

Využití systému g - h rozdělení při analýze dat

Dipl. Ing. Milirka, VÚTZ, Čvár Králové n.L.

1. Úvod

Při průzkumové analýze dat je účelem kvantifikovat jejich statistické zvláštnosti, které se týkají zejména tvaru výběrového rozdělení. Běžně se pro tyto účely používá kvantilových charakteristik šíkmosti a špičatosti [1]. Další možností je použití vhodného systému empirických frekvenčních funkcí, pomocí kterého se approximuje výběrové rozdělení. Oba tyto přístupy je možné spojit s využitím takových systémů frekvenčních funkcí, které vycházejí z kvantilové funkce (jako je např. Tukeyův systém lambda rozdělení). Vzhledem k tomu, že při analýze dat je referenční rozdělení normální (a odchylky od normality jsou obvyklejší píspisovány transformací původně normálně rozdělených dat), je vhodné, aby empirický systém rozdělení byl generován jako transformace normálně rozdělených náhodných veličin (jako je Johnsonův systém).

Syntézou těchto požadavků je tzv. systém g-h rozdělení, jehož parametry jsou přímo kvantilové míry šíkmosti a špičatosti. Lze ho tedy konstruovat velmi jednoduše přímo z experimentálních dat. Navíc je "automaticky" zajištěna robustnost resp. lokální charakterizovatelnost zvláštností rozdělení dat.

Systém g-h rozdělení lze kromě průzkumové analýzy dat uplatnit také pro generaci náhodných veličin z rozdělení s definovanými šíkmostmi a špičatostmi (které mohou být závislé na vzdálenosti od mediánu). Tyto náhodné veličiny lze s výhodou použít při různých simulacích studiích resp. ověřování robustních metod.

V tomto příspěvku jsou přehledně uvedeny základní vlastnosti g-h rozdělení, které se dostalo i do statistické encyklopédie [2].

Je popsána struktura části programu EXGR, která využívá g-h rozdělení pro účely průzkumové analýzy.

2. DEFINICE g-h ROZDĚLENÍ

Systém g-h empirických rozdělení je založen na monotónní transformaci standardizované náhodné veličiny s normálním rozdělením. Jde o jeden z translačních systémů rozdělení, kdy se nejprve provádí standardizace původních náhodných proměnných (kvantilů) až na normalizované kvantily

$$z = (x - \tilde{x}_{0.5}) / R \quad (1)$$

kde $\tilde{x}_{0.5}$ je parametr (medián) a R je parametr měřítka.

Pro třídu g-rozdělení (pouze sešikmených) se volí obecná transformace typu

$$Q_{g,0}(z) = G(z) \cdot z \quad (2)$$

kde $G(z)$ je lichá funkce, pro kterou musí platit, že

$$Q_{g,0}(0) = 0 \quad \lim_{z \rightarrow 0} Q_{g,0}(z) \approx z$$

$G(z)$ je vlastně operátor šíkmosti, který závisí na parametru šíkmosti g. Je vhodné, aby pro $g > 0$ vycházela rozdělení sešikmené vpravo, pro $g < 0$ rozdělení sešikmené vlevo a pro $g = 0$ normální rozdělení. Všem těmto požadavkům vyhovuje jednoduchá funkce

$$G(z) = [\exp(g \cdot z) - 1] / (g \cdot z) \quad (3)$$

Lze jednoduše nalézt, že frekvenční funkce g-rozdělení $f_g(x)$ je dána vztahem

$$f_g(x) = 1 / [2\pi \cdot |(x - \tilde{x}_{0.5}) \cdot g + R|] \cdot \exp(-[(x - \tilde{x}_{0.5}) \cdot g / R + 1]^2 / (2g^2))$$

Jde tedy o lognormální rozdělení (rozdělení typu S₁ v Johnsonově systému frekvenčních funkcí).

Pro třídu h-rozdělení (symetrických s různou špičatostí) se volí transformace typu

$$Q_{0,h}(z) = H(z) \cdot z \quad (4)$$

Operator špičatosti $H(\zeta)$ musí být rostoucí kladná funkce závislá na parametru špičatos-
ti h . Pro $h > 0$ je dle o normální rozdělení a tím je $h > 0$ větší, tím má odpovídající rozděle-
ní delší konc. Tento výhoduje volba

$$H(\zeta) = \exp(h\zeta^2/2) \quad (5)$$

Z rov. (5) plyne, že parametr h může být i záporný. Podmínka monotonnosti funkce $Q_{g,h}$ je
však nerušena, pokud je $\zeta^2 > -1/h$. Třída h -rozdělení definované rov. (4) a (5) se
v oblasti konců chová jako Paretovo rozdělení. Volba faktoru $1/2$ v rov. (5) také zajistíuje,
že pro $h \approx 1$ je $Q_{g,h}$ rozdělení blízké Cauchyho rozdělení [3].

Frekvenční funkci h -rozdělení lze však získat pouze numericky, protože nelze nalézt analy-
tický transformaci inverzní k rov. (4).

Pro obecnou třídu g-h rozdělení se volí dvouparametrická transformace

$$Q_{g,h}(\zeta) = G(\zeta) \cdot H(\zeta) \cdot \zeta \quad (6)$$

kde $G(\zeta)$ je voleno podle rov. (3) a $H(\zeta)$ je definováno rov. (5). Také v tomto případě je
třeba odpovídající frekvenční funkci počítat numericky.

Pro generaci náhodných čísel z g-h rozdělení se zvolenými parametry g a h postačuje genera-
vat kvantily ζ_p standardního normálního rozdělení (zde p určuje 100 P %ní kvantil) a dosa-
zovat do rov. (6) za $\zeta = \zeta_p$.

3. MOMENTY

I když se ani při generaci náhodných průměrných ani průzkumové analýze dat momenty nepouží-
vají, je pro účely porovnání s ostatními systémy empirických rozdělení vhodné znát výrazy pro
první čtyři momenty. Ty zde budou funkciemi parametrů g , h . Pro standardizované náhodné pro-
měnné $Y = Q_{g,h}(\zeta)$ z g-h rozdělení je střední hodnota definována vztahem

$$E(Y) = \frac{1}{g^{1/2}-1} \left[\exp\left(\frac{g^2}{2(1-h)}\right) - 1 \right] \quad 0 \leq h \leq 1 \quad (7)$$

a pro rozptyl platí

$$D(Y) = \frac{1}{g^{2/2-2h}} \left[\exp(2g^2) - 2\exp(g^2) + 1 \right] - \frac{1}{g^{2(1-h)}} \left[\exp\left(\frac{g^2}{2}\right) - 1 \right]^2 \quad 0 \leq h < 0,5 \quad (8)$$

kde $w = g^2/(1-2h)$.

Pro případ g-rozdělení je $h = 0$ a pak z rov. (7)

$$E(Y) = [\exp(g^2/2) - 1] / g \quad (7a)$$

resp. z rov. (8)

$$D(Y) = \exp(g^2) \cdot [\exp(g^2) - 1] / g^2 \quad (8a)$$

Sikost g_1 je u g-rozdělení vyjádřitelná ve tvaru

$$g_1 = [\exp(3g^2) - 3\exp(g^2) + 1] / \sqrt{[\exp(g^2) - 1]^3} \quad (9)$$

a pro špičatost g_2 platí

$$g_2 = [\exp(6g^2) - 4\exp(3g^2) + 6\exp(g^2) - 3] / (\exp(g^2) - 1)^2 \quad (10)$$

Pro případ h-rozdělení je $g = 0$, takže $E(Y) = g_1 = 0$. Pro rozptyl zde platí

$$D(Y) = 1 / \sqrt{(1-2h)^3} \quad 0 \leq h < 0,5 \quad (8b)$$

a špičatost je ve tvaru

$$g_2 = 3(1-2h)^3 / \sqrt{(1-4h)^6} \quad (10a)$$

Při výprážtu vyšších momentů g-h rozdělení je možné použít Martínezova vztahu

$$E(Y^n) = \frac{1}{g^n \sqrt{1-nh}} \sum_{i=0}^n (-1)^i \binom{n}{i} \exp\left\{\frac{[(n-i)g]^2}{2(1-nh)}\right\} \quad (11)$$

který platí pro $g \neq 0$ a $0 \leq h \leq 1/n/3$.

4. ODHADY PARAMETRU g-h ROZDĚLENÍ

S ohledem na návaznost na průzkumovou analýzu dat se pro odhad parametrů g, h resp. obou používá výběrového mediánu $\tilde{x}_{0,5}$ a dvojic kvantilů $\tilde{x}_p, \tilde{x}_{1-p}$ pro vhodné $0 < p < 0,5$. U malých výběrů se pro tyto účely využívá všech pořadkových statistik $x_{(i)}$. U větších výběrů se používá tzv. písmenových hodnot, kde $p = 2^{-i}$, $i = 2,3$ (odpovídají kvartilem, oktildům, sedecilem atd.)

g - Rozdělení

Pro zvolené p lze určit parametr g (a také R) dvou typů

$$\tilde{x}_p = \tilde{x}_{q5} + R \cdot Y = \tilde{x}_{0,5} + R \left[\frac{\exp(g \cdot \tilde{x}_p) - 1}{g} \right] \quad (12a)$$

$$\tilde{x}_{1-p} = \tilde{x}_{q5} + R \left[\frac{\exp(-g \cdot \tilde{x}_p) - 1}{g} \right] \quad (12b)$$

kde \tilde{x}_p jsou standardizované kvantily normálního rozdělení (vzhled k symetrii je $\tilde{x}_p = -\tilde{x}_{1-p}$). Po jednoduchých úpravách vyjde

$$g_p = -\frac{1}{2p} \cdot \ln \left[\frac{\tilde{x}_{1-p} - \tilde{x}_{q5}}{\tilde{x}_{q5} - \tilde{x}_p} \right] \quad (13)$$

Index p zde ukazuje, že parametr šíkosti obecně může záviset na poloze kvantilů vzhledem k mediánu.

Pokud je g_p přibližně konstantní (resp. není funkci p) postupuje se tak, že se určí medián \tilde{g} ze všech g_p . Jednoduše je pak možné odhadnout parametr R jako směrnici v Q-Q grafu, kde se vynášejí hodnoty \tilde{x}_p ve. $Q_{g,0}(\tilde{x}_p)$. Přímka v tomto grafu indikuje, že lze použít g-rozdělení s konstantní hodnotou g. Pro stejný účel je možné použít také graf symetrie, kdy se vynášejí polosumy $0,5(\tilde{x}_p + \tilde{x}_{1-p})$ ve. $\tilde{x}_p/2$. S využitím rov. (12a) a (12b) snadno určíme teoretickou závislosti

$$0,5(\tilde{x}_p + \tilde{x}_{1-p}) = \tilde{x}_{q5} + 0,5 \frac{R}{g} [\exp(g \cdot \tilde{x}_p) + \exp(-g \cdot \tilde{x}_p) - 2] \approx \tilde{x}_{q5} + R \cdot g \cdot \tilde{x}_p^2 / 2 \quad (14)$$

Plati-li tedy předpoklad konstantnosti parametru se šíkmení g, vyjde v grafu symetrie přibližně lineární závislost.

Pokud není g_p konstantní, předpokládá se obyčejně, že je to jednoduchý polynom vzhledem k normalizovaným kvantilům \tilde{x}_p resp. jejich čtverců. Jednoduché je záocení typu /3/

$$g_p = g_0 + g_1 \tilde{x}_p^2 \quad (15)$$

Rov. (15) lze snadno ověřit na základě grafu šíkosti, kdy se vynáší g_p ve. \tilde{x}_p^2 . Pokud vyjde v tomto grafu přibližně lineární závislost, je třeba uvažovat obecnější rozdělení s nekonstantním g_p vyjádřeným rov. (15). Parametry g_0, g_1 v rov. (15) lze snadno určit z úseku a směrnice v grafu šíkosti (doporučuje se použití robustní regrese).

h - Rozdělení

Při znaloosti parametru polohy $\tilde{x}_{0,5}$ a parametru měřítka R lze určit pro každé \tilde{x}_p odpovídající hodnotu h_p přímo z definičního vztahu (4). Vyjde

$$h_p = [2 \ln \left(\frac{(\tilde{x}_p - \tilde{x}_{q5})/R}{\tilde{x}_p} \right)] / \tilde{x}_p^2 \quad (16)$$

v případě, že je vyběrové rozdělení symetrické, musí pochopitelně být $h_p = h_{1-p}$. Jako vhodný parametr měřítka se doporučuje interkvartilový odhad směrodatné odchylky mediánu /3/

$$R = 0.926 \cdot (\tilde{x}_{0.75} - \tilde{x}_{0.25}) / \sqrt{n} \quad (17)$$

kde n je rozsah výběru.

Za předpokladu symetrického rozdělení lze přímo z definičních vztahů

$$\hat{x}_p = \tilde{x}_{0.5} + R \cdot \hat{\sigma}_p \cdot \exp(h \cdot \hat{\sigma}_p^2/2) \quad (18a)$$

$$\hat{h}_{1-p} = \tilde{x}_{0.5} - R \cdot \hat{\sigma}_p \cdot \exp(h \cdot \hat{\sigma}_p^2/2) \quad (18b)$$

(zde $\hat{\sigma}_p < 0$ a $\hat{h}_{1-p} = -\hat{h}_p$) dospět k lineární funkci vzhledem k h

$$\ln \left[\frac{\hat{x}_{1-p} - \hat{x}_p}{-2\hat{\sigma}_p} \right] = \ln R + h \cdot \hat{\sigma}_p^2/2 \quad (19)$$

Z rov. (19) je patrné, že vynesením $\ln[(\hat{x}_{1-p} - \hat{x}_p)/(-2\hat{\sigma}_p)]$ ve. $\hat{\sigma}_p^2/2$ vyjde v případě platnosti h -rozdělení s konstantním parametrem špičatosti h přibližně lineární závislost. Tato závislost se označuje jako graf pseudo-sigma a umožňuje odhad h i R ze směrnice resp. úseku regresní přímky.

Pro normální rozdělení je graf pseudosigma prakticky horizontální přímka s nulovou směrnicí a úsekem $\ln R$. Pro nesymetrická rozdělení vycházejí grafy pseudosigma nelineární.

Pokud není h_p konstantní (ale $h_p \neq h_{1-p}$) předpokládá se, že jde opět o jednoduchou závislost typu rov. (15). Tedy

$$h_{1-p} = h_0 + h_1 \hat{\sigma}_p^2 \quad (20)$$

Rov. (20) lze jednoduše ověřit na základě grafu špičatosti, kdy se vynese h_{1-p} v závislosti na $\hat{\sigma}_p^2$. Pokud vyjde v tomto grafu přibližně lineární závislost, znamená to, že platí obecnější rozdělení s nekonstantním parametrem špičatosti vyjádřeným rov. (20).

g-h rozdělení

To je nejčastější případ, kdy je rozdělení dat nesymetrické a ani po symetrizační transformaci neodpovídá délce koncov normálnímu rozdělení. Vzhledem k volbě $Qgh(\xi)$ ve tvaru rov. (6) a funkcím $G(\xi), H(\xi)$, lze snadno určit, že rov. (13) platí nezávisle na velikosti h . To znamená, že lze nejprve nalézt odhad parametru špičatosti g (stejně jako u "čistého" g-rozdělení) a pak provést symetizační transformaci dat před odhadem h .

Při znalosti g můžeme snadno určit s využitím z rov. (6) polorozdílu

$$\tilde{x}_{1-p} - \tilde{x}_{0.5} = \frac{R}{g} (\exp(-g\hat{\sigma}_p) - 1) \cdot \exp(h \cdot \hat{\sigma}_p^2/2) \quad (21)$$

Po úpravě vyjde lineární závislost vzhledem k h

$$y^* = \ln \left[\frac{g(\tilde{x}_{1-p} - \tilde{x}_{0.5})}{\exp(-g\hat{\sigma}_p) - 1} \right] = \ln R + h \cdot \hat{\sigma}_p^2/2 \quad (22)$$

Vynesením y^* ve. $\hat{\sigma}_p^2/2$ (modifikovaný graf pseudosigma) vyjde v případě platnosti g-h rozdělení s konstantními g , h přibližně lineární závislost.

Jednoduše lze také postupovat v případě, že $\hat{\sigma}_p$ je vyjádřeno rov. (15) a h_p je vyjádřeno rov. (20). Protože je odhad $\hat{\sigma}_p$ nezávislý na h_p , lze stejně jako u "čistého" g-rozdělení nalézt parametry g_0, g_1 . Pro odhad parametrů h_0, h_1 však již nelze použít rov. (20), ale je třeba provést symetizační transformaci (stejně jako u rov. (21)).

Po úpravách vyjde vztah

$$y^* = \ln R + \frac{h_0}{2} \hat{\sigma}_p^2 + \frac{h_1}{2} \hat{\sigma}_p^4 \quad (23)$$

$$g_p^* = \ln \left[\frac{(\tilde{x}_{p/2} - \tilde{x}_{0,5}) \tilde{g}_p}{\exp(-\tilde{g}_p \tilde{x}_p) - 1} \right]$$

se určuje při znalosti \tilde{g}_p , \tilde{x}_p , pro každé p. Závislost \tilde{g}_p^* ve. $\tilde{x}_p/2$ se označuje jako zobecněný graf šířitosti. Výjde-li parabolický, znamená to, že je nutné použít zobecněné g-h rozdělení s nekonstantními parametry g, h.

Pro ověřování platnosti různých typů g-h rozdělení lze pochopitelně použít formální aparát lineární regrese a testovat významnost směrnic, resp. úseků ve výše uvedených grafech. Pro účely předběžné analýzy dat však běžně postačuje posouzení vlastních grafů.

5. PROGRAM EXGR

Program EXGR byl vytvořen v jazyce MPL pro stolní počítač HP 9825 vybavený plotterem. Skládá se z řady různých grafů pro posouzení statistických zvláštností dat. Kromě řady dalších se pro posouzení šířitosti používá grafy

- symetrie, t.j. závislosti $0,5(\tilde{x}_p + \tilde{x}_{1-p})$ ve. $\tilde{x}_p/2$
- šířitosti, t.j. závislosti \tilde{g}_p ve. $\tilde{x}_p/2$

Pro posouzení šířitosti se používá grafy:

- pseudosigma, t.j. závislosti $\ln[(\tilde{x}_{p/2} - \tilde{x}_p)/(-2\tilde{g}_p)]$ ve. $\tilde{x}_p/2$
- šířitosti, t.j. závislosti \tilde{h}_p ve. $\tilde{x}_p/2$.

V případě, že výjde výrazná šířitost vyjádřená parametrem g (resp. \tilde{g}_p), provádí se před posouzením šířitosti symetrizační transformace kvantilů

$$\tilde{x}_p^* = \tilde{g}_p \tilde{x}_p / [\exp(\tilde{g}_p \tilde{x}_p) - 1] \quad (24)$$

Vzhledem k účelu použití se pouze kreslí jednotlivé grafy spolu s odpovídajícím robustním (mediánovým) odhadem regresní přímky.

6. ZÁVĚR

Z výše uvedeného je patrné, že manipulace s g-h rozdělením je velmi snadná. Navíc je možné nejdříve určit šířitost (část g-rozdělení) a pak zpracovávat šířitost (část h-rozdělení) po symetrizační transformaci. Pro účely generace náhodných čísel je výhodné, že je g-h rozdělení definované přes kvantilové funkce. Toho se také s výhodou používá při odhadech parametrů. Pro účely předběžné analýzy dat je zase výhodné, že lze jednoduše popsat i případy nekonstantní šířitosti resp. šířitosti v závislosti na vzdálenosti od mediánu.

Také grafy symetrie a pseudosigma lze interpretovat s ohledem na jejich význam vzhledem ke g-h rozdělení.

Při použití g-h rozdělení jako systému empirických rozdělení je však nevýhodou, že nelze obecně analyticky nalézt odpovídající frekvenční funkce. Pokud je to účelem, bude zřejmě výhodnější použití jiného translačního systému empirických rozdělení, jako je Johnsonův, Shapiro-Wilkův atp.

7. LITERATURA

- /1/ Milítký J.: Zpracování experimentálních dat I, předběžná analýza, Skripta DT Ostrava 1986
- /2/ Encyclopedia of Statistical Sciences, vol 4, str. 298-301, J.Wiley New York 1982
- /3/ Hoaglin D.C., Mosteller F., Tukey J.W. Eds. : Exploring data tables trends and shapes, J.Wiley New York 1985, kap. 11