	šikmost	špičatost
výška	28	168	192	181	181,07	5,894	-0,249	0,089
hmotnost	28	61	90	73	73,36	7,670	0,172	-0,533

Komentář: Směrodatná odchylka a průměr mohou být silně ovlivněny extrémními hodnotami. Je-li šikmost kladná, rozložení dat má prodloužený pravý konec, mluvíme o kladně zešikmeném rozložení. V záporně zešikmeném rozložení má naopak rozložení prodloužený levý konec. Je-li špičatost záporná, jedná se o ploché rozložení dat, je-li kladná, jde o strmé rozložení dat. V případě normálního rozložení je pak špičatost rovna 0.

Úkol 5.

Vypočítejte a interpretujte Pearsonův koeficient korelace proměnných Výška a Hmotnost

- pro celý soubor
- pro ženy
- pro muže.

Řešení:

I. Pearsonův koeficient korelace proměnných Výška a Hmotnost pro celý soubor

$$r_{12} = 0,739702$$

Komentář: Existuje silná kladná korelace mezi proměnnými Výška a Hmotnost pro celý soubor. Můžeme tedy říci, že čím vyšší jsou hodnoty jedné proměnné, tím vyšší jsou hodnoty druhé a naopak.

II. Pearsonův koeficient korelace proměnných Výška a Hmotnost pro ženy

$$r_{12} = 0,5082377$$

Komentář: Existuje středně silná kladná korelace mezi proměnnými Výška a Hmotnost pro ženy. Můžeme tedy říci, že čím vyšší jsou hodnoty jedné proměnné, tím vyšší jsou hodnoty druhé a naopak.

III. Pearsonův koeficient korelace proměnných Výška a Hmotnost pro muže

$$r_{12} = 0,487723$$

Komentář: Existuje středně silná kladná korelace mezi proměnnými Výška a Hmotnost pro muže. Můžeme tedy říci, že čím vyšší jsou hodnoty jedné proměnné, tím vyšší jsou hodnoty druhé a naopak.

Úkol 6.

Najděte rovnici regresní přímky vyjadřující závislost proměnné Hmotnost na proměnné Výška. Jaký je index determinace a co vyjadřuje? Jaká je predikovaná hodnota hmotnosti pro výšku 175 cm? Nalezenou regresní přímku zakreslete do dvourozměrného tečkového diagramu.

Řešení:

I. Rovnice regresní přímky

Rovnice: $\text{hmotnost} = -99,502 + 0,9405 \cdot \text{výška}$

Komentář: Zvětší-li se výška o 1 bod, hmotnost se zvětší v průměru o 0,9405 bodu.

II. Index determinace

Index determinace je **0,5472**

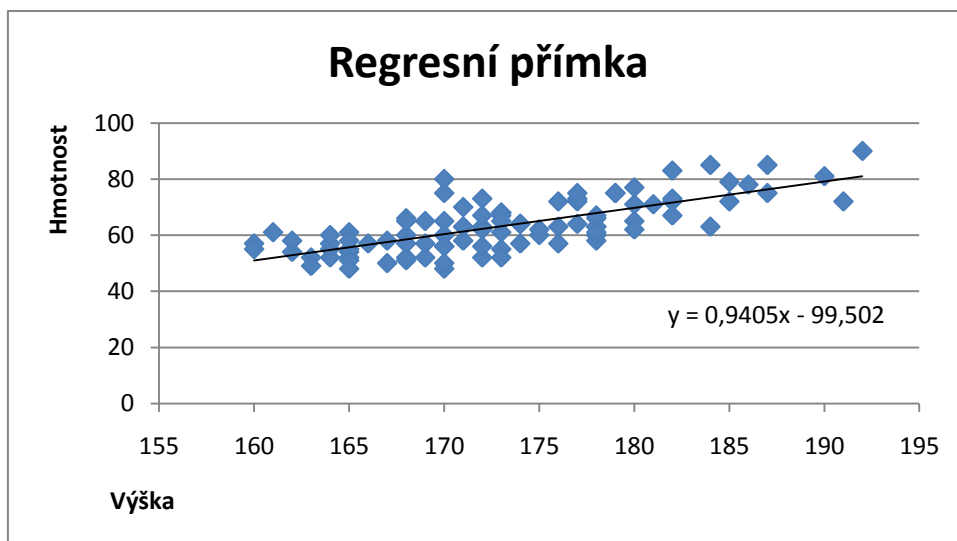
Komentář: Model regresní přímky vysvětluje variabilitu proměnné hmotnost z 54%.

III. Predikovaná hodnota hmotnosti pro výšku 175

Vypočteme, že predikovaná hodnota hmotnosti pro výšku 175 je:

$\text{hmotnost} = -99,502 + 0,9405 \cdot 175 = \mathbf{65,0855}$

IV. Dvourozměrný tečkový diagram



Komentář: Z grafu je patrné, že regresní přímka je vhodná na modelování dané závislosti, jelikož body jsou přibližně rozmístěny kolem regresní přímky.