

# 11 Kontingenční tabulky

## 11.1 Test nezávislosti

Nechť náhodný vektor  $\mathbf{X} = (Y, Z)'$  má diskrétní rozdělení, přičemž veličina  $Y$  nabývá hodnot  $1, \dots, r$  a veličina  $Z$  hodnot  $1, \dots, c$ . Označme

$$p_{ij} = P(Y = i, Z = j), \quad p_{i.} = \sum_j p_{ij}, \quad p_{.j} = \sum_i p_{ij}.$$

Předpokládejme, že se uskutečnil výběr o rozsahu  $n$  z tohoto rozdělení. Nechť  $n_{ij}$  je počet těch případů, kdy se ve výběru vyskytla dvojice  $(i, j)$ . Náhodné veličiny  $n_{ij}$  mají pak sdružené *multinomické rozdělení* s parametrem  $n$  a s pravděpodobnostmi  $p_{ij}$ . Proti odstavci 10.1 máme tu nepodstatnou odlišnost, že pravděpodobnosti  $p_{ij}$  i empirické četnosti  $n_{ij}$  píšeme ve tvaru matic místo ve tvaru vektorů. Matici  $(n_{ij})$  se říká *kontingenční tabulka*. Podobně jako v dřívějších analogických situacích píšeme

$$n_{i.} = \sum_j n_{ij}, \quad n_{.j} = \sum_i n_{ij}.$$

Samozřejmě platí

$$n = \sum_i n_{i.} = \sum_j n_{.j} = \sum_i \sum_j n_{ij}.$$

Číslům  $p_{i.}$  a  $p_{.j}$  se říká *marginální pravděpodobnosti* a hodnotám  $n_{i.}$  a  $n_{.j}$  *marginální četnosti*. Matice pravděpodobností  $(p_{ij})$  a kontingenční tabulka  $(n_{ij})$  jsou uvedeny v tab. 11.1.

V praxi vzniká kontingenční tabulka tak, že se na statistických jednotkách sledují dva znaky. Tyto znaky mohou být svou povahou diskrétní a nabývat jen konečně mnoha hodnot (např. muž — žena, krevní skupina 0 — A — B — AB), nebo sami záměrně vytvoříme jen konečně mnoho kategorií (místo uvedení přesného reakčního času zkoumaných osob na nějaký podnět se zaznamená jako výsledek jen to, zda reakce byla rychlá či pomalá). Někdy může mít znak objektivně danou číselnou hodnotu (např. počet dětí), jindy přiřazená číselná hodnota zachovává uspořádání, ale nevypovídá o skutečné velikosti znaku (např. barvy uspořádáme do skupin podle rostoucí vlnové délky; přiřazeno však je nakonec jen pořadí barev,

Tabulka 11.1: Matice pravděpodobností a kontingenční tabulka

Matice pravděpodobností

Y	Z			$\Sigma$
	1	$\cdots$	c	
1	$p_{11}$	$\cdots$	$p_{1c}$	$p_{1.}$
$\cdots$	$\cdots$	$\cdots$	$\cdots$	$\cdots$
r	$p_{r1}$	$\cdots$	$p_{rc}$	$p_{r.}$
$\Sigma$	$p_{.1}$	$\cdots$	$p_{.c}$	1

Kontingenční tabulka

Y	Z			$\Sigma$
	1	$\cdots$	c	
1	$n_{11}$	$\cdots$	$n_{1c}$	$n_{1.}$
$\cdots$	$\cdots$	$\cdots$	$\cdots$	$\cdots$
r	$n_{r1}$	$\cdots$	$n_{rc}$	$n_{r.}$
$\Sigma$	$n_{.1}$	$\cdots$	$n_{.c}$	n

ne vlnová délka). V mnoha případech přiřazujeme čísla  $1, 2, \dots$  jen jako označení, ale nejde ve skutečnosti ani o přesné hodnoty ani o objektivní pořadí (např. uvedená čísla slouží jako kódy pro různé dělnické profese). Zde se omezíme na základní popis obecných kontingenčních tabulek a nebudeme se zabývat speciálními metodami, jimiž lze vytěžit další informaci v případě uspořádaných hodnot znaku nebo v případě znaků přímo vyjádřených číselně.

Nejčastější úlohou při rozboru kontingenčních tabulek je provedení testu hypotézy, že veličiny  $Y$  a  $Z$  jsou na sobě nezávislé. Než přistoupíme k odvození tohoto testu nezávislosti, dokážeme jedno přípravné tvrzení.

**Věta 11.1** *Veličiny  $Y$  a  $Z$  jsou nezávislé tehdy a jen tehdy, platí-li  $p_{ij} = p_{i.}p_{.j}$  pro všechny dvojice  $(i, j)$ .*

*Důkaz.* Podle definice uváděné v teorii pravděpodobnosti jsou  $Y$  a  $Z$  nezávislé právě tehdy, platí-li  $P(Y \in A, Z \in B) = P(Y \in A)P(Z \in B)$  pro kterékoli množiny  $A \subset \{1, \dots, r\}$ ,  $B \subset \{1, \dots, c\}$ . Zvolíme-li  $A = \{i\}$ ,  $B = \{j\}$ , pak odtud plyne, že v případě nezávislosti veličin  $Y$  a  $Z$  musí platit  $p_{ij} = p_{i.}p_{.j}$ . Obrácenou implikaci dokážeme pro  $A = \{1, 2\}$ ,  $B = \{1, 2, 3\}$ . Obecný případ se dokazuje úplně stejně, jen zápis je o něco složitější. Máme

$$\begin{aligned} P(Y \in A, Z \in B) &= \sum_{i=1}^2 \sum_{j=1}^3 p_{ij} = \sum_{i=1}^2 \sum_{j=1}^3 p_{i.}p_{.j} \\ &= \left( \sum_{i=1}^2 p_{i.} \right) \left( \sum_{j=1}^3 p_{.j} \right) = P(Y \in A)P(Z \in B). \quad \square \end{aligned}$$

Proto hypotézu nezávislosti  $H_0$  můžeme psát ve tvaru

$$H_0 : p_{ij} = p_{i.}p_{.j}, \quad i = 1, \dots, r; \quad j = 1, \dots, c.$$

To znamená, že pravděpodobnosti  $p_{ij}$  multinomického rozdělení jsou funkcemi menšího počtu neznámých parametrů, jimiž jsou marginální pravděpodobnosti  $p_i$  a  $p_j$ . Ale tyto marginální pravděpodobnosti nejsou nezávislé, neboť

$$\sum_i p_i = \sum_j p_j = 1.$$

Máme-li splnit podmínku (4) věty 10.4, pak do neznámých parametrů nemůžeme počítat pravděpodobnosti  $p_r$  a  $p_c$ , neboť ty lze už z ostatních marginálních pravděpodobností vypočítat. Počet neznámých parametrů je tudíž  $m = r - 1 + c - 1 = r + c - 2$ . Přitom se samozřejmě omezuje jen na ty situace, kdy všechny marginální pravděpodobnosti jsou kladné. Kdyby tomu tak nebylo, prostě by se některé řádky nebo některé sloupce vynechaly.

K odhadu neznámých parametrů  $p_1, \dots, p_{r-1}$  a  $p_{.1}, \dots, p_{.c-1}$  použijeme soustavu rovnic (10.14). Místo veličin  $X_i$  zde samozřejmě máme veličiny  $n_{ij}$ . Dostáváme

$$\sum_{j=1}^c \left( \frac{n_{ij}}{p_i} - \frac{n_{rj}}{p_r} \right) = 0, \quad i = 1, \dots, r-1 \quad (11.1)$$

a

$$\sum_{i=1}^r \left( \frac{n_{ij}}{p_j} - \frac{n_{ic}}{p_c} \right) = 0, \quad j = 1, \dots, c-1, \quad (11.2)$$

protože pro  $h = 1, \dots, r-1$  platí

$$\frac{\partial p_i p_j}{\partial p_h} = \begin{cases} p_j & \text{pro } i = h, \\ -p_j & \text{pro } i = r, \\ 0 & \text{v ostatních případech.} \end{cases}$$

Podobný výsledek dostaneme pro parciální derivace podle  $p_l$ . Vzorec (11.1) však platí i pro  $i = r$  a (11.2) platí i pro  $j = c$ . Proto místo (11.1) můžeme psát

$$\frac{n_{i.}}{p_i} - \frac{n_{r.}}{p_r} = 0, \quad i = 1, \dots, r. \quad (11.3)$$

Odtud

$$n_i = \frac{n_{r.}}{p_r} p_i, \quad i = 1, \dots, r.$$

Sečtením přes všechna  $i$  dostaneme  $n = n_{r.}/p_r$ , takže odtud máme pro  $p_r$  odhad  $\hat{p}_r = n_{r.}/n$ . Dosazením do (11.3) nakonec získáme odhady

$$\hat{p}_i = \frac{n_{i.}}{n}, \quad i = 1, \dots, r.$$

Podobně se vyřeší i (11.2) a dostane se

$$\hat{p}_{.j} = \frac{n_{.j}}{n}, \quad j = 1, \dots, c.$$

Podle věty 10.4 má pak veličina

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{\left(n_{ij} - \frac{n_{i.}n_{.j}}{n}\right)^2}{\frac{n_{i.}n_{.j}}{n}} \quad (11.4)$$

asymptoticky rozdělení  $\chi^2$ , jehož počet stupňů volnosti je roven

$$rc - (r + c - 2) - 1 = (r - 1)(c - 1).$$

Vzhledem k (10.8) můžeme  $\chi^2$  také počítat podle vzorce

$$\chi^2 = n \sum_{i=1}^r \sum_{j=1}^c \frac{n_{ij}^2}{n_{i.}n_{.j}} - n. \quad (11.5)$$

Když vyjde  $\chi^2 \geq \chi_{(r-1)(c-1)}^2(\alpha)$ , zamítáme hypotézu  $H_0$  o nezávislosti veličin  $Y$  a  $Z$ . Ke shodě s limitním rozdělením se opět vyžaduje, aby všechny teoretické četnosti  $n_{i.}n_{.j}/n$  byly větší než 5. Není-li tato podmínka splněna, spojují se některé řádky nebo některé sloupce.

**Příklad 11.2** U 6800 mužů byla zjišťována barva očí a barva vlasů (viz Yule a Kendall 1950). Výsledky jsou uvedeny v tab. 11.2.

Tabulka 11.2: Barva očí a vlasů

Barva očí	Barva vlasů				Celkem
	Světlá	Kaštanová	Černá	Zrzavá	
Světle modrá	1 768	807	189	47	2 811
Šedá nebo zelená	946	1 387	746	53	3 132
Tmavohnědá	115	438	288	16	857
Celkem	2 829	2 632	1 223	116	6 800

Podle vzorce (11.4) nebo (11.5) vypočteme  $\chi^2 = 1073,51$ . Vzhledem k tomu, že  $\chi^2 \geq \chi_6^2(0,05) = 12,59$ , zamítneme hypotézu, že barva očí a barva vlasů u mužů jsou nezávislé znaky.  $\diamond$

## 11.2 Test homogenity multinomických rozdělení

Někdy se stává, že marginální řádkové četnosti  $n_{i.}$  jsou předem stanoveny. Pak ale  $i$ -tý řádek kontingenční tabulky  $n_{i1} \dots n_{ic}$  má *multinomické rozdělení* s parametry  $n_{i.}, q_{i1}, \dots, q_{ic}$ , kde  $q_{i1}, \dots, q_{ic}$  jsou nějaké pravděpodobnosti splňující podmínku  $q_{i1} + \dots + q_{ic} = 1$ . Většinou je pak třeba testovat *hypotézu homogenity*, která říká, že pravděpodobnosti  $q_{i1}, \dots, q_{ic}$  nezávisí na řádkovém indexu  $i$  [takže všechny řádky matice  $(q_{ij})$  jsou stejné]. Na základě jistého zobecnění věty 10.4 lze dokázat (viz Cramér 1946), že i v tomto případě za platnosti hypotézy homogenity má veličina  $\chi^2$  počítaná podle vzorce (11.4) nebo (11.5) asymptoticky rozdělení  $\chi^2_{(r-1)(c-1)}$ . Hypotéza homogenity se zamítá v případě  $\chi^2 \geq \chi^2_{(r-1)(c-1)}(\alpha)$ .

**Příklad 11.3** V severozápadním Skotsku byla provedena studie, která měla prokázat, zda je procentuelní zastoupení krevních skupin na celém území homogenní či nikoli. V oblasti Eskdale bylo náhodně vybráno 100 osob, v oblasti Annandale 125 osob a v oblasti Nithsdale 253 osob. Výsledky převzaté ze sbírky Osborn (1979) jsou uvedeny v tab. 11.3.

Tabulka 11.3: Krevní skupiny

Oblast	A	B	0	AB	Celkem
Eskdale	33	6	56	5	100
Annandale	54	14	52	5	125
Nithsdale	98	35	115	5	253
Celkem	185	55	223	15	478

Podle (11.5) dostaneme  $\chi^2 = 10,45$ . Toto číslo je menší než kritická hodnota  $\chi^2_6(0,05) = 12,59$ . Nemůžeme tudíž zamítnout hypotézu, že rozdělení krevních skupin je ve všech třech oblastech stejné.  $\diamond$

Všimněme si podrobněji testu homogenity v případě  $c = 2$ .

**Věta 11.4** *Je-li  $c = 2$ , pak*

$$\chi^2 = \frac{n^2}{n_{.1}n_{.2}} \sum_{i=1}^r \frac{n_{i1}^2}{n_{i.}} - n \frac{n_{.1}}{n_{.2}}. \quad (11.6)$$