

# ASTAc/02 Biostatistika

## 6. cvičení



**Opakování**  
**Analýza kontingenčních tabulek**  
**Základy korelační analýzy**

# Co byste měli umět z minula:



1. Určit, kdy je vhodné použít pro testování hypotéz parametrické a neparametrické testy – ověřování předpokladů.
2. Vybrat typ neparametrického testu – jednovýběrový, párový nebo dvouvýběrový?
3. Provést testování v softwaru Statistica – Wilcoxonův test, znaménkový test, Mannův-Whitneyho test, Kruskalův-Wallisův test, mediánový test.
4. Interpretovat výsledky testování.

# Analýza kontingenčních tabulek



**Kontingenční tabulky**

**Pearsonův chí-kvadrát test (test dobré shody)**

**Fisherův exaktní test**

**McNemarův test**

# Kontingenční tabulka - opakování



- Frekvenční sumarizace dvou kategoriálních proměnných (binárních, nominálních nebo ordinálních proměnných).
- Obecně: **R x C kontingenční tabulka** (R – počet kategorií jedné proměnné, C – počet kategorií druhé proměnné).
- Speciální případ: 2 x 2 tabulka = čtyřpolní tabulka.
- Kontingenční tabulky: **absolutních četností, celkových procent, řádkových/sloupcových četností**
- Příklad: Sumarizace vyšetřených osob podle pohlaví a výsledku diagnostického testu.

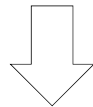
Pohlaví	Výsledek vyšetření		Celkem
	Nemocný	Zdravý	
Muž	45	11	56
Žena	25	6	31
<b>Celkem</b>	70	17	87

# Ukázka kontingenční tabulky



- **Vztah pohlaví a výskytu onemocnění** (pozor na hodnocení nesmyslného vztahu)

	Nemocný	Zdravý	Celkem	
Muž	a	b	a + b	→ Marginální absolutní četnost
Žena	c	d	c + d	
Celkem	a + c	b + d	a + b + c + d = N	→ Celkový počet hodnot Simultánní absolutní četnost



	Nemocný	Zdravý	Celkem
Muž	45	11	56
Žena	25	6	31
Celkem	70	17	87



**Jsou více nemocní muži nebo ženy?**

# Co analyzujeme u kontingenčních tabulek?



- Analýza kontingenčních tabulek umožňuje analyzovat **vazbu mezi dvěma kategoriálními proměnnými**. Základním způsobem testování je tzv. chí-kvadrát test, který **srovnává pozorované četnosti kombinací kategorií oproti očekávaným četnostem**, které vychází z teoretické situace, kdy je vztah mezi proměnnými náhodný.
- Test dobré shody je využíván také pro **srovnání pozorovaných četností proti očekávaným četnostem daných určitým pravidlem** (typickým příkladem je Hardy-Weinbergova rovnováha v genetice).
- Specifickým typem výstupů odvozených z kontingenčních tabulek jsou tzv. **poměry šancí a relativní rizika**, využívaná často v medicíně pro identifikaci a popis rizikových skupin pacientů.

# Test dobré shody - základní teorie

Testová statistika:

$$\chi^2 = \sum \frac{\left[ \begin{array}{c} \text{pozorovaná} - \text{očekávaná} \\ \text{četnost} \quad \text{četnost} \end{array} \right]^2}{\text{očekávaná četnost}}$$

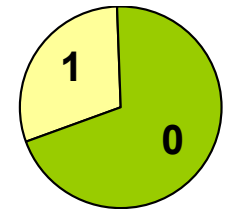
$$\chi^2 = \underbrace{\frac{\left[ \begin{array}{c} \text{pozorovaná} - \text{očekávaná} \\ \text{četnost} \quad \text{četnost} \end{array} \right]^2}{\text{očekávaná četnost}}}_{\text{1. jev}} + \underbrace{\frac{\left[ \begin{array}{c} \text{pozorovaná} - \text{očekávaná} \\ \text{četnost} \quad \text{četnost} \end{array} \right]^2}{\text{očekávaná četnost}}}_{\text{2. jev}} + \dots$$

$$\chi^2 > \chi^2_{(1-\alpha)} (s.v.) \quad \dots \text{ zamítáme } H_0$$

1 - hladina významnosti

stupně volnosti

# Test dobré shody: příklad I



Binomické jevy (1/0)

$$\chi^2_{(1)} = \frac{\left[ \begin{array}{c} \text{pozorovaná} \\ \text{četnost} \end{array} - \begin{array}{c} \text{očekávaná} \\ \text{četnost} \end{array} \right]^2}{\underbrace{\text{očekávaná četnost}}_{\text{I. jev 1}}} + \frac{\left[ \begin{array}{c} \text{pozorovaná} \\ \text{četnost} \end{array} - \begin{array}{c} \text{očekávaná} \\ \text{četnost} \end{array} \right]^2}{\underbrace{\text{očekávaná četnost}}_{\text{II. jev 2}}}$$

## Příklad



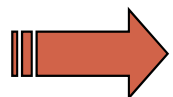
10 000 lidí hází mincí → rub: 4 000 případů (R)  
→ líc: 6 000 případů (L)



Lze výsledek považovat za statisticky významně odlišný (nebo neodlišný) od očekávaného poměru R : L = 1 : 1 (tzn. že je výsledek hodu mincí náhodný)?

$$\chi^2 = \frac{(4000 - 5000)^2}{5000} + \frac{(6000 - 5000)^2}{5000} = 400$$

Tabulková hodnota:  $\chi^2_{(0,95)} (v = k - 1 = 1) = \underline{3,84}$  ( $0,95 = 1 - \alpha$ )



**Rozdíl je vysoce statisticky významný ( $p < 0,001$ )**



# Test dobré shody: příklad II



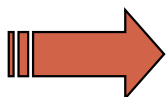
Celkem bylo zkoumáno 250 semen určitého druhu rostliny a roztríděno do následujících kategorií: žluté/hladké; žluté/vrásčité; zelené/hladké; zelené/vrásčité. Předpokládaný poměr výskytu těchto kategorií v populaci je **9 : 3 : 3 : 1**. Následující tabulka obsahuje původní data z pozorování a dále postup při testování  $H_0$ .

	žluté/hladké	žluté/vrásčité	zelené/hladké	zelené/vrásčité	n
$f_{\text{poz.}}$	152	39	53	6	250
$f_{\text{oček.}}$	140,6250	46,8750	46,8750	15,6250	

$$\nu = k - 1 = 3$$

$$\chi^2 = \frac{11,3750^2}{140,6250} + \frac{7,8750^2}{46,8750} + \frac{6,1250^2}{46,8750} + \frac{9,6250^2}{15,6250} = 8,972$$

Tabulková hodnota:  $\chi^2_{(0,95)} (\nu = k - 1 = 3) = \underline{7,81} \quad (0,95 = 1 - \alpha)$



**Zamítáme hypotézu shody pozorovaných četností s očekávanými**

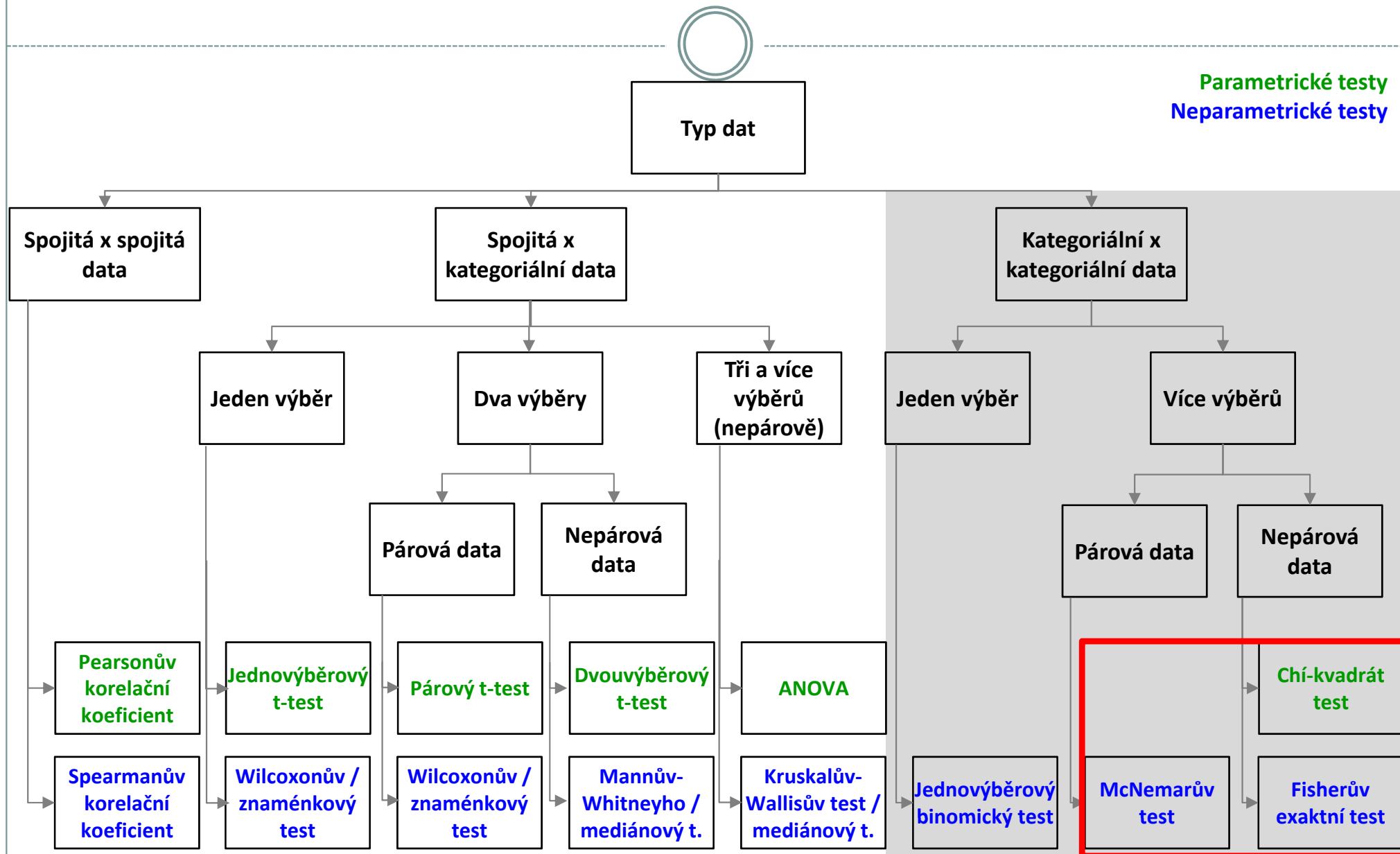
# Kontingenční tabulka - hypotézy



- **NEZÁVISLOST** (Pearsonův chí-kvadrát test, Fisherův exaktní test)
  - Jeden výběr, 2 charakteristiky – obdoba nepárového uspořádání
  - Např.: existence vztahu mezi barvou očí a známkou z biostatistiky u studentů
- **SHODA STRUKTURY** (Pearsonův chí-kvadrát test, Fisherův exaktní test)
  - Tzv. test homogenity
  - Více výběrů, jedna charakteristika – obdoba nepárového uspořádání
  - Např.: věková struktura pacientů s diabetem v  $K$  nemocnicích (tj.  $K$  výběrů)
- **SYMETRIE** (McNemarův test)
  - Jeden výběr, opakovaně jedna charakteristika – obdoba párového uspořádání
  - Např.: posouzení stavu stromů ve dvou sezónách

# Základní rozhodování o výběru statistických testů

## - analýza kontingenčních tabulek



# Kontingenční tabulka - obecně



- Máme dvě nominální veličiny, X (má r variant) a Y (má s variant)
- Kontingenční tabulka typu r x s

$x_{[j]} \backslash y_{[k]}$	$y_{[1]}$	.....	.....	$y_{[s]}$	$n_{j.}$
$x_{[1]}$	$n_{11}$	.....	.....	$n_{1s}$	$n_{1.}$
.	.	.....	.....	.	.
.	.	.....	.....	.	.
$x_{[r]}$	$n_{r1}$	.....	.....	$n_{rs}$	$n_{r.}$
$n_{.k}$	$n_{.1}$	.	.	$n_{.s}$	$n$

Marginální absolutní četnost

Marginální absolutní četnost

Simultánní absolutní četnost

- Označení:  
 $n_{jk}$  - simultánní absolutní četnost,  
 $n_{j.}$  - marginální absolutní četnost

# Testování nezávislosti – Pearsonův chí-kvadrát test



- Souvisí spolu výskyt dvou nominálních znaků měřených na jediném výběru?
- Příklad: Barva očí (modrá, zelená, hnědá) a barva vlasů (hnědá, černá, blond) u vybraných 30 studentů jsou nezávislé.
- **Nulová hypotéza:** Znaky X a Y jsou nezávislé náhodné veličiny.
- **Alternativní hypotéza:** Znaky X a Y jsou závislé náhodné veličiny.
- Test: **Pearsonův chí-kvadrát**

$$K = \sum_{j=1}^r \sum_{k=1}^s \frac{(n_{jk} - e_{jk})^2}{e_{jk}} \stackrel{\text{H}_0 \text{ platí}}{\approx} \chi^2((r-1)(s-1))$$

Očekávané (teoretické) četnosti  $e_{jk}$ :  $e_{jk} = \frac{n_{j.} \cdot n_{.k}}{n}$

- $H_0$  zamítáme na hladině významnosti  $\alpha$ , pokud  $K \geq \chi_{1-\alpha}^2((r-1)(s-1))$
- **Předpoklady testu ?**

# Testování nezávislosti – Pearsonův chí-kvadrát test



- **Předpoklady Pearsonova chí-kvadrát testu:**

1. **Jednotlivá pozorování** shrnutá v kontingenční tabulce **jsou nezávislá**, tj. každý prvek patří jen do jedné buňky kont. tabulky, nemůže zároveň patřit do dvou.
2. **Podmínky dobré aproximace:** Očekávané (teoretické) četnosti jsou aspoň v 80 % případů větší nebo rovné 5 a ve 100 % případů nesmí být pod 2 (pokud není tento předpoklad splněn, je vhodné sloučit kategorie s nízkými četnostmi).

- **Měření síly závislosti:**

Cramérův koeficient:  $V = \sqrt{\frac{K}{n(m-1)}}$ , kde  $m = \min\{r, s\}$ ,  $V$  je z intervalu (0,1)

Význam hodnot: 0-0,1....zanedbatelná závislost

0,1-0,3...slabá závislost

0,3-0,7...střední závislost

0,7-1 silná závislost

# Kontingenční tabulky

$H_0$  : Nezávislost dvou jevů A a B



**Kontingenční  
tabulka  
2 x 2**

$\begin{array}{c} \rightarrow \\ \downarrow \end{array} \begin{array}{c} B \\ A \end{array}$	+	-	Podíl (+)
+	a	b	$\frac{a}{(a+b)}$ $p_1$
-	c	d	$\frac{c}{(c+d)}$ $p_2$
Podíl (+)	$\frac{a}{(a+c)}$	$\frac{b}{(b+d)}$	

$$N = a + b + c + d$$

$$P(B^+) = \frac{(a+b)}{N}$$

$$P(B^-) = \frac{(c+d)}{N}$$

## Očekávané četnosti:

$$F_{(A)} = \frac{(a+b)(a+c)}{N}$$

$$F_{(C)} = \frac{(a+c)(d+c)}{N}$$

$$\chi^2_{\nu=1} = \sum_{i=1}^4 \frac{(f_i - F_i)^2}{F_i}$$

$$F_{(B)} = \frac{(a+b)(b+d)}{N}$$

$$F_{(D)} = \frac{(b+d)(c+d)}{N}$$

$$s.v. = 1 = (r-1) * (c-1)$$

# Kontingenční tabulky: příklad



gen \ †	Ano	Ne	Σ
Ano	20	82	102
Ne	10	54	64
Σ	30	136	166

$$F_A = 102 * 30 / 166 = 18,43$$

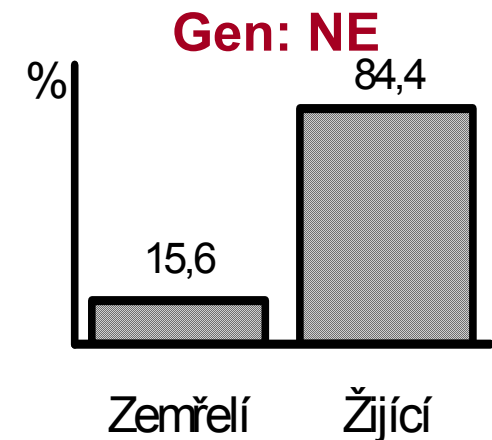
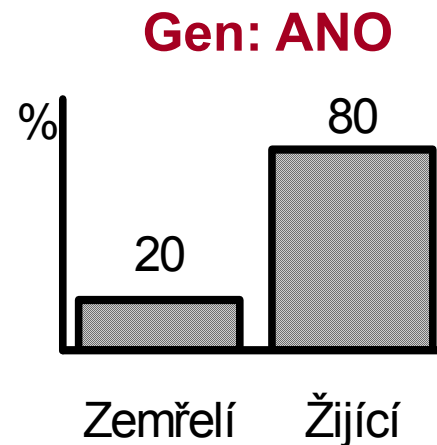
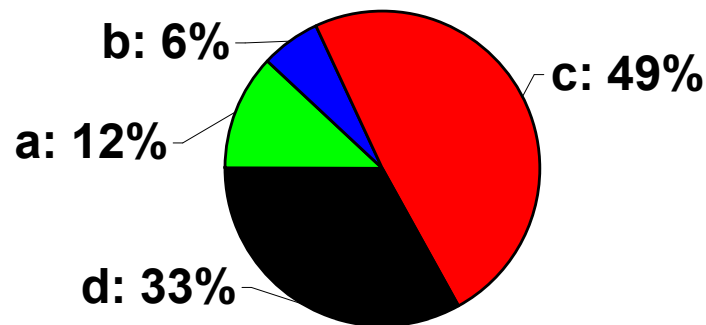
$$F_B = 102 * 136 / 166 = 83,57$$

$$F_C = 11,57$$

$$F_D = 52,43$$

$$\chi^2_{(1)} = \frac{(20-18,43)^2}{18,43} + \frac{(82-83,57)^2}{83,57} + \frac{(10-11,57)^2}{11,57} + \frac{(54-52,43)^2}{52,43} = 0,423 \quad 0,423 < \chi^2_{0,95}^{(1)} = 3,84$$

## Kontingenční tabulka v obrázku





# Řešení v softwaru Statistica



- Datový soubor může být zadán 2 způsoby:
  - **Původní data** (co řádek, to subjekt charakterizovaný danými kategoriálními proměnnými),
  - **Agregovaná data** (kontingenční tabulka, četnosti všech kombinací kategorií 2 kategoriálních proměnných) – analýza agregovaných dat možná i pomocí webových kalkulačtorů.

# Způsob 1: Řešení v softwaru Statistica I

- Na hladině významnosti 0,05 testujte hypotézu o nezávislosti genu a stavu pacienta. Simultánní četnosti znázorněte graficky.

- **Původní datový soubor**  
(co řádek, to subjekt)

- V menu **Statistics** zvolíme **Basic statistics**,  
Vybereme **Tables and banners**  
(v češtině **Kontingenční tabulky**)

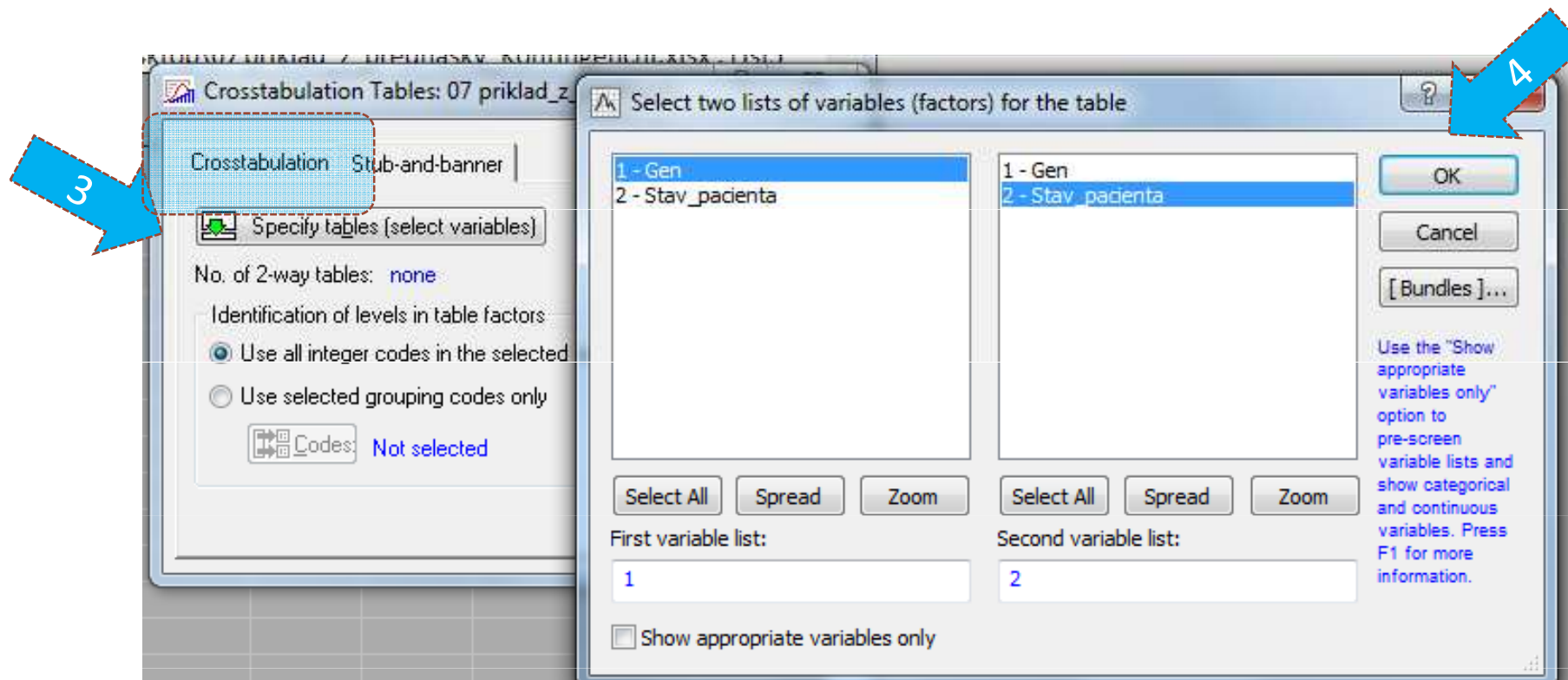
The screenshot shows the Statistica software interface. The 'Statistics' menu is open, and the 'Basic statistics' option is highlighted with a red dashed box and a blue arrow labeled '1'. Below the menu, the 'Basic Statistics and Tables' dialog box is open, and the 'Tables and banners' option is selected with a blue arrow labeled '2'. The background shows a data table with columns 'Gen' and 'Stav\_p' and rows 1 through 12.

	1 Gen	2 Stav_p
1	přítomer	úmrtí
2	přítomer	úmrtí
3	přítomer	úmrtí
4	přítomer	úmrtí
5	přítomer	úmrtí
6	přítomer	úmrtí
7	přítomer	úmrtí
8	přítomer	úmrtí
9	přítomer	úmrtí
10	přítomer	úmrtí
11	přítomer	úmrtí
12	přítomer	úmrtí

# Způsob 1: Řešení v softwaru Statistica II

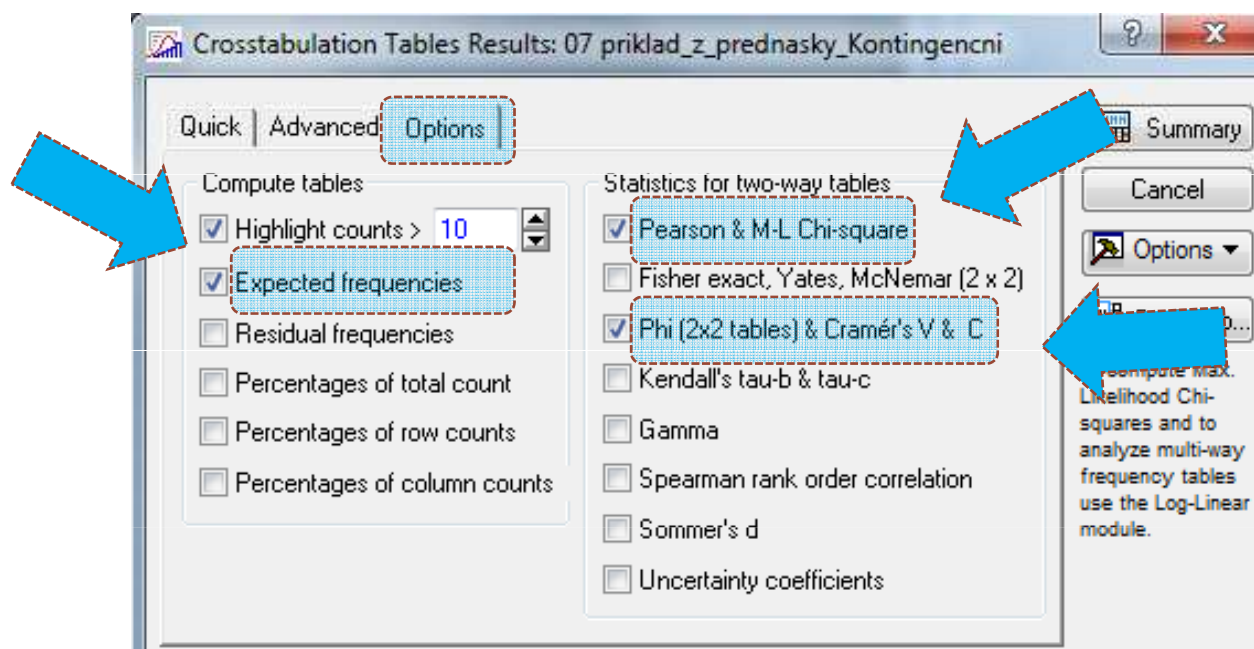


- Vybereme proměnné, které chceme testovat



# Způsob 1: Řešení v softwaru Statistica III

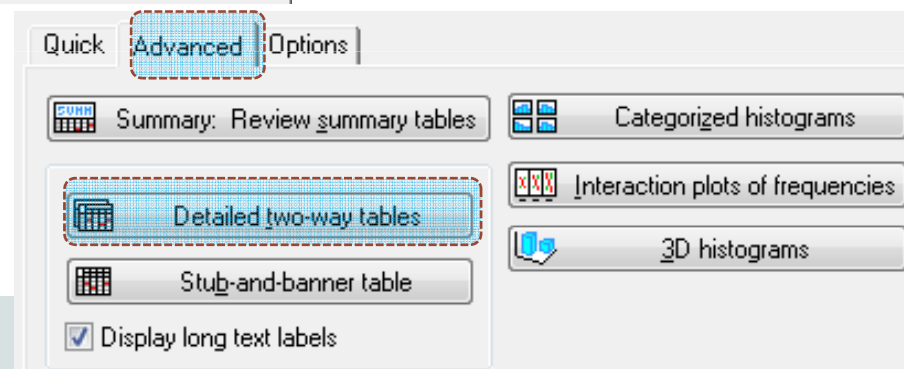
- Na záložce **Options** zaškrtneme **Expected frequencies (Očekávané četnosti)** (k ověření podmínek dobré aproximace)



- Zaškrtneme Pearsonův chí-kvadrát

- Pokud chceme vypočítat i Cramérův koeficient zaškrtneme Phi & Cramer's V

- Poté se vrátíme na záložku **Advanced**, kde a zvolíme **Detailed two-way tables**



# Způsob 1: Řešení v softwaru Statistica IV

**Tab.1: Pozorované četnosti**

Summary Frequency Table (07 priklad\_z\_prednasky\_K  
Marked cells have counts > 10  
(Marginal summaries are not marked)

Gen	Stav_pacienta úmrtí	Stav_pacienta žijící	Row Totals
přítomen	20	82	102
nepřítomen	10	54	64
All Grps	30	136	166

**Tab. 2: Očekávané četnosti**

Summary Table: Expected Frequencies (07 priklad\_z\_pre  
Marked cells have counts > 10  
Pearson Chi-square: ,421322, df=1, p=,516278

Gen	Stav_pacienta úmrtí	Stav_pacienta žijící	Row Totals
přítomen	18,43373	83,5663	102,0000
nepřítomen	11,56627	52,4337	64,0000
All Grps	30,00000	136,0000	166,0000



**Jsou splněny podmínky dobré aproximace?**

**Tab. 3: Paersonův chí-kvadrát**

Hodnota testové statistiky      Počet stupňů volnosti      p- hodnota

Statistics: Gen(2) x Stav\_pacienta(2)

Statistic	Chi-square	df	p
Pearson Chi-square	4213223	df=1	p=,51628
M-L Chi-square	,4277117	df=1	p=,51311
Phi for 2 x 2 tables	,0503794		
Tetrachoric correlation	,0949754		
Contingency coefficient	,0503156		

# Způsob 2: Řešení v softwaru Statistica I



- Na hladině významnosti 0,05 testujte hypotézu o nezávislosti genu a stavu pacienta. Simultánní četnosti znázorněte graficky.

- **Agregovaný datový soubor**

- V menu **Statistics** zvolíme **Basic statistics**, vybereme **Tables and banners** (v češtině **Kontingenční tabulky**)

	1 Gen	2 Stav_pacienta	3 Četnost
1	přítomen	úmrti	20
2	přítomen	žijící	82
3	nepřítomen	úmrti	10
4	nepřítomen	žijící	54

Basic Statistics and Tables: Spreadsheet11

Quick

- Descriptive statistics
- Correlation matrices
- t-test, independent, by groups
- t-test, independent, by variables
- t-test, dependent samples
- t-test, single sample
- Breakdown & one-way ANOVA
- Breakdown; non-factorial tables
- Frequency tables
- Tables and banners**
- Multiple response tables
- Difference tests: r, %, means
- Probability calculator

OK

Cancel

Options

Open Data

SELECT CASES

10 W

# Způsob 2: Řešení v softwaru Statistica II



- Vybereme proměnné, které chceme testovat

The image shows two overlapping dialog boxes from the Statistica software. The background dialog is titled 'Crosstabulation Tables: 07 priklad\_...' and has tabs for 'Crosstabulation' and 'Stub-and-banner'. It includes a 'Specify tables (select variables)' button, a 'No. of 2-way tables: none' field, and radio buttons for 'Use all integer codes in the selected' (selected) and 'Use selected grouping codes only'. There is also a 'Codes: Not selected' field. The foreground dialog is titled 'Select up to 6 lists of grouping variables:' and contains six columns, each with a list of variables: '1 - Gen', '2 - Stav\_pacienta', and '3 - Četnost'. Below these lists are 'Spread' and 'Zoom' buttons for each column. At the bottom, there are input fields for 'List1:' through 'List6:', with 'List1' containing the number '1' and 'List2' containing '2'. A checkbox 'Show appropriate variables only' is at the bottom left. On the right side of the foreground dialog, there are 'OK', 'Cancel', and '[Bundles]...' buttons, along with a help text: 'Use the "Show appropriate variables only" option to pre-screen variable lists and show categorical and continuous variables. Press F1 for more information.'

# Způsob 2: Řešení v softwaru Statistica III



- Zapneme **váhy** (vpravo ikonka černých vah **w**), jako váhy vybereme proměnnou **četnost** (tj. proměnnou, ve které jsou uvedeny počty případů jednotlivých kombinací kategorií)

	1 Gen	2 Stav_pacienta	3 Četnost
1	přítomen	úmrtí	20
2	přítomen	žijící	82
3	nepřítomen	úmrtí	10
4	nepřítomen	žijící	54

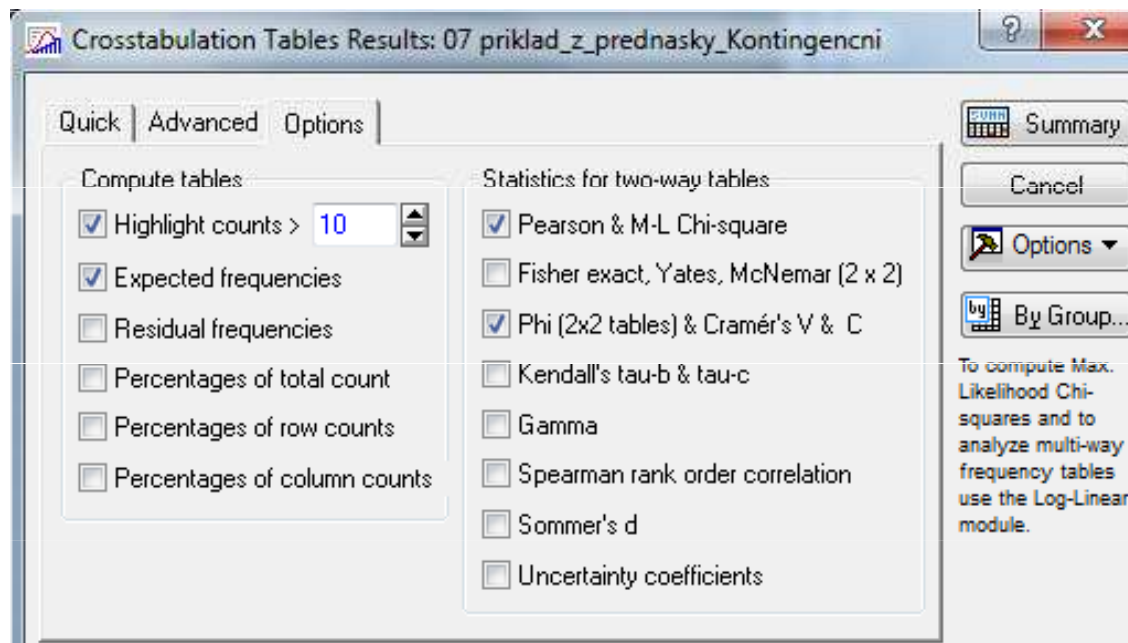
The image shows two dialog boxes overlaid on a spreadsheet. The first dialog box, titled 'Crosstabulation Tables: Spreadsheet11', has the 'Weighted moments' checkbox checked. A blue arrow labeled '1' points to this checkbox. The second dialog box, titled 'Analysis/Graph Case Weights', has 'Use Spreadsheet weights' selected and 'Četnost' entered in the 'Weight variable' field. A blue arrow labeled '2' points to the 'Četnost' field. A third blue arrow labeled '3' points to the 'OK' button in the second dialog box.



# Způsob 2: Řešení v softwaru Statistica IV



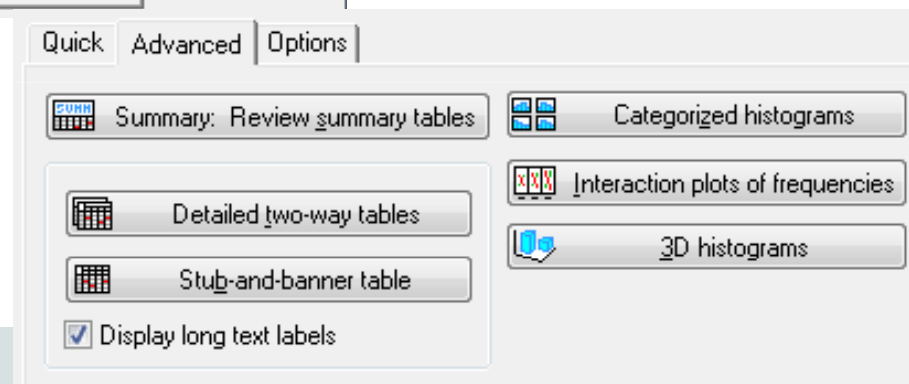
- Na záložce **Options** zaškrtneme **Expected frequencies (Očekávané četnosti)** (k ověření podmínek dobré aproximace)



- Zaškrtneme Pearsonův chí-kvadrát

- Pokud chceme vypočítat i Cramérův koeficient zaškrtneme Phi & Cramer's V

- Poté se vrátíme na záložku **Advanced**, kde a zvolíme **Detailed two-way tables**



# Testování homogenity (shody struktury)

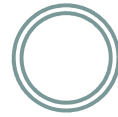


- Motivace: Zajímá nás výskyt nominálního znaku u  $r$  nezávislých výběrů z  $r$  různých populací.
- Příklad: Je zájem o sport stejný u děvčat jako u chlapců?
- Nulová hypotéza: pravděpodobnostní rozdělení kategoriální proměnné je stejné v různých populací
- Test: **Pearsonův chí-kvadrát**

		Dívky	Chlapci	
Zájem o sport	Ano	$a$	$b$	$a+b$
	Ne	$c$	$d$	$c+d$
		$a+c$	$b+d$	$n$

*Některé marginální četnosti (buď sloupcové nebo řádkové) jsou předem pevně stanoveny*

# Fisherův exaktní test



- Využití ve čtyřpolní tabulce s nízkými četnostmi, které znemožňují použití Pearsonova chí-kvadrát testu.
- Patří mezi **neparametrické testy** pracující s daty na nominální škále, v nejjednodušší podobě ve dvou třídách: pozitivní/negativní, úspěch/neúspěch apod.
- Nulová hypotéza předpokládá rovnoměrné zastoupení sledovaného znaku u dvou nezávislých souborů.
- Slovo exaktní (přímý) znamená, že se přímo vypočítává pravděpodobnost odmítnutí, resp. platnosti nulové hypotézy.

# Fisherův exaktní test



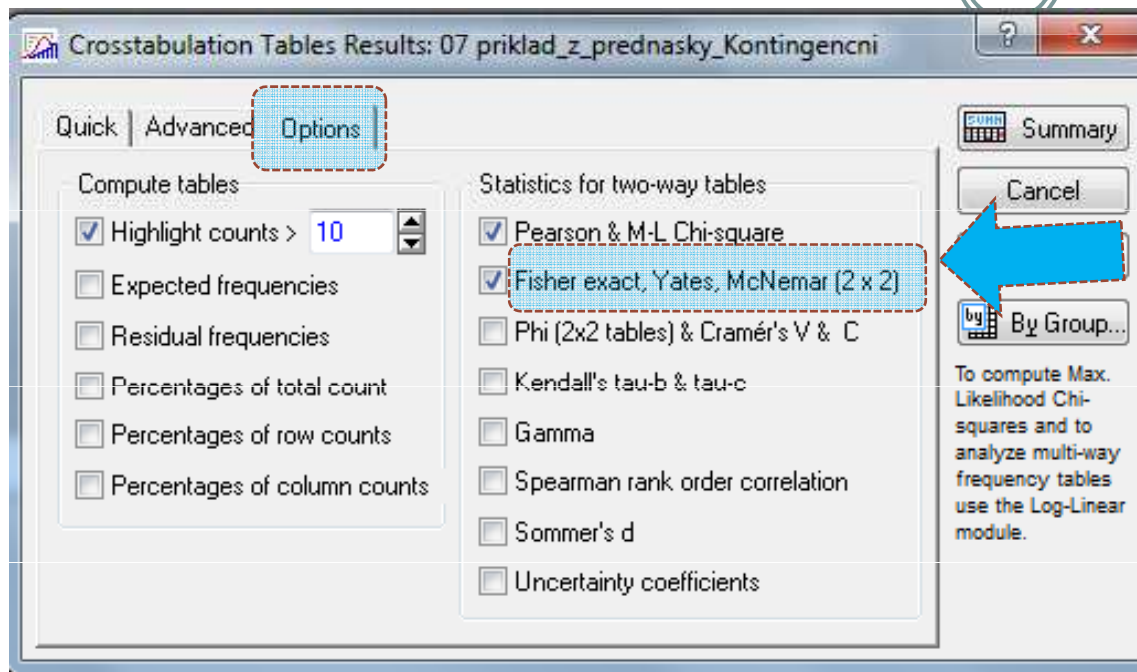
- Výpočet „přesné“ p-hodnoty, která zde hraje roli testové statistiky:
  - spočítá se parciální pravděpodobnost čtyřpolní tabulky  $p_1$ :

Sledovaný jev	Skupina		Celkem
	Experimentální	Kontrolní	
Ano	$a$	$b$	$a + b$
Ne	$c$	$d$	$c + d$
Celkem	$a + c$	$b + d$	$n$

$$p_1 = \frac{(a+b)! * (c+d)! * (a+c)! * (b+d)!}{N! * a! * b! * c! * d!}$$

- Spočítá se  $p_a$  všech možných tabulek při zachování marginálních četností (řádkové a sloupcové součty) a výsledná p-hodnota je součtem  $p_a$  menších nebo stejných jako  $p_1$ , která přísluší pozorované tabulce.

# Řešení v softwaru Statistica: Fisherův exaktní test



- Na záložce **Options** zaškrtneme **Fisher exact**

- Výstupní tabulka

Statistic	Statistics: Gen(2) x Stav_paci		
	Chi-square	df	p
Pearson Chi-square	,4213223	df=1	p=,51628
M-L Chi-square	,4277117	df=1	p=,51311
Yates Chi-square	,1952605	df=1	p=,65857
Fisher exact, one-tailed			p=,33259
two-tailed			p=,54314
McNemar Chi-square (A/D)	14,71622	df=1	p=,00012
(B/C)	54,79348	df=1	p=,00000

Pro jednostranný test

Pro oboustranný test



# Test hypotézy o symetrii (McNemarův test pro čtyřpolní tabulku)



- Motivace: Na osobách sledujeme binární proměnnou před pokusem a po něm, cílem je zjistit, zda došlo ke změně v rozdělení této proměnné.
- **Analýza párových dichotomických proměnných**

Četnostní tabulka

		po		$n_{j.}$
		+	-	
před	+	$a$	$b$	$a+b$
	-	$c$	$d$	$c+d$
$n_{.k}$		$a+c$	$b+d$	$n$

Tabulka teoretických pravděpodobností

		po		
		+	-	
před	+	$p_{11}$	$p_{12}$	$p_{1.}$
	-	$p_{21}$	$p_{22}$	$p_{2.}$
		$p_{.1}$	$p_{.2}$	

- Nulová hypotéza:  $p_{ij} = p_{ji}$ , pokus nemá vliv na výskyt daného znaku
- Testová statistika:  $\chi^2 = \frac{(|b - c| - 1)^2}{b + c}$  pokud je větší než kritická hodnota  $\chi^2$  rozdělení o jednom stupni volnosti (vhodné pro počty údajů  $b+c > 8$ ), pak nulovou hypotézu zamítáme

# McNemarův test: příklad I



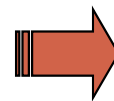
Zjistěte, zda výuka o pozitivním působení sportu na zdraví vede ke změně postojů žáků ke sportování.

**Nulová hypotéza:** Počet žáků, kteří změní svůj postoj pozitivním směrem, je pouze náhodně odlišný od počtu žáků, kteří změní svůj postoj negativním směrem.

		Postoj po výuce		
		+	-	
Postoj před výukou	+	5	3	8
	-	16	2	18
		21	5	26

$$\chi^2 = \frac{(|3 - 16| - 1)^2}{3 + 16} = 7,58$$

Tabulky:  $\chi^2_{1-\alpha}(v = k(k-1)/2 = 1) = 3,84$  *Stupně volnosti*



**H<sub>0</sub> zamítnuta**

Závěr: Výuka má pozitivní vliv na postoj žáků vzhledem k provozování sportu.

# Řešení v softwaru Statistica: McNemarův test



## Datový soubor

	1 postoj_pred_vyukou	2 postoj_po_vyuce	3 cetnost
1	kladný	kladný	5
2	záporný	kladný	16
3	kladný	záporný	3
4	záporný	záporný	2

## Výstupní kontingenční tabulka

2-rozměrná tabulka: Pozorované četnosti (příklad\_pos  
Četnost označených buněk > 10

	postoj_po_vyuce kladný	postoj_po_vyuce záporný	Řádk. součty
postoj_pred_vyukou kladný	5	3	8
postoj_pred_vyukou záporný	16	2	10
Celk.	21	5	26

Crosstabulation Tables Results: 07 priklad\_z\_prednasky\_Kontingencni

Quick | Advanced | Options

Compute tables

- Highlight counts > 10
- Expected frequencies
- Residual frequencies
- Percentages of total count
- Percentages of row counts
- Percentages of column counts

Statistics for two-way tables

- Pearson & M-L Chi-square
- Fisher exact, Yates, McNemar (2 x 2)
- Phi (2x2 tables) & Cramér's V & C
- Kendall's tau-b & tau-c
- Gamma
- Spearman rank order correlation
- Sommer's d
- Uncertainty coefficients

Buttons: Summary, Cancel, Options, By Group...

To compute Max. Likelihood Chi-squares and to analyze multi-way frequency tables use the Log-Linear module.

- Na záložce **Options** zaškrtneme **McNemar (2x2)**

- Výstupní tabulka

Statist.	Chi-kvadr.	sv	p
Yatesův chí-kv.	1.074735	df=1	p=.29988
Fisherův přesný, 1-str.			p=.15026
Fisherův přesný, 2-str.			p=.28051
McNemarův chí-kv. (A/D)	57.14286	df=1	p=.44969
McNemarův chí-kv. (B/C)	7.578948	df=1	p=.00591

- ↪ 2 hodnoty testových statistik a p-hodnoty, podle toho, kde jsou ve výstupní kontingenční tabulce uloženy četnosti, u kterých jsme při opakovaném měření zaznamenali rozdílné výsledky (A/D nebo B/C)



# Analýza kontingenčních tabulek na webu



- 2x2 tabulky: <http://graphpad.com/quickcalcs/contingency1/>
- 2x3 tabulky: <http://www.vassarstats.net/fisher2x3.html>
- 2x5 (nebo menší) tabulky:  
<http://www.quantitativeskills.com/sisa/statistics/fiveby2.htm>
- 3x3 tabulky: <http://vassarstats.net/fisher3x3.html>

# Společný příklad – testování homogenity



Očkování proti chřipce se zúčastnilo 460 dospělých, z nichž 240 dostalo očkovací látku proti chřipce a 220 dostalo placebo. Na konci experimentu onemocnělo 100 lidí chřipkou, 20 z nich bylo z očkované skupiny a 80 z kontrolní skupiny. Je to dostatečný důkaz, že očkovací látka byla účinná?

**Nulová hypotéza:** Procento výskytu chřipky je v očkované a kontrolní skupině stejné.

1. *Vytvořte si na základě zadání datový soubor v softwaru STATISTICA (agregovaná data ve formě kontingenční tabulky).*
2. *Testujte platnost nulové hypotézy pomocí Pearsonova chí-kvadrát testu.*
3. *Testujte platnost nulové hypotézy pomocí Fisherova exaktního testu.*
4. *Který z testů je vhodné použít a proč?*

# Základy korelační analýzy



**Korelace a regrese**  
**Pearsonův korelační koeficient**  
**Spearmanův korelační koeficient**

# Proč hodnotit vztah dvou spojitých veličin?



- Vztah mezi dvěma spojitými veličinami v jedné skupině:
  1. Chceme zjistit, jestli mezi nimi **existuje vztah** – např. jestli vyšší hodnoty jedné veličiny znamenají nižší hodnoty jiné veličiny,
  2. Chceme **predikovat hodnoty** jedné veličiny na základě znalosti hodnot jiných veličin,
  3. Chceme **kvantifikovat vztah** mezi dvěma spojitými veličinami – např. pro použití jedné veličiny na místo druhé veličiny.

# Korelační a regresní analýza



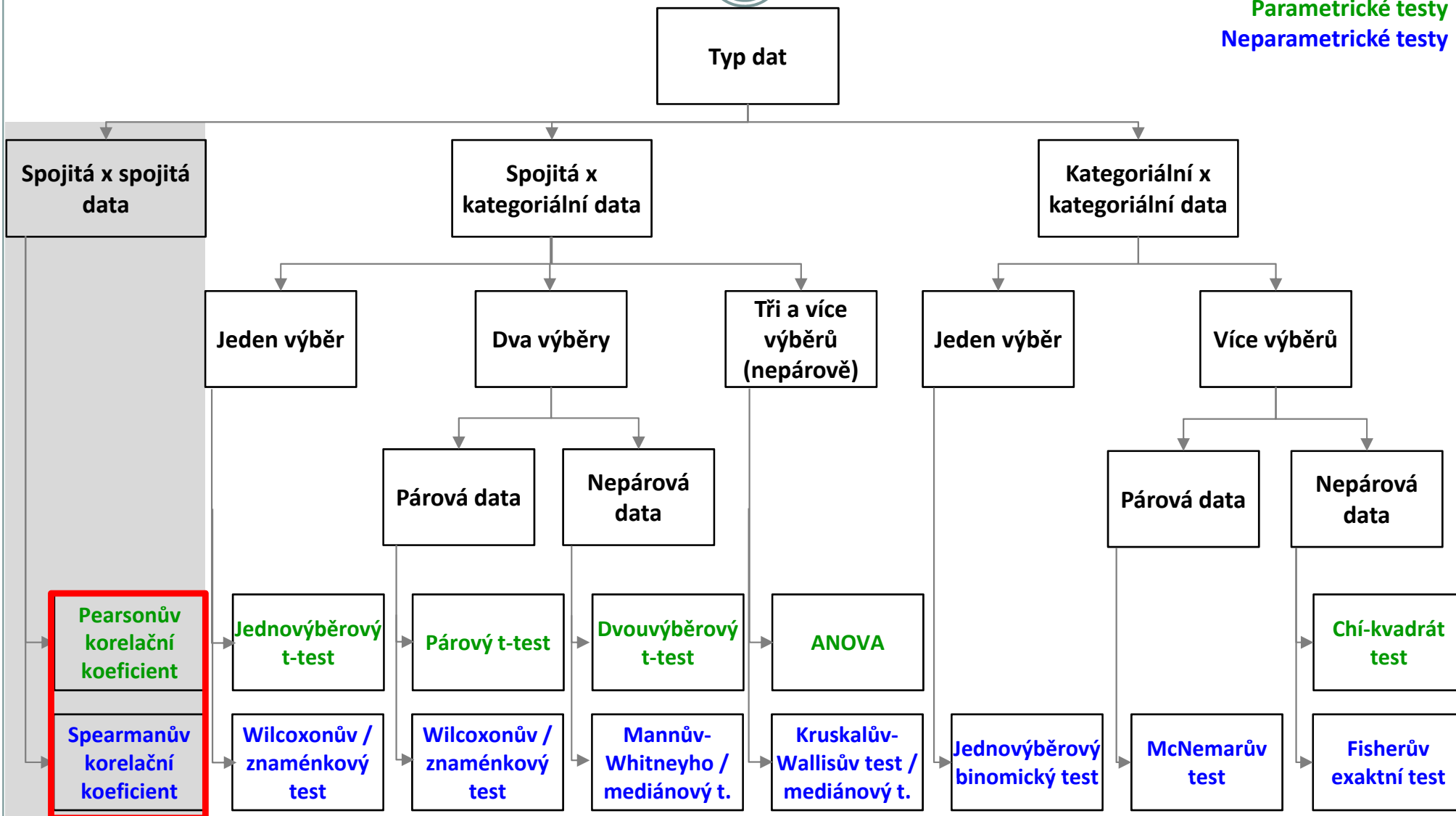
- **Korelační analýza** je využívána pro vyhodnocení míry vztahu dvou spojitých proměnných. Obdobně jako jiné statistické metody, i korelace mohou být parametrické nebo neparametrické.
- **Regresní analýza** vytváří model vztahu dvou nebo více proměnných, tedy jakým způsobem jedna proměnná (vysvětlovaná) závisí na jiných proměnných (prediktorech). Regresní analýza je obdobně jako ANOVA nástrojem pro vysvětlení variability hodnocené proměnné.

# Základní rozhodování o výběru statistických testů

## - korelační analýza



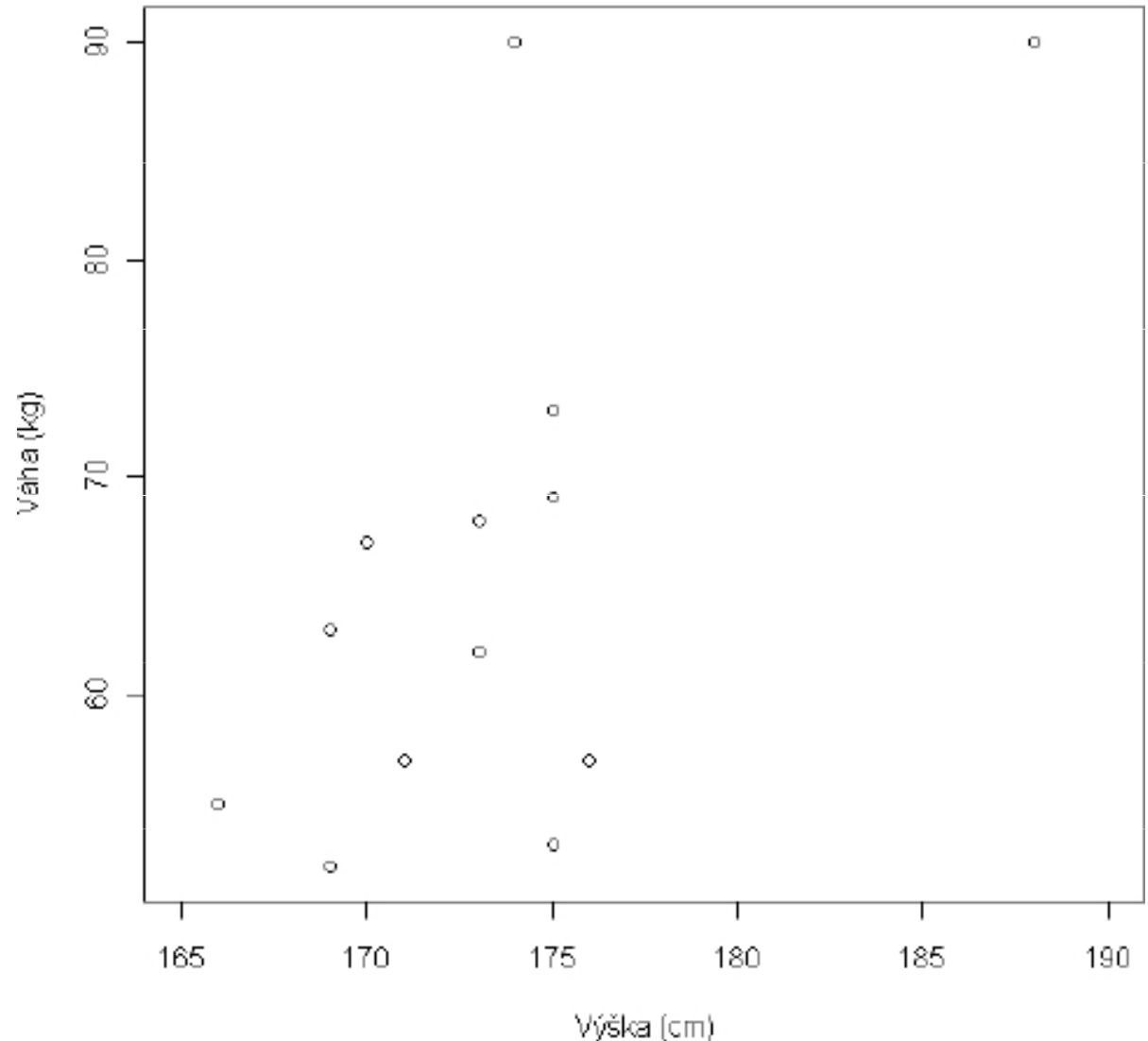
Parametrické testy  
Neparametrické testy



# Vizuální hodnocení vztahu dvou proměnných



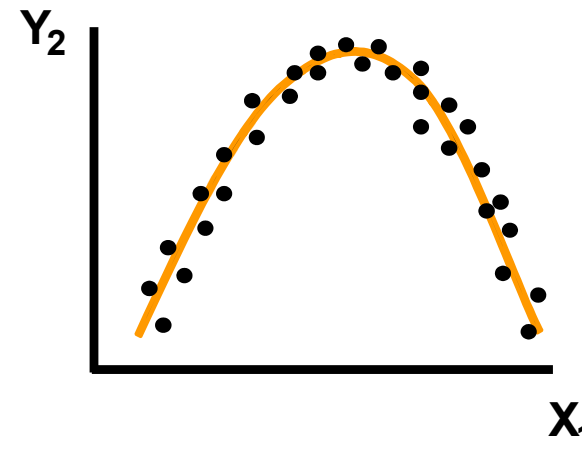
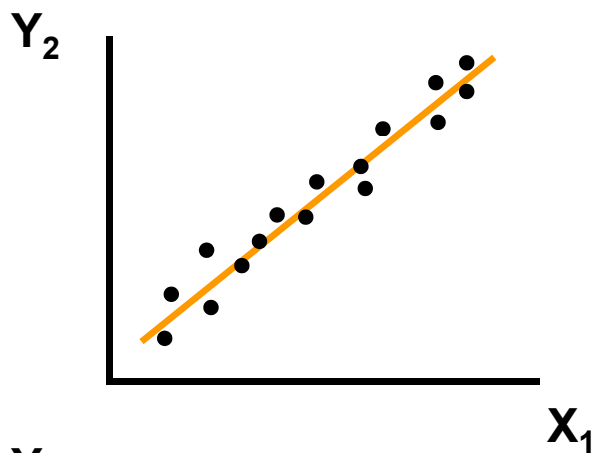
- Nejjednodušší formou je **bodový graf** (x-y graf), tzv. scatterplot.
- Vztah výšky a váhy studentů  
Biostatistiky pro  
matematické biologie  
– jaro 2010:



# Korelace



**Korelace – vztah (závislost) dvou znaků (parametrů)**



$X_2 \backslash X_1$	ANO	NE
ANO	a	b
NE	c	d

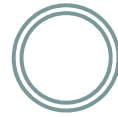


# Korelační koeficienty



- **Korelační koeficient ( $r$ )** – kvantifikuje míru vztahu mezi dvěma spojitými veličinami ( $X$  a  $Y$ ).
  - **Pearsonův korelační koeficient** – parametrický, hodnotí míru lineární závislosti mezi 2 spojitými proměnnými,
  - **Spearmanův korelační koeficient** – neparametrický, hodnotí míru pořadové závislosti mezi 2 spojitými proměnnými.
  - Hodnota  $r$  je kladná, když vyšší hodnoty  $X$  souvisí s vyššími hodnotami  $Y$ , naopak hodnota  $r$  je záporná, když nižší hodnoty  $X$  souvisí s vyššími hodnotami  $Y$ .
  - Nabývá hodnot od -1 do 1:
    - $r = 0 \rightarrow$  nekorelované
    - $r > 0 \rightarrow$  kladně korelované
    - $r < 0 \rightarrow$  záporně korelované

# Test hypotézy $H_0: r = 0$

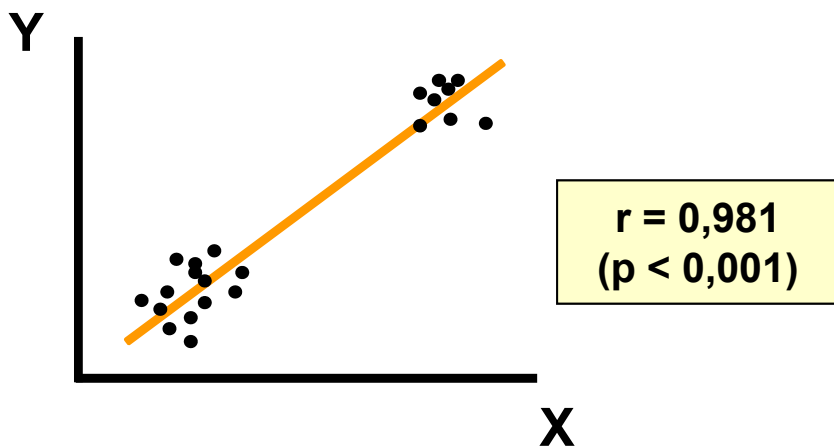


- K měření těsnosti lineárního vztahu 2 spojitých proměnných
  - $r = 0 \rightarrow$  nekorelované
  - $r > 0 \rightarrow$  kladně korelované
  - $r < 0 \rightarrow$  záporně korelované
- $H_0$ : proměnné  $X, Y$  jsou stochasticky nezávislé náhodné veličiny  
( $r = 0$ )  
 $H_A$ : proměnné  $X, Y$  nejsou stochasticky nezávislé náhodné veličiny  
( $r \neq 0$ )
- Testování pomocí intervalu spolehlivosti nebo výpočet testové statistiky (srovnání s kritickou hodnotou nebo výpočet p-hodnoty)

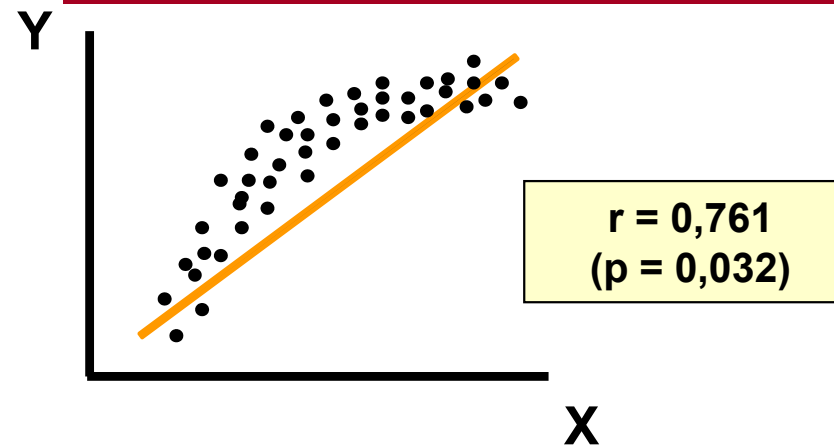
# Problémy s výpočtem $r$



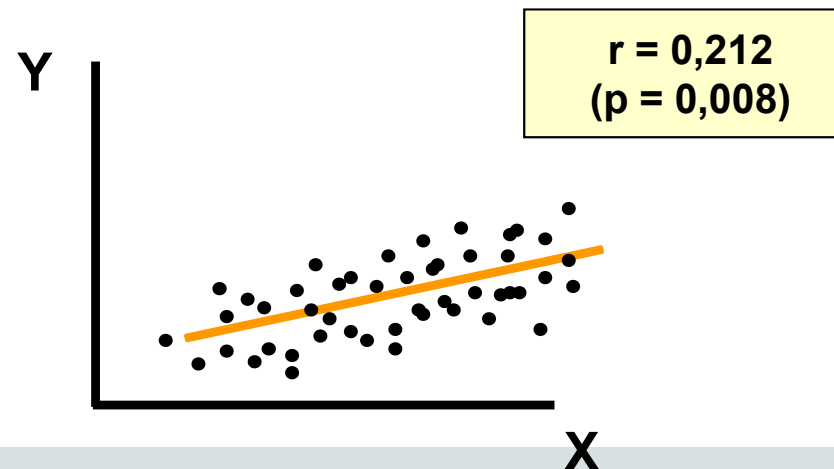
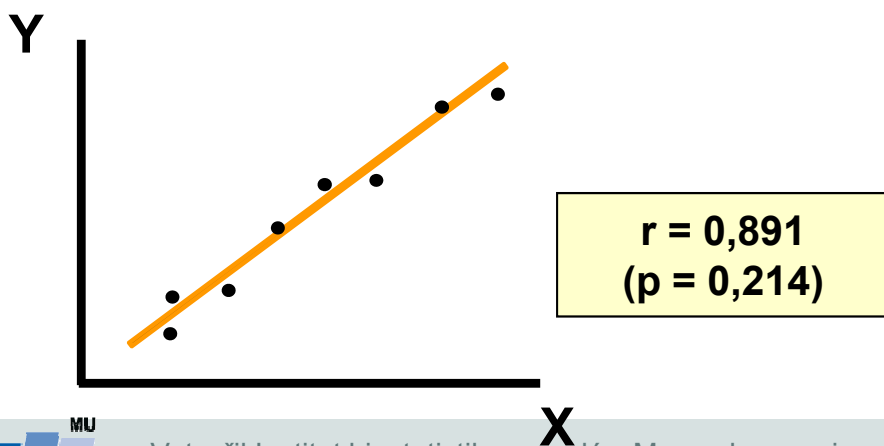
## Problém více skupin



## Nelineární vztah



## Problém velikosti výběru



# Řešení v softwaru Statistica: Pearsonův korelační koeficient I

Prozkoumejte lineární vztah mezi výškou a váhou u 13 studentů. Testujte hypotézu, že jsou tyto proměnné nezávislé.

1. Záložka **Statistics**
2. **Basic Statistics**
3. **Correlation matrices**
4. Potvrdíme: **OK**

The screenshot shows the Statistica software interface. The 'Statistics' ribbon is active, with 'Basic Statistics' selected. A data table is visible with columns '1 vyska' and '2 vaha', and rows 1 through 13. The 'Basic Statistics and Tables: priklad\_vyska...' dialog box is open, with 'Correlation matrices' selected in the 'Quick' list. The 'OK' button is highlighted. Blue arrows with numbers 1 through 4 point to the Statistics ribbon, Basic Statistics, Correlation matrices, and the OK button respectively.

	1 vyska	2 vaha
1	175	69
2	166	55
3	170	67
4	169	52
5	188	90
6	175	53
7	176	
8	171	
9	173	68
10	175	73
11	173	62
12	174	90
13	169	63

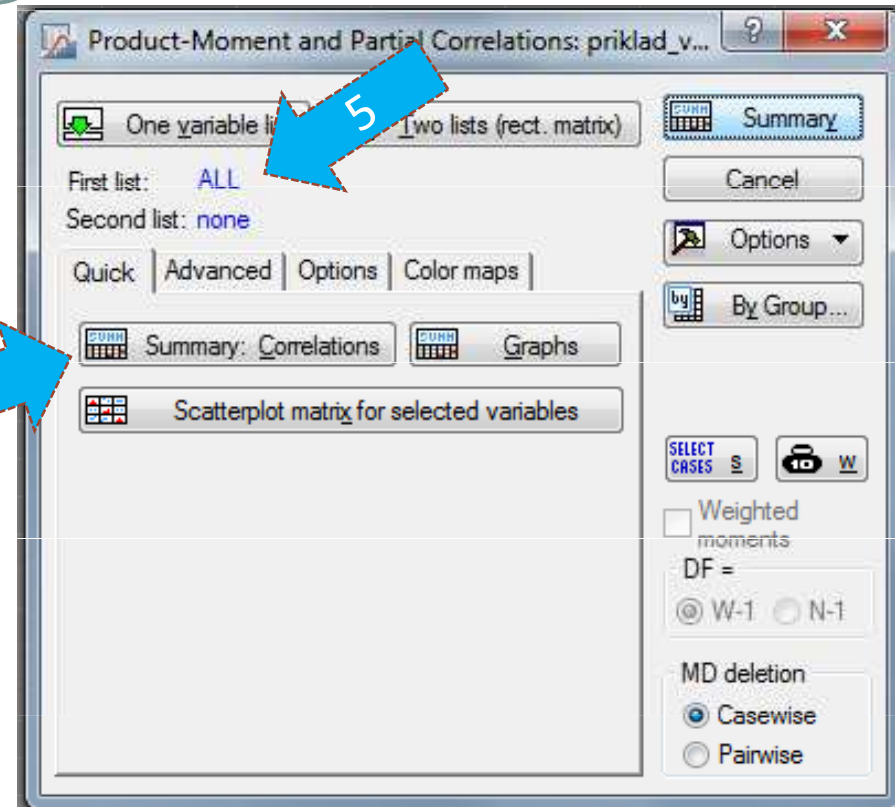
# Řešení v softwaru Statistica: Pearsonův korelační koeficient II

5. Vybereme spojité proměnné pro hodnocení vztahu (váha a výška).

Na záložce **Options** můžeme vybrat formu výstupu (pouze p-hodnoty, matice korelačních koeficientů a p-hodnot ap.).

## 6. **Summary: Correlations**

Jedna z možných výstupních tabulek:



Correlations (priklad_vyska_vaha.sta)				
Marked correlations are significant at $p < ,05000$				
N=13 (Casewise deletion of missing data)				
Variable	Means	Std.Dev.	vyska	vaha
vyska	173,3846	5,31568	1,000000	0,639675
vaha	65,8462	12,54224	0,639675	1,000000

p-hodnota  $< 0,05$  - test hypotézy  $H_0: r = 0$ , lze vypsát i konkrétní hodnotu (změna formy výstupu na záložce **Options**)

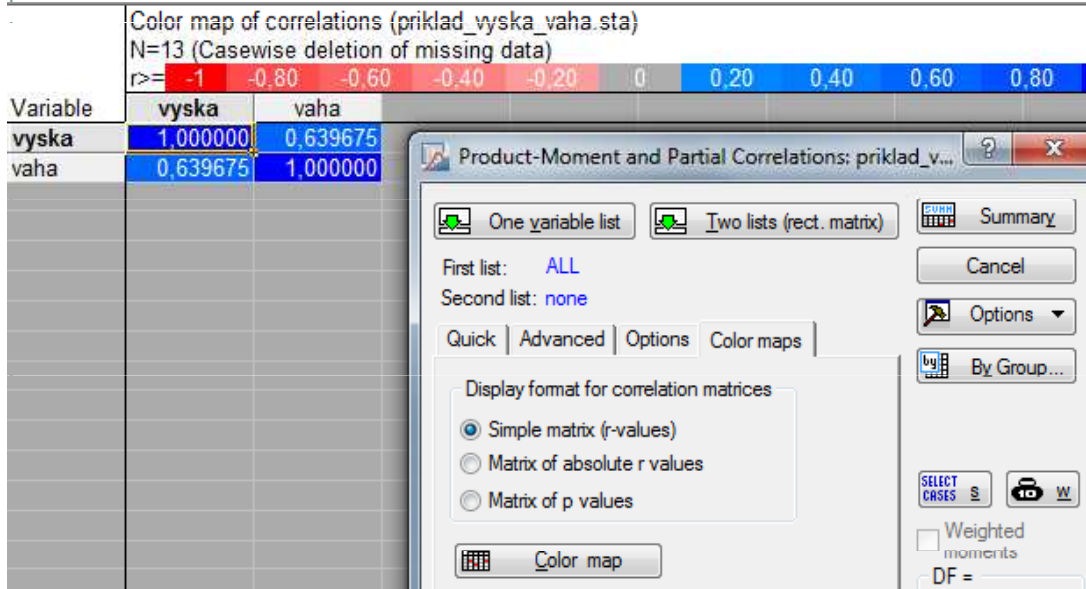
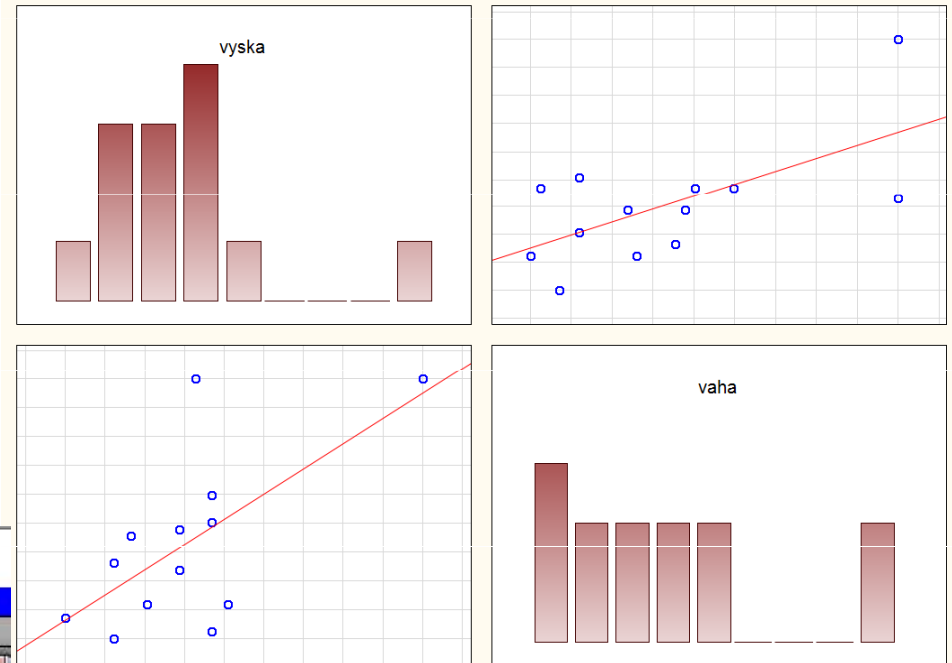
**Pearsonovy korelační koeficienty**

# Řešení v softwaru Statistica: Pearsonův korelační koeficient III

Záložka **Quick / Advanced** umožňuje vykreslit různé druhy grafů (2D, 3D v případě více proměnných, matice bodových grafů s histogramy na diagonále ap.).

*Jsou v daném případě splněny předpoklady (dvourozměrné normální rozdělení, absence odlehlých pozorování, lineární vztah)?*

Correlations (priklad\_vyska\_vaha.sta 2v\*13c)



Na záložce **Color maps** můžeme získat matici korelačních koeficientů (nebo příslušných p-hodnot) obarvenou dle odpovídající barevné škály. Vhodné zejména při zkoumání vztahů mezi více spojitými proměnnými.

# Řešení v softwaru Statistica: Spearmanův korelační koeficient I

Prozkoumejte pořadový vztah mezi výškou a váhou u 13 studentů. Testujte hypotézu, že jsou tyto proměnné nezávislé.

1. Záložka **Statistics**
2. **Nonparametrics**
3. **Correlations**
4. Potvrdíme: **OK**

The screenshot shows the Statistica software interface. The ribbon is set to 'Statistics'. The data table contains the following information:

	1	2
	vyska	vaha
1	175	69
2	166	55
3	170	67
4	169	52
5	188	90
6	175	53
7	176	57
8	171	57
9	173	68
10	175	62
11	173	62
12	174	90
13	169	63

The 'Nonparametric Statistics: prikad\_vyska\_vaha.sta' dialog box is open, showing the 'Quick' tab. The 'Correlations (Spearman, Kendall tau, gamma)' option is selected. The 'OK' button is highlighted.

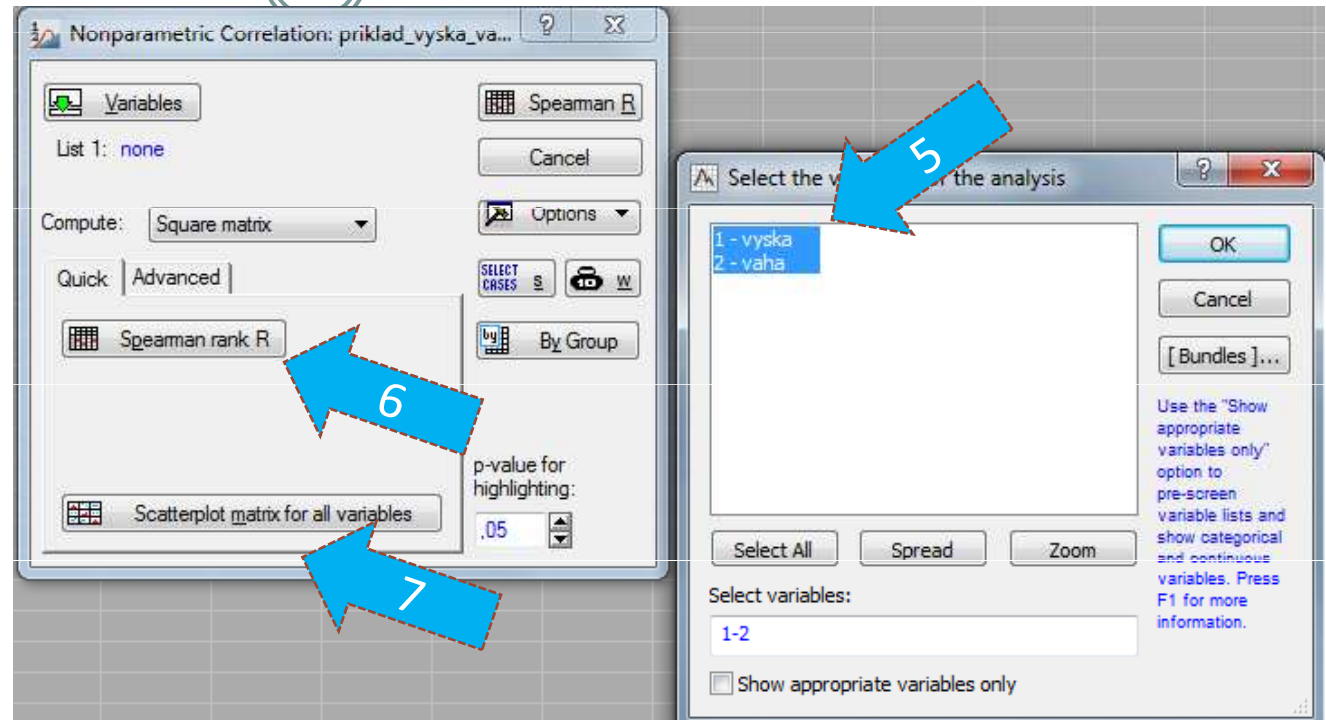
# Řešení v softwaru Statistica: Spearmanův korelační koeficient II

5. Výběr proměnných –  
**Variables – Select variables**  
(vyska, vaha) – **OK**

6. Pod možností Compute  
můžeme vybrat formu  
výstupu (čtvercová matice -  
**Square matrix**, příp. detailní  
výsledky).

7. Lze vykreslit i matici  
bodových grafů s histogramy  
na diagonále (**Scatterplot  
matrix for all variables**).

Jedna z forem výstupní  
tabulky:



p-hodnota  $< 0,05$  - test hypotézy  $H_0: r = 0$ ,  
lze vypsát i konkrétní hodnotu

Spearman Rank Order Correlations (priklad_vyska_vaha.sta)		
MD pairwise deleted		
Marked correlations are significant at: $p < .05000$		
Variable	vyska	vaha
vyska	1.000000	0.469452
vaha	0.469452	1.000000

**Spearmanovy  
korelační  
koefficienty**