

Cvičení 2.: Intervalové rozložení četností, výpočet číselných charakteristik nominálních a ordinálních znaků

Úkol 1.: Datový soubor vysvah.sta obsahuje údaje o hmotnosti (znak X, v kg), výšce (znak Y, v cm) a pohlaví (znak Z, 0 – žena, 1 – muž) 50 náhodně vybraných studentů. Načtete tento soubor do systému STATISTICA. Proměnným X, Y, Z vytvořte návěští „hmotnost“, „výška“ a „pohlaví“. Popište, co u znaku Z znamenají varianty 0, 1. Podle Sturgesova pravidla najděte optimální počet třídících intervalů pro znaky X a Y a vhodně stanovíte meze třídících intervalů.

Návod: Soubor – Otevřít – vybereme příslušný adresář se souborem vysvah.sta – Otevřít. Kurzor nastavíme na X – 2x klikneme myší – Dlouhé jméno hmotnost – OK, kurzor nastavíme na Y – 2x klikneme myší – Dlouhé jméno výška – OK, kurzor nastavíme na Z – 2x klikneme myší – Dlouhé jméno pohlaví, Text. hodnoty – 0 žena, 1 – muž - OK. Protože případů je 50, podle Sturgesova pravidla je optimální počet třídících intervalů 7. Musíme zjistit minimum a maximum, abychom vhodně stanovili třídící intervaly: Statistika - Základní statistiky/tabulky – Popisné statistiky - OK - Proměnné X,Y – OK – Detailní výsledky – ponecháme zaškrtnuté pouze Minimum&maximum – Výpočet.

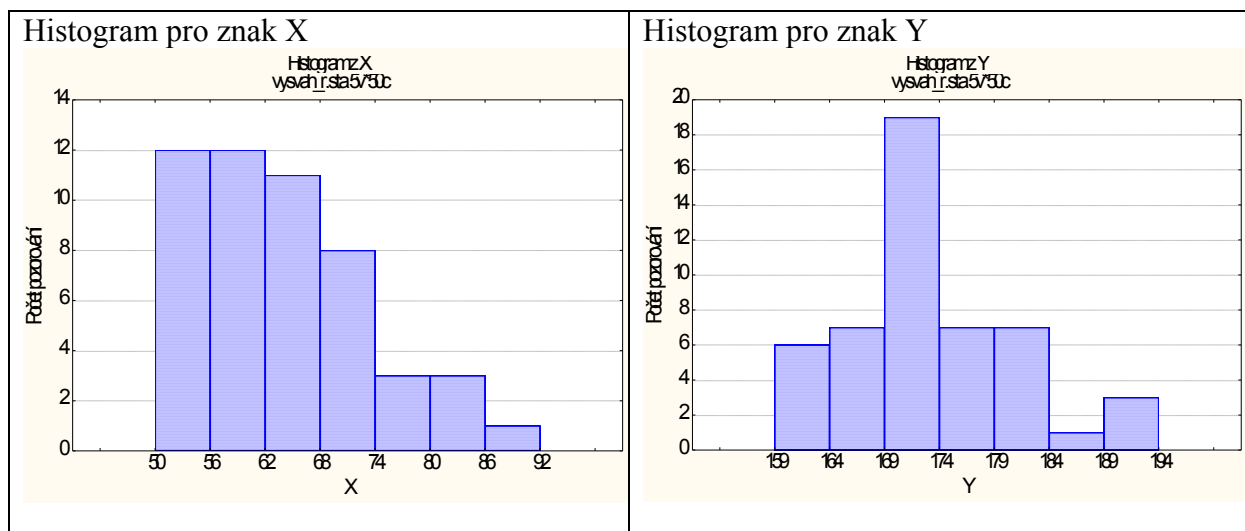
Promě	Popisné statistiky (vys		
	N platn	Minim	Maxim
X	50	51,00	90,00
Y	50	160,0	192,0

Pro X je minimum 51 a maximum 90, tedy dolní mez prvního třídícího intervalu volíme 50, horní mez posledního třídícího intervalu 92. Celkem tedy třídící intervaly pro znak X budou: (50,56>, (56,62>, (62,68>, (68,74>, (74,80>, (80,86>, (86,92>.

Pro Y je minimum 160 a maximum 192, tedy dolní mez prvního třídícího intervalu volíme 159, horní mez posledního třídícího intervalu 194. Celkem tedy třídící intervaly pro znak Y budou: (159,164>, (164,169>, (169,174>, (174,179>, (179,184>, (184,189>, (189,194>.

Úkol 2.: Vytvořte histogram pro X a pro Y.

Návod: Grafy – Histogramy – Proměnné X – vypneme Normální proložení – Detaily – zaškrtneme Hranice – Určit hranice – zvolíme Zadejte hraniční rozmezí – Minimum = 50, Krok = 6, Maximum = 92 - OK – OK. Po vykreslení histogramu lze 2 x klepnout na pozadí grafu a ve volbě Všechny možnosti měnit různé vlastnosti grafu. Analogicky pro Y.



Úkol 3.: Proved'te zakódování hodnot proměnných X a Y do příslušných třídících intervalů. Všem hodnotám proměnné X, které leží v intervalu (50,56>, přiřaďte hodnotu 53 atd. až všem hodnotám proměnné X, které leží v intervalu (86,92>, přiřaďte hodnotu 89. Analogicky pro proměnnou Y, tj. všem hodnotám výšky, které leží v intervalu (159,164>, přiřaďte hodnotu 161,5 atd. až všem hodnotám výšky, které leží v intervalu (189,194> přiřaďte hodnotu 191,5.

Návod: Vytvoříme dvě nové proměnné: Vložit – Přidat proměnné – 2 – Za Y – OK – přejmenujeme je na RX a RY. Nastavíme se kurzorem na RX – Data – Překódovat - vyplníme podmínky pro všech 7 kategorií. (Pozor – podmínky píšeme ve tvaru $X > 50$ and $X \leq 56$ atd.). Pak klepneme na OK.

Analogicky překódujeme hodnoty proměnné Y do proměnné RY.

Úkol 4.: Sestavte kontingenční tabulky absolutních četností (relativních četností, sloupcově a řádkově podmíněných relativních četností) dvourozměrných třídících intervalů pro (X,Y). Graficky znázorněte simultánní absolutní četnosti.

Návod: Při tvorbě kontingenčních tabulek musí být proměnné celočíselné. Proto proměnnou RY vynásobíme 10 (do jejího Dlouhého jména napíšeme =10*RY). Tím vlastně dostaneme středy třídících intervalů pro výšku vyjádřenou v mm.

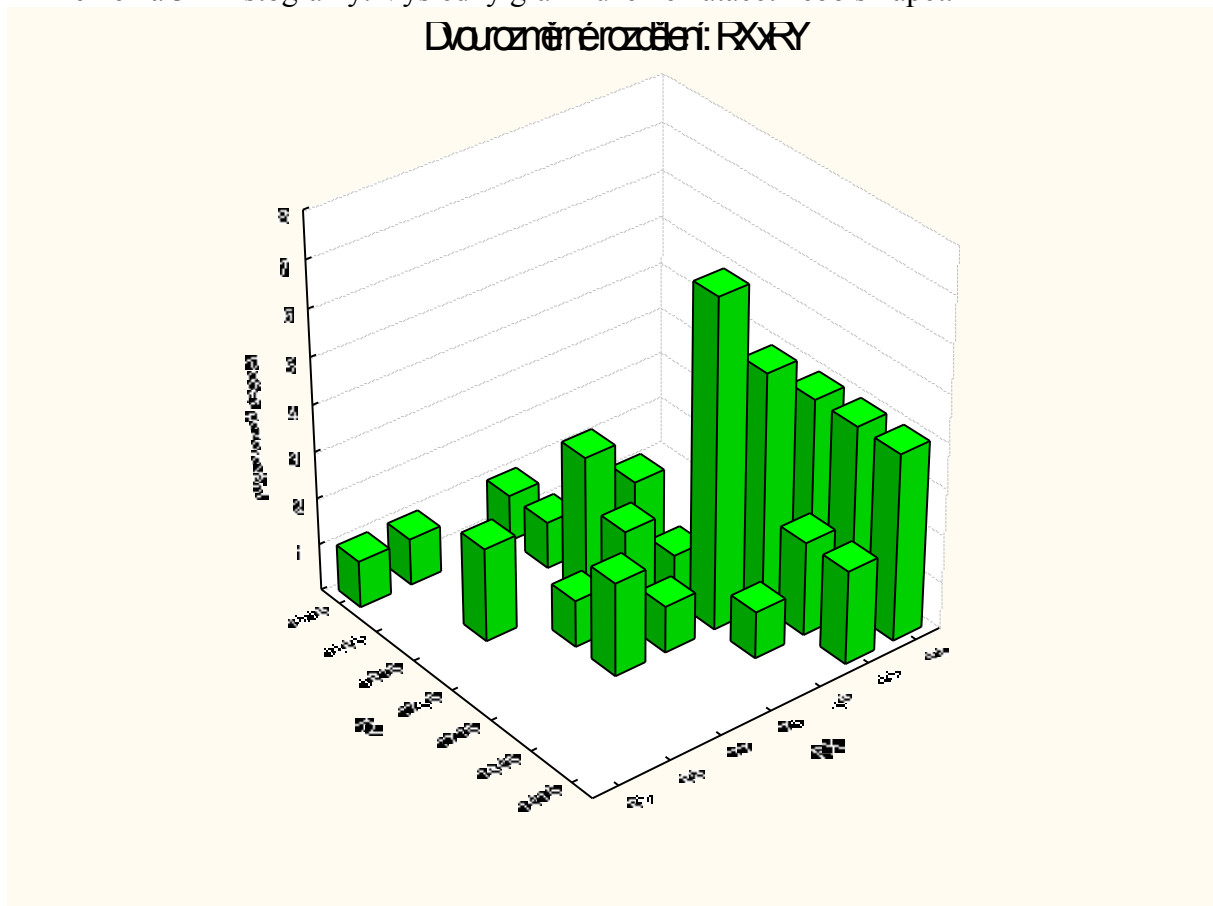
Statistiky – Základní statistiky/tabulky – OK - Kontingenční tabulky – OK – Specif. tabulky - List 1 RX, List 2 RY, OK, Výpočet.

Kontingenční tabulka absolutních četností:

RX	RY 16'	RY 16'	RY 17'	RY 17'	RY 18'	RY 18'	RY 19'	Rad souč
53	4	4	4	0	0	0	0	12
59	2	2	5	3	0	0	0	12
65	0	1	1	1	2	0	0	5
71	0	0	1	2	3	1	0	7
77	0	0	2	1	0	0	0	3
83	0	0	0	0	2	0	1	3
89	0	0	0	0	0	0	1	1
VS.SKI	6	7	11	7	7	1	3	51

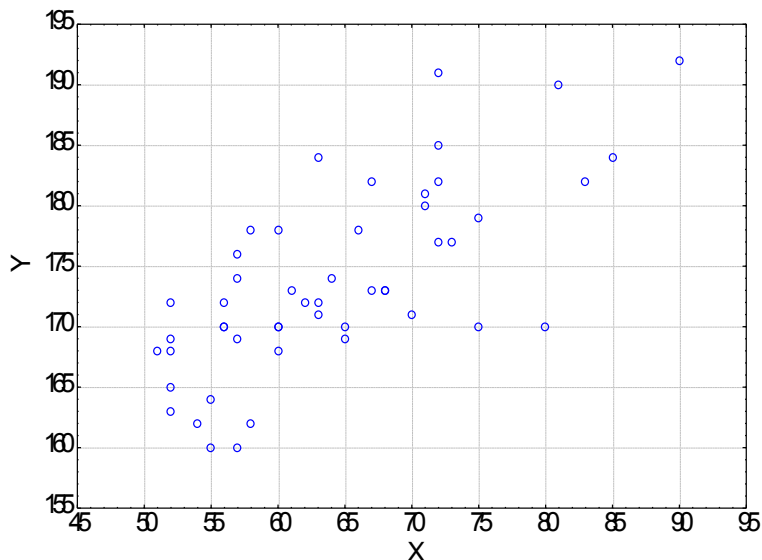
Chceme-li získat kontingenční tabulku relativních četností, resp. sloupcově či řádkově podmíněných relativních četností, na záložce Možnosti zaškrtneme Procenta celkového počtu resp. Procenta z počtu ve slouci či Procenta z počtu v řádku.

Simultánní absolutní četnosti graficky znázorníme tak, že na záložce Detailní výsledky klikneme na 3D histogramy. Výsledný graf můžeme natáčet nebo sklápět.



Úkol 5.: Nakreslete dvourozměrný tečkový diagram pro (X,Y).

Návod: Grafy – Bodové grafy – Proměnné X,Y – OK - vypneme Lineární proložení – OK.



Vidíme, že mezi oběma proměnnými existuje určitý stupeň přímé lineární závislosti – s růstem hmotnosti vesměs rostou hodnoty výšky a naopak.

Samostatná práce: úkoly 1 až 5 proveďte zvlášť pro muže a zvlášť pro ženy.

Úkol 6.: U 100 náhodně vybraných domácností byl zjišťován způsob zásobování bramborami (znak X, varianty 1 = vlastní sklep, 2 = jinde, 3 = nákup) a bydliště (znak Y, varianty 1 = velké město, 2 = malé město, 3 = vesnice).

způsob zásobování	bydliště		
	velké město	malé město	vesnice
vlastní sklep	13	15	14
jinde	11	7	2
nákup	19	9	10

a) Pro oba znaky určíme modus.

b) Vypočteme Cramérův koeficient znaků X, Y.

Návod: Otevřeme nový datový soubor se třemi proměnnými X, Y, četnost a devíti případy. Do proměnné X napíšeme 3 jedničky, 3 dvojky a 3 trojky, do proměnné Y napíšeme 3 krát pod sebe 1, 2, 3 a do proměnné četnost napíšeme odpovídající simultánní absolutní četnosti dvojic variant (X, Y), tj. 13, 15, 14, 11, 7, 2, 19, 9, 10. Proměnným vytvoříme návěští a popíšeme význam jednotlivých variant.

ad a) Výpočet modu: Statistiky – Základní statistiky/tabulky – Popisné statistiky – OK – klikneme na tlačítko se závažím – zaškrtneme Stav zapnuto, vybereme proměnnou vah četnost – OK - Proměnné X, Y – OK – Detailní výsledky – zaškrtneme Modus.

Promě	Popisné statistiky	
	Modu	Četno modu
X	1,000	4,000
Y	1,000	4,000

Proměnná X má modus 1, tj. nejvíce domácností skladuje brambory ve vlastním sklepě a proměnná Y má také modus 1, tj. nejvíce domácností bydlí ve velkém městě.

ad b) Výpočet Cramérova koeficientu: Statistika – Základní statistiky/tabulky – Kontingenční tabulky – OK – Specif. tabulky - List 1 X, List 2 Y - OK – na záložce Možnosti ve Statistikách 2 rozměrných tabulek zaškrtneme F_i (tabulky 2x2) & Cramérovo V & C – přejdeme na záložku Detailní výsledky – Detailní 2-rozm. tabulky.

Statist.	Statist. : X(3) x Y(3)		
	Chi-kv	sv	p
Pearsonův chi-k	6,420	df=	p=,16
M-V chi-kvadr.	7,075	df=	p=,13
FI	,2533		
Kontingenční ko	,2456		
Cramér. V	,1791		

Na posledním řádku najdeme, že Cramérův koeficient nabývá hodnoty 0,179, tedy mezi způsobem zásobování bramborami a bydlištěm domácnosti existuje jen slabá závislost – viz následující tabulka:

Cramérův koeficient	interpretace
mezi 0 až 0,1	zanedbatelná závislost
mezi 0,1 až 0,3	slabá závislost
mezi 0,3 až 0,7	střední závislost
mezi 0,7 až 1	silná závislost

Úkol 7.: Datový soubor znamky.sta obsahuje údaje o 20 studentech 1. ročníku ekonomicky zaměřené vysoké školy. Znak X – známka z matematiky v 1. zkušebním termínu (má varianty 1, 2, 3, 4), znak Y – známka z angličtiny v 1. zkušebním termínu (má rovněž varianty 1, 2, 3, 4), znak Z – pohlaví studenta (0 – žena, 1 – muž).

Otevřeme datový soubor znamky.sta.

a) Pro známky z matematiky a angličtiny vypočteme medián, dolní a horní kvartil, kvartilovou odchylku a vytvoříme krabicový diagram.

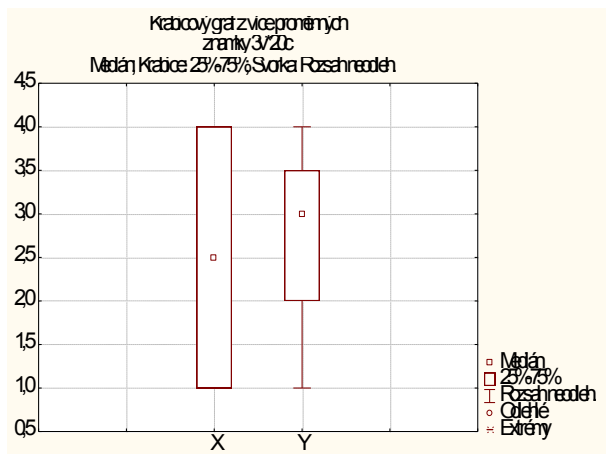
b) Vypočteme Spearmanův korelační koeficient známek z matematiky a angličtiny pro všechny studenty, pak zvlášť pro muže a zvlášť pro ženy. Získané výsledky budeme interpretovat.

Návod:

ad a) Statistika – Základní statistiky/tabulky – Popisné statistiky – OK – Proměnné X, Y – OK – Detailní výsledky - zaškrtneme Medián, Dolní & horní kvartily, Kvartil. rozpětí – Výpočet.

Promě.	Popisné statistiky (znamky)			
	Medián	Spod. kvartil	Horní kvartil	Kvartil. rozpětí
X	2,500	1,000	4,000	3,000
Y	3,000	2,000	3,500	1,500

Vytvoření krabicového diagramu: Grafy – 2D Grafy – Krabicové grafy – vybereme Vícenásobný – Proměnné X, Y – OK.



ad b) Statistika – Neparametrická statistika – Korelace – OK – Proměnné X, Y – OK – Spearman R.

Pro všechny:

Spearmanovy korelace (znamí ChD vynechány párově, Označ. korelace jsou významné)		
Promě	X	Y
X	1,000	0,688
Y	0,688	1,000

Počítáme-li Spearmanův korelační koeficient pro ženy (resp. pro muže), použijeme filtr: tlačítko Select Cases – Zapnout filtr – včetně případů – některé, vybrané pomocí výrazu Z=0 (resp. Z=1).

Pro ženy:

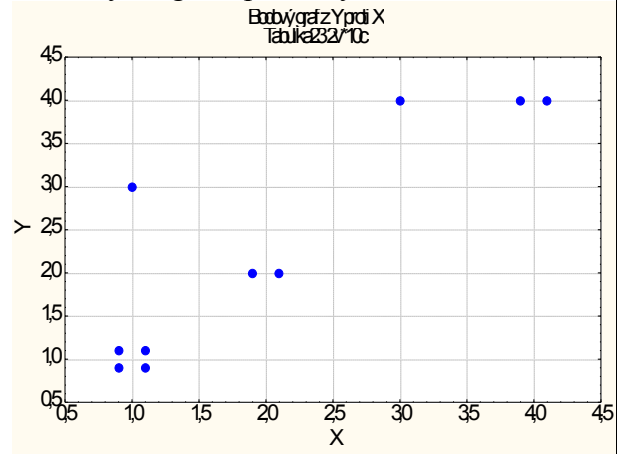
Spearmanovy korelace (znamí ChD vynechány párově, Označ. korelace jsou významné, Zhrnout podmínku: Z=0)		
Promě	X	Y
X	1,000	0,860
Y	0,860	1,000

Pro muže:

Spearmanovy korelace (znamí ChD vynechány párově, Označ. korelace jsou významné, Zhrnout podmínku: Z=1)		
Promě	X	Y
X	1,000	0,373
Y	0,373	1,000

Vidíme, že nejsilnější přímá pořadová závislost mezi známkami z matematiky a angličtiny je u žen, $r_s = 0,86$. U mužů je tato závislost mnohem slabší, $r_s = 0,37$. U žen tedy dochází k tomu, že se sdružují podobné známky z obou předmětů, zatímco u mužů se projevuje spíše tendence k různým známkám. Je to zřetelně vidět na dvourozměrných tečkových diagramech.

Tečkový diagram pro ženy



Tečkový diagram pro muže

