

OVĚŘENÍ NORMÁLNÍHO ROZDĚLENÍ - MĚNĚNÍ BÍŽNÉ TESTY

Jana Kubanová

Ústav matematiky, FES, Univerzita Pardubice

Abstract: *The normal distribution of population is presumption of many statistical tests. The acceptance of this presumption is most frequently tested by chi-square test and by one-sample Kolmogorov–Smirnov’s test. Some other tests are described in this article, serving for confirmation of this hypothesis, such as the tests based on skewness and curtosis, the Jarque – Berra’s test, the Shapiro-Wilk’s normality test and the D’Agostin test.*

Při statistických analýzách se velmi často používají testy, které vycházejí z předpokladu normálního rozložení základního souboru. V některých případech můžeme být o splnění tohoto předpokladu přesvědčeni na základě dřívějších zkušeností nebo z výzkumů realizovaných na stejných nebo podobných souborech. Pokud si nejsme jisti splněním předpokladu normality, je třeba rozhodnout, kterým ze statistických testů. Nejčastěji používané jsou χ^2 test a jednovýběrový Kolmogorovův–Smirnovův test (Anděl, 1993).

Nyní vzpomeneme stručně ještě několik dalších testů, na jejichž základě lze ověřit normalitu rozložení základního souboru.

a) Testy založené na šikmosti a špičatosti

Pomocí těchto testů testujeme nulovou hypotézu H_0 , že náhodný výběr pochází ze základního souboru X s normálním rozložením pravděpodobností $N(\mu, \sigma)$.

O normálním rozdělení je známo, že jeho koeficienty šikmosti a špičatosti (asymetrie a excesu) jsou rovny nule. Toho se využije v následujícím testu.

Je známo, že výběrové koeficienty šikmosti a špičatosti (asymetrie a excesu) jsou definovány následujícím způsobem:

$$\text{Výběrový koeficient šikmosti } S_k = \frac{M_3}{s^3} = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{x})^3}{\left(\sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{x})^2} \right)^3}$$

$$\text{Výběrový koeficient špičatosti } E_k = \frac{M_4}{s^4} - 3 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{x})^4}{\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{x})^2 \right)^2} - 3$$

Koeficient šikmosti i špičatosti normálního rozdělení je roven nule, proto by měly i tyto odhady být blízké nule.

Lze dokázat, že při náhodném výběru dostatečného rozsahu z normálního rozdělení pravděpodobností $N(\mu, \sigma)$ má náhodná veličina S_k přibližně $N(ES_k, \sqrt{DS_k})$ rozdělení a náhodná veličina E_k přibližně $N(EE_k, \sqrt{DE_k})$ rozdělení pravděpodobností, kde pro střední hodnotu a disperzi uvedených koeficientů platí (Anděl, 1978):

$$ES_k = 0 \qquad DS_k = \frac{6(n-2)}{(n+1)(n+3)}$$

$$EE_k = \frac{-6}{n+1} \qquad DE_k = \frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}$$

Při testu založeném na šikmosti zamítneme hypotézu o normálním rozložení pravděpodobností základního souboru tehdy, když

$$\frac{|S_k|}{\sqrt{DS_k}} \geq z_\alpha,$$

kde $z_\alpha = \Phi^{-1}\left(\frac{2-\alpha}{2}\right)$. Tuto hodnotu hledáme v tabulkách $N(0, 1)$ rozdělení pravděpodobností.

Při testu založeném na špičatosti hypotézu o normalitě rozložení zamítneme, je-li

$$\frac{|E_k - EE_k|}{\sqrt{DE_k}} \geq z_\alpha$$

V nichž případech jsou testy založené na šikmosti a špičatosti citlivější na porušení normality než uvedený χ^2 test.

b) Test kombinace výběrové šikmosti a špičatosti (Jarque – Berra)

Testujeme opět nulovou hypotézu H_0 , že náhodný výběr pochází ze základního souboru X s normálním rozložením pravděpodobností $N(\mu, \sigma)$.

K testování nulové hypotézy používáme testovací kritérium tvaru

$$\chi = \frac{S_k^2}{DS_k} + \frac{(E_k - EE_k)^2}{DE_k},$$

kde S_k je výběrový koeficient šikmosti a E_k výběrový koeficient špičatosti. Za předpokladu platnosti H_0 má náhodná veličina χ asymptoticky χ^2 rozložení pravděpodobností se 2 stupni volnosti. Hypotézu H_0 zamítáme v případě, kdy $\chi > \chi^2_{\alpha,2}$.

c) Shapiro-Wilkův test normality

Dalším z testů, který může sloužit k ověření normality je test odvozený Shapiro a Wilkem. Test je vhodný pro výběry rozsahu $7 \leq n \leq 30$. V testu se používá realizace $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ uspořádaného náhodného výběru $X_{(1)}, \dots, X_{(n)}$.

Testujeme nulovou hypotézu

H_0 : X má $N(\mu, \sigma)$ rozdělení pravděpodobností proti alternativní hypotéze

H_1 : X nemá $N(\mu, \sigma)$ rozdělení pravděpodobností.

Testovací kritérium má tvar:
$$SW = \frac{\left(\sum_{i=1}^m a_i(n)(X_{(n-i+1)} - X_{(i)}) \right)^2}{\sum_{i=1}^n (X_i - \bar{x})^2},$$

kde $a_i(n)$ jsou tabelované konstanty, $m = \frac{n}{2}$, je-li n sudé číslo a $m = \frac{n-1}{2}$, je-li n liché číslo, $X_{(n-i+1)}$ a $X_{(i)}$ jsou pořádkové statistiky vytvořené z náhodného výběru X_1, \dots, X_n jeho uspořádáním do neklesající posloupnosti.

Konstanty $a_i(n)$ souvisí s očekávanými pořádkovými statistikami normálního rozložení, tzn., že má-li náhodná veličina X rozdělení $N(\mu, \sigma)$, budou body $(a_i(n), X_i)$ seskupeny až na náhodné odchylky kolem regresní přímky se směrnicí σ .

Princip testu je v tom, že se odhadne parametr σ náhodnou veličinou $s^* = \sum_{i=1}^n a_i X_i$ a jeho odhad se porovná s odhadem založeným na náhodné veličině $\sum_{i=1}^n (X_i - \bar{x})^2$.

Hypotézu o normalitě rozložení zamítáme, je-li $SW < SW^*(n, \alpha)$. Hodnoty $SW^*(n, \alpha)$ hledáme ve statistických tabulkách (Palenčár a kol., 2001).

V následující tabulce 1 jsou uvedeny koeficienty $a_i(n)$ Shapiro-Wilkova testu a 95%-ní kvantily pro četnosti $n = 5$ až $n = 10$ náhodného výběru.

Vzhledem k symetrii kvantilů normálního rozdělení platí: $a_i = -a_{n-i+1}$.

Tab.1: Koeficienty $a_i(n)$ Shapiro-Wilkova testu a 95%-ní kvantily pro $n = 5$ až $n = 10$ náhodného výběru.

| i | $n = 5$ | $n = 6$ | $n = 7$ | $n = 8$ | $n = 9$ | $n = 10$ |
|---------------|---------|---------|---------|---------|---------|----------|
| 1 | 0,6646 | 0,6431 | 0,6233 | 0,6052 | 0,5888 | 0,5739 |
| 2 | 0,2413 | 0,2816 | 0,3031 | 0,3164 | 0,3244 | 0,3291 |
| 3 | 0 | 0,0875 | 0,1401 | 0,1743 | 0,1976 | 0,2141 |
| 4 | - | - | 0 | 0,0561 | 0,0947 | 0,1224 |
| $SW^*_{0,95}$ | 0,762 | 0,788 | 0,803 | 0,818 | 0,829 | 0,842 |

Pro větší četnosti náhodného výběru se hodnoty koeficientů a_i i hodnoty testovacího kritéria $W^*_{1-\alpha}$ počítají podle Roystonova postupu (Royston, 1982).

Konstanty a_i jsou tabelované např. v (Owen, 1962), (Sarhan, Greenberg, 1962) ale i v ČSN 01 0225.

Praktické pokusy ukázaly, že tento test je vzhledem k odchylkám od normality citlivější, než ostatní známé testy dobré shody.

d) D'Agostinův test

Tento test se používá k testování hypotézy o normalitě pro výběry o rozsahu $30 \leq n \leq 100$, tzn. že opět testujeme nulovou hypotézu

$H_0: X$ má $N(\mu, \sigma)$ rozdělení pravděpodobností proti alternativní hypotéze

$H_1: X$ nemá $N(\mu, \sigma)$ rozdělení pravděpodobností.

Při testování postupujeme tak, že hodnoty realizace náhodného výběru uspořádáme do neklesající posloupnosti $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ a definujeme náhodnou veličinu tvaru

$$D = \frac{\sqrt{n}(A - 0,2820948)}{0,0299860}$$

$$\text{kde } A = \frac{\sum_{i=1}^n \left(\frac{n+1}{2} - i \right) (x_{(n-i+1)} - x_{(i)})}{\sqrt{n^3 \sum_{i=1}^n (x_{(i)} - \bar{x})^2}}$$

a $x_{(n-i+1)}, x_{(i)}$ jsou hodnoty pořádkových statistik z náhodného výběru X_1, \dots, X_n .

Testovanou hypotézu zamítneme, platí-li pro hodnotu náhodné veličiny D :

$$D < D\left(n, \frac{\alpha}{2}\right) \text{ nebo } D > D\left(n, 1 - \frac{\alpha}{2}\right).$$

Pokud se snažíme hodnotit simultánně výsledky různých testů normality, vyvstávají problémy s určením chyby prvního druhu. Proto je vhodné se na tyto výsledky testů dívat s určitým nadhledem a považovat je pouze za orientační kritéria při vytváření konečného názoru.

V některých případech může k posouzení normality dobře posloužit grafické znázornění hodnot statistického znaku. Uvádíme pro orientaci některé možnosti:

1. *Histogram rozdělení četností statistického znaku.*

Histogram nám umožní posoudit, zda hodnoty statistického znaku připomínají průběh Gaussovy křivky.

2. *Graf empirické distribuční funkce.*

Sledujeme, zda je graf empirické distribuční funkce podobný esovité křivce distribuční funkce normálního rozložení. Před sestavením grafu je vhodné data normovat.

3. *Q-Q diagram.*

Na svislé ose y nanášíme výběrové kvantily, na vodorovné ose x jim odpovídající kvantily normálního rozložení pravděpodobností. Pokud výběr odpovídá normálnímu rozložení pravděpodobností, pak body, odpovídající jednotlivým pozorováním, budou ležet kolem přímky $y = x$.

Literatura:

1. Anděl, J.: *Statistické metody*. Matfyzpress, Praha 1993
2. Kubanová, J.: *Statistické metody pro ekonomickou a technickou praxi*, 2003 (v tisku)
3. Owen, D.B.: *Handbook of Statistical Tables*. Massachusetts, Adison-Wesley, 1962
4. Palenčár, R., Ruiz, J.M., Janiga, I., Horníková, A.: *Štatistické metódy v metrologických a skúšobných laboratóriach*. Bratislava 2001, ISBN 80-968449-3-8
5. Royston, J.P.: *Appl. Statist.* 31, 115, 1982
6. Sarhan, A.E., Greenberg, B.G.: *Contribution to Order Statistics*. New York, J. Wiley 1962

Kontaktní adresa:

Paed. Dr. Jana Kubanová, CSc., Ústav matematiky, FES, Univerzita Pardubice
Studentská 84, 532 10 Pardubice
tel.: 466 036 046
e-mail: Jana.Kubanova@upce.cz

Recenzoval: doc. RNDr. Bohdan Linda, CSc., vedoucí Ústavu matematiky, FES,
Univerzita Pardubice