

Kvantilio vertinimas baigtinėje populiacijoje

Danutė KRAPAVICKAITĖ (MII, VGTU)

el. paštas: krapav@ktl.mii.lt

1. Įvadas

Nagrinėjama baigtinė populiacija $\mathcal{U} = \{1, 2, \dots, N\}$ su tyrimo kintamuoju y , kurio reikšmės $\{y_1, y_2, \dots, y_N\}$ nežinomos ir gali būti sužinotos imties elementams. Populiacijos parametras $\theta = \theta(y_1, y_2, \dots, y_N)$ taip pat nežinomas. Jį reikia ivertinti. Tuo tikslu renkama tikimybinė imtis $i \subset \mathcal{U}$, kurios elementus paprastumo dėlei žymėsime tiesiog $i = \{1, 2, \dots, n\}$. Dažniausiai vertinamas populiacijos parametras yra suma

$$t = t_y = y_1 + y_2 + \dots + y_N.$$

Norint ją ivertinti, turint tikimybinę imtį i , galima pasinaudoti Horvitz–Thompson sumos ivertiniu ([2]):

$$\hat{t}_y = \sum_{k \in i} \frac{y_k}{\pi_k}.$$

Jis yra nepaslinktasis. Jo dispersija

$$D\hat{t}_y = \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}} (\pi_{kl} - \pi_k \pi_l) \frac{y_k y_l}{\pi_k \pi_l}$$

gali būti vertinama nepaslinktuoju ivertiniu

$$\widehat{D}\hat{t}_y = \sum_{k \in i} \sum_{l \in i} \left(1 - \frac{\pi_k \pi_l}{\pi_{kl}}\right) \frac{y_k y_l}{\pi_k \pi_l},$$

pažymėjus elementų priklausymo imčiai tikimybes $\pi_k = P\{k \in i\}$, $\pi_{kl} = P\{k \in i, l \in i\}$.

Literatūroje retai nagrinėjamas, bet praktikoje naudojamas populiacijos parametras yra tyrimo kintamojo y pasiskirstymo F bet kurio q lygio, $0 < q < 1$, kvantilis $Q(q)$. Atskiras kvantilio – medianos ($q = 1/2$) ir jos pasikliautinojo intervalo vertinimo atvejis išnagrinėtas [2]. Kai kurios kvantilių funkcijų vertinimo galimybės nagrinėtos [1].

Šiame darbe siūlomas kvantilio – dviejų sumų santykio ivertinys ir kvantilio pasikliautinojo intervalo ivertinys. Šie ivertiniai gali būti taikomi Lietuvos oficialiosios statistikos vykdomame namų ūkių biudžetų tyriame.

2. Kvantilis

Kvantiliu vadinamas toks skaičius $Q(q)$, $0 < q < 1$, kuris tenkina nelygybes

$$F(Q(q)) \leq q \leq F(Q(q) + 0).$$

Kvantilis gali būti lygties

$$F(Q(q)) = q \tag{1}$$

sprendiniu. Tačiau tyrimo kintamojo reikšmių pasiskirstymo funkcija baigtinėje populiacijoje yra laiptuota, ir (1) lygtis turi arba be galio daug sprendinių, arba jų visai neturi. Tyrimo kintamojo reikšmės populiacijoje surūšiuojamos didėjimo tvarka: $y_1^* \leq y_2^* \leq \dots \leq y_N^*$. Jei egzistuoja toks $k_0 \in \mathcal{U}$, kad $F(y_{k_0}^*) = q$, tai kvantilis $Q(q) = (y_{k_0+1}^* + y_{k_0}^*)/2$. Jei $F(y_k^*) \neq q$ su visais $k \in \mathcal{U}$, tai atsiras toks k_0 , kuriam $F(y_{k_0-1}^*) < q$, tačiau $F(y_{k_0}^*) > q$. Tada kvantilis $Q(q) = y_{k_0}^*$.

Norint gauti kvantilio įvertinį $\widehat{Q}(q)$, reikia nagrinėti populiacijos pasiskirstymo funkcijos $F(x)$ įvertinį $\widehat{F}(x)$, tenkinantį nelygybes

$$\widehat{F}(\widehat{Q}(q)) \leq q \leq \widehat{F}(\widehat{Q}(q) + 0).$$

Imties reikšmės rūšiuojamos didėjimo tvarka: $y_1^* \leq y_2^* \leq \dots \leq y_n^*$. Kvantilio įvertinys bus

$$\widehat{Q}(q) = \begin{cases} (y_{k_0+1}^* + y_{k_0}^*)/2, & \text{jei yra } k_0 \in \mathbf{i}: \widehat{F}(y_{k_0}^*) = q, \\ y_{k_0}^*, & \text{jei } \widehat{F}(y_{k_0-1}^*) < q, \widehat{F}(y_{k_0}^*) > q \end{cases}$$

kokiam nors $k_0 \in \mathbf{i}$.

Pasiskirstymo funkciją galima užrašyti

$$F(x) = \#A_x/N, \tag{2}$$

čia $A_x = \{k \in \mathcal{U}: y_k \leq x\}$, ir $\#A_x$ reiškia aibės A_x elementų skaičių. Norint įvertinti kvantili, reikia įvertinti pasiskirstymo funkciją, o ji gali būti užrašyta ir kitaip. Iveskime papildomą kintamąjį $z(x) = \{z_1(x), \dots, z_N(x)\}$, čia

$$z_k(x) = \begin{cases} 1, & \text{jei } y_k \leq x, \\ 0, & \text{jei } y_k > x, \end{cases}$$

$k = 1, 2, \dots, N$. Tada

$$F(x) = t_{z(x)}/N, \tag{3}$$

ir pasiskirstymo funkcijos (2) vertinimo uždavinyssusiveda į dviejų sumų santykio vertinimo uždavini:

$$\widehat{F}(x) = \frac{\widehat{t}_{z(x)}}{\widehat{N}}, \quad \widehat{t}_{z(x)} = \sum_{k \in i} \frac{z_k(x)}{\pi_k}, \quad \widehat{N} = \sum_{k \in i} \frac{1}{\pi_k}.$$

3. Kvantilio pasikliautinasis intervalas

Radus konstantas $c_1 > 0$ ir $c_2 > 0$ tokias, kad

$$P\{c_1 \leq \widehat{F}(Q(q)) \leq c_2\} \geq 1 - \alpha,$$

intervalas $(\widehat{F}^{-1}(c_1), \widehat{F}^{-1}(c_2))$ bus α lygio ($0 < \alpha < 1$) pasikliautinasis intervalas kvantiliui $Q(q)$. Jei santykio įvertinys $\widehat{F}(Q(q))$ turi normalujį skirstinį su vidurkiu q , tai

$$\begin{aligned} c_1 &= q - Q_{N(0,1)}(\alpha/2) \sqrt{\mathbf{D}\widehat{F}(Q(q))}, \\ c_2 &= q + Q_{N(0,1)}(\alpha/2) \sqrt{\mathbf{D}\widehat{F}(Q(q))}, \end{aligned} \tag{4}$$

čia $Q_{N(0,1)}(\alpha/2)$ yra standartinio normaliojo skirstinio $1 - \alpha/2$ lygio kvantilis.

$\mathbf{D}\widehat{F}(Q(q))$ apytikslę dispersiją užrašome, naudojantis santykio įvertinio apytikslės dispersijos formule, gaunama, skleidžiant $\mathbf{D}\widehat{F}(Q(q))$ Teiloro eilute:

$$\mathbf{A}\mathbf{D}\widehat{F}(Q(q)) = \frac{1}{N^2} \mathbf{D}(\widehat{t}_{z(Q(q))} - q\widehat{N}). \tag{5}$$

Praktiškai šios apytikslės dispersijos panaudoti negalime, nes ji priklauso nuo visų populiacijos reikšmių, todėl ją dar reikia įvertinti. Iverčiu naudojama

$$\widehat{\mathbf{D}}\widehat{F}(Q(q)) = \frac{1}{N^2} \widehat{\mathbf{D}}(\widehat{t}_{z(\widehat{Q}(q))} - q\widehat{N}).$$

Tokiu būdu iš (4) gauname q lygio kvantilio $Q(q)$ pasikliautinąjį intervalą (α lygio)

$$\begin{aligned} &\left(\widehat{F}^{-1}\left(q - Q_{N(0,1)}(\alpha/2) \sqrt{\widehat{\mathbf{D}}\widehat{F}(Q(q))}\right), \right. \\ &\left. \widehat{F}^{-1}\left(q + Q_{N(0,1)}(\alpha/2) \sqrt{\widehat{\mathbf{D}}\widehat{F}(Q(q))}\right) \right). \end{aligned} \tag{6}$$

4. Pavyzdys

Pajamų vienam namų ūkio nariui pasiskirstymo kvantilių ir jų pasikliautinujų intervalų vertinimas namų ūkių biudžetu tyrime

Šiame tyriame naudojama sluoksninė imtis su dvieju etapu lizdine imtimi sluoksniuose. Pirmajį sluoksnį sudaro trijų didžiųjų miestų: Vilniaus, Kauno, Klaipėdos, Šiaulių ir Panevėžio namų ūkiai. Antrajam priklauso vidutinių ir mažųjų miestų gyventojų namų ūkiai, ir trečiam – kaimo gyventojų namų ūkiai. Tarkime, kad visi gyventojai yra ištraukti iš gyventojų registrą. 10 680 namų ūkių imtis paskirstoma į sluoksnius proporcingai jų dydžiui. Imčių planai sluoksniuose yra tokie:

1-asis sluoksnis. Iš gyventojų registro renkama paprastoji atsitiktinė 4 476 asmenų imtis. Tiriami išrinktuju asmenų namų ūkiai.

2-asis sluoksnis. Vidutinių ir mažųjų miestelių gyventojai apjungiami į lizdus, turinčius nuo 1 000 iki 10 000 gyventojų. 1-ajame etape su tikimybėmis, proporcingomis lizdo dydžiui, renkama grąžintinė 20 lizdų imtis iš 140 populiacijoje esančių lizdų. 2-ajame etape kiekviename iš išrinktuju lizdų iš gyventojų registro renkama paprastoji atsitiktinė 132 asmenų imtis. Jų namų ūkiai laikomi išrinktais į imtį.

3-asis sluoksnis. Suformuojami 463 lizdai, kuriuose yra nuo 300 iki 2 000 namų ūkių. 1-ajame etape su tikimybėmis, proporcingomis lizdo dydžiui, išrenkama grąžintinė 33 lizdų imtis. 2-ajame etape išrinktuosiuose lizduose iš gyventojų registro renkama paprastoji atsitiktinė 108 asmenų imtis. Jų namų ūkiai laikomi išrinktais į imtį.

Tyrimo kintamasis y – asmens pajamos. Jų pasiskirstymo funkcija $F(x)$, pasinaudojus (3), gali būti užrašyta tokiu būdu:

$$F(x) = \frac{t_{z(x)}}{N} = \sum_{h=1}^3 \frac{N_h}{N} \frac{t_{z(x)}^{(h)}}{N_h} = \sum_{h=1}^3 \frac{N_h}{N} F_h(x),$$

čia N – populiacijos dydis ēmimo sąraše (gyventojų registre), N_h – sluoksnio dydis ēmimo sąraše, $F_h(x)$ – pajamų pasiskirstymo funkcija h -ajame sluoksnuje, $t_{z(x)}^{(h)}$ – kintamojo $z(x)$ reikšmių suma h -ajame sluoksnuje. Kvantislis vertinamas, iprastiniu būdu nagrinėjant funkciją $F(x)$.

Kvantislis, būdamas (1) lygties sprendiniu, tuo pačiu yra lygties

$$\sum_{h=1}^3 \frac{N_h}{N} \widehat{F}_h(Q(q)) = q$$

sprendiniu. Norint rasti kvantilio pasikliautinająjį intervalą (6), reikia skaičiuoti funkcijos $\widehat{F}(Q(q)) = \sum_{h=1}^3 \frac{N_h}{N} \widehat{F}_h(Q(q))$ apytiksle dispersiją

$$AD\widehat{F}(Q(q)) = AD \left(\sum_{h=1}^3 \frac{N_h}{N} \widehat{F}_h(Q(q)) \right) = \sum_{h=1}^3 \frac{N_h^2}{N^2} AD\widehat{F}_h(Q(q)). \quad (7)$$

Apytikslės dispersijos $AD\widehat{F}(Q(q))$ skaičiavimas

1-ajame sluoksnje elemento tikimybė priklausyti imčiai $\pi_{1k} = n_1 m_{1k} / N_1$, čia n_1 – 1-ojo sluoksnio imties dydis, m_k – k-ojo namų ūkio registruotų asmenų skaičius, $k \in \mathcal{U}_1$, \mathcal{U}_1 – 1-ojo sluoksnio gyventojų populiacija. Tikimybė, kad k-sis ir l-sis namų ūkiai kartu priklausys imčiai, yra

$$\pi_{1kl} = \frac{n_1(n_1 - 1)m_{1k}m_{1l}}{N_1(N_1 - 1)}.$$

Tada, pasinaudojus Horvitz–Thompson sumos įvertiniu ir (5), gauname

$$\begin{aligned} AD\widehat{F}_1(Q(q)) &= \frac{1}{N_1^2} D(\widehat{t}_{z(Q(q))}^{(1)} - F_1(Q(q))\widehat{N}_1) \\ &= \frac{1}{N_1^2} \sum_{k \in \mathcal{U}_1} \sum_{l \in \mathcal{U}_1} (\pi_{1kl} - \pi_{1k}\pi_{1l}) \frac{z_{1k}(Q(q)) - F_1(Q(q))}{\pi_{1k}} \cdot \frac{z_{1l}(Q(q)) - F_1(Q(q))}{\pi_{1l}} \\ &= -\frac{1}{N_1^2} \left(1 + \frac{1}{N_1 - 1}\right) \left(1 - \frac{n_1}{N_1}\right) \frac{1}{n_1} \sum_{k \in \mathcal{U}_1} \sum_{l \in \mathcal{U}_1} u_{1k}u_{1l}, \end{aligned} \quad (8)$$

čia $u_{1k} = z_k(Q(q)) - F_1(Q(q))$ yra kintamojo $u_1 = z(Q(q)) - F_1(Q(q))$ reikšmė.

2-jame ir 3-jame sluoksniuose turima dvių etapų lizdinė imtis. I etape – grąžintinė lizdinė imtis su lizdų išrinkimo tikimybėmis $p_{hi} = L_{hi}/\sum_{j \in \mathcal{U}_h} L_{hj}$, imties dydis n_h , todėl lizdo priklausymo imčiai tikimybė

$$\pi_{hi} = \frac{n_h L_{hi}}{\sum_{j \in \mathcal{U}_h} L_{hj}},$$

\mathcal{U}_h – lizdų populiacija h-ajame sluoksnje, $h = 2, 3$, L_{hi} – registruotų asmenų skaičius lizde.

II etapo lizdas – namų ūkis. h-ojo sluoksnio i-ojo lizdo k-ojo namų ūkio tikimybė priklausyti imčiai yra

$$\pi_{hik} = \frac{\tilde{n}_{hi} m_{hik}}{L_{hi}},$$

čia $h = 2, 3$, $i \in \mathcal{U}_h$, $h = 2, 3$, \tilde{n}_{hi} – II etapo imties dydis h-ojo sluoksnio i-ajame lizde, $k = 1, 2, \dots, m_{hik}$, m_{hik} – asmenų skaičius i-ojo lizdo k-ajame namų ūkyje.

Pažymėjus \mathbf{i}_h I etapo imtį h-ajame sluoksnje, parametru θ įvertiniui $\widehat{\theta}$ galioja lygybė

$$D\widehat{\theta} = D(E(\widehat{\theta}|\mathbf{i}_h)) + E(D(\widehat{\theta}|\mathbf{i}_h)).$$

Ja naudojamas, skaičiuojant dispersiją 2-ajame ir 3-ajame sluoksniuose:

$$\begin{aligned} D(\widehat{t}_{z(Q(q))}^{(h)} - F_h(Q(q))\widehat{N}_h) \\ = D(E(\widehat{t}_{z(Q(q))}^{(h)} - F_h(Q(q))\widehat{N}_h|\mathbf{i}_h)) + E(D(\widehat{t}_{z(Q(q))}^{(h)} - F_h(Q(q))\widehat{N}_h|\mathbf{i}_h)). \end{aligned}$$

2-ojo ir 3-ojo sluoksnį I etape imti laikysime grąžtine, išrinkta su tikimybėmis, proporciniomis lizdo dydžiui. Tada sumai vertinti I etape galima naudotis Hansen–Hurwitz sumos įvertiniu ([2]). II etape galima naudotis tokiu pat Horvitz–Thompson sumos įvertiniu, kaip ir 1-ajame sluoksnje. Sukombinavus šiuos įvertinius, gaunama pasiskirstymo funkcijos įvertinio apytiksles dispersijos sluoksnje išraiška

$$\begin{aligned} ADF_h(Q(q)) &= \frac{1}{N_h^2} D(\hat{t}_{zQ(q)}^{(h)} - F_h(Q(q))\hat{N}_h) \\ &= \frac{1}{N_h^2 n_h} \sum_{i \in \mathcal{U}_h} \left[p_{hi} \left(\frac{\hat{t}_{ui}^{(h)}}{p_{hi}} - t_u^{(h)} \right)^2 \right. \\ &\quad \left. - \frac{1}{\tilde{n}_{hi} L_{hi}^2 p_{hi}} \left(1 - \frac{\tilde{n}_{hi}}{L_{hi}} \right) \left(1 + \frac{1}{L_{hi} - 1} \right) \sum_{k \in \mathcal{U}_{hi}} \sum_{l \in \mathcal{U}_{hi}} u_{hik} u_{hil} \right], \quad (9) \end{aligned}$$

$p_{hi} = L_{hi}/\sum_{j \in U_h} L_{hj}$, $t_u^{(h)} = \sum_{i,k} u_{hik}$, $u_{hik} = z_{hik}(Q(q)) - F_h(Q(q))$, $\hat{t}_{ui}^{(h)} = \sum_{k \in \mathcal{U}_{hi}} u_{hik}/\pi_{hik}$, $h = 2, 3$, \mathbf{i}_{hi} – imtis h -ojo sluoksnio i -ajame lizde. \mathcal{U}_{hi} – h -ojo sluoksnio i -ojo lizdo namų ūkių populiacija, z_{hik} – kintamojo $z(x)$ reikšmė h -ojo sluoksnio i -ojo lizdo k -ajam asmeniui.

Apytiksle pasiskirstymo funkcijos įvertinio dispersija, ištačius dispersiją išraiškas sluoksniuose (8) ir (9) i (7), tampa

$$\begin{aligned} ADF(Q(q)) &= \frac{1}{N^2} \left\{ -\frac{1}{N_1^2} \left(1 + \frac{1}{N_1 - 1} \right) \left(1 - \frac{n_1}{N_1} \right) \frac{1}{n_1} \sum_{k \in \mathcal{U}_1} \sum_{l \in \mathcal{U}_1} u_{1k} u_{1l} \right. \\ &\quad + \sum_{h=2}^3 \frac{1}{n_h} \sum_{i \in \mathcal{U}_h} \left[p_{hi} \left(\frac{\hat{t}_{uh_i}}{p_{hi}} - t_{uh_i} \right)^2 \right. \\ &\quad \left. \left. - \frac{1}{\tilde{n}_{hi} L_{hi}^2 p_{hi}} \left(1 - \frac{\tilde{n}_{hi}}{L_{hi}} \right) \left(1 + \frac{1}{L_{hi} - 1} \right) \sum_{k \in \mathcal{U}_{hi}} \sum_{l \in \mathcal{U}_{hi}} u_{hik} u_{hil} \right] \right\}. \end{aligned}$$

Pasiskirstymo funkcijos įvertinio dispersijos įvertinys

$$\begin{aligned} \widehat{D}\hat{F}(Q(q)) &= \frac{1}{N^2} \left\{ -\frac{1}{N_1^2} \left(1 + \frac{1}{N_1 - 1} \right) \left(1 - \frac{n_1}{N_1} \right) \frac{1}{n_1} \sum_{k \in \mathbf{i}_1} \sum_{l \in \mathbf{i}_1} \frac{1}{\pi_{kl}} \widehat{u}_{1k} \widehat{u}_{1l} \right. \\ &\quad + \sum_{h=2}^3 \frac{1}{n_h(n_h - 1)} \sum_{i \in \mathbf{i}_h} \left[\left(\frac{\widehat{t}_{uh_i}^{(h)}}{p_{hi}} - \widehat{t}_u^{(h)} \right)^2 \right. \\ &\quad \left. \left. - \frac{1}{\tilde{n}_{hi} L_{hi}^2 p_{hi}} \left(1 - \frac{\tilde{n}_{hi}}{L_{hi}} \right) \left(1 + \frac{1}{L_{hi} - 1} \right) \sum_{k \in \mathcal{U}_{hi}} \sum_{l \in \mathcal{U}_{hi}} \frac{1}{\pi_{kl}} \widehat{u}_{hik} \widehat{u}_{hil} \right] \right\}, \end{aligned}$$

$\widehat{u}_{hik} = z_{hik}(\widehat{Q}(q)) - \widehat{F}_h(\widehat{Q}(q))$, $h = 2, 3$, $\widehat{u}_{1k} = z_k(\widehat{Q}(q)) - \widehat{F}_1(\widehat{Q}(q))$, išstatytas i (5), duoda kvantilio apytiksli pasikliautinaji intervalą.

Autorė dėkoja VGTU studentei Jurgitai Čimžaitei, dalyvavusiai ruošiant ši darbą.

Literatūra

- [1] T. Orusild, Confidence intervals for functions of quantiles under finite population sampling, *Central Statistical Bureau of Latvia, Workshop for the Baltic Countries on Survey Sampling Theory and Methodology*, Jurmala, Latvia (1998), pp. 51–55.
- [2] C.-E. Sarndal, B. Swensson, J. Wretman, *Model Assisted Survey Sampling*, Springer-Verlag, New-York (1992).

Estimation of a quantile in finite population

D. Krapavickaitė

An estimator of the quantile of a study variable distribution function in finite population and its confidence interval is proposed. The estimator is adapted to the two-stage cluster sample which is used in a household budget survey carried out in Lithuania.