

## Adaptyvus glodinimo pločio parinkimas pasiskirstymo tankio statistiniame įvertyme

M. Kavaliauskas (KTU, VDU)

Paskutiniu metu, vystantis kompiuterinei technikai, suintensyvėjo sudėtingesnių statistinės analizės metodų plėtojimas. Statistinės analizės paketais pradeda naudotis ne tik matematikai, bet ir kitų sričių specialistai: medikai, sociologai, ekonomistai, etc., todėl atsiranda automatizuotų procedūrų poreikis. Statistinės procedūros dirbančios be žmogaus pagalbos (ivairių parametrų nustatymo ir pan.) yra ypač reikalingos, jei jos yra sudėtingesnio statistinės analizės paketo sudėtine dalimi. Šis darbas yra bandymas gauti pakankamai tikslų ir automatizuotą neparametrinį pasiskirstymo tankio įvertį, bei ištirti jo savybes Monte-Karlo būdu.

Nagrinėsime branduolinį pasiskirstymo tankio įvertį:

$$\hat{f}(x) = \hat{f}_h(x) = \frac{1}{n} \sum_{j=1}^n W(x - X_j), \quad (1)$$

čia  $W(x) = \frac{1}{h} \psi\left(\frac{x}{h}\right)$ ,  $h$  – branduolio plotis,  $\psi(x)$  – branduolio funkcija, tenkinanti sąlygą:

$$\int_{-\infty}^{\infty} \psi(x) dx = 1. \quad (2)$$

Šiame darbe modeliavimas buvo atliekamas naudojant Jepaničnikovo branduoli:

$$\psi(x) \begin{cases} \frac{3}{4}(1-x^2), & |x| \leq 1, \\ 0, & |x| > 1. \end{cases}$$

Vertinant pasiskirstymo tankio reikšmę taške  $x$ , natūralu parinkti glodinimo plotį atsižvelgiant ne tik į imties tūrį  $n$ , bet ir į tankio savybes taško  $x$  aplinkoje. Tuo būdu  $h = h(x, n)$ . Remsimės asimptotinėmis įverčių savybėmis, kai  $n \rightarrow \infty$ ,  $h = h(n) \rightarrow 0$ . Nuostolių funkcija laikysime  $E(f(x) - \hat{f}(x))^2$ . Jei apsiribosime įverčiais su neneigiamais lyginiais branduoliais, kaip paprastai daroma praktikoje, tai įverčio poslinkis dukart tolygiai diferencijuojamo tankio atveju lygus:

$$b_h(x) = c_1 \frac{f''(x)h^2}{2} + o(h^2), \quad (4)$$

kur

$$c_1 = \int y^2 \psi(y) dy. \quad (5)$$

Taip pat gerai žinoma asimptotinė įverčio dispersijos išraiška:

$$\sigma_h^2 = \frac{c_2 f(x) - h f^2(x)}{nh} + O\left(\frac{h^2}{n}\right) = \frac{c_2 f(x)}{nh} + O\left(\frac{1}{n}\right). \quad (6)$$

Čia

$$c_2 = \int \psi^2(y) dy. \quad (7)$$

Iš šių formulų seka, kad jei taške  $x$   $f''(x) \neq 0$ , tai asimptotiškai optimalus branduolio plotis  $h(x, n)$  yra lygus:

$$h(x, n) = \left( \frac{c_2 f(x)}{c_1^2 (f''(x))^2 n} \right)^{1/5}. \quad (8)$$

Betarpškai pasinaudoti (8) negalima, nes  $f(x)$  ir jos išvestinės nėra žinomos. Be to reikia atsižvelgti į galimą situaciją, kai  $f''(x)$  artima nuliui. Todėl konstruosime įvertį tokiu būdu:

$$h = h(x, n): \hat{b}_h^2(x) + \hat{\sigma}_h^2(x) \rightarrow \min. \quad (9)$$

*Pastaba.* Žymėjimų paprastumo sumetimais, tiek optimalų branduolio ploti, tiek jo įvertį žymėsime  $h$  tikėdamiesi, kad tai nejės painiavos.

Poslinkio ir dispersijos įverčius nusakykime lygybėmis:

$$\hat{b}_n(x) = \frac{c_1 \lambda(x, h) h^2}{2}, \quad (10)$$

kur  $\lambda(x, h) = \hat{f}_h''(x)$ ,

$$\hat{\sigma}_h^2(x) = \frac{c_2 \hat{f}_h(x)}{nh}. \quad (11)$$

I poslinkio įvertį įeinantis pasiskirstymo tankio antrosios išvestinės įvertis  $\lambda(x, h)$  randamas tokiu būdu. Pažymėkime:

$$Q(y) = F(x + y) - F(x - y) - \frac{y}{h} (F(x + h) - F(x - h)). \quad (12)$$

Skleisdami pasiskirstymo funkciją Teiloro eilute taško  $x$  aplinkoje gauname:

$$Q(y) = \frac{f''(x)(y^3 - yh^2)}{3} + o(yh^3). \quad (13)$$

Todėl  $\lambda(x, h)$  apibrėžkime taip:

$$\lambda = \lambda(x, h): \max_{0 \leq y \leq h} \left| \hat{Q}(y) - \frac{\lambda(y^3 - yh^2)}{3} \right| \rightarrow \min, \quad (14)$$

čia  $\hat{Q}(y)$  apibrėžiamas (12) lygybe, keičiant  $F(x)$  į empirinę pasiskirstymo funkciją  $\hat{F}(x)$ .

Tokiu būdu gauto pasiskirstymo tankio įverčio testavimas parodė, kad ji reikia pataisyti. Norint įvertinti pasiskirstymo tankio antrają išvestinę, reikia imti platesnę taško  $x$  aplinką, negu vertinant pačią tankio funkciją. Todėl įvertis pataisomas įvedant papildomą koeficientą  $c > 1$ :

$$\hat{b}_h(x) = \frac{c_1 \lambda(x, ch) h^2}{2}. \quad (15)$$

Pasiskirstymo tankio įvertis apibrėžiamas (1), (9), (11), (15) lygbybėmis buvo tiriamas Monte-Karlio būdu. Tyrimui naudojami Gauso skirstinių mišiniai, nes būtent mišinių atveju, kai tankio glodumo savybės stipriai skiriasi įvairiomis  $x$  reikšmėms, efektyvu naudoti branduolio plotį priklausomą nuo  $x$ . Be to, neparametrinis Gauso skirstinių tankio įvertis reikalingas netgi norint gauti geros kokybės parametrinius įverčius (žr. [2]).

Aprašytas įvertis buvo lygintas su:

**I.** Asimptotiskai optimaliu pseudoįverčiu, kurio parametras  $h$  gaunamas analogiškai kaip ir (8) formulėje, bet jis nepriklauso nuo  $x$ :

$$h = \left( \frac{c_2}{c_1^2 n \int (f''(x))^2 dx} \right)^{1/5}. \quad (16)$$

**II.** Statistiniu įverčiu, kurio  $h$  nepriklauso nuo  $x$ , bet parenkamas pagal imtį, minimizuojant nuostoliu  $E \|\hat{f} - f\|_2^2$  įvertį.

**III.** Adaptyviu pseudoįverčiu, kurio  $h$  optimaliai parenkamas pagal  $f(x)$  lokalias savybes taško  $x$  aplinkoje.

Trumpai aprašysime **II** ir **III** statistikų konstrukcijas.

**II įvertis.** Šis įvertis aprašytas [1].

Apskaičiuokime įverčio vidutinę kvadratinę paklaidą:

$$\|\hat{f}(x) - f(x)\|_2^2 = \int \hat{f}^2(x) dx - 2 \int \hat{f}(x) f(x) dx + \int f^2(x) dx. \quad (17)$$

Kadangi paskutinis narys nepriklauso nuo įverčio, tai apibrėžkime  $h$  tokiu būdu:

$$h: \int \hat{f}^2(x) dx - 2 \int \hat{f}(x) f(x) dx \rightarrow \min, \quad (18)$$

$$\begin{aligned} \int \hat{f}^2(x) dx &= \int \left( \frac{1}{nh} \sum_{i=1}^n \psi \left( \frac{x - X_i}{h} \right) \right)^2 dx = \frac{1}{n^2 h} \sum_{i,j=1}^n K \left( \frac{X_i - X_j}{h} \right) \\ &= \frac{2}{n^2 h} \sum_{i=1}^n \sum_{j=i+1}^n K \left( \frac{X_i - X_j}{h} \right) + \frac{K(0)}{nh}, \end{aligned} \quad (19)$$

čia  $K(x)$  yra branduolio  $\psi(x)$  sąsuka su pačiu savimi.

$$\int \hat{f}(x) f(x) dx \approx \int \hat{f}(x) dF_n(x) \approx \frac{1}{n^2 h} \sum_{i=1}^n \sum_{j=1, j \neq i}^n \psi\left(\frac{X_i - X_j}{h}\right). \quad (20)$$

Pastebėsime, kad dėmenų su  $j = i$  išmetimas sumažina poslinkį.

Taigi iš (19) ir (20) gauname:

$$h: \frac{2}{n^2 h} \sum_{i=1}^n \sum_{j=i+1}^n \left( K\left(\frac{X_i - X_j}{h}\right) - 2\psi\left(\frac{X_i - X_j}{h}\right) \right) + \frac{K(0)}{nh} \rightarrow \min. \quad (21)$$

**III pseudojvertis.** Natūralu pseudojvertį apibrėžti naudojant tikrąjį poslinkį ir dispersiją

$$h = h(x, n): b_h^2(x) + \sigma_h^2(x) \rightarrow \min, \quad (22)$$

tačiau taip apibrėžtas pseudojvertis priklauso ne tik nuo lokalių tankio savybių taško  $x$  aplinkoje, bet ir nuo tankio reikšmių daug nutolusių nuo  $x$ . Todėl pseudojvertis modifikuotas tokiu būdu. Kadangi:

$$\begin{aligned} b_h(x) &= \int_{-1}^1 \psi(y) dy \int_0^{hy} f'(x+z) dz = \int_0^1 \psi(y) dy \int_0^{hy} dz \int_{-z}^z f''(x+u) du \\ &= \int_{-h}^h f''(x+u) \left( \int_{|u|}^h (z-|u|) W(z) dz \right) du \stackrel{\text{def}}{=} \int_{-h}^h f''(x+u) G_h(u) du, \end{aligned} \quad (23)$$

tai pseudojvertį  $h(x)$  apibrėžkime:

$$h(x): (b_h^*(x))^2 + \sigma_h^2(x) \rightarrow \min, \quad (24)$$

kur

$$|b_h^*(x)| = \int_{-h}^h |f''(x+y)| G_h(y) dy. \quad (25)$$

**Tyrimo schema.** Tyrimas atliktas Monte-Karlio metodu. Buvo generuojami vienmačių Gauso skirstinių mišinių su žinomais parametrais imtys. Jų tankis yra:

$$f(x) = \sum_{i=1}^k p_i \varphi(x, m_i, \sigma_i). \quad (26)$$

Gautomis imtims skaičiuojami aukšciau minėti pasiskirstymo tankio funkcijos įverčiai ir pseudojverčiai. Šie įverčiai lyginami su teorine pasiskirstymo tankio funkcija.

Įverčio tikslumas vertintas  $L_1$ ,  $L_2$  ir  $L_\infty$  metrikose:

$$\begin{aligned}\varepsilon_1 &= \int |\hat{f}(x) - f(x)| dx, \\ \varepsilon_2 &= \int (\hat{f}(x) - f(x))^2 dx, \\ \varepsilon_\infty &= \max_x |\hat{f}(x) - f(x)|.\end{aligned}\quad (27)$$

Eksperimentų rezultatų lentelė.

Gauso mišinių sudėtis:

- 1)  $n = 500$ ;  $p_1 = 0.7$ ,  $p_2 = 0.3$ ;  $m_1 = 0$ ,  $m_2 = 3$ ;  $\sigma_1 = 0.4$ ,  $\sigma_2 = 1$ ;
- 2)  $n = 2000$ ;  $p_1 = 0.7$ ,  $p_2 = 0.3$ ;  $m_1 = 0$ ,  $m_2 = 3$ ;  $\sigma_1 = 0.4$ ,  $\sigma_2 = 1$ ;
- 3)  $n = 2000$ ;  $p_1 = 0.7$ ,  $p_2 = 0.3$ ;  $m_1 = 0$ ,  $m_2 = 6$ ;  $\sigma_1 = 0.4$ ,  $\sigma_2 = 2$ .

Gauso mišiniai	Metrika	Siūlomas įvertis	Asimptotinis pseudoįvertis (I)	Fiksuoto pločio įvertis (II)	Optimalus pseudoįvertis (III)
1	$L_1$	0.130268	0.132292	0.093799	0.094259
	$L_2$	0.065963	0.075425	0.048306	0.049532
	$L_\infty$	0.095950	0.086234	0.074546	0.051814
2	$L_1$	0.065705	0.077963	0.059230	0.055269
	$L_2$	0.033966	0.041728	0.030135	0.029183
	$L_\infty$	0.037409	0.084537	0.037334	0.035365
3	$L_1$	0.064688	0.100593	0.081787	0.056844
	$L_2$	0.032094	0.042985	0.032597	0.028600
	$L_\infty$	0.038153	0.081336	0.040551	0.036456

Preliminarūs tyrimo rezultatai rodo, kad siūlomas įvertis duoda mažesnes paklaidas negu asimptotinis pseudoįvertis I, tačiau nevisuomet vertina geriau negu fiksuooto pločio įvertis II. Iš pateiktos lentelės matome, kad esant mažam imties kiekiui, geresnis yra įvertis II, didesnės imties atveju (antras mišinys) įverčio ir fiksuooto pločio įverčio II kokybės yra panašios, tačiau pakeitus mišinio formą taip, kad glodumas įvairiose tankio dalyse dar labiau skirtuosi (trečias mišinys), geresnius rezultatus duoda įvertis apibrėžiamas (1), (9), (11), (15).

#### Išvados.

Įverčių palyginamoji analizė parodė, kad įverti galima tobulinti. Vienas iš neišspręstų klausimų yra konstantos  $c$  įeinančios į (15) išraišką parinkimas. Bandymai rodo, kad  $c$  turėtų būti  $x$  funkcija, nes esant mažiemis  $h$  nepakanka dvigubo pločio antrajai pasiskirstymo tankio išvestiniai įvertinti. Taip pat pastebėtas minimumo apibrėžiamo (9) nestabilumas. Esant mažiemis  $h$  funkcija  $\hat{b}_h^2(x) + \hat{\sigma}_h^2(x)$  yra „šokinėjanti“, o pernelyg dideliems  $h$  galima rasti kitų šios funkcijos minimumų, todėl minimumo reikia ieškoti optimaliai parinktame intervale  $[h_{min}; h_{max}]$ . Manoma, kad įvertis geriau vertins tankį, jei  $h$  įverti lokalai suglodinsime.

Iverčio savybių išsamesnio tyrimo modeliavimo būdu rezultatai bus pateikti kitame straipsnyje.

Dėkoju prof. R.Rudzkiui už pagalbą rašant šį straipsnį.

## LITERATŪRA

- [1] B. W. Silverman, *Density Estimation for Statistics Data Analysis*, Chapman and Hall, 1986. ISBN 0-412-24620-1.
- [2] M. Radavičius and R. Rudzkis, Statistical estimation of a mixture of Gaussian distribution, *Acta Applicandae Mathematicae*, **38** (1995), 37–54.

### An adaptive choice of smoothing width in statistical estimation of distribution density

*M. Kavaliauskas*

This article illustrates the adaptive kernel probability density function estimation method, shows encountered problems, gives solutions and suggestions. The simulation results are discussed.