

11 Testování nezávislosti v kontingenčních tabulkách

11.1 Kontingenční tabulky

- $(X_1, Y_1), \dots, (X_n, Y_n)$... dvouozměrný náhodný výběr rozsahu n
- X, Y ... nominální znaky:
- znak X ... r variant: $x_{[1]}, \dots, x_{[r]}$; znak Y ... s variant: $y_{[1]}, \dots, y_{[s]}$

Kontingenční tabulka (KT)

- Absolutní četnosti v KT
 - n_{jk} ... absolutní simultánní četnosti j -té variandy znaku X a k -té variandy znaku Y
 - $n_{j\cdot} = n_{j1} + \dots + n_{js}$... absolutní marginální četnosti j -té variandy znaku X
 - $n_{\cdot k} = n_{1k} + \dots + n_{rk}$... absolutní marginální četnosti k -té variandy znaku Y

Pearsonův χ^2 test

- asymptotický test
 - musíme ověřit podmínu dobré approximace
 - `chisq.test(data, correct=F)$expected`
 - alespoň 80 % případů musí být ≥ 5 a zbylých 20 % nesmí klesnout pod 2.
- $H_0 : X, Y$ jsou stochasticky nezávislé.
- $H_1 : X, Y$ nejsou stochasticky nezávislé.
- porovnáváme pozorované četnosti n_{jk} a teoretické četnosti $\frac{n_{j\cdot}n_{\cdot k}}{n}$ dvojice variant $(x_{[j]}, y_{[k]})$
- za platnost H_0 si jsou n_{jk} a $\frac{n_{j\cdot}n_{\cdot k}}{n}$ podobné
- Testovací statistika:
$$K = \sum_{j=1}^r \sum_{k=1}^s \frac{\left(n_{jk} - \frac{n_{j\cdot}n_{\cdot k}}{n}\right)^2}{\frac{n_{j\cdot}n_{\cdot k}}{n}}$$
- Kritický obor: $W = \langle \chi^2_{1-\alpha}((r-1)(s-1)), \infty \rangle$
- `chisq.test(data, correct=F)`

Měření závislosti, Cramérův koeficient

- Cramérův koeficient

$$V = \sqrt{\frac{K}{n(m-1)}},$$

kde $m = \min\{r, s\}$.

Cramérův koeficient	interpretace
0 – 0.1	zanedbatelná závislost
0.1 – 0.3	slabá závislost
0.3 – 0.7	střední závislost
0.7 – 1	silná závislost

- `cramersV(data)` z knihovny `lsr`

11.2 Čtyřpolní kontingenční tabulky

- náhodné veličiny X, Y mají pouze 2 varianty \rightarrow čtyřpolní kontingenční tabulka
- značení: $n_{11}a =, n_{12} = b, n_{21} = c, n_{22} = d$

11.2.1 Pearsonův χ^2 test

- viz výše; asymptotický test
- kritický obor: $W = \langle \chi^2_{1-\alpha}(1), \infty \rangle$
- `chisq.test(data, correct=F)`

11.2.2 Fisherův faktoriálový test

- přesný test
- `fisher.test(data)`

Podíl šancí ve čtyřpolní KT

- pokus se provádí za dvojích různých okolností a může skončit buď úspěchem nebo neúspěchem

- 1.okolnost: podíl počtu úspěchů ku počtu neúspěchů: $\frac{a}{c}$
- 2.okolnost: podíl počtu úspěchů ku počtu neúspěchů: $\frac{b}{d}$
- $o\rho \dots$ teoretický podíl šancí
 - X, Y nezávislé \rightarrow potom $o\rho = 1$

- $OR \dots$ výběrový podíl šancí

$$OR = \frac{\frac{a}{c}}{\frac{b}{d}} = \frac{ad}{bc}$$

- Závislost X, Y je tím silnější, čím více se OR ($o\rho$) liší od 1.
- OR resp. $o\rho \in \langle 0; \infty \rangle \rightarrow$ preferujeme logaritmus podílu šancí
- $\ln(OR)$ resp. $\ln(o\rho) \in \langle -\infty; \infty \rangle$

Test podílem šancí

- $H_0 : X, Y$ jsou stochasticky nezávislé $\dots \ln o\rho = 0$
- $H_1 : X, Y$ nejsou stochasticky nezávislé $\dots \ln o\rho \neq 0$.
- Testová statistika

$$T_0 = \frac{\ln OR}{\sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}}$$

- Kritický obor: $W = (-\infty; -u_{1-\alpha/2}) \cup (u_{1-\alpha/2}; \infty)$
- $100(1 - \alpha)\%$ **asymptotický** interval spolehlivosti

$$(d, h) = \left(\ln OR - \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}} u_{1-\alpha/2}; \ln OR - \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}} u_{\alpha/2} \right).$$