

PSY117

Statistická analýza dat v psychologii

**Přednáška 11 2016**

---

# **TESTY PRO KATEGORICKÉ PROMĚNNÉ – NEPARAMETRICKÉ METODY**

... a to mělo, jak sám vidíte, nedozírné následky.

*Smrt'*

# Analýza četností hodnot kategorických (=O, N) proměnných

---

Výzkumné otázky...

- Liší se preference politických stran?
- Liší se poměrné zastoupení kuřáků mezi ženami a muži?
- Souvisí nějak individuální volební preference s odhadem měsíčního příjmu respondenta?
- Otázky směřují
  - buď k rozdílu četností různých jevů v rámci jedné proměnné (četnost různých jevů v populaci),
  - k rozdílu četností jevu mezi různými proměnnými (četnost jevu v různých populacích),
  - Nebo k pravděpodobnosti výskytu dvou (či více) jevů současně.

# $\chi^2$ test dobré shody

---

- Liší se empirické četnosti nějakých jevů od teoreticky očekávaných četností?
  - Preference politických stran ve volbách...
  - Tedy jedna nominální proměnná, jeden výběr

- Testujeme  $p$  rozdílů mezi empirickými-pozorovanými ( $f_o$ ) a očekávanými ( $f_e$ ) četnostmi

- Mírou rozdílu je hodnota  $\chi^2$ , která má rozložení  $\chi^2$  s  $\nu=k-1$  stupni volnosti a průměrem  $\nu$

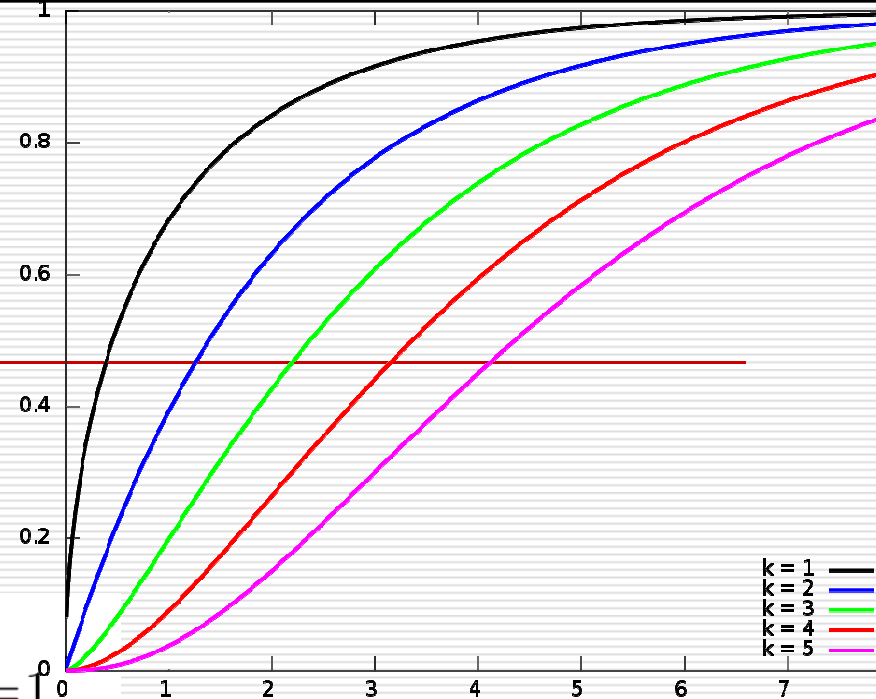
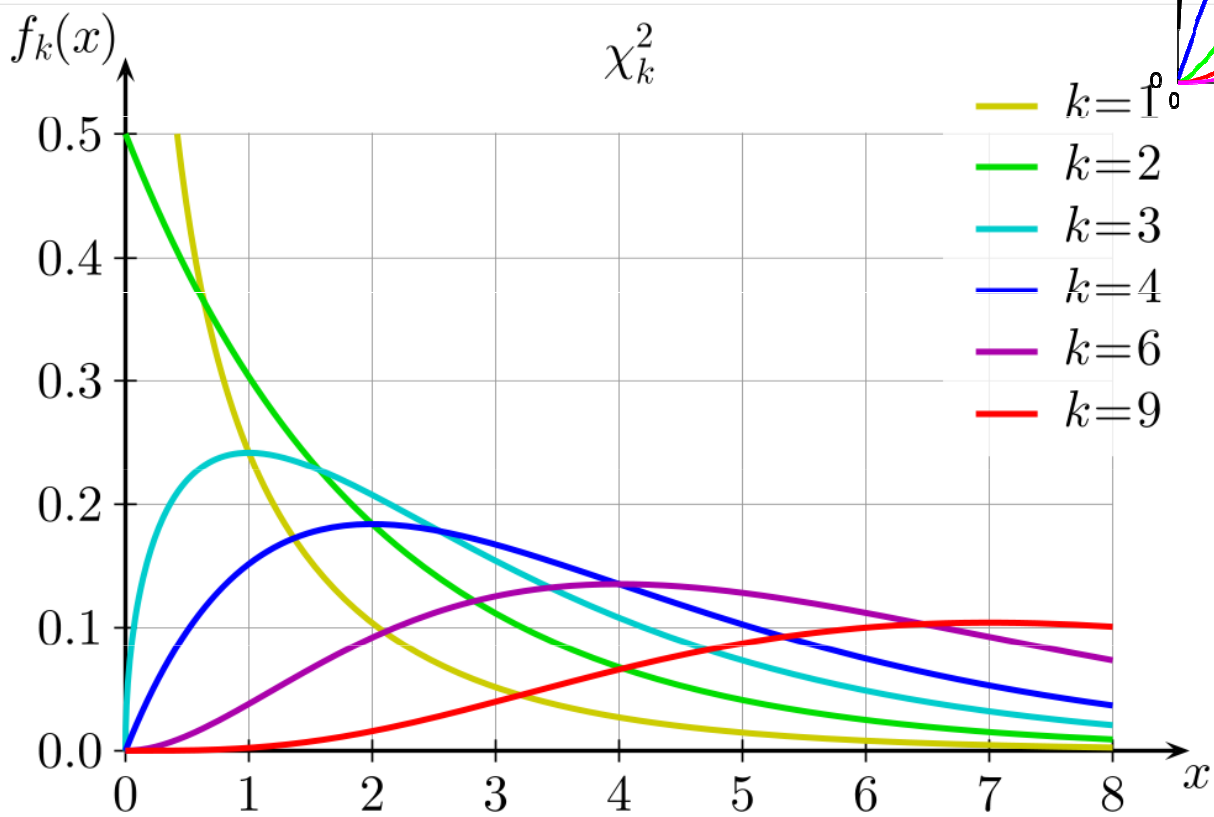
$$\chi^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i} = \sum_{i=1}^k \frac{(fo_i - fe_i)^2}{fe_i}$$

- $H_0: \chi^2 = \nu$  vs.  $H_1: \chi^2 > \nu$

$k$  je počet kategorií,  $n$  velikost vzorku,  $n_i$  četnost kat.  $i$ ,  
 $p_i$  teoretická  $p$ -nost jevu v kategorii  $i$ ;

- Pro získání pravděpodobnosti  $\chi^2$  CHISQ.DIST( $\chi^2$ ; df; 1); CHISQ.INV( $p$ ; df)
- Očekávané četnosti stanovujeme na základě teoretického předpokladu.
- $n_i$  a  $np_i$  vždy jako **četnosti**; nikdy ne procenta (ztráta informace o velikosti vzorku)

# Rozdělení $\chi^2$



# Ve kterém městě byste žili nejraději?

---

Kategorie	n	p	np	(n-np)^2/np
Paříž	28	0,2	28	0
New York	28	0,2	28	0
Londýn	28	0,2	28	0
L.A.	28	0,2	28	0
Tokio	28	0,2	28	0
Celkem	140	1	140	0
Chi <sup>2</sup>				<b>0</b>

$$\chi^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i}$$

$$P(\chi^2 > 0 \mid \chi^2 = 4) \approx 1$$

---

# Ve kterém městě byste žili nejraději?

---

Kategorie	n	p	np	(n-np)^2/np
Paříž	38	0,2	28	3,57
New York	37	0,2	28	2,89
Londýn	22	0,2	28	1,29
L.A.	25	0,2	28	0,32
Tokio	18	0,2	28	3,57
<b>Celkem</b>	140	1	140	11,64
<b>Chi2</b>				<b>11,64</b>

$$\chi^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i}$$

$$P(\chi^2 > 11,64 \mid \chi^2 = 4) = 1 - \text{CHISQ.DIST}(11,64; 4; 1) = 0,02$$

---

# Závislost kategorických proměnných

- ❑ Jaká je souvislost preference politické strany a úrovně hrubého příjmu voliče?
- ❑ Jaká je pravděpodobnost společného výskytu dvou jevů  $x$  a  $y$  možných?
- ❑ Kontingenční tabulka ... řádky  $\times$  sloupce =  $r \times s$ ;  $i \times j$
- ❑ Ve těle tabulky jsou četnosti jednotlivých kombinací, v okrajích tzv. **marginální četnosti** – sumy sloupců nebo řádků. Tedy  $n_{12}$  znamená počet osob ve druhém sloupci prvního řádku; počet osob, u nichž nastal jev  $A_1$  a současně  $B_2$ .

Kategorie	$B_1$	$B_2$	...	$B_s$	Řádkové součty
$A_1$	$n_{11}$	$n_{12}$	...	$n_{1s}$	$n_{1.}$
$A_2$	$n_{21}$	$n_{22}$	...	$n_{2s}$	$n_{2.}$
...	...	...	...	...	...
$A_r$	$n_{i1}$	$n_{i2}$	...	$n_{ij}$	$n_{i.}$
<b>Sloupcové součty</b>	<b><math>n_{.1}</math></b>	<b><math>n_{.2}</math></b>	...	<b><math>n_{.j}</math></b>	<b><math>n</math></b>

# Závislost kategorických proměnných

- ❑  $\chi^2$  test nezávislosti(homogeneity)
- ❑ Očekávané četnosti  $f_e: m_{ij}$  (očekávaná četnost v  $i$ - $j$ -té buňce)( $i$  – řádky,  $j$  – sloupce)
- ❑ Testová statistika je  $\chi^2$
- ❑ Stupně volnosti:  $df = (i-1)*(j-1)$

$$f_e = m_{ij} = \frac{n_{i.} \cdot n_{.j}}{n} \quad \chi^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(n_{ij} - m_{ij})^2}{m_{ij}} = \sum_{i=1}^r \sum_{j=1}^s \frac{(fo_{ij} - fe_{ij})^2}{fe_{ij}}$$

Kategorie	B <sub>1</sub>	B <sub>2</sub>	...	B <sub>s</sub>	Řádkové součty
A <sub>1</sub>	n <sub>11</sub>	n <sub>12</sub>	...	n <sub>1s</sub>	n <sub>1.</sub>
A <sub>2</sub>	n <sub>21</sub>	n <sub>22</sub>	...	n <sub>2s</sub>	n <sub>2.</sub>
...	...	...	...	...	...
A <sub>r</sub>	n <sub>i1</sub>	n <sub>i2</sub>	...	n <sub>ij</sub>	n <sub>i.</sub>
<b>Sloupcové součty</b>	<b>n<sub>.1</sub></b>	<b>n<sub>.2</sub></b>	<b>...</b>	<b>n<sub>.j</sub></b>	<b>n</b>



# Př. $\chi^2$ test nezávislosti(homogeneity)

## Vztah bydliště a počtu holínek

Pozorované Řádková %	0	1	>2	Řádkové součty
<b>Velkoměsto</b>	10 67%	1 7%	4 27%	<b>15</b>
<b>Maloměsto</b>	15 43%	19 54%	1 3%	<b>35</b>
<b>Vesnice</b>	15 30%	20 40%	15 30%	<b>50</b>
<b>Sloupcové součty</b>	<b>40</b>	<b>40</b>	<b>20</b>	<b>100</b>

Očekávané/ dílčí $\chi^2$	0	1	>2	Řádkové součty
<b>Velkoměsto</b>	6/ 2,7	6/ 4,2	3/ 0,3	<b>15</b>
<b>Maloměsto</b>	14/ 0,1	14/ 1,8	7/ 5,1	<b>35</b>
<b>Vesnice</b>	20/ 1,3	20/ 0	10/ 2,5	<b>50</b>
<b>Sloupcové součty</b>	<b>40</b>	<b>40</b>	<b>20</b>	<b>100</b>

$$\chi^2 = 17,9 \quad df = (3-1) * (3-1) = 4 \quad P(\chi^2 > 17,9 \mid \chi^2 = 4) = 0,001$$

# Síla vztahu v kontingenční tabulce

---

- ❑ Pro tabulky 2x2 **Phi**  $\phi = \sqrt{\frac{\chi^2}{n}}$
- ❑ Pro tabulky 3x3 a více **koeficient kontingence** (Pearson)  $C = \sqrt{\frac{\chi^2}{\chi^2+n}}$
- ❑ Pro tabulky  $r \times s$  **Cramerovo V**  $V = \sqrt{\frac{\chi^2}{n(k-1)}}$  kde  $k$  je menší z  $r$  a  $s$
- ❑ Všechny koeficienty v intervalu  $<0;1>$ .
- ❑ Pro tabulky větší než 2x2 je často třeba identifikovat buňku(y), kde jsou největší odchylky od očekávaných četností
  - Skrze výpočet **reziduí**, tj. rozdílů mezi pozorovanou a očekávanou četností:  $n_{ij} - m_{ij} = res_{ij}$ 
    - ❑ tyto „zbytkové“ hodnoty lokalizují odchylky od pravděpodobnostního rozdělení
    - ❑ Součet residuí v tabulce je vždy nula
  - **Standardizovaná rezidua** (Pearsonova):  $R = (n_{ij} - m_{ij})/\sqrt{m_{ij}}$ 
    - ❑ rozdělení standardizovaných residuí je normální s průměrem 0 a sm. odchylkou 1; tedy  $R \geq \pm 1,96$  jsou „zajímavá“ pro interpretaci, významně přispívají k signifikanci  $\chi^2$ .
- ❑ Analýza tabulky skrze  $\chi^2$  je nespolehlivá, je-li  $\min(m_{ij}) < 5$ . *I řídké jevy musí mít šanci ☺*
  - Pro tuto situaci máme tzv. permutační testy (v SPSS „exact“)
- ❑ Hendl str. 297 – 313.

kontingenční koeficient  $C = \sqrt{(17,9/(17,9+100))}=0,4$   
 Cramérovo  $V = \sqrt{(17,9/(100*2))}=0,3$

Pozorované Řádková % St. rezidua	0	1	>2	Řádkové součty
<b>Velkoměsto</b>	10 67% <b>1,6</b>	1 7% <b>-2,0</b>	4 27% <b>0,6</b>	<b>15</b>
<b>Maloměsto</b>	15 43% <b>0,3</b>	19 54% <b>1,3</b>	1 3% <b>-2,3</b>	<b>35</b>
<b>Vesnice</b>	15 30% <b>-1,1</b>	20 40% <b>0</b>	15 30% <b>1,6</b>	<b>50</b>
<b>Sloupcové součty</b>	<b>40</b>	<b>40</b>	<b>20</b>	<b>100</b>

# Testy středních hodnot pro ordinální proměnné – neparametrické metody

---

- Metody užívající *parametrů* normálního rozložení ( $m, s$ ) mají svá omezení, když...
  - data pochází z rozložení, které se od normálního výrazně liší (tvar, či odlehlé hodnoty)
  - data mají spíše ordinální charakter; nebo se jedná o krátké intervalové škály
- *Neparametrické* metody
  - jsou *robustní* vůči rozložení dat...
  - mají nižší sílu testu (tj. vyšší požadavky na velikost vzorku)
- Testy pro mediány
  - Pro jeden výběr: znaménkový test, Wilcoxonův test
  - Pro párové srovnání: Wilcoxonův test
  - Pro 2 nezávislé výběry: Mann-Whitney U, Kolmogorov-Smirnov Z

# Jeden výběr, znaménkový test

---

- $H$ : Je medián roven  $k$ ?  $H_0: Md = k; H_1: Md \neq k$
  - Platí-li  $H_0$ , mělo by nad i pod hypotetizovaným mediánem být stejné množství případů
  - Asymptotický test pomocí normálního rozložení:
    - $Z^+$  ( $Z^-$ ) je počet hodnot vyšších (nižších) než hypotetizovaný medián
    - Hodnoty rovné mediánu ignorujeme a odečítáme z  $n$
    - Platí-li  $H_0$ ,  $Z^+ = Z^-$  a  $Z^+ + Z^- = n$ .
    - Testová statistika  $z = (2Z_+ - n)/\sqrt{n}$  má asymptoticky normální rozložení, (přesně má binomické rozložení).
    - $P = 2 * (1 - \text{NORM.S.DIST}(z))$
  - Jedná se tedy o alternativu  $t$ -testu pro jeden výběr;
  - Pro závislé výběry (=párové srovnání) spočítáme  $d_i = x_i - y_i$  a znaménkovým testem testujeme  $H_0: Md_d = 0$ .
-

# Neparametrické testy pro nezávislé výběry

---

## □ Mediánový test

- Je-li společný medián dvou výběrů shodný, leží na jedné straně  $Md$  50% každého výběru.
- Určíme  $Md$  pro celý soubor; pokud platí  $H_0$ , četnosti hodnot ležících nad i pod  $Md$  by měly být stejné pro  $x$  i  $y$ .
- Pokud  $H_0$  neplatí, budou četnosti výrazně asymetrické, v „diagonále“.
- Při  $n > 30$  lze užít asymptoticky normálně rozloženou testovou statistiku  $z$ :

$$z = \frac{(ad - bc)\sqrt{n}}{\sqrt{(a+b)(b+d)(a+c)(c+d)}}$$

	Sk A	Sk B	$\Sigma$
<Md	$a$	$b$	$a+b$
>Md	$c$	$d$	$c+d$
$\Sigma$	$a+c$	$b+d$	$n$

Silnější alternativou je Wilcoxonův test pro nezávislé výběry nebo Mann-Whitney U, popřípadě další.

---

# Shrnutí

---

- Pro nominální data máme testy založené na chí-kvadrátu
    - Test dobré shody
    - Test nezávislosti/homogeneity
  - Pro ordinální data a výrazně nenormálně rozložená intervalová máme „neparametrické“ testy
    - Jejich primitivní verze jsem si ukázali
    - „Pojmenované“ testy je zpřesňují
-