

Gnostická teorie dat je pokusem o matematický model poznávání kvantitativních vlastností světa za působení neurčitosti, o které není známo, zda vyhovuje statistickým předpokladům. Z jediného axiomu se dokazují významné zákonitosti jednotlivých dat, z nichž lze za použití druhého axiomu odvodit vlastnosti datových souborů a jejich charakteristiky použitelné k robustnímu odhadování polohy, měřítka i distribučních funkcí datových souborů, k testování homogenity souborů apod.

MATEMATICKÁ DATA VERSUS REÁLNÁ DATA

Statistická teorie je vybudována pro data generovaná myšlenými matematickostatistickými modely, nebo pro data podle těchto modelů vypočtená, či fyzikálně sestavená tak, aby vlastnosti dat dobře aproximovaly vlastnosti matematicky definovaných dat. Praktická použitelnost statistických metod je pak podmíněna shodou vlastností skutečných dat s teoretickými. Takovou shodu často nelze zaručit, dokonce někdy ani rozpoznat pro nedostatečné množství dat, proměnnost situací, nestacionárnost atp. Nemůžeme-li ale přisoudit datům statistické vlastnosti a přesto potřebujeme vybudovat teorii vedoucí k metodám použitelným pro zpracování dat, musíme najít jejich jinou důležitou vlastnost, z níž lze teorii odvodit: Reálná data přijímaná za základ gnostické teorie jsou výsledky buď měření nebo čítání reálných kvantit. Z teorie měření vyplývá, že tato data jsou striktně kladná a konečná. V praxi se setkáváme i s daty, která byla získána dodatečnou transformací výsledků měření nebo čítání, mají např. posunutou nulu ap. Ta musí být před dosazením do gnostických vzorců transformována do požadované primární formy.

GNOSTICKÁ TEORIE JEDNOTLIVÝCH DAT

První axiom a jeho důsledky pro teorii kvantifikace

Shrňme stručně gnostickou teorii jednotlivých dat dle /1/, kde jsou uvedeny přesnější formulace tvrzení i s důkazy. Nechť "z" je skutečný výsledek kvantifikace, tj. měření nebo čítání, vyjádřený číslem a nechť "z₀" je (ideální) hodnota, kterou bychom dostali jako výsledek kvantifikace, kdyby byl proces kvantifikace zbaven jakékoliv neurčitosti. První axiom gnostické teorie pak charakterizuje model reálných dat takto:

$$z = z_0 \xi \quad (z, z_0, \xi \in \mathbb{R}_+) \quad /1/$$

Snadno zjistíme, že množina transformací zobrazujících jakoukoliv datovou položku tvaru /1/ na sebe či na jinou položku dat, lišící se pouze parametrem ξ charakterizujícím vliv neurčitosti tvoří komutativní grupu, jejíž strukturální operací je násobení operátorů. Veličina "z" dle /1/ je však bilineární funkcí proměnných z_0 a ξ . Záměna proměnných

$$x = z_0 \cdot \text{ch} \Omega \quad y = z_0 \cdot \text{sh} \Omega \quad / \text{kde } \Omega = \ln \xi / \quad /2/$$

pak dává možnost přiřadit každé položce dat (z) typu /1/ dvourozměrový vektor $\underline{u}^T := (x, y)^T$ a také $\underline{u}^T := (y, x)^T$, přičemž $z = x + y$. Pro neznamlost z_0 ani ξ nedovedeme takový rozklad dat uskutečnit, ale jeho teoretický tvar nám dovoluje odvozovat důležité závěry. Především zjišťujeme, že množina matic \underline{K} zobrazujících vektory typu \underline{u} na sebe nebo na jiné \underline{u}' , tvoří komutativní grupu, jejíž strukturální operací je běžné maticové násobení. Tato grupa je izomorfní jednak s grupou operátorů ξ na \mathbb{R}_+ a jednak s grupou Lorentzových transformací na dvourozměrné varietě vektorů. Lze dokázat, že metriku na této varietě nemůžeme libovolně zadat, ta je již určena, neboť jediná metrika invariantní k Lorentzovým transformacím je pseudoeuklidovská metrika. Jako důsledek prvního axiomu tedy dostáváme, že vliv kvantifikační neurčitosti na reálná data lze matematicky modelovat jako ortogonální rotaci vektorů v Minkowského rovině,

při níž se body reprezentující reálná data pohybují po pseudoeuklidovské kružnici, tj. po euklidovské hyperbole. Relativní délka dráhy tohoto pohybu může tedy být snadno vyčíslena, je rovna veličině $|\Omega|$ o níž se snadno přesvědčíme, že vyhovuje jednoduchému variačnímu principu: je to maximum z možných relativních délek drah pohybů mezi zadanými koncevními body po různých křivkách. V tomto smyslu tedy lze hovořit o tom, že vliv kvantifikační neurčitosti maximalizuje znehodnocení reálných dat.

Teorie estimačních transformací a ideální gnostický cyklus

Nás ale zajímá i další transformace dat, estimační. Ta má výsledku kvantifikace (s) přiřadit odhad \hat{s}_0 . V dvourozměrné reprezentaci je stanovení estimační transformace ekvivalentní nalezení dráhy po níž se má pohybovat konec vektoru \underline{y} nebo \underline{u} , aby se nakonec dostal co nejbližší k bodu $(s_0, 0)^T$ nebo $(0, s_0)^T$. Abychom stanovili nejlepší dráhu estimační transformace, studujeme dále kvantifikační transformace Lorentzova typu, o nichž si dokážeme, že splňují podmínky homogenity, dvojitě symetrie a rovnoměrné regularity. Matematický zápis těchto vlastností je však hyperbolickou verzí známých podmínek analytičnosti funkce komplexní proměnné. Estimační transformace plně duální ke kvantifikační je proto daná euklidovskými ortogonálními otáčením vektorů v Gaussově rovině. Vyhovuje podmínkám analytičnosti Cauchy-Riemanna a příslušná dráha odpovídá opačnému variačnímu principu: její relativní délka je minimální, duálně k (2) platí pro estimační dráhu

$$\begin{aligned} x &= r \cos \omega & y &= -r \sin \omega & /3/ \\ \text{kde} & & r &= \sqrt{x^2 + y^2} & /4/ \\ & & \operatorname{tg} \omega &= -\operatorname{th} \Omega & /5/ \end{aligned}$$

Z obou drah, kvantifikační a estimační, lze sestavit uzavřený cyklus transformací, který nazveme ideálním gnostickým cyklem. Tento cyklus hraje významnou úlohu jako vzor pro postup skutečného procesu odhadování, lze dokázat jeho důležité extrémální vlastnosti.

Charakteristiky nepodobnosti gnostických událostí

Gnostickou událostí je skutečnost, že byl získán výsledek kvantifikace tvaru /1/. Za gnostické události můžeme proto přijmout i vektor \underline{u} a \underline{u}' mající složky x, y tvaru /2/ nebo /3/, jejichž součet dá datovou položku tvaru /1/. V předchozím odstavci jsme se přesvědčili o plodnosti geometrického přístupu, který dovolil odvodit důležité vlastnosti matematického modelu neurčitosti s úvah o metrice. Zůstaneme geometrickému přístupu věrni při posuzování vztahů mezi dvěma událostmi s využitím pojmu podobnosti, který však oproti běžnému zobecníme tak, aby chom vyčerpali všechny aspekty podobnosti. Dvě gnostické události \underline{u} a \underline{u}' mající složky x, y a x', y' budeme považovat za $\hat{\epsilon}, \hat{\xi}$ -podobné, jestliže platí

$$x/y = \hat{\epsilon} \cdot (x'/y')^{\hat{\xi}} \quad / \text{kde } \hat{\epsilon} = \pm 1 \text{ a } \hat{\xi} = \pm 1 / \quad /6/$$

Za číselnou charakteristiku $\hat{\epsilon}, \hat{\xi}$ -nepodobnosti pak přijmeme rozdíl obou těchto poměrů. Dosazením z /2/ se pak přesvědčíme, že charakteristikami nepodobnosti v případě kvantifikace jsou hyperbolické sinusy a cosinusy součtů a rozdílů parametrů neurčitosti Ω obou událostí. S použitím /3/ dostaneme pro estimaci další čtyři charakteristiky nepodobnosti, obyčejné sinusy a cosinusy součtů a rozdílů parametrů ω spjatých s parametry Ω vztahem /5/. Obvyklá "trejúhelníková" podobnost odpovídající hodnotám $\hat{\epsilon} = 1$ a $\hat{\xi} = 1$ nás na tomto místě zajímá nejméně, protože pouze ostatní tři kombinace hodnot $\hat{\epsilon}$ a $\hat{\xi}$ umožňují kvantitativně charakterizovat kvalitu gnostické události novým způsobem, charakteristikami "nepodobnosti sobě". Pro $\underline{u}' = \underline{u}$ jsou totiž numericky charakteristikami nepodobnosti mezi \underline{u} a \underline{u}' (tedy nepodobnosti mezi \underline{u} a \underline{u}) hyperbolické sinusy a cosinusy dvojnásobného argumentu Ω a obyčejné sinusy a cosinusy dvojnásobného argumentu ω . Zapišme explicitní výrazy pro tyto charakteristiky i snadno odvoditelné vztahy mezi nimi vyplývající z /5/:

$$\begin{aligned} f_1 &:= \cos 2\omega = 1/\operatorname{ch} 2\Omega = (x^2 - y^2)/(x^2 + y^2) = 2/(\xi^2 + \xi^{-2}) & /7/ \\ h_0 &:= \sin 2\omega = -\operatorname{th} 2\Omega = 2xy/(x^2 - y^2) = (\xi^2 - \xi^{-2})/2 = -h_0^* & /8/ \\ h_0 &:= \sin 2\omega = -\operatorname{th} 2\Omega = -2xy/(x^2 + y^2) = -(\xi^2 - \xi^{-2})/(\xi^2 + \xi^{-2}) = -h_0^* & /9/ \\ \text{kde} & & \xi &:= x/s_0 = (x+y)/s_0 & /10/ \end{aligned}$$

V předcházejícím odstavci jsme se již setkali s maticemi K typu

$$K_q(\Omega) = \begin{pmatrix} \cosh \Omega & \sinh \Omega \\ \sinh \Omega & \cosh \Omega \end{pmatrix} \quad /11/$$

zobrazujícími jednu gnostickou událost tvaru /2/ na jinou. Taková matice je prvkem komutativní grupy izomorfní s Lorentzovou grupou. Duální k /11/ je estimační operátor

$$K_e(\omega) = \begin{pmatrix} \cosh \omega & -\sinh \omega \\ \sinh \omega & \cosh \omega \end{pmatrix} \quad /12/$$

kteřý je také prvkem komutativní grupy. Ta však již izomorfní s Lorentzovou grupou není. Součiny matic z obou grup již jednoznačně spjatý nejsou i když prvky obou grup jsou vzájemně jednoznačně přiřazeny vztahem /5/. Uvedené veličiny a vztahy nám postačují k vyslovení (a podle /1/ i k důkazům) těchto důležitých tvrzení:

- 1) Vliv neurčitosti na data lze modelovat jako ortogonální otáčení vektoru v rovině (operátorem $K_q(\Omega)$ nad Minkowského rovinou nebo operátorem $K_e(\omega)$ nad Euklidovou rovinou).
- 2) Veličiny $K_q^2(\Omega) = K_q(2\Omega)$ a $K_e^2(\omega) = K_e(2\omega)$ se při takových transformacích chovají jako tensory.
- 3) Veličiny f, f^{-1}, h_e a h_q /7/-/9/, číselné charakteristiky nepodobnosti gnostických událostí sobě se shodují se složkami těchto tenzorů.

Platí tedy identity

$$K_q^2(\Omega) = \begin{pmatrix} f^{-1} & h_q \\ h_q & f \end{pmatrix} \quad K_e^2(\omega) = \begin{pmatrix} f & -h_e \\ -h_e & f \end{pmatrix} \quad /13/$$

a my máme důvod nazvat tyto veličiny tensory nepodobnosti. Vidíme, že charakterizují nekvalitu dat způsobenou neurčitostí.

Základy informace jednotlivé datové položky způsobené neurčitostí

Složky tenzorů nepodobnosti dále hrají významnou úlohu, stojí proto za to je pojmenovat:

- f^{-1} nevěrnost
- f věrnost
- h_q kvantifikační irrelevance
- h_e estimační irrelevance

Při kvantifikaci se kvantifikační irrelevance mění od 0 do hodnoty h_q . Vypočítáme střední hodnotu úhlu 2ω na tomto intervalu, veličinu

$$I_q = \int_0^{h_q} 2\omega dh_q \quad /14/$$

Hodnotě h_q dosažené na konci kvantifikace odpovídá hodnota estimační irrelevance h_e počátku estimace. Při estimaci se pak tato irrelevance bude měnit od h_e do nuly. Vypočítáme střední hodnotu úhlu 2Ω na tomto intervalu:

$$I_e = \int_{h_e}^0 2\Omega dh_e \quad /15/$$

Použijeme-li pro nějaké komplexní p ($p \neq 0, p \neq 1$) funkci

$$H_p = -p \ln p - (1-p) \ln(1-p) \quad /16/$$

kde pod $n(\cdot)$ rozumíme hlavní hodnotu a zavedeme-li označení

$$p_q = (1 + \sqrt{-1} h_q) / 2 \quad p_e = (1 + h_e) / 2 \quad /17/$$

dostaneme integrály /14/ a /15/ ve tvaru

$$I_q = H_{1/2} - H_{p_q} \quad /18/$$

$$I_e = H_{1/2} - H_{p_e} \quad /19/$$

Veličiny I_q a I_e interpretujeme jako změny informace v kvantifikační a estimační části ideálního gnostického cyklu. K takové interpretaci je několik závažných důvodů přesto, že jsme tyto veličiny odvodili pro jedinou položku dat a není dán žádný obvyklý pravděpodobnostní model:

- 1) Obě veličiny I_q i I_e jsme odvodili jako rozdíly funkce H mající formu shodnou se statistickou entropií termodynamického systému a tudíž i s Shannonovou informací binárního pravděpodobnostního systému.
- 2) Při vstřelování veličiny dat z od ideální hodnoty I_q klesá od nuly monotonně a při měření veličiny z od sašené (poserované) hodnoty I_e roste monotonně od nuly. Pro každou posero-

vanou hodnotu $z_1 \neq z_0$ však platí

$$I_q(z_1/z_0) + I_g(z_1/z_0) < 0 \quad (0 < z_1/z_0 < \infty) \quad /20/$$

cež interpretujeme jako nemožnost "napravit" estimací zcela ztrátu informace způsobenou vlivem neurčitosti, a to ani v případě ideálního gnostického cyklu.

- 3) Parametry funkcí H podobně jako pravděpodobnost v případě Shannonovy informace jsou veličiny p_q a p_g použitelné jako číselné charakteristiky očekávání, že neznámá veličina z_0 má hodnotu převyšující pozorovanou hodnotu.
- 4) Zřídla polí veličin I_q a I_g jsou úměrná sřídům polí změn termodynamické entropie. V gnostickém cyklu tedy dochází ke změně informace na úkor změny entropie a naopak. Hypotéza o existenci konverse entropie na informaci a naopak byla vyslovena velice dávno, zde máme matematický model této konverse pro obecný případ neurčitosti každé jednotlivé datové položky.

Po přijetí informační interpretace funkcí I_q a I_g lze přistoupit k formulaci významného tvrzení o optimalitě ideálního gnostického cyklu.

Ideální gnostický cyklus a jeho informační optimalita

Ideální hodnotě z_0 můžeme přiřadit dva vsájemně duální vektory $\underline{u}_0^T = (z_0, 0)$ a $\underline{u}_0 = (0, z_0)$. Byl-li výsledek skutečné kvantifikace z_g , jsou příslušné gnostické události v dvourozměrovém tvaru $\underline{u}_g^T = (x_g, y_g)$ a $\underline{u}_g = (y_g, x_g)$, přičemž složky x_g, y_g mají tvar /2/ při kvantifikaci s úhlem Ω_g a při estimaci tvar /3/ s úhlem ω_g a amplitudou $r_g = \sqrt{x_g^2 + y_g^2}$. Kvantifikační část dráhy ideálního gnostického cyklu (IGC) představuje oblouk Minkowského kružnice vycházející z bodu \underline{u}_0 a končící v bodě \underline{u}_g (s hlediska euklidovského zobrazení je to oblouk rovnosečné hyperboly). Z bodu \underline{u}_g pak IGC pokračuje podél oblouku obyčejné kružnice o poloměru r_g až do bodu $\underline{u}_g^{-T} = (r_g, 0)$ a do původního východního bodu \underline{u}_0 se cyklus uzavírá pohybem představujícím bodu podél úsečky $\underline{u}_g \underline{u}_0$. Ke každému takovému IGC graficky znázorněnému jako

$$\text{IGC} := \overset{\frown}{\underline{u}_0 \underline{u}_g} || \overset{\frown}{\underline{u}_g \underline{u}_g^{-T}} || \overset{\frown}{\underline{u}_g^{-T} \underline{u}_0} \quad /21/$$

existuje ovšem duální cyklus

$${}^c\text{IGC} := \overset{\frown}{\underline{u}_0 \underline{u}_g} || \overset{\frown}{\underline{u}_g \underline{u}_g^{-T}} || \overset{\frown}{\underline{u}_g^{-T} \underline{u}_0} \quad /22/$$

který při zobrazení v rovině (x, y) je s IGC symetrický vzhledem k ose kvadrantu. Třetí, uzavírající přímkovou část cyklů nazýváme atenuací, protože pohyb podél přímky procházející počátkem souřadnic lze představit jako důsledek změny měřítka či "zeslabení" veličin nějakým stenuátorem. Cyklus je tedy definován jako posloupnost tří na sebe navazujících myšlených pohybů v rovině:

$$\text{IGC} := \text{Kvantifikace} || \text{Estimace} || \text{Atenuace}$$

probíhajících po již stanovených drahách. Představme si však jiný uzavřený cyklus VGC vsaklý s IGC malou variací při safixovaných koncových bodech z_0 a z_g . Pro každý takový cyklus je v /1/ dokázáno, že platí

$$\oint_{\text{VGC}} dI \leq \oint_{\text{IGC}} dI < 0 \quad /23/$$

kde I značí postupně při vyčíslení integrálu I_q , I_g a I_a , přičemž I_a je změna informace při atenuaci (ta je nulová, protože změna měřítka nemění poměr $y(x)$). Analogická nerovnost platí i pro duální cykly ${}^c\text{IGC}$ a ${}^c\text{VGC}$. Odtud plyne, že IGC je optimální v tom smyslu, že minimalizuje ztrátu informace způsobené neurčitostí. V tomto smyslu je IGC vzorem pro praktické odhadování, k němuž se budeme snažit co nejvíce přiblížit.

GNOSTICKÁ TEORIE DATOVÉHO SOUBORU

Charakteristiky datového souboru by měly reprezentovat společné vlastnosti jednotlivých dat. Mají být určeny tak, aby jejich neurčitost byla menší než "průměrná" neurčitost jednotlivých dat. K tomu je třeba data složit podle nějakého kompozičního zákona. Prosté sečítání dat běžné při vytváření průměrů se neosvědčilo v případech, kdy se požaduje robustnost odhadu vůči odlehlým pozorováním. V robustní statistice

dostává každá položka dat při vytváření odhadu parametru polohy či měřítka svou individuální váhu závislou na vzdálenosti této položky od hodnoty parametru polohy. V gnostické teorii /2/ je taková individuální váha určena druhým axiomem této teorie, který stanovuje kompoziční zákon pro skládání dat.

Druhý axiom gnostické teorie

Necht $Z(z_0, n)$ je datový soubor sestavený z dat z_1, \dots, z_n majících tvar /1/. Charakteristikou souboru Z nazveme funkci všech těchto dat a z_0 . Za charakteristickou událost souboru Z přijmeme událost $u_0^T = (z_0 \operatorname{ch} \Omega_0, z_0 \operatorname{sh} \Omega_0)$ v případě kvantifikace a $u_0^T = (r_0 \cos \omega_0, -r_0 \sin \omega_0)$ v případě estimace, přičemž veličiny Ω_0, ω_0 jsou charakteristiky souboru Z .

Druhý axiom gnostické teorie říká, jak určit charakteristiky Ω_0 a ω_0 :

$$\underline{\kappa}_q^2(\Omega_c) = \frac{1}{w_q^2} \sum_1^n \underline{\kappa}_q^2(\Omega_i) \quad \underline{\kappa}_\omega^2(\omega_c) = \frac{1}{w_\omega^2} \sum_1^n \underline{\kappa}_\omega^2(\omega_i) \quad /24/$$

kde

$$w_q^2 = \operatorname{Det} \left(\sum_1^n \underline{\kappa}_q^2(\Omega_i) \right) \quad w_\omega^2 = \operatorname{Det} \left(\sum_1^n \underline{\kappa}_\omega^2(\omega_i) \right) \quad /25/$$

jsou normalizační váhy. Gnostický kompoziční zákon tedy skládá vlivy jednotlivých dat souboru silně nelineárně tak, aby tenzory nepodobnosti charakteristické události souboru byly normovanými součty tenzorů nepodobnosti jednotlivých dat. Tenzorový charakter sčítaných veličin zajišťuje formální správnost takového skládání, dále však uvidíme, že pro právě tento způsob skládání je i závažný přímý důvod - zajištění konsistence gnostické teorie s relativistickou fyzikou. Takový výběr axiomu skládání dat se ovšem podstatně liší od často heuristických a formálních důvodů vedoucích k pravidlům pro vážení dat přijímaným v mnoha desítkách metod robustní statistiky.

Důsledky druhého axiomu

Ze vztahů /17/ snadno získáme charakteristiky Ω_0 a ω_0 :

$$2\Omega_0 = \operatorname{arcth} \left(\bar{h}_q / \bar{r}^{-1} \right) \quad 2\omega_0 = \operatorname{arctg} \left(\bar{h}_\omega / \bar{r} \right) \quad /26/$$

přičemž symbol \bar{q} značí zde i dále aritmetický průměr hodnot q_1, \dots, q_n . Veličina Ω_0 má povahu charakteristické chyby dat hodnocené z hlediska kvantifikačního procesu, zatímco ω_0 odpovídá chybě dat měřené v estimační metrice. Lze ukázat, že při slabé neurčitosti dat (malé relativní chyby všech dat souboru) se všechny charakteristiky $\Omega_0, -\omega_0, \bar{h}_q$ i \bar{h}_ω liší od průměrné relativní chyby dat pouze o veličinu $O(\epsilon^2)$, kde ϵ je největší absolutní hodnota relativní chyby dat. Uvedené čtyři veličiny tedy představují různá gnostická zobecnění relativních chyb dat platná i pro případy silné neurčitosti. Jejich užitečnost je v tom, že se liší stupněm robustnosti vůči odlehlym pozorováním.

Ve vztazích /19/ vystupují i dvě další důležité veličiny, průměrná nevěrnost \bar{r}^{-1} a průměrná věrnost \bar{r} . Při výpočtu normalizačních vah podle /18/ se čtverce těchto veličin mohou rozepsat pomocí veličin \bar{r}^{-2} a \bar{r}^2 a pomocí průměrů součinů irrelevantní. Poslední z těchto veličin nabízí možnost zobecnění korelačních koeficientů, zatímco odchylky prvních čtyř veličin od 1 představují gnostická zobecnění výběrových rozptylů. Lze se přesvědčit, že všechny čtyři veličiny $(\bar{r}^{-1} - 1)/2, (1 - \bar{r})/2, (\bar{r}^{-2} - 1)/4$ a $(1 - \bar{r}^2)/4$ se při omezení všech relativních chyb hodnotou ϵ liší od relativní hodnoty výběrového rozptylu o veličinu $O(\epsilon^3)$. Všechny čtyři veličiny tedy jsou gnostickými zobecněními výběrového rozptylu. Lze ukázat, že se vzájemně liší stupněm robustnosti vůči odlehlym pozorováním.

Gnostická teorie datových souborů nám tedy nabízí řadu charakteristik různé robustnosti. Závislost těchto veličin na neznámé ideální hodnotě z_0 umožňuje formulovat a řešit problém odhadování ideální hodnoty jako úlohu extremalizace vybrané charakteristiky. V dalším příspěvku /4/ je uvedeno souhrnné řešení takového problému.

O KONSISTENCI GNOSTICKÉ TEORIE

Za konsistentní považujeme teorii odvozenou pomocí nějakého jazyka či kalkulu formálně správně z nerozporného a úplného systému axiomů. Takovou vlastnost "vnitřní" či vertikální konsistence zatím můžeme gnostické teorii přisoudit (dokud nebude dokázán opak). U teorie, která má poskytovat model nějakého výseku skutečnosti je však důležitá i jiná, "vnější" či "horizontální" konsistence: vzhledem k "nyní" nepochybným teoriím "sousedních" výseků skutečnosti. Lze ukázat /3/, že gnostická teorie je v souladu s teorií měření, s relativistickou fyzikou a s klasickou termodynamikou a že za určitých zvláštních podmínek, které umožňuje stanovit dává gnostická teorie výsledky konvergující k výsledkům klasické (nerobustní) statistické teorie a teorie informace. Lze rovněž ukázat, že tak rozsáhlou návaznost na jiné teorie existující statistická teorie neprokazuje.

Literatura:

- /1/ Kovanic P., Gnostical theory of individual data, v tisku
- /2/ Kovanic P., Gnostical theory of small samples of real data, v tisku
- /3/ Kovanic P., On relations between information and physics, v tisku (články vyjdou v časopise Problems of control and information theory, č. 4, 5 a 6, r. 1984)
- /4/ Kovanic P., Gnostické algoritmy zpracování dat, Sborník ROBUST '84, MFF-UK (1984)

GNOSTICKÉ ALGORITMY ZPRACOVÁNÍ DAT

Pavel Kovanic

SOUHRN

Příspěvek podává stručnou informaci o současném stavu algoritmizace robustních estimačních metod založených na gnostické teorii dat. Jde zejména o odhadování parametrů polohy jednorozměrových datových souborů, jejich měřítka, distribuční funkce a hustoty pravděpodobnosti. Souhrnně je popsán systém pro interaktivní gnostickou analýzu datových souborů na osobním počítači.

GNOSTICKÉ PARAMETRY POLOHY A MĚŘÍTKA

Gnostická teorie stručně vyložená v jiném příspěvku v tomto sborníku /1/ nepoužívá statistický model. Přesto však nabízí charakteristiky datových souborů blízké se svým smyslem tomu, k čemu statistika používá některé parametry svých modelů, k vyšetření polohy a variability datového souboru. Gnostický model dat doplněný oproti /1/ o parametr měřítka má tvar

$$z_i = z_0 \exp(s \Omega_i) \quad /1/$$

kde z_i je i -tá položka datového souboru $Z(n, z_0, s)$ obsahujícího n výsledků kvantifikace (tj. měření nebo čítání) kvantity, která by se při ideálních podmínkách kvantifikace zobrazila jako číslo z_0 , tzv. ideální hodnota. Veličina Ω_i charakterizuje vliv neurčitosti na i -tou položku dat a veličina " s " ^{měřítka} Především tyto dvě veličiny budeme odhadovat při zpracování datových souborů.

GNOSTICKÁ OPTIMALITA ODHADŮ POLOHY

V předchozím příspěvku /1/ je uvedena dvojice gnostických charakteristik E_g a E_e , které mají smysl sobecné "střední" chyby dat. Jsou to kvantifikační a estimační irrelevance, složky tensorů nepodobnosti. Tamtéž jsou uvedeny čtyři gnostická sobecná poměrného výběrového rozptylu určená odchylkami průměrné věrnosti \bar{F} , průměrné nevěrnosti \bar{f}^{-1} a průměrů čtverců těchto veličin od jedničky. Za nejlepší gnostický odhad J -tého typu můžeme přijmout veličinu s_j extremalizující některou ze šestí

uvedených charakteristik datového souboru $Z(n, z_0, s)$. Každý z takto získaných odhadů má důležitou interpretaci jak je vidět z Tab. 1:

Typ odhadu	Podmínka	Interpretace optimality odhadu
z_{qI}	$df^{-2}/dz_0=0$	Minimalizace intenzity ztráty informace při kvantifikaci
z_{qP}	$df^{-1}/dz_0=0$	Minimální zvýšení termodynamické entropie při kvantifikaci
z_{qS}	$\bar{H}_q=0$	1) Minimalizace celkové kvantifikační ztráty informace 2) Kvantifikační symetrie dat
z_{eS}	$\bar{H}_e=0$	1) Maximalizace celkového estimačního zvýšení informace 2) Estimační symetrie dat
z_{eP}	$d\bar{f}/dz_0=0$	Maximální snížení entropie při estimaci
z_{eI}	$d\bar{f}/dz_0=0$	1) Maximalizace intenzity zvyšování informace při estimaci 2) Maximalizace hustoty pravděpodobnosti

Tab. 1: Přehled optimalizačních podmínek pro stanovení gnostických odhadů parametru polohy a jejich teoretická interpretace. Symbol \bar{q} značí aritmetický průměr n -tice veličin q_1, \dots, q_n .

ROVNICE PRO ODHAD GNOSTICKÉHO PARAMETRU POLOHY

Dosažením do podmínek optimality dostáváme pro všech 6 případů z Tab. 1 jedinou rovnici

$$\sum_{i=1}^n ((z_0/z_i)^{2/s} - (z_i/z_0)^{2/s}) / ((z_0/z_i)^{2/s} + (z_i/z_0)^{2/s})^M = 0 \quad /2/$$

kde exponent M je určen Tab. 2:

Odhad	z_{qI}	z_{qP}	z_{qS}	z_{eS}	z_{eP}	z_{eI}
M	-1	0	0	1	2	3

Tab. 2: Exponent M pro rovnici /1/ odhadu parametru polohy z_0 .

Důsledkem volby exponentu M je stupeň a druh robustnosti odhadu (blíže viz /2/). Pro $M=3$ je odhad maximálně robustní a pro $M=-1$ citlivý k odlehlým pozorováním.

ODHAD MĚŘÍTKA A HUSTOTY DATOVÉHO SOUBORU

Odhad parametru měřítka je spjat s hustotou dat interpretovatelnou jako gnostický odhad hustoty pravděpodobnosti datového souboru. Ten je dán vzorcem

$$\frac{dP}{d} = \frac{1}{n} \sum_i 4 / ((z_i/z)^{2/s} + (z/z_i)^{2/s})^2 \quad /3/$$

kde má smysl průměrné pravděpodobnosti, že neznámá veličina z_0 je větší než z , podmíněné tím, že data mají hodnoty z_1, \dots, z_n . Známe-li parametr měřítka, můžeme tedy z dat sestavit odhad hustoty pravděpodobnosti, aniž bychom potřebovali statistický model. Vzorec /3/ můžeme však také použít k získání algoritmu pro odhadování parametru měřítka. Teoretický prostředek k tomu najdeme opět v gnostické teorii: za nejlepší odhad parametru měřítka přijmeme tu jeho hodnotu, která zaručuje shodu mezi průměrnou nepodobností dat a mezi střední nepodobností vypočtenou pomocí odhadnuté hustoty pravděpodobnosti jako funkci parametru měřítka.

SYSTEM PROGRAMŮ PRO GNOSTICKOU ANALÝZU DAT

V současné době již existuje programová realizace algoritmů vyplývajících z gnostické teorie, umožňujících detailní analýzu jednorozměrových datových souborů. Je vytvořena jako interaktivní systém programů v jazyce BASIC pro osobní počítač SINCLAIR ZX-81 s pamětí 16K. Má následující funkce:

- 1/ Vstup dat z klávesnice. Data mohou odpovídat buď multiplikativnímu nebo aditivnímu modelu vlivu neurčitosti.
- 2/ Výběr typu operace na datech
- 3/ Odhadování parametru polohy pro volitelnou robustnost (parametr M , viz Tab. 2)
- 4/ Odhadování parametru měřítka s volitelnou robustností
- 5/ Odhadování hranic datového souboru
- 6/ Odhadování pravděpodobnosti $P(z_0 > z)$
- 7/ Odhadování hustoty pravděpodobnosti dP/dz
- 8/ Testování hypotéz o homogenitě datového souboru a určování počtu shluků
- 9/ Určování vah jednotlivých dat
- 10/ Uspořádávání dat
- 11/ Odhadování základních statistických charakteristik datového souboru
- 12/ Numerický i grafický výstup výsledků

Tento programový systém byl vyzkoušen na různých aplikacích a prokázal výhodné vlastnosti gnostických algoritmů. Výsledky dávají naději, že touto cestou získají vědecká pracoviště účinný a flexibilní nástroj pro zkoumání vlastností malých datových souborů.

Literatura:

- /1/ Kovanic P., Základy gnostické teorie dat, Sborník ROBUST '84, MFF-UK (1984).