

GEOMETRICKÉ MOMENTY

Zdeněk FABIÁN

ÚIVT AV ČR

Abstract: Geometric moments are defined as expectations of the k -th power of a newly introduced "geometric function" of distribution. They are alternative numerical characteristics of distributions which generally exist and are expressed by means of simple formulas. Parametric estimates based on empirical geometric moments are discussed.

Резюме: Введена "геометрическая" функция распределения. Математические ожидания ее k -той степени мы называем геометрическими моментами. Показано, что они существуют и даны простыми формулами. Введены и обсуждены оценки основанные на геометрических моментах.

I. ÚVOD.

Bud' $T \subset R$ otevřený interval v reálné přímce a \mathcal{B}_T jeho borelovská σ -algebra. Bud' Π_T množina všech absolutně spojitých rozdělení P_T na \mathcal{B}_T s hustotami p , které pro jednoduchost předpokládáme na T spojitě diferencovatelné. Budu nepořádně říkat "rozdělení $p \in \Pi_T$ " nebo i "rozdělení p " místo "rozdělení $P_T \in \Pi_T$ s hustotou p ". Bud' $\Theta \subset R^m$ konvexní množina. Budeme uvažovat obvyklý parametrický model náhodného experimentu

$$\mathcal{P}_T^\theta = \{T, \mathcal{B}_T, p(u|\theta) : \theta \in \Theta\},$$

regulární v Cramer-Raově smyslu. Funkce θ pro pevné $u \in T$,

$$s_j(u|\theta) = \frac{\partial \log p(u|\theta)}{\partial \theta_j}, \quad j = 1, \dots, m, \quad (1)$$

nazýváme *věrohodnostními skóry* parametru θ_j .

V modelu $\mathcal{P}_R^\mu = \{R, \mathcal{B}_R, p(x - \mu) : \mu \in R\}$ je věrohodnostní skór parametru polohy μ v bodě $\mu = 0$ totožný pro pevné x s hodnotou skórové funkce

$$s(x) = -\frac{p'(x)}{p(x)}. \quad (2)$$

Podobně Fisherova informace o parametru μ , $E_p s_1^2$, je v bodě $\mu = 0$ totožná s Fisherovou informací *rozdělení*, definovanou např. v (Cover, Thomas 1991) jakožto

$$I_p = E_p s^2 = \int_R s^2(x)p(x)dx. \quad (3)$$

(2) a (3) charakterizují rozdělení v počátku parametrického prostoru a je možno je považovat za jisté prototypy věrohodnostních skóru parametrů (souřadnicových funkcí statistické variety \mathcal{P}_T^θ) a Fisherovy informační matice $I_p(\theta) = (E_p s_j s_k)_1^m$ (metrického tenzoru variety). Protože platí $\int_R s'(x)p(x)dx = I_p$, je funkce $s(x)/I_p$ totožná s vlivovou (influenční) funkcí maximálně věrohodného estimátoru parametru μ v bodě $\mu = 0$. Funkce $s^2(x)$ je nezáporná, je rovna nule v bodě $x = 0$ s minimální informací (kde $p'(0) = 0$) a její střední hodnota má význam informace. Lze ji tedy patrně interpretovat jako *relativní informaci* přiřazenou bodu $x \in R$ rozdělením p .

Skórová funkce je vlastně jakýmsi zárodkem nejdůležitějších funkcí a veličin studovaných v teorii odhadu. Zdálo by se tedy celkem přirozené přiřadit výběrovému prostoru $T = R$ geometrii ve smyslu pseudometriky

$$d(x_1, x_2) = |s(x_2) - s(x_1)|. \quad (4)$$

Buď $X = (x_1, \dots, x_n)$ náhodný výběr z rozdělení p . Pišme (zde i v dalším textu) \sum místo $\sum_{i=1}^n$. Pseudometrika (4) "doporučuje" hledat "výběrové těžiště" souboru X jakožto $x^* : n^{-1} \sum s(x_i - x^*) = 0$, definovat výběrovou informaci obsaženou v X jako $\bar{I}_p = n^{-1} \sum s^2(x_i - x^*)$ nebo výběrovou vzájemnou informaci X a Y jako $\bar{\rho}_{XY} = n^{-1} \sum s_x(x_i - x^*)s_y(y_i - y^*)$. Charakteristiky datových souborů, v nichž namísto "syrových" dat x_i vystupují (třeba i latentní) hodnoty skórové funkce, $s(x_i)$ nebo $s(x_i|\theta)$, nazveme *geometrické charakteristiky* náhodných výběrů.

Je zřejmé, že geometrické charakteristiky zohledňují charakter rozdělení. Je-li např. $p \sim e^{-|x|}$ (rozdělení s těžkými chvosty), je pro $|x| \rightarrow \infty$ $|s| = O(1)$, takže geometrické charakteristiky jsou robustní. Pro rozdělení s malou pravděpodobností výskytu odlehklých hodnot (*ostré* rozdělení), na př. $p \sim e^{-x^2}$, je pro $|x| \rightarrow \infty$ $s(x) = O(x)$ a geometrické charakteristiky jsou k případným odlehklým hodnotám citlivé. To je právě případ normálního rozdělení, pro něž $s(x) = x$ a geometrické charakteristiky se shodují s tradičními elementárními statistickými charakteristikami (výběrová informace je zde rovna rozptylu). Viděno z opačného konce, tradiční charakteristiky náhodných výběrů jsou šité na míru výběru z normálního rozdělení, nepřihlížejí k povaze rozdělení skutečného či předpokládaného a proto často nedávají použitelné výsledky a je nutno užívat různých umělých postupů (s výjimkami jako metoda maximální věrohodnosti).

Důvod, proč se ve statistice geometrické charakteristiky dat explicitně nepoužívají je samozřejmě pádny: vše co bylo řečeno o modelu \mathcal{P}_R^μ skutečně funguje jen v případě rozdělení s kladnou pravděpodobnostní mírou na celé reálné přímce. Je-li $T \neq R$, nejsou hodnoty $s(x)$ v žádném bodě $\theta \in \Theta$ shodné s hodnotami jakéhokoli věrohodnostního skóru, (3) nemusí být (a není) informací a geometrické charakteristiky dávají nesmyslné hodnoty.

Tento problém je snad vyřešen v následující kapitole, kde je definována

funkce, kterou je třeba pro konstrukci geometrických charakteristik použít. Nazýváme ji *geometrickou funkcí rozdělení*. V případě $T = R$ je to skórová funkce, případě $T \neq R$ je to funkce, která má vlastnosti obdobné vlastnostem skórové funkce na $T = R$. Uvidíme, že to je funkce pro většinu parametrických rozdělení vlastně známá a prakticky používaná. Její analytické vyjádření pro *holé* rozdělení bez parametrů však donedávna (Fabián 1994) známo nebylo, snad pro svou koncepční nezvyklost (je závislé na T). Vhodnější název geometrické funkce by asi byl *vlivová* funkce rozdělení, termín je však již obsazen.

2. GEOMETRICKÁ FUNKCE ROZDĚLENÍ.

Množinu Π_T všech přípustných rozdělení na T budeme považovat za obraz množiny Π_R zprostředkovaný nějakým zobrazením $\varphi : R \rightarrow T$, monotónním a dostatečně hladkým. Každé rozdělení $p \in \Pi_T$ má tak v Π_R svůj *vzor*, který označíme p_R . Pro příslušné distribuční funkce platí $F(u) = F_R(\varphi^{-1}(u))$ a pro odpovídající hustoty, položíme-li $x = \varphi^{-1}(u)$, platí

$$p(u) = \frac{dF(u)}{du} = \frac{dF_R(x)}{dx} \frac{dx}{du} = \frac{p_R(x)}{L(u)}, \quad (5)$$

kde $L(u) = \varphi'(\varphi^{-1}(u))$ je jakobián transformace $\varphi : R \rightarrow T$.

Definice. Buď $\varphi : R \rightarrow T$ vhodné zobrazení, buď rozdělení $p \in \Pi_T$ a s skórová funkce jeho vzoru $p_R \in \Pi_R$. Geometrickou funkci rozdělení q definujeme jako obraz skórové funkce vzoru při zobrazení φ , t.j. vztahem

$$q(u) = s(\varphi^{-1}(u)). \quad (6)$$

Vyjádření q pomocí p a φ podává

Věta 1.

$$q(u) = \frac{1}{p(u)} \frac{d}{du} (-L(u)p(u)). \quad (7)$$

Důkaz. Podle (2) a (5) je

$$q(u) = -\frac{1}{p_R(x)} \frac{d}{dx} (p_R(x)) = -\frac{1}{L(u)p(u)} \frac{d}{du} (L(u)p(u)) \cdot L(u).$$

□

Věta 2. $E_p q = 0$.

Důkaz. Označme $c_1 = \inf\{u : u \in T\}$, $c_2 = \sup\{u : u \in T\}$. Podle (7) a (5),

$$E_p q = \int_{c_1}^{c_2} q(u)p(u) du = -L(u)p(u)|_{c_1}^{c_2} = p_R(x)|_{-\infty}^{\infty} = 0.$$

Geometrickou funkci parametrického rozdělení definujeme jako

$$q(u|\theta) = \frac{1}{p(u|\theta)} \frac{d}{du} (-L(u)p(u|\theta)). \quad (8)$$

Lze snadno ukázat, že pro (8) platí vztah (6) (s parametrickou skórovou funkcí $s(x|\theta) = -p_T^{-1}(x|\theta)\partial p_T(x|\theta)/\partial x$) a Věta 2.

Pro $\theta \in \Theta \subset R^m$ máme teď namísto m věrohodnostních skórovů jedinou funkci charakterizující poměry ve výběrovém prostoru za vlády rozdělení p . To může být výhoda při studiu vzdálenosti v T v bodě $p(\cdot|\theta)$ variety P_T^θ , kterou lze definovat obdobně jako v (4), t.j. vztahem $d(x_1, x_2|\theta) = |s(x_2|\theta) - s(x_1|\theta)|$ (porovnejte s obdubnou konstrukcí vzdálenosti pomocí věrohodnostních skórovů v (Oller 1989)), ale je to nevýhoda při konstrukci odhadů pro $m > 1$. Uvidíme, že tato nevýhoda může být někdy vyvážena užitím lineárně nezávislých funkcí $q^k(u|\theta)$, $k = 1, \dots, m$.

Zavedeme *transformovaný parametr polohy* vztahem $\nu = \varphi(\mu)$. Pak platí

Věta 3.

$$q(u|\nu) = L(\nu)s_1(u|\nu)$$

where s_1 is given by (1). *Důkaz.* Položme $r = x - \mu$. Platí

$$\begin{aligned} s_1(u|\nu) &= \frac{1}{p(u|\nu)} \frac{dp(u|\nu)}{d\nu} = \frac{L(u)}{p_R(r)} \frac{d(L^{-1}(u)p_R(r))}{d(r)} \frac{d(x - \varphi^{-1}(\nu))}{d\nu} \\ &= \frac{p'_R(r)}{p_R(r)} \frac{-1}{\varphi'(\nu)} = s(x - \mu)L^{-1}(\nu) = L^{-1}(\nu)q(u|\nu) \end{aligned}$$

podle (5) a věty o derivaci inverzní funkce. \square

Podle Věty 3 je tedy geometrická funkce rozdělení pro většinu parametrických rozdělení věrohodnostním skórem nejzajímavějšího parametru. Uvážíme-li i Větu 2, má q na T vlastnosti obdobné vlastnostem s na R . Pro různá $\varphi : R \rightarrow T$ dostáváme ovšem různá q , lišící se jakobiánem zobrazení. Sortiment φ však není moc široký. Pro $\varphi_R : R \rightarrow R$ je vhodné zvolit identické zobrazení (geometrická funkce parametrického rozdělení je pak parametrická skórová funkce). Pro $\varphi_{R^+} : R \rightarrow (0, \infty)$ je také asi jediná rozumná volba,

$$z = \varphi_{R^+}(x) = e^x, \quad (9)$$

kteřé lognormálnímu rozdělení přiřazuje jako vzor rozdělení normální a log-Cauchyovu rozdělení Cauchyovo. Pak je $\varphi^{-1}(z) = \ln z$, $L(z) = z$ a geometrická funkce rozdělení je dána vztahem

$$q(z) = zs(z) - 1, \quad (10)$$

kde $s = -p'/p$. Zvolíme-li ovšem pevně φ_{R^+} , je geometrická funkce rozdělení na (a, b) už jednoznačně určena. Pro zobrazení $\varphi_{ab} : R \rightarrow (a, b)$ by mělo totiž platit

$$\lim_{\substack{a \rightarrow 0 \\ b \rightarrow \infty}} \varphi_{ab}^{-1}(w) = \varphi_{R^+}^{-1}(z). \quad (11)$$

Obecné zobrazení splňující (11) při φ_{R^+} zvoleném v (9) je

$$\varphi_{ab}^{-1}(w) = \ln \frac{(b-c)(w-a)}{(c-a)(b-w)}, \quad (12)$$

kde $a < c < b$ a kde případně $c = c(a, b)$ za podmínky $\lim_{\substack{a \rightarrow 0 \\ b \rightarrow \infty}} c(a, b) = 1$. Ježto však

$$\frac{d}{dw}(\varphi_{ab}^{-1}(w)) = L^{-1}(w) = \frac{b-a}{(w-a)(b-w)},$$

nezávisí L na c a libovolná z transformací (12) definuje tutéž geometrickou funkci rozdělení ve tvaru

$$q(w) = [(w-a)(b-w)s(w) + 2w - (b+a)]/(b-a). \quad (13)$$

Pro $T = (0, 1)$ se (13) redukuje na

$$q(w) = w(1-w)s(w) + 2w - 1. \quad (14)$$

Zobrazení $\varphi_{01}^{-1}(w) = \ln(w/(1-w))$ je logit transformace.

Poznamenejme, že za uvažovaných předpokladů je vztah mezi hustotou a geometrickou funkcí vzájemně jednoznačný.

3. GEOMETRICKÉ MOMENTY.

3.1. Holá rozdělení.

Definice. Momentům geometrické funkce q rozdělení p ,

$$M_k = \int_T q^k(u)p(u)du, \quad (15)$$

budeme říkat *geometrické momenty*.

Věta 4. Geometrické momenty rozdělení na $T \neq R$ jsou shodné s geometrickými momenty vzoru.

Důkaz. Podle (6) a (5)

$$\int_{c_1}^{c_2} q^k(u)p(u)du = \int_{c_1}^{c_2} s^k(\varphi^{-1}(u))p_R(\varphi^{-1}(u))L^{-1}(u)du = \int_R s^k(x)p_R(x)dx.$$

Je známo, že obvyklé (*Euklidovské*) momenty $m_k = E_p x^k$ rozdělení s těžkými chvosty často neexistují (z rozdělení dále zmiňovaných je to Cauchyovo, extrémní hodnoty II, log-logistické a log-Cauchyovo, viz příklad 1). Naproti tomu všechny geometrické momenty za rozumných předpokladů existují. To je evidentní právě v případech rozdělení s těžkými chvosty, která mají omezenou geometrickou funkci. Tvzení pro rozdělení s neomezenou q stačí na základě Věty 4 vyslovit pro případ $T = R$.

Věta 5. Buď s skórová funkce rozdělení p_R . Buď $|s(x)| > |x|^\alpha$, $\alpha > 0$ pro $|x| > 1$. Nechť existuje takové x_1 , že pro $|x| > x_1$ platí $s^2(x) > k s'(x)$. Pak integrál $M_k = E_{p_R} s^k$ konverguje.

Důkaz. Zřejmě platí $Q(x) = \int s(x)dx > x/(\alpha + 1)$. Podle (2) je

$$g_k(x) = s^k(x)p_R(x) = cs^k(x)e^{-Q(x)}.$$

Stačí uvažovat $\int_a^\infty g_k(x)dx$. Pro $x > x_1$ je g_k klesající, neboť

$$g'_k(x) = ce^{-Q(x)}s^{k-1}(x)[ks'(x) - s^2(x)] < 0.$$

Zřejmě tedy lze nalézt bod $x_2 > 0$, pro nějž $g_k(x_2) \leq k^k$. Zvolme $a = \max(1, x_1, x_2)$. Protože pak

$$\frac{1}{ck^k}g_k(x) = \left(-\frac{d}{dx}e^{-\frac{1}{k}Q(x)}\right)^k \leq -\frac{d}{dx}e^{-\frac{1}{k}Q(x)},$$

integrál $\int_a^\infty g_k(x)dx \leq -ck^k e^{-\frac{1}{k}Q(x)}|_a^\infty < \infty$. □

Podmínka Věty 5 je pro hladká rozdělení typu $|s(x)| > |x|^\alpha$, $\alpha > 0$ splněna.

$M_1 = 0$ podle Věty 2. Hodnota $u^* : q(u^*) = 0$ je tak, vedle střední hodnoty, mediánu a módu, alternativním "těžištěm" rozdělení vzhledem ke geometrii generované rozdělení. Lze tedy "geometrické výběrové těžiště" souboru $U = (u_1, \dots, u_n)$ odhadovat jako $u^* : n^{-1} \sum q(u_i|u^*) = 0$. Konkrétně, pro $T = R$ je $x^* = n^{-1} \sum s(u_i)$, pro $T = (0, \infty)$ je

$$z^* : \quad n^{-1} \sum q(z_i|z^*) = n^{-1} \sum s(\ln(z_i/z^*)) = 0$$

a pro $T = (0, 1)$

$$w^* : \quad n^{-1} \sum s\left(\ln \frac{w_i(1-w^*)}{(1-w_i)w^*}\right) = 0.$$

s je skórová funkce vzoru. V případě lognormálního rozdělení je z^* obvyklý geometrický průměr. Podle Věty 3 je "geometrické těžiště" totožné s maximálně věrohodným odhadem zobecněného parametru polohy.

Překvapivě, $u^* : q(u^*) = 0$ je zároveň bodem, v němž je relativní informace o rozdělení p minimální (ačkoli pro $T \neq R$ není $p(u^*)$ maximem $p(u)$). Podle (5) je totiž $p(u)$ podílem dvou členů, z nichž jakobián zobrazení φ je společný všem rozdělením na T a nenese tudíž žádnou informaci o p . Nejméně informativním bodem rozdělení je tedy bod, v němž je člen $p_R(\varphi^{-1}(u))$ maximální. Podle (5) a (7)

$$\frac{d}{du}p_R(\varphi^{-1}(u)) = \frac{d}{du}(L(u)p(u)) = -q(u)p(u),$$

takže to je právě bod $u^* : q(u^*) = 0$. Funkce q^2 má tedy minimum v bodě s minimální informací, je nezáporná a její střední hodnota má význam informace. $q^2(u)$ je tedy asi možno interpretovat jako funkci popisující relativní informaci obsaženou v pozorování $u \in T$. $M_2 = E_p q^2$ lze pak zřejmě považovat za střední informaci spojitého rozdělení.

V případě diskretního rozdělení $p = \{p_i\}$ je neurčitost náhodného experimentu vyjádřena Shannonovu entropií $H_p = \sum -\ln p_i \cdot p_i$. Její obdobu pro absolutně spojitá rozdělení, diferenciální entropii $h_S = \int_T -\ln p(u) p(u) du$, podobně interpretovat nelze, protože může být záporná. Z tabulky příkladu 2 je patrné, že rozdělení s těžkými konci (jistě více neurčitá než ostrá rozdělení) mají menší hodnoty M_2 než ostrá rozdělení. Domnívám se, že převrácená hodnota střední informace rozdělení

$$h_F = M_2^{-1}, \quad (16)$$

kteří v příkladě 4 říká *Fisherova entropie*, by mohla být chápána jako míra neurčitosti spojitých rozdělení.

Řekneme, že rozdělení na T je φ -symetrické, je-li jeho vzor symetrické rozdělení. Pro φ -symetrické rozdělení platí $M_3 = 0$. Na $T = (0, \infty)$ je φ -symetrické rozdělení dáno vztahem

$$p(1/z) = z^{-1}p_R(-x) = z^{-1}p_R(x) = z^{-2}p(z).$$

Konečně, moment M_4 má přibližně opačný smysl než m_4 .

3.2. Parametrická rozdělení.

Geometrickými momenty rozdělení $p(u|\theta)$ jsou v parametrickém případě integrály

$$M_k(\theta) = \int_T q^k(u|\theta)p(u|\theta) du. \quad (17)$$

Překvapivým výsledkem výpočtů geometrických momentů pro různá rozdělení je fakt, že výsledné vzorce obsahují pouze parametry, a ne různé neelementární funkce parametrů, jak tomu často bývá v případě Euklidovských

momentů. Nejen z důvodu, že existují, se tedy geometrické momenty zdají být přirozenějšími numerickými charakteristikami rozdělení než Euklidovské momenty (viz příklad 5).

Empirických geometrických momentů lze použít pro odhady parametrů rozdělení na základě náhodného výběru U z rozdělení $p(u|\theta) \in \Pi_T$. Odhady $\hat{\theta}_G$ určené z rovnic

$$\hat{\theta}_G : \quad n^{-1} \sum_{i=1}^n q^k(u_i|\hat{\theta}_G) = M_k(\hat{\theta}_G), \quad k = 1, \dots, m$$

budeme nazývat *G-odhady*. Dá se očekávat a je možná dokázáno (Fabián 1996), že při existenci konečných $r_{jk}(\theta) = E_\theta[\partial(q_k(u|\theta) - M_k(\theta))/\partial\theta_j]$ a platnosti podmínky

$$\text{Det}(r_{jk}(\theta))_1^m \neq 0 \quad (18)$$

pro každé $\theta \in \Theta$ jsou G-odhady silně konzistentní a asymptoticky normální. V citované práci byla studována asymptotická vydatnost G-odhadů. Ukázalo se, že závisí na charakteru geometrické funkce rozdělení. Pro $q \sim \varphi(x^\alpha)$, $\alpha \geq 1$ jsou vydatnosti G-odhadů nižší než vydatnosti Euklidových momentových (EM) odhadů. Označme N neomezenou a O omezenou geometrickou funkci. Pro typy $O-N$ (zleva O , zprava N) a $N-O$ je zpravidla asymptotická relativní vydatnost (ARV) odhadu zobecněného parametru polohy poměrně blízká k jedné, ARV ostatních parametrů je nízká. Rozdělení typu $O-O$ vykazují ARV G-odhadů hodnoty blízké k jedné. G-odhady jsou v tomto případě jednodušší alternativou maximálně věrohodných (ML) odhadů. ML i G-odhady jsou pro tuto třídu rozdělení robustní vzhledem k parametru polohy; G-odhady jsou však robustní i vzhledem k parametru měřítka a je docela možné, že jim pro tuto třídu rozdělení lze před odhady typu ML dát přednost (příklad 7). Poznamenejme, že tento závěr neplatí pro rozdělení s redescenční geometrickou funkcí: pro Cauchyovo rozdělení není splněna podmínka (18).

4. NĚKOLIK PŘÍKLADŮ.

Příklad 1. Zobecněné skórové funkce některých holých rozdělení jsou uvedeny v tabulkách 1-3 pro $T = R$, $(0, \infty)$ a $(0, 1)$.

TABULKA 1.
Skórové funkce a hustoty některých rozdělení na R

typ	$s(x)$	$p_R(x)$	rozdělení
$N-N$	$\sinh x$	$\frac{1}{2K_0(1)}e^{-\cosh x}$	
$N-N$	x	$\frac{1}{\sqrt{2\pi}}e^{-x^2/2}$	normální
$O-N$	$e^x - 1$	$e^x e^{-e^x}$	dvojitě exponenciální
$N-O$	$1 - e^{-x}$	$e^{-x} e^{-e^{-x}}$	extremní hodnoty I
$O-O$	$\operatorname{tgh}(x/2)$	$\frac{e^x}{(1+e^x)^2}$	logistické
$O-O$	$2x/(1+x^2)$	$\frac{1}{\pi}(1+x^2)^{-1}$	Cauchyho

kde K_0 je Besselova funkce III. druhu.

TABULKA 2.
Geometrické funkce a hustoty rozdělení na $(0, \infty)$ se vzory v Tab. 1

typ	$q(z)$	$p(z)$	rozdělení
$N-N$	$\frac{1}{2}(z - 1/z)$	$\frac{1}{2K_0(1)z}e^{-\frac{1}{2}(z+1/z)}$	Waldovo
$N-N$	$\ln z$	$\frac{1}{\sqrt{2\pi}z}e^{-\frac{1}{2}\ln^2 z}$	lognormální
$O-N$	$z - 1$	e^{-z}	exponenciální
$N-O$	$1 - 1/z$	$z^{-2}e^{-1/z}$	extr. hodn. II
$O-O$	$(z - 1)/(z + 1)$	$1/(z + 1)^2$	log-logistické
$O-O$	$2 \ln z/(1 + \ln^2 z)$	$\frac{1}{\pi z}(1 + \ln^2 z)^{-1}$	log-Cauchyho

TABULKA 3.
Geometrické funkce a hustoty rozdělení na $(0, 1)$ se vzory v Tab. 1.

typ	$q(w)$	$p(w)$
$N-N$	$\frac{1}{2} \frac{(2w-1)}{w(1-w)}$	$\frac{1}{2K_0(1)w(1-w)}e^{-\frac{w^2-w+1/2}{w(1-w)}}$
$N-N$	$\ln \frac{w}{1-w}$	$\frac{1}{\sqrt{2\pi}w(1-w)}e^{-\frac{1}{2}\ln^2 \frac{w}{1-w}}$
$O-N$	$\frac{2w-1}{1-w}$	$\frac{1}{(1-w)^2}e^{-\frac{w}{1-w}}$
$N-O$	$\frac{2w-1}{w}$	$\frac{1}{w^2}e^{-\frac{1-w}{w}}$
$O-O$	$2w - 1$	1
$O-O$	$2 \ln \frac{w}{1-w}/(1 + \ln^2 \frac{w}{1-w})$	$\frac{1}{\pi w(1-w)} \left(1 + \ln^2 \frac{w}{1-w}\right)^{-1}$

V tabulce 3 je druhé z rozdělení typu $N-N$ Johnsonovo U_B rozdělení (Patil 1984) a první z $O-O$ rozdělení rovnoměrné. Ostatní nejsou myslím známa. Z tabulek je patrné, že popis rozdělení pomocí geometrické funkce je většinou jednodušší než popis pomocí hustoty.

Příklad 2. Geometrické momenty rozdělení z tabulek 1-3 jsou uvedeny v tabulce 4.

TABULKA 4.

typ	M_2	M_3	M_4
$N-N$	1.430	0	11.578
$N-N$	1	0	3
$O-N$	1	2	9
$N-O$	1	-2	9
$O-O$	1/3	0	1/5
$O-O$	1/2	0	3/8

Příklad 3. Rozdělení extrémní hodnoty II (Tab.1). Střední hodnota rozdělení neexistuje, hodnota mediánu je 0.693 a módu 0.5. Z rovnice $q(z^*) = 0$ dostáváme "geometrické těžiště rozdělení" v bodě $z^* = 1$.

Příklad 4. Rovnoměrné rozdělení na $T = [0, b]$. Diferenciální entropie rozdělení je $h_S(b) = \log b$ a její hodnota a dokonce znaménko závisí na zvolené délkové jednotce. Geometrická funkce rozdělení je $q(w|b) = 2w/b - 1$ a $M_2(b) = b^{-3} \int_0^b (2w - b)^2 dw = 1/3$. Fisherova entropie (16) rozdělení je $h_F = 3$ nezávisle na délce intervalu.

Rovnoměrné rozdělení je rozdělením s maximální diferenciální entropií na intervalu; diferenciální entropie libovolného jiného rozdělení na $T = (0, 1)$ je tedy záporná. Uvažujme rozdělení beta s hustotou

$$p(w|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} w^{\alpha-1} (1-w)^{\beta-1}, \quad w \in [0, 1], \alpha, \beta \geq 0,$$

kde B je beta funkce. Toto rozdělení nemá transformovaný parametr polohy, takže příslušná geometrická funkce $q(w|\alpha, \beta) = (\alpha + \beta)w - \alpha$ není úměrná žádnému z věrohodnostních skóru. Platí

$$M_2(\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \int_0^1 [(\alpha + \beta)w - \alpha]^2 w^{\alpha-1} (1-w)^{\beta-1} dw = \frac{\alpha\beta}{\alpha + \beta + 1}.$$

Hodnoty Fisherovy entropie rozdělení beta zachycuje následující tabulka.

β	4	2	1	0.5	0.25
$h_F(1, \beta)$	1.5	2	3	5	9

Fisherova entropie je ovšem kladná, z tabulky je však patrné, že rovnoměrné rozdělení není rozdělením s maximální neurčitostí, měřeno pomocí (16). To ale může být rozumný výsledek. Vztah $h_F > 3$ platí pro malé hodnoty parametrů, kdy je hustota rozdělení beta antimodální. Výsledek pozorování je tedy ve Fisherově smyslu více neurčitý, když je očekáván v jedné ze dvou téměř oddělených oblastí, než když je očekáván se stejnou relativní pravděpodobností na celém intervalu.

Příklad 5. Mějme model

$$\mathcal{P}_{0,\infty}^\theta = \{(0, \infty), \mathcal{B}_{(0,\infty)}, p(z|\nu, \beta, \lambda) : \nu, \beta, \lambda \in (0, \infty)\}$$

kde

$$p(z|\nu, \beta, \lambda) = \frac{\beta\lambda^\lambda}{z\Gamma(\lambda)} \left(\frac{z}{\nu}\right)^{\beta\lambda} e^{-\lambda(z/\nu)^\beta}. \quad (19)$$

Rodina (19) vznikla zavedením možných parametrů do geometrické funkce $q(z) = z - 1$ exponenciálního rozdělení, s výsledkem

$$q(z|\nu, \beta, \lambda) = \lambda\beta[(z/\nu)^\beta - 1],$$

který byl použit při integraci rovnice (10). Členy rodiny (19) jsou např. rozdělení Weibullovo, Maxwellovo, gamma, Erlangovo či chi-kvadrát. Při výpočtu obou typů momentů se použije integrál

$$\int_0^\infty z^{\alpha-1} e^{-rz^\beta} dz = \frac{1}{\beta} r^{-\alpha/\beta} \Gamma(\alpha/\beta).$$

Označíme-li $u = z/\nu$, platí

$$\begin{aligned} m_k &= \frac{\nu^k \lambda^\lambda \beta}{\Gamma(\lambda)} \int_0^\infty u^{k+\beta\lambda-1} e^{-\lambda u^\beta} du = \left(\frac{\nu}{\lambda^{1/\beta}}\right)^k \frac{\Gamma(\lambda + k/\beta)}{\Gamma(\lambda)}, \\ M_k &= \frac{\lambda^{k+\lambda} \beta^{k+1}}{\Gamma(\lambda)} \int_0^\infty (u^\beta - 1)^k u^{\beta\lambda-1} e^{-\lambda u^\beta} du \\ &= \beta^k \sum_{j=0}^k (-\lambda)^j \binom{k}{j} \frac{\Gamma[((k-j)\beta + \lambda\beta)/\beta]}{\Gamma(\lambda)}. \end{aligned}$$

V případě geometrických momentů se v argumentu gamma funkce zkrátí β , takže výsledek je neuvěřitelně jednoduchý:

$$M_1 = 0, \quad M_2 = \lambda\beta^2, \quad M_3 = 2\lambda\beta^3, \quad M_4 = 3\lambda(\lambda + 2)\beta^4.$$

Obecně, geometrické momenty nezávisí na transformovaném parametru polohy. Podle Věty 3 je Fisherova informace o parametru ν rovna $I_p(\nu) = M_2/\nu^2$. Pro zajímavost, diferenciální entropie rozdělení (19) je dána výrazem

$$h_S(\nu, \beta, \lambda) = -\log \beta - \lambda \log \lambda + \log \Gamma(\lambda) + \log \nu + (\lambda - 1/\beta)(\log \lambda - \psi(\lambda)) + \lambda.$$

Poznamenejme ještě, že rozdělení φ -symetrické k (19) podle osy $z = 1$ (jeho vzor je symetrický podle osy $x = 0$ se vzorem (19)), má geometrické momenty $M_k = (-1)^k M_k$. Jeho Euklidovské momenty jsou dány vzorcem

$$\tilde{m}_k = (\nu \lambda^{1/\beta})^k \frac{\Gamma(\lambda - k/\beta)}{\Gamma(\lambda)}$$

a existují jen pro $k < \beta\lambda$.

Příklad 6. Mějme speciální případ Lomaxova rozdělení s hustotou a geometrickou funkcí

$$p_L(z|\lambda) = \frac{\lambda}{(z+1)^{1+\lambda}}, \quad q_L(z|\lambda) = \frac{\lambda z - 1}{z+1}$$

a Gumbelovo zobecnění logistického rozdělení

$$p_G(z|\lambda) = \frac{1}{B(\lambda, \lambda)} \frac{z^{\lambda-1}}{(1+z)^{2\lambda}}, \quad q_G(z|\lambda) = \lambda \frac{z-1}{z+1}.$$

Rovnice pro EM, ML a G-odhady jsou

	Lomax	Gumbel
EM :	$\frac{1}{n} \sum z_i = 1/(\hat{\lambda}_E - 1)$	$\frac{1}{n} \sum z_i = \hat{\lambda}_E/(\hat{\lambda}_E - 1)$
ML :	$\frac{1}{n} \sum \log(1+z_i) = 1/\hat{\lambda}_L$	$\frac{1}{n} \sum \log \frac{z_i}{1+z_i} = 2[\psi(\hat{\lambda}_L) - \psi(2\hat{\lambda}_L)]$
G :	$\frac{1}{n} \sum \frac{\lambda_G z_i - 1}{z_i + 1} = 0$	$\frac{1}{n} \sum \left(\frac{z_i - 1}{z_i + 1}\right)^2 = 1/(2\hat{\lambda}_G + 1)$

kde $\psi = \Gamma'/\Gamma$. EM platí pro $\lambda > 1$. V případě Gumbelova rozdělení bylo nutno použít 2. geometrický moment. Odvozené vzorce pro asymptotické rozptyly odhadů jsou

	Lomax	Gumbel
$Var \hat{\lambda}_E$	$\frac{\lambda(\lambda-1)^2}{\lambda-2}$	$\frac{\lambda(\lambda-1)^2(2\lambda-1)}{\lambda-2}$
$Var \hat{\lambda}_L$	λ^2	$2(\psi'(\lambda) - 2\psi'(2\lambda))$
$Var \hat{\lambda}_G$	$\frac{\lambda(\lambda+1)^2}{\lambda+2}$	$\frac{\lambda(2\lambda+1)^2}{2\lambda+3}$

a platí v případě EM odhadu pro $\lambda > 2$. V následující tabulce jsou uvedeny asymptotické relativní vydatnosti $e_{\hat{\lambda}_E} = (Var \hat{\lambda}_E / Var \hat{\lambda}_L)^{1/2}$ odhadů typu EM a $e_{\hat{\lambda}_G} = (Var \hat{\lambda}_G / Var \hat{\lambda}_L)^{1/2}$ G-odhadů, které jsou blízké jedné.

TABULKA 5.

λ	Lomax		Gumbel	
	$e_{\hat{\lambda}_E}$	$e_{\hat{\lambda}_G}$	$e_{\hat{\lambda}_E}$	$e_{\hat{\lambda}_G}$
1	–	0.750	–	0.782
2	–	0.889	–	0.906
3	0.8	0.938	0.258	0.948

Příklad 7. Bylo generováno 30 náhodných výběrů o délce 200 vzorků z logistického rozdělení $l(0, 1)$, kde

$$l(\mu, \sigma) = p(x|\mu, \sigma) = \frac{e^{-\frac{x-\mu}{\sigma}}}{(e^{-\frac{x-\mu}{\sigma}} + 1)^2},$$

a z kontaminovaného rozdělení $l_{\varepsilon c} = (1 - \varepsilon)l(0, 1) + \varepsilon l(0, c)$.

Označme $y = (x - \mu)/\sigma$. Skórová funkce rozdělení $l(\mu, \sigma)$ je $s(x|\mu, \sigma) = \sigma^{-1}(e^y - 1)/(e^y + 1)$. Rovnice pro výpočet ML a G-odhadů parametrů jsou

$$\begin{array}{ll} \text{ML} & \text{G} \\ \sum s(\hat{y}_{Li}) = 0 & \sum s(\hat{y}_{Gi}) = 0 \\ \sum (\hat{y}_{Li}s(\hat{y}_{Li}) - 1/\hat{\sigma}_L^2) = 0 & \sum (s^2(\hat{y}_{Gi}) - 1/3\hat{\sigma}_G^2) = 0, \end{array}$$

kde $\hat{y}_{Li} = (x_i - \hat{\mu}_L)/\hat{\sigma}_L$ a podobně pro G . Střední hodnoty odhadnutých parametrů $\bar{\theta} = \sum_{j=1}^{30} \hat{\theta}_j/30$ uvádíme v následující tabulce.

TABULKA 6.

ε	c	$\bar{\mu}_L$	$\bar{\mu}_G$	$\bar{\sigma}_L$	$\bar{\sigma}_G$
0		0.0846	0.0849	1.0111	1.0160
0.05	5	0.0538	0.0558	1.2318	1.0930
0.05	10	0.0632	0.0802	1.5402	1.1234

Z tabulky je patrné že, na rozdíl od ML odhadu, je G-odhad parametru σ málo ovlivněn odlehlými hodnotami v generovaných datech.

PODĚKOVÁNÍ.

Děkuji všem, kteří jste dočetli až sem, zejména pak těm z vás, kteří mi sdělí svůj názor a zvláště těm, kteří při tom dodrží zásadu fandí, ale zůstan člověkem.

LITERATURA

- Cover, T.M., Thomas, J.A., *Elements of Information Theory*, J. Wiley, New York, 1991.
- Fabián, Z., Generalized score function and its use, *Transactions of 12-th Prague Conference on Information Theory*, pp.67-72, 1994.
- Fabián, Z., Geometric moments and geometric moment estimates, *výzk.zpr. V-694, ÚIVT AV ČR, Praha 1996*.
- Hampel, F.R., Rousseeuw, P.J., Ronchetti, E.M., Stahel, W.A., *Robust Statistic. The Approach Based on Influence Functions*, J. Wiley, New York, 1987.
- Lehmann, E.J., *Theory of point estimation*, J. Wiley, New York, 1983.
- Oller, J.M., Some geometrical aspects of data analysis and statistics. *Statistical Data Analysis and Inference*, Y. Dodge (ed.), Elsevier, North-Holland, 1989.
- Patil, G.P., Boswell, M.T., Ratnaparkhi, M.V., *Dictionary and Classified Bibliography of Statistical Distributions in Scientific Work*. Int. cooperative Publ. House, Maryland, 1984.