

# Neparametrické testy

---

- parametrické a neparametrické testy
  - pořadové neparametrické testy
  - test Chí-kvadrát
    - test nezávislosti proměnných
    - test dobré shody
-

# Parametrické testy

---

- ❑ t-testy a analýza rozptylu jsou tzv. parametrické testy
  - ❑ parametr = charakteristika populace (průměr, rozptyl)
  - ❑ parametrické testy používají při výpočtech charakteristiky populace (parametry)
-

# Parametrické testy

---

- parametrické testy pracují s předpoklady o charakteristikách populace
  - např. u t-testu předpokládáme, že směrodatné odchyly výběrů mohou posloužit jako odhad pro směrodatnou odchylku populace
  - podobně počítají s normálním rozdělením měřeného znaku
-

# Parametrické testy

---

- pokud nejsou tyto předpoklady splněny, můžeme dojít k nepřesným výsledkům

# Neparametrické testy

---

- neparametrické testy nezávisí na charakteristikách populace ani o nich nečiní žádné závěry
  - není vyžadováno normální rozdělení znaku
  - proto jsou tyto testy označovány také jako „distribution-free“ testy
-

# Neparametrické testy

---

- proč potom vůbec používat parametrické testy?
    - mnoho parametrických testů je poměrně „odolných“ (tzv. robustních) vůči narušení předpokladů testu (např. menší odchylky od normálního rozdělení výsledky nezkreslí)
    - parametrické testy mají větší statistickou sílu než neparametrické (větší pravděpodobnost zjištění rozdílu, pokud skutečně existuje)
    - pro některé typy analýz neparametrické metody nejsou (např. neexistuje obecně přijímaná neparametrická faktoriální ANOVA)
-

# Neparametrické testy

---

- hlavní **výhody** neparametrických testů
    - nejsou omezeny předpokladem normálního rozdělení
    - jsou často založeny na pořadí, dají se použít i pro ordinální data (kde můžeme spočítat pouze průměr, nikoli medián) i pro nominální (test Chí-kvadrát)
    - nejsou citlivé na extrémní hodnoty (jsou většinou založeny na mediánu)
-

# Neparametrické testy

---

## □ hlavní **nevýhody** neparametrických testů

- menší statistická síla
  - pro složitější analýzy často není neparametrická varianta metody k dispozici
-



# Neparametrické testy

---

- přehled neparametrických ekvivalentů parametrických testů
    - t-test pro nezávislé výběry – Mann-Whitney U test
    - t-test pro závislé výběry – Wilcoxon test
    - analýza rozptylu – Kruskal-Wallis test
    - opakovaná měření (ANOVA) – Friedman Rank Test
-

# Test Chí-kvadrát

---

- chí-kvadrát může být použit
    - pro testování rozdělení jedné proměnné (test dobré shody)
    - testování nezávislosti dvou proměnných
-

# Test Chí-kvadrát

---

- chí-kvadrát pro testování nezávislosti proměnných se používá pro nominální nebo ordinální proměnné
  - data jsou uspořádána do tzv. kontingenční tabulky (viz příklad)
-

# Příklad

---

- zajímá nás, jak souvisí model manželství s jeho vydařeností
    - model manželství má kategorie: dominance žena, dominance muž, kooperace
    - vydařenost má 3 kategorie – vydařené, průměrné, nevydařené
  - pozn.: jde o manželství rodičů respondentů, tak jak je posuzují oni (zdroj dat – výzkum doc. Plaňavy)
-

# Příklad

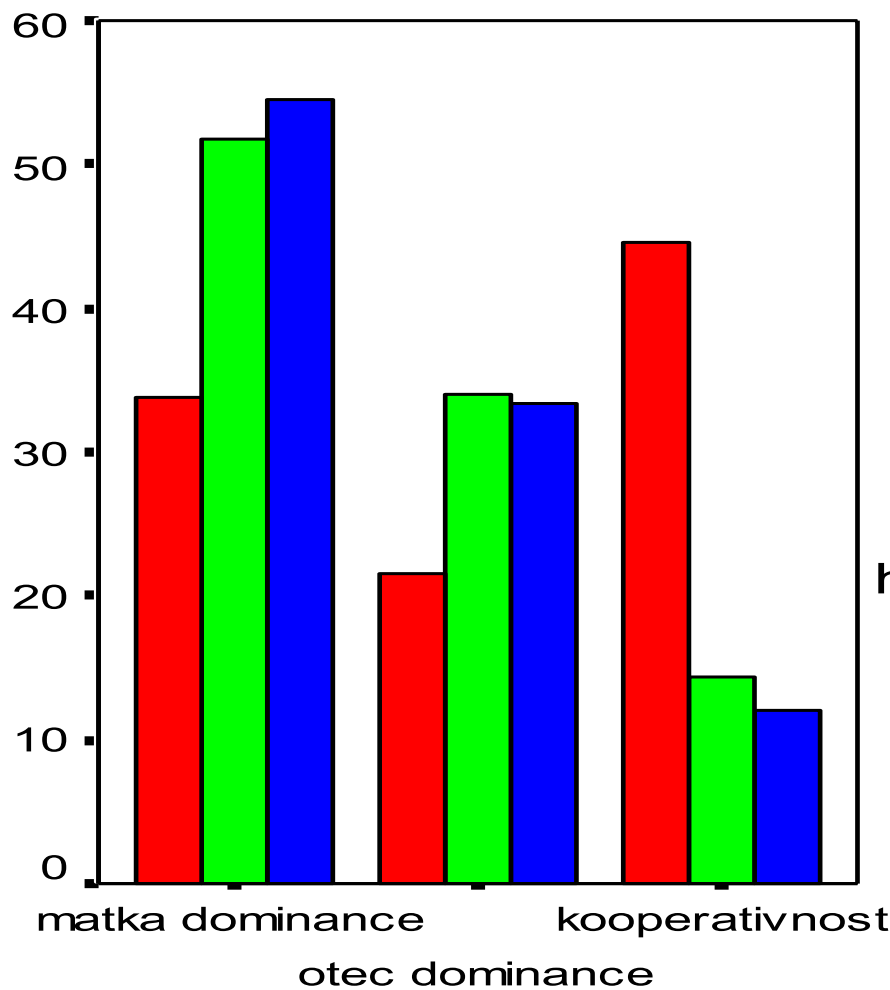
---

- otázka zní: liší se podíl vydařených, průměrných a nevydařených manželství u rodin, kde dominovala matka, rodin, kde dominoval otec a u rodin, kde nedominoval ani jeden z nich?
-

# Kontingenční tabulka (SPSS)

rodiny - muž \* hodnocení manželství rodiců - muž  
 Count

	hodnocení manželství rodiců - muž			Total
	vydarene	umerne	nevdydarene	
rodiny - matka dominanta	22	29	18	69
rodiny - otec dominanta	14	19	11	44
rodiny - kooperativne	29	8	4	41
Total	65	56	33	154



hodnoceni manzelstvi

vydarene

prumerne

nevydarene

model rodic. rodiny - muz

# Test Chí-kvadrát

---

- chí-kvadrát porovnává očekávané a pozorované četnosti
  - očekávané jsou četnosti za předpokladu, že proměnné jsou nezávislé
-



# rodic. rodiny - muz \* hodnoceni manzelstvi rodicu - muz Crosstab

		oceni manzelstvi rodicu -			Total
		ydarenerumerne	nydaren		
model roc matka domin rodiny - r	Count	22	29	18	69
	% within model rodic. rodiny	31,9%	42,0%	26,1%	100,0%
otec domina	Count	14	19	11	44
	% within model rodic. rodiny	31,8%	43,2%	25,0%	100,0%
kooperativno	Count	29	8	4	41
	% within model rodic. rodiny	70,7%	19,5%	9,8%	100,0%
Total	Count	65	56	33	154
	% within model rodic. rodiny	42,2%	36,4%	21,4%	100,0%



# Příklad

---

- v našem příkladu bylo 42,2% vydařených manželství
  - pokud by proměnné (model a vydařenost manželství) byly vzájemně nezávislé, poměr vydařených manželství v jednotlivých modelech manželství by měl být přibližně stejný (a odrážet celkový podíl) – 42%
  - podobně ostatní kategorie...
-

# Test Chí-kvadrát

---

- očekávané četnosti – výpočet:

$$O_{ij} = (r_i s_j) / N$$

(pro každé políčko tabulky se vynásobí celkové četnosti z příslušného řádku se sloupcovými četnostmi a vydělí celkovým počtem osob)

---

# Příklad

---

rodiny - muž \* hodnocení manželství rodiců - muž  
 Count

	hodnocení manželství rodiců - muž			Total
	matka dominantní	otec dominantní	kooperativní	
rodiny - muž	22	29	18	69
rodiny - žena	14	19	11	44
Total	29	8	4	41
Total	65	56	33	154

# Příklad

---

- pro první políčko tabulky (vydařená manželství s dominantní matkou) je očekávaná četnost

$$O_{ij} = (r_i s_j) / N$$

$$O_{11} = (r_1 s_1) / N$$

$$O_{11} = (69 * 65) / 154$$

$$\underline{O_{11}} = 29,12$$

---

# Očekávané četnosti

rodic. rodiny - muz \* hodnoceni manzelstvi rodicu - muz Crossta

			oceni manzelstvi rodicu -			Total
			vydarene	umerne	evydarene	
model rodiny - r	matka domina	Count	22	29	18	69
		Expected C	29,1	25,1	14,8	69,0
	otec domina	Count	14	19	11	44
		Expected C	18,6	16,0	9,4	44,0
	kooperativnc	Count	29	8	4	41
		Expected C	17,3	14,9	8,8	41,0
Total		Count	65	56	33	154
		Expected C	65,0	56,0	33,0	154,0

# Test Chí-kvadrát

---

- chí-kvadrát porovná očekávané četnosti s pozorovanými

$$\chi^2 = \sum [( \text{pozor. četnosti} - \text{oček.} )^2 / \text{oček.}]$$

---

# Příklad

---

$$\chi^2 = \sum [(\text{pozor. četnosti} - \text{oček.})^2 / \text{oček.}]$$

$$\chi^2 = (-7,1)^2/29,1 + 3,9^2/25,1 + 3,2^2/14,8 +$$
$$(-4,6)^2/18,6 + 3^2/16 + 1,6^2/9,4 +$$
$$11,7^2/17,3 + (-6,9)^2/14,9 + (-4,8)^2/8,8$$

$$\chi^2 = \mathbf{18,71}$$

---



# Test Chí-kvadrát

---

- pro vyhledání kritické hodnoty  $\chi^2$  v tabulce musíme vypočítat ještě počet stupňů volnosti (df)
- **df = (ř-1) (s-1)**

(tj. počet řádků -1 krát počet sloupců -1)

---

# Příklad

---

□  $df = (ř-1) (s-1)$

$df = (3-1) * (3-1)$

$df = 4$

□ v tabulkách vyhledáme kritickou hodnotu  $\chi^2$  pro  $df = 4$  a 5% hladinu významnosti

□  $\chi^2_{\text{krit}} = \mathbf{9,49}$

---

# Příklad

---

□  $\chi^2_{\text{krit}} = 9,49$

□  $\chi^2 = 18,71$

- **závěr:** vypočítaná hodnota je větší než kritická hodnota - očekávané a pozorované četnosti se liší na 5% hladině významnosti (tj. je malá pravděpodobnost, že proměnné jsou nezávislé)
-

# Test Chí-kvadrát v SPSS

---

## Chi-Square Tests

	Value	df	symp. Sig. (2-sided)
Pearson Chi-Square	18,712 <sup>a</sup>	4	,001
Likelihood Ratio	18,837	4	,001
Linear-by-Linear Association	11,482	1	,001
N of Valid Cases	154		

a.0 cells (.0%) have expected count less than minimum expected count is 8,79.

# Chí-kvadrát pro 1 proměnnou

---

- tzv. test dobré shody (goodness-of-fit test)
  - opět porovnává očekávané a pozorované četnosti
  - předpokladem očekávaných četností není tentokrát nezávislost proměnných (máme jen 1)
-

# Test dobré shody

---

- jak určíme očekávané četnosti?
  - 2 způsoby:
    - předpoklad vyplývá z teorie (např. u genetických dat – poměr osob s projevem dominantní a recesivní alely)
    - nebo můžeme předpokládat náhodné rozdělení do kategorií
-

# Příklad

---

- je počet sebevražd stejný každý den v týdnu?
  - zjistíme data pro rok 2000 (ČR)
-

# Příklad

---

pondělí	255
úterý	247
středa	240
čtvrtek	206
pátek	236
sobota	192
neděle	226

---



# Příklad

---

## □ očekávané četnosti

- stejný počet sebevražd pro každý den v týdnu
  - celkem 1602 sebevražd
  - očekávaná četnost pro každý den je 228,9
-

# Příklad

---

## sebevraždy - den v týdnu

	Observed	Expected	Residuals
pondělí	255	228,9	26,1
úterý	247	228,9	18,1
středa	240	228,9	11,1
čtvrtek	206	228,9	-22,9
pátek	236	228,9	7,1
sobota	192	228,9	-36,9
neděle	226	228,9	-2,9
Total	1602		

# Příklad

---

- vzorec pro výpočet je stejný
  - $\chi^2 = 13,44$
  - $df = k - 1$  (počet kategorií - 1)
  - $df = 6$
  - pro  $df = 6$  a 5% hladinu významnosti je  $\chi^2_{\text{krit}} = 12,59$
  - **rozdíl je statisticky významný**
-

# Výstup v SPSS

---

## Test Statistics

	nebevrazdy en v tydnu
Chi-Square	13,444
df	6
Asymp.	,036

a.0 cells (,0%) have expected frequency less than 5. The minimum expected cell frequency is

# Omezení Chí-kvadrátu

---

- 2 potenciální problémy:
    - malý počet osob – pokud má velké % políček tabulky očekávanou četnost menší než 5 (v ideálním případě by všechna měla mít oček. četnost nejméně 5 osob)
    - příliš velký počet osob – čím vyšší  $N$ , tím vyšší  $\chi^2$  (vyjdou významné i malé rozdíly)
-

# Míry asociace

---

- obecná definice – síla a směr vztahu
  - míry asociace pro nominální data
    - založené na chí-kvadrátu
    - PRE míry
  - míry asociace pro ordinální data
-

# Míry asociace

---

- míry asociace vyjadřují **těsnost vztahu proměnných** (a případně **směr** vztahu)
  - z chí-kvadrátu se dozvíme pouze, **zda nějaký vztah mezi proměnnými existuje** (tj. zda se liší četnosti pozorované a četnosti očekávané za předpokladu, že proměnné jsou nezávislé)
-

# Míry asociace

---

- **těsnost (síla) vztahu** – vyjádřena absolutní hodnotou koeficientu
  - není shoda v tom, od jaké hodnoty je vztah považován za těsný (někdy uváděno  $>0.70$ , jindy  $>0.30$ ), středně těsný či slabý
-



# Míry asociace

---

- **směr vztahu** – pouze u ordinálních a kardinálních proměnných
  - **pozitivní vztah** – čím vyšší hodnoty jedné proměnné, tím vyšší hodnoty druhé proměnné
  - **negativní vztah** - čím vyšší hodnoty jedné proměnné, tím nižší hodnoty druhé proměnné
-

# Míry asociace pro nominální data

---

- míry asociace pro nominální data ukazují pouze sílu vztahu dvou proměnných, nikoli směr či jiné informace o povaze vztahu
  - rozlišujeme míry založené na chí-kvadrátu a míry PRE
-

# Míry založené na chí-kvadrátu

---

- velikost hodnoty chí-kvadrát je ovlivněna velikostí výběru a počtem kategorií tabulky
  - účelem koeficientů založených na chí-kvadrátu je eliminovat tyto vlivy
-

# Míry založené na chí-kvadrátu

---

- rozsah koeficientů je obvykle mezi 0 a 1
    - čím vyšší hodnota, tím těsnější vztah
    - 0 – žádný vztah
    - 1 – absolutní vztah (z hodnot jedné proměnné můžeme předpovědět hodnoty druhé proměnné)
  - pro koeficienty je možno spočítat statistickou významnost
-

# Míry založené na chí-kvadrátu

---

- mezi nejčastěji užívané míry asociace založené na chí-kvadrátu patří koeficienty
    - $\Phi$  (Phi)
    - Cramerovo  $V$  (Cramer's  $V$ )
    - koeficient kontingence (Contingency Coefficient)
-

# Míry založené na chí-kvadrátu

---

- **Fí koeficient** - užívá se pro tabulky 2x2 (tj. pro dichotomické proměnné, např. pohlaví)
  - vypočte se tak, že se hodnota chí-kvadrátu vydělí počtem osob a výsledek se odmocní
-

# Míry založené na chí-kvadrátu

---

- koeficient kontingence – užívá se někdy místo  $F$  pro tabulky větší než  $2 \times 2$
  - bohužel jeho max. hodnota je nižší než 1 (závisí na počtu políček tabulky)
  - neužívá se proto příliš často
-

# Míry založené na chí-kvadrátu

---

- Cramerovo  $V$  – podobný výpočet jako  $F_i$ ; počet osob se navíc násobí počtem řádků - 1
    - (pokud je počet řádků menší než počet sloupců, jinak počtem sloupců - 1)
  - používá se pro tabulky větší než 2x2
-



# Příklad

---

- příklad z minulé přednášky - jak souvisí model manželství s jeho vydařeností
  - Chí-kvadrát = 18.71
  - počet osob  $N = 154$
  - $m = \text{počet řádků} - 1 = 3 - 1 = 2$
-

# Kontingenční tabulka (SPSS)

rodiny - muž \* hodnocení manželství rodičů - muž Cross  
 Count

	hodnocení manželství rodičů - muž			Total
	vydarene	umerne	nevdydareni	
rodiny - matka dominantni	22	29	18	69
rodiny - otec dominantni	14	19	11	44
rodiny - kooperativni	29	8	4	41
Total	65	56	33	154

# Příklad

---

□ tabulka 3x3 – použijeme Cramerovo V

$$\square V = \sqrt{\chi^2 / (N * m)}$$

$$\square V = \sqrt{18.71 / (154 * 2)}$$

$$\square \mathbf{V = 0,246}$$

---

# Příklad

---

- **interpretace:** hodnota 0,246 je poměrně nízká – vztah mezi modelem manželství a jeho vydařeností není příliš těsný (i když statisticky významný – viz výstup v SPSS)
  - v SPSS jsou uvedeny oba koeficienty (F i V), je třeba zvolit ten správný pro každou tabulku
-

# Výstup v SPSS

---

## Symmetric Measures

	Value	Approx. Sig.
Nominal by Phi	,349	,001
Nominal Cramer's V	,246	,001
N of Valid Cases	154	

a. Not assuming the null hypothesis.

b. Using the asymptotic standard error assuming the null hypothesis.

# Míry PRE

---

- ❑ **PRE** je zkratka pro **Proportional Reduction in Error** (poměrná redukce chyby odhadu)
  - ❑ princip PRE: porovnání odhadu hodnot závislé proměnné bez znalosti hodnot nezávislé proměnné a s její znalostí (o kolik se sníží chyba odhadu?)
-

# Míry PRE

---

- **příklad** – jaký je vztah mezi pohlavím a užíváním rtěnky?\*
  - vypočítáme koeficient **lambda**
  - pokud bychom měli odhadnout, zda náhodně vybraný respondent používá rtěnku: jaká je pravděpodobnost chybného odhadu?
- 
- \*převzat z Disman: Jak se vyrábí sociologická znalost
-

# Míry PRE

---

- můžeme očekávat, že více lidí rtěnku nepoužívá než používá (naprostá většina mužů + některé ženy)
  - takže bude výhodnější odhadnout, že náhodně vybraný respondent rtěnku nepoužívá
  - pravděpodobnost chyby závisí na podílu lidí užívajících rtěnku
-



# Míry PRE

---

## RTENKA

	Frequency	Percent
Valid nepoužíva	97	60,6
používa	63	39,4
Total	160	100,0

---

# Míry PRE

---

- při tomto podílu osob je pravděpodobnost chyby asi 40% (když budeme odhadovat, že náhodný respondent rtěnku neužívá)
  - ze 160 případů bychom se zmýlili 63x
-

# Míry PRE

---

- o kolik by se chyba zmenšila, kdybychom znali pohlaví respondenta?**
  - pro muže bychom odhadovali, že rtěnku nepoužívá, pro ženu naopak - že ji používá
-

# Míry PRE

---

## POHLAVI \* RTENKA Crosstabulation

Count

	RTENKA		Total
	nepoužívá	používá	
POHLAVI muži	78	2	80
ženy	19	61	80
Total	97	63	160

---

# Míry PRE

---

- pokud bychom znali pohlaví respondenta, zmýlili bychom se ve svém odhadu 21x (2 x u muže a 19x u ženy)
  - **o kolik by se náš odhad zlepšil?**
-

# Míry PRE

---

- chyby předtím – chyby teď  
=  $63 - 21 = 42$
  - poměrná redukce chyby (tj. vzhledem k předchozím chybám) = **lambda** =  $42/63 = \mathbf{0,667}$
  - **chyba v odhadu užívání rtěnky se sníží asi o 67%, pokud známe pohlaví respondenta**
-

# Míry PRE

---

- rozsah koeficientu  $\lambda$  je od 0 do 1
  - **0** znamená, že znalost hodnoty nezávislé proměnné vůbec nesníží chybu v odhadu hodnot závislé proměnné; **proměnné jsou vzájemně nezávislé**
  - čím blíže **1**, tím lépe můžeme z hodnot nezávislé proměnné předpovědět hodnoty závislé proměnné
-

# Míry PRE

---

- v SPSS jsou počítány 3 varianty koeficientu lambda
    - symetrická – není určeno, co je závislá a co nezávislá proměnná
    - 2 asymetrické – pro proměnnou 1 jako závislou a pro proměnnou 2 jako závislou
-



# Výstup v SPSS

---

## Directional Measures

	Value	Asymp. Std. Error	Approx. T	Approx. Sig.
Nominál Lambd Symmetric	,706	,063	7,120	,000
by Nominál POHLAVI Dep	,738	,051	9,187	,000
Nominál RTENKA Dep	,667	,082	5,057	,000

# Míry PRE pro nominální data

---

- kromě koeficientu lambda se užívají také
    - Goodmanovo a Kruskalovo **tau**  
(nevyužívá při predikci nejčastější kategorii závislé proměnné jako lambda, ale rozdělení ve všech kategoriích závisle proměnné)
    - Cohenova **Kappa** – pro měření **shody dvou posuzovatelů**
-

# Míry asociace pro ordinální data

---

- u ordinálních dat je výpočet založen na poměru souhlasných a nesouhlasných párů případů
  - **souhlasný** pár případů – hodnota obou proměnných je vyšší (nebo nižší) u jednoho člena páru
  - **nesouhlasný** pár případů – hodnota jedné proměnné je u jednoho člena páru vyšší a hodnota druhé proměnné je nižší
-

# Míry asociace pro ordinální data

---

- pokud je většina párů souhlasných, je hodnota gamma kladná – tj. **pozitivní vztah** (až +1)
  - pokud je většina párů nesouhlasných, je hodnota gamma záporná – tj. **negativní vztah** (až -1)
  - pokud je počet souhlasných a nesouhlasných párů vyrovnán – gamma kolem 0
-

# Míry asociace pro ordinální data

---

- gamma je symetrická míra – nedělá rozdíly mezi závislou a nezávislou proměnnou
  - asymetrická varianta koeficientu gamma – **Somersovo D**
  - **Kendalovo tau b** – bere v úvahu i nerozhodné páry (tzv. ties); ale hodnoty v rozsahu -1 až +1 mohou být získány pouze pro čtvercové tabulky (tj. stejný počet kategorií obou proměnných)
-

# Shrnutí

---

- u nominálních dat hodnota míry asociace proměnných indikuje sílu vztahu – rozsah od 0 do 1
    - nejužívanější  $F_i$  nebo Cramerovo  $V$ ; když víme, která proměnná nezávislá - lambda
  - u ordinálních dat míry asociace indikují jak sílu vztahu (abs. hodnota koeficientu), tak směr vztahu
-

# Kontrolní otázky

---

- ❑ hlavní rozdíl mezi parametrickými a neparametrickými testy
  - ❑ výhody a nevýhody neparametrických testů
  - ❑ kdy je možno využít chí-kvadrát jako test nezávislosti proměnných? (pro jaké typy proměnných?)
  - ❑ kdy se chí-kvadrát využívá jako test dobré shody?
-

# Kontrolní otázky

---

- k čemu slouží míry asociace proměnných?
  - rozdíl mezi mírami založenými na chí-kvadrátu a mírami PRE
  - nejužívanější míry pro nominální data
  - nejužívanější míry pro ordinální data
-