

# Testy dobré shody při neznámých parametrech, ověřování exponenciálního rozdělení

Pavel Hellebrand  
20.března 2008

# Obsah

- Multinomické rozdělení
- Testy dobré shody při známých parametrech
- Testy dobré shody při neznámých parametrech
- Ověřování exponenciálního rozdělení

# Multinomické rozdělení 1

- Mějme urnu a v ní kuličky  $k$  různých barev. Nechť pravděpodobnost vytažení kuličky  $i$ -té barvy je rovna  $p_i$   $i = 1, 2, \dots, k$ , přičemž

$$(1) \quad 0 < p_i < 1 \quad p_1 + \dots + p_k = 1.$$

- Za těchto podmínek  $n$ -krát nezávisle na sobě vybereme (s vracením) po jedné kuličce. Označme  $X_i$  počet kuliček  $i$ -té barvy, které takto byly vybrány. Je zřejmé, že sdružené rozdělení pravděpodobnosti náhodných veličin  $X_1, \dots, X_k$  je dáno vzorcem

$$(2) \quad P(X_1 = x_1, \dots, X_k = x_k) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}$$

pro  $x_i = 0, 1, \dots, n$  ( $i = 1, 2, \dots, k$ ),  $x_1 + \dots + x_k = n$

- Rozdělení dané vzorcem (2) se nazývá *multinomické s parametry*  $n, p_1, \dots, p_k$

# Multinomické rozdělení 2

- Pro multinomické rozdělení platí, že všechna jeho marginální podmíněná rozdělení jsou opět multinomická.
- Všechna jednorozměrná marginální rozdělení jsou *binomická*. ( $X_i$  má binomické rozdělení s parametry  $n$  a  $p_i$ )
- $EX_i = np_i$ ,  $\text{var } X_i = np_i(1 - p_i)$ ,  $1 \leq i \leq k$ ,  
 $\text{cov}(X_i, X_j) = -np_i p_j$ ,  $1 \leq i \neq j \leq k$ .
- Jestliže  $\mathbf{X} = (X_1, \dots, X_k)'$  má multinomické rozdělení (2), pak náhodná veličina

$$(3) \quad \chi^2 = \sum_{i=1}^k \frac{(X_i - np_i)^2}{np_i}$$

má při  $n \rightarrow \infty$  asymptoticky rozdělení  $\chi_{k-1}^2$

# Multinomické rozdělení 3

- Vzorec (3) lze snadno upravit na tvar

$$\chi^2 = \sum_{i=1}^k \frac{X_i^2}{np_i} - n$$

který je vhodnější pro výpočet, avšak v praxi se dává přednost vzorci (3), protože je z něj lépe vidět jakou měrou přispívá každý sčítanec k celkovému součtu  $\chi^2$ .

# Test dobré shody při známých parametrech 1

- Předpokládejme, že výsledky pozorování byly uspořádány do  $k$  tříd s četnostmi  $X_1, \dots, X_k$  (empirické četnosti)
- Dále předpokládejme, že teoretické modelové rozdělení četností je reprezentováno četnostmi  $np_i$  (očekávané četnosti)
- Potom shodu mezi empirickým a teoretickým rozdělením posuzujeme pomocí testovacího kritéria

$$(3) \quad \chi^2 = \sum_{i=1}^k \frac{(X_i - np_i)^2}{np_i}$$

- Tzv. Pearsonův Chí-kvadrát test
- Pro to aby se tato veličina řídila asymptoticky chí-kvadrát rozdělením o  $k-1$  stupni volnosti je žádoucí aby  $n > 50$ .

# Test dobré shody při známých parametrech 2

- Dále je nutné, aby pro teoretické četnosti platilo  $np_i \geq 5$
- Nevyhovují-li některé četnosti této podmínce, lze dosáhnout jejího splnění sloučením několika sousedních tříd. Tím se sníží počet stupňů volnosti, neboť  $k$  je rovno počtu tříd po sloučení.
- Pokud pro hodnotu testovací statistiky platí  $\chi^2 < \chi^2_{\alpha}(k-1)$   
pak testovanou hypotézu nezamítáme na hladině významnosti alpha.

# Chí-kvadrát test dobré shody – příklad se zmrzlinou

- *Řetězec cukráren, který nabízí 4 druhy zmrzliny otevřel provozovnu v nové lokalitě. Ve stávajících provozovnách řetězce byla dosud struktura prodeje podle druhů zmrzliny následující: vanilková 62%, čokoládová 18%, jahodová 12%, pistáciová 8%. Po otevření provozovny v nové lokalitě máme záznam o následujícím prodeji: vanilková 120, čokoládová 40 jahodová 18, pistáciová 22.*
- *Vyjádřete se pomocí statistického testu ke shodě či odlišnosti struktury prodeje v nové lokalitě oproti dosavadním prodejům řetězce.*



# Příklad se zmrzlinou - řešení

- Pro získání očekávaných četností u prodeje zmrzliny (při platnosti stávající struktury prodeje pro novou lokalitu) aplikujeme dosavadní strukturu prodeje na celkové prodané množství v nové lokalitě (kde je prodáno celkem 200 kusů zmrzliny):
- Např. u vanilkové: očekávaná četnost při 200 prodaných kusech = 62% \* 200 = 134 kusů
- Tyto očekávané četnosti konfrontujeme se skutečně pozorovanými (chí-kvadrát test dobré shody), výpočet testového kritéria:

$$\chi^2 = \sum_{i=1}^k \frac{(X_i - np_i)^2}{np_i} = \dots = 4,32$$

# Příklad se zmrzlinou – řešení

Výpočet je naznačen v následující tabulce:

	vanilková	čokoládová	jahodová	pistáciová	$\Sigma$
struktura prodeje	62%	18%	12%	8%	100%
nová provozovna	120	40	18	22	200
oč. při stejné struktuře	124	36	24	16	200
chi-square:	0,13	0,44	1,50	2,25	<b>4,32</b>

# Příklad se zmrzlinou - řešení

- Spočtenou hodnotu testového kritéria porovnáme s příslušným kvantilem rozdělení  $\chi^2$  s (k-1), tedy se 3 stupni volnosti. Pro 5% hladinu významnosti půjde o kvantil  $\chi^2_{(1-\alpha)}$ , tedy o kvantil  $\chi^2_{0,95} = 7,82$
- Spočtená hodnota testového kritéria (4,32) nepřekračuje mez vymežující kritický obor (7,82), nachází se v oboru přijetí a na zvolené 5%ní hladině významnosti hypotézu o shodě struktury prodeje nezamítáme.

# Testy dobré shody při neznámých parametrech

- V praxi se často stává, že pravděpodobnosti  $p_1, \dots, p_k$  uvažovaného multinomického rozdělení závisejí na nějakých neznámých parametrech  $a_1, \dots, a_m$ . Vzniká pak problém jak pozměnit testovací kritérium, aby se hodilo i na tento příklad.
- Nabízí se možnost tyto parametry odhadnout, odtud získat i odhady pro pravděpodobnosti  $p_1, \dots, p_k$  a do vzorce (3) pak dosadit tyto odhady.
- Lze očekávat, že pak rozdělení chí-kvadrát bude mít o tolik stupňů volnosti méně, kolik parametrů jsme museli odhadovat.

# Testy dobré shody při neznámých parametrech

- Označme  $\mathbf{a} = (a_1, \dots, a_m)'$ . Předpokládejme, že  $p_1 = p_1(\mathbf{a}), \dots, p_k = p_k(\mathbf{a})$  jsou dostatečně hladké funkce proměnné  $\mathbf{a}$ . Protože platí pro každé  $\mathbf{a}$

$$(4) \quad p_1(\mathbf{a}) + \dots + p_k(\mathbf{a}) = 1,$$

derivováním dostaneme

$$(5) \quad \frac{\partial p_1(\mathbf{a})}{\partial a_j} + \dots + \frac{\partial p_k(\mathbf{a})}{\partial a_j} = 0 \quad j = 1, 2, \dots, m.$$

- Nyní vzniká problém, jakým způsobem pořídit odhady parametru  $\mathbf{a}$ . Jedna možnost spočívá v tom, že se za odhad vezme ta hodnota  $\mathbf{a}$ , která při daných veličinách  $X_1, \dots, X_k$  minimalizuje  $\chi^2$  ve vzorci (3)
- Jedná se o jakousi analogii metody nejmenších čtverců. Říkáme, že jde o odhad parametru  $\mathbf{a}$  pořízený metodou minimálního  $\chi^2$

# Testy dobré shody při neznámých parametrech

- Po derivaci vzorce (3) dostaneme tuto soustavu rovnic

$$(6) \quad -\frac{1}{2} \frac{\partial \chi^2}{\partial a_j} = \sum_{i=1}^k \left\{ \frac{X_i - np_i(\mathbf{a})}{p_i(\mathbf{a})} + \frac{[X_i - np_i(\mathbf{a})]^2}{2np_i^2(\mathbf{a})} \right\} \frac{\partial p_i(\mathbf{a})}{\partial a_j} = 0, \quad j=1, 2, \dots, m.$$

- Ukazuje se, že vliv druhého členu v (6) při velkém  $n$  není příliš podstatný, takže řešení soustavy rovnic

$$(7) \quad \sum_{i=1}^k \frac{X_i - np_i(\mathbf{a})}{p_i(\mathbf{a})} \frac{\partial p_i(\mathbf{a})}{\partial a_j} = 0, \quad j=1, 2, \dots, m,$$

se příliš neliší od řešení soustavy (6). Vzhledem k (5) lze soustavu (7) upravit na tvar

$$(8) \quad \sum_{i=1}^k \frac{X_i}{p_i(\mathbf{a})} \frac{\partial p_i(\mathbf{a})}{\partial a_j} = 0, \quad j=1, 2, \dots, m.$$

# Testy dobré shody při neznámých parametrech

- Řešením soustavy (8) je tzv. odhad parametru  $\mathbf{a}$  modifikovanou metodou minimálního  $\chi^2$
- **Věta:** *Budiž dáno  $k$  funkcí  $p_1(\mathbf{a}), \dots, p_k(\mathbf{a})$ , kde  $\mathbf{a} = (a_1, \dots, a_m)'$ . Předpokládejme, že  $m < k - 1$ . Necht' pro všechny body  $\mathbf{a}$  nedegenerovaného konečného uzavřeného intervalu  $A$  z  $R_m$  platí:*

1.  $p_1(\mathbf{a}) + \dots + p_k(\mathbf{a}) = 1$

2. *Existuje takové  $c > 0$ , že  $p_i(\mathbf{a}) > c$  pro  $i = 1, 2, \dots, k$ .*

3. *Každá funkce  $p_i(\mathbf{a})$  má spojité derivace*

$$\frac{\partial p_i(\mathbf{a})}{\partial a_j} \quad \text{a} \quad \frac{\partial^2 p_i(\mathbf{a})}{\partial a_j \partial a_s} \quad (j, s = 1, 2, \dots, m)$$

4. *Matice  $(\frac{\partial p_i(\mathbf{a})}{\partial a_j})_{i=1, j=1}^{k, m}$  má hodnotu  $m$ .*

*Necht'  $\mathbf{a}^0$  je vnitřním bodem  $A$ . Označme  $p_i^0 = p_i(\mathbf{a}^0)$ .*

# Testy dobré shody při neznámých parametrech

*Nechť  $\mathbf{X} = (X_1, \dots, X_k)'$  má multinomické rozdělení s parametry  $n, p_1^0, \dots, p_k^0$ .*

*Pak soustava rovnic (8) má právě jeden kořen  $\mathbf{a}$  takový, že  $\mathbf{a}$  konverguje k  $\mathbf{a}^0$  podle pravděpodobnosti při  $n \rightarrow \infty$ . Dosadíme-li tento kořen do výrazu*

$$\chi^2 = \sum_{i=1}^k \frac{[X_i - np_i(\mathbf{a})]^2}{np_i(\mathbf{a})},$$

*má veličina  $\chi^2$  při  $n \rightarrow \infty$  asymptoticky rozdělení chí-kvadrát s  $k - m - 1$  stupni volnosti.*

Důkaz: například v Andělovi (Matematická statistika)

- Tato věta se používá hlavně k ověřování typu rozdělení a při hodnocení kontingenčních tabulek.



# Ověřování exponenciálního rozdělení

- Chceme ověřit, zda daný výběr pochází z exponenciálního rozdělení s hustotou

$$(9) \quad f(x) = \frac{1}{\lambda} e^{-x/\lambda} \quad \text{pro } x > 0,$$

kde  $\lambda > 0$  je neznámý parametr. Položme  $b_i = ih$ ,  $i = 0, 1, \dots, k-1$  kde  $h > 0$  je vhodně zvolená délka třídy. Necht'  $b_k = \infty$ . Máme třídy

$$(10) \quad \langle 0, b_1 \rangle, \langle b_1, b_2 \rangle, \dots, \langle b_{k-2}, b_{k-1} \rangle, \langle b_{k-1}, \infty \rangle$$

jejichž četnosti necht' jsou  $X_1, \dots, X_k$ ; opět položíme  $X_1 + \dots + X_k = n$ .

- Pravděpodobnost, že jednotka padne do  $i$ -té třídy, je rovna

$$(11) \quad p_i = e^{-b_{i-1}/\lambda} - e^{-b_i/\lambda}, \quad i = 1, 2, \dots, k.$$

# Ověřování exponenciálního rozdělení

- Položíme-li  $\infty e^{-\infty} = 0$ , máme

$$\frac{d p_i}{d \lambda} = \frac{1}{\lambda^2} (b_{i-1} e^{-b_{i-1}/\lambda} - b_i e^{-b_i/\lambda}), \quad i=1, 2, \dots, k.$$

- Rovnice (8) má pak tvar

$$\frac{1}{\lambda^2} \sum_{i=1}^k X_i \frac{b_{i-1} e^{-b_{i-1}/\lambda} - b_i e^{-b_i/\lambda}}{e^{-b_{i-1}/\lambda} - e^{-b_i/\lambda}} = 0.$$

Odtud

$$\sum_{i=1}^k X_i \frac{b_{i-1} - b_i e^{(-b_{i-1} - b_i)/\lambda}}{1 - e^{(b_{i-1} - b_i)/\lambda}} = 0.$$

# Ověřování exponenciálního rozdělení

- Dosazením za  $b_i$  dostaneme

$$e^{-h/\lambda} = \frac{hX_2 + 2hX_3 + \dots + (k-1)hX_k}{hX_1 + 2hX_2 + \dots + khX_k - hXk}.$$

- Označme

$$(12) \quad \bar{X} = (hX_1 + 2hX_2 + \dots + khX_k) / n.$$

- Pak máme

$$e^{-h/\lambda} = \frac{n\bar{X} - nh}{n\bar{X} - hX_k},$$

- Takže

$$(13) \quad \lambda = -h / \ln \frac{n\bar{X} - nh}{n\bar{X} - hX_k}.$$

# Ověřování exponenciálního rozdělení

- Hodnotu  $\lambda$  z (13) dosadíme do (11) a vypočteme

$$(14) \quad \chi^2 = \sum_{i=1}^k \frac{(X_i - np_i)^2}{np_i}$$

- V případě  $\chi^2 \geq \chi_{k-2}^2(\alpha)$  zamítneme hypotézu o exponenciálním rozdělení na hladině, která je asymptoticky rovna  $\alpha$

# Literatura

- Anděl J., Matematická statistika, SNTL, Praha, 1978
- Zvára K., Štěpán J., Pravděpodobnost a matematická statistika, MATFYZPRESS, Praha, 2001