

Apie tiesinės regresijos identifikavimo metodus

Romualdas SALĖTIS* (VGTU, MII)

el. paštas: romas@goda.vtu.lt

1. Įvadas

Mažiausių kvadratų metodas yra populiariausias taikomasis statistinis metodas identifikuojant tiesinės regresijos modelius. Galiojant tam tikroms sąlygoms, mažiausių kvadratų metodu gautas įvertis yra efektyvus (Gauso-Markovo teorema). Tačiau, praktikoje jos nėra pilnai patenkinamos, todėl yra tikslinga naudoti robastines įvertinimo procedūras, kurios yra nejautrios nedideliems nuokrypiams nuo prielaidų apie modelį, o jų pagalba gautas įvertis yra beveik efektyvus.

Šiame darbe yra tiriami mažiausių kvadratų ir medianinio įvertinimo metodai bei įvairios jų modifikacijos, kurios yra nejautrios nedidelės dalies duomenų užteršimui grubiai klaidingais duomenimis – išsiskiriančiais stebėjimais, bei kurie įvertina ir eliminuoja modelio liekanų koreliuotumo ir heteroskedastiškumo (dispersijos priklausomybės nuo laiko) įtaką.

Siekiant palyginti nagrinėjamų metodų įvertinimo tikslumą, esant aukščiau išvardintoms situacijoms, imitacinio modeliavimo būdu buvo atlikta serija eksperimentų, kurių aprašymas ir rezultatai yra pateikiami šio darbo eksperimentinėje dalyje.

2. Tiesinės regresijos modelis ir jo identifikavimo metodai

Tarkime, kad laike stebime atsitiktinius dydžius $Y(t)$, $t = 1, 2, \dots, n$, kurie yra aprašomi regresijos lygtimi:

$$Y(t) = \boldsymbol{\theta}^T \mathbf{x}(t) + \xi(t), \quad t = 1, 2, \dots, n. \quad (1)$$

čia: $\boldsymbol{\theta}$ – nežinomas daugiamatis parametras, $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)^T$;

$\mathbf{x}(t)$ – determinuotas vektorius, $\mathbf{x}(t) = (x_1(t), \dots, x_d(t))^T$, kur $x_i(t)$, $i = 1, 2, \dots, d$ – žinomo pavidalo funkcijos;

$\xi(t)$ – atsitiktiniai dydžiai, modelyje vadinami liekanomis.

Tiesinės regresijos modelyje (1) dažnai yra priimanamos šios prielaidos.

1. Liekanų vidurkis yra lygus nuliui: $\mathbf{E}\xi(t) = 0$.
2. Liekanų $\xi(t)$ pasiskirstymas nepriklauso nuo $\mathbf{x}(t)$.
3. Liekanų dispersija yra pastovi: $\mathbf{D}\xi(t) = \sigma^2$.
4. Liekanos yra nekoreliuotų atsitiktinių dydžių seka: $\text{cov}(\xi(t), \xi(s)) = 0$, kai $t \neq s$.

*Darbą remia Lietuvos valstybinis mokslo ir studijų fondas, Grantas K-014.

5. Liekanos $\xi(t)$ yra pasiskirstę pagal Gauso dėsnį.

Parametro θ mažiausių kvadratų metodo įvertis yra gaunamas minimizuojant liekanų kvadratų sumą:

$$\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^n (Y(t) - \theta^T x(t))^2 = \arg \min_{\theta} (Y - X\theta)^T (Y - X\theta) \quad (2)$$

čia: Y – Y stebėjimų vektorius, $Y = (Y(1), Y(2), \dots, Y(n))^T$;

X – $x_j(t)$ funkcijų reikšmių matrica, $j = 1, 2, \dots, d$,

$$X = \begin{pmatrix} x_1(1) & x_2(1) & \dots & x_d(1) \\ x_1(2) & x_2(1) & \dots & x_d(2) \\ \dots & \dots & \dots & \dots \\ x_1(n) & x_2(n) & \dots & x_d(n) \end{pmatrix}.$$

Jei matrica $X^T X$ yra neišsigimusi ($\det X^T X \neq 0$), tai parametro θ įvertis gaunamas iš formulės:

$$\hat{\theta} = (X^T X)^{-1} X^T Y.$$

Praktikoje trečioji ir ketvirtoji modelio prielaidos dažnai yra netenkinamos ir tuomet mažiausių kvadratų metodu gautas įvertis nėra efektyvus. Šiuo atveju efektyvus įvertis yra gaunamas apibendrintu mažiausių kvadratų metodu:

$$\hat{\theta} = (X^T Q^{-1} X)^{-1} X^T Q^{-1} Y.$$

kur Q – liekanų kovariacijos matrica, (prie sąlygos kad $\det X^T Q^{-1} X \neq 0$).

Suprantama, jei liekanų seka būtų stacionari plačiąja prasme, tai kovariacijos matricą Q vienareikšmiškai apibrėžtų šių liekanų kovariacijos funkcija. Tačiau dėl liekanų dispersijos kitimo laike – heteroskedastiškumo ($D\xi(t) = \sigma^2(t)$), dažnai taip nėra. Šioje situacijoje paprastai daroma prielaida apie normuotų liekanų stacionarumą. Tegul seka $\varepsilon(t) = \xi(t)/\sigma(t)$ yra stacionari atsitiktinių dydžių seka su nuliniu vidurkiu ir pastovia dispersija ($E\varepsilon(t) = 0$, $D\varepsilon(t) = \sigma_0^2$). Šiuo atveju θ įvertis gaunamas iš formulės:

$$\hat{\theta} = (XS^{1/2}R^{-1}S^{-1/2}X)^{-1}X^T S^{-1/2}R^{-1}S^{-1/2}Y. \quad (3)$$

Čia: $S = \text{diag}(\sigma^2(1), \sigma^2(2), \dots, \sigma^2(n))$; R – stacionarių liekanų $\varepsilon(t)$ kovariacijos matrica. Praktikoje matricos S ir R , dažnai yra nežinomos, todėl formulėje (3) keičiamos statistiniais įverčiais.

Praktikoje dažnos yra situacijos, kai imtis yra užteršta išsiskiriančiais stebėjimais. (Išsiskiriančiais stebėjimais šiuo atveju yra vadinami stebėjimai tam tikrais laiko momentais, kai aiškiai

negalioja priimtas modelis.) Todėl modelį (3) reikia apibendrinti, stacionarioms liekanoms $\xi(t)$ suteikiant svorius $0 \leq \psi(\xi(t)) \leq 1$:

$$\hat{\theta} = (X^T \hat{\Psi}^{1/2} \hat{S}^{-1/2} \hat{R}^{-1} \hat{S}^{-1/2} \hat{\Psi}^{1/2} X)^{-1} X \hat{\Psi}^{1/2} \hat{S}^{-1/2} \hat{R}^{-1} \hat{S}^{-1/2} \hat{\Psi}^{1/2} Y, \quad (4)$$

kur $\Psi = \text{diag}(\psi(\varepsilon(1)), \psi(\varepsilon(2)), \dots, \psi(\varepsilon(n)))$. $\hat{\Psi}$ gauname, nestebimą seką $\xi(t)$ keičiant jos įverčiais: $\hat{\varepsilon}(t) = Y(t) - \tilde{\theta}^T x(t)$, kur $\tilde{\theta}$ – pradinis θ įvertis.

Svoriai $\psi(\varepsilon(t))$ gali būti suteikti remiantis taip vadinama „ 3σ “ taisykle. Ši taisyklė yra ypač taikoma Gauso sekoms su nuliniu vidurkiu ir pastovia dispersija. Pagal ją iš imties yra pašalinamas stebėjimas, kuris modulių yra didesnis už tris vidutinius standartinius nuokrypius $|\varepsilon(t)| > 3\sigma_0$.

Įvertį (4) galima dar labiau apibendrinti, aprioriai įvedant papildomus svorius $0 \leq w(t) \leq 1$, dėl kitų priežasčių, pvz. dėl modelio parametrų lėto kitimo laike, t.y. dėl duomenų senėjimo. Šiuo atveju θ įvertis gaunamas iš formulės:

$$\hat{\theta} = (X^T \Lambda^{1/2} \hat{R}^{-1} \hat{\Lambda}^{1/2} X)^{-1} X^T \Lambda^{1/2} \hat{R}^{-1} \Lambda^{1/2} Y, \quad (5)$$

kur $\Lambda = W \Psi S^{-1}$; $W = \text{diag}(w(1), w(2), \dots, w(n))$.

Metodas (5) yra mažiausių kvadratų metodo apibendrinimas susijęs su regresijos modelio sąlygų silpninimu, tačiau šį metodą galima apibendrinti ir kitaip – nagrinėjant bendresnę nei kvadratinę nuostolių funkciją. Priminsime, kad klasikinio mažiausių kvadratų metodo atveju įverčiai gaunami minimizuojant liekanų kvadratų sumą (2), o siekiant robastiškumo galima minimizuoti sumą $\sum_{t=1}^n \rho((Y(t) - \theta^T x(t)))$, kur ρ yra funkcija auganti lėčiau nei kvadratinė. Bendru atveju šio metodo modifikacija yra pavidalo:

$$\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^n \rho(Y(t) - \theta^T x(t)) \psi(\varepsilon(t)) w(t).$$

Tokio robastinio metodo pavyzdžiu būtų modifikuotas medianinis įvertinimo metodas:

$$\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^n |Y(t) - \theta^T x(t)| \psi(\varepsilon(t)) w(t) / \sigma(t).$$

Atveju $\psi(\cdot) = w(\cdot) = s(\cdot) = 1$ gautume įprastinį medianinį parametro θ įvertį. Vienamačiu atveju ($d = 1$) šis įvertis reikštų, kad liekanų $\varepsilon(t)$ empirinė mediana prilyginama nuliui.

3. Eksperimentų aprašymas

Tiriant metodų įvertinimo efektyvumą eksperimentiškai yra generuojamos atsitiktinės sekos su žinomais parametrais, o po to šių parametrų įverčiai lyginami su jų tikromis reikšmėmis. Metodo tikslumo matu yra naudojamas parametro įvertinimo paklaidos modulio vidurkis:

$$\Delta = \Delta(\theta, \hat{\theta}) = E \sqrt{(\theta - \hat{\theta})^T (\theta - \hat{\theta})}.$$

Kadangi matematinę paklaidos modulio vidurkį sunku apskaičiuoti, šiame darbe jis keičiamas aritmetiniu:

$$\bar{\Delta} = \frac{1}{N} \sum_{i=1}^N \Delta_i,$$

čia N bandymų skaičius, Δ_i žymi paklaidos modulį i -tos sekos atveju.

Eksperimentiškai prie įvairių sąlygų buvo tiriamas tiesinio trendo parametrų įverčių tikslumas. Tegul stebima seka tenkina lygybę: $Y(t) = a + bt + \xi(t)$, $t = 1, 2, \dots, n$, kur $E\xi(t) = 0$. Visuose toliau aprašytuose eksperimentuose buvo generuojama po 100 skirtingų nepriklausomų atsitiktinių imčių realizacijų, o parametrai $(a, b) = \theta^T$ įvertinti šiais dešimčia metodų:

1. $\hat{\theta} = (X^T X)^{-1} X^T Y$ – mažiausių kvadratų metodo įvertis.
2. $\hat{\theta} = (X^T \hat{S}^{-1} X)^{-1} X^T \hat{S}^{-1} Y$.
3. $\hat{\theta} = (X^T \hat{\Psi} X)^{-1} X^T \hat{\Psi} Y$.
4. $\hat{\theta} = (X^T \hat{\Psi} \hat{S}^{-1} X)^{-1} X^T \hat{\Psi} \hat{S}^{-1} Y$.
5. $\hat{\theta} = (X^T \hat{R}^{-1} X)^{-1} X^T \hat{R}^{-1} Y$.
6. $\hat{\theta} = (X^T \hat{\Psi}^{1/2} \hat{S}^{1/2} \hat{R}^{-1} \hat{S}^{-1/2} \hat{\Psi}^{1/2} X)^{-1} X^T \hat{\Psi}^{1/2} \hat{S}^{-1/2} \hat{R}^{-1} \hat{S}^{-1/2} \hat{\Psi}^{1/2} Y$.
7. $\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^n |Y(t) - \theta^T x(t)|$ – medianinio įvertinimo metodo įvertis.
8. $\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^n |Y(t) - \theta^T x(t)| / \hat{\sigma}(t)$.
9. $\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^n |Y(t) - \theta^T x(t)| \Psi(\hat{\varepsilon}(t))$.
10. $\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^n |Y(t) - \theta^T x(t)| \Psi(\hat{\varepsilon}(t)) / \hat{\sigma}(t)$.

Pastaba. Metoduose 2, 4, 6, 8, 10 kintanti dispersija yra aprašoma funkcija: $\sigma^2(t) = \exp(\beta_0 + \beta_1 t)$, kur parametrai β_0 ir β_1 yra vertinami mažiausių kvadratų metodu. Metoduose 3, 4, 5, 8, 9 svoriai $\Psi(\hat{\varepsilon}(t))$ yra suteikiami remiantis „3 σ “ taisykle. O metoduse 5 ir 6 liekanos $\varepsilon(t)$ yra laikomos $ARMA(1, 1)$ procesu.

1 eksperimentas. Šio eksperimento metu buvo gauti parametro θ įverčiai, kai visos modelio (1) prielaidos buvo tenkinamos. Atsitiktinės imtys buvo generuojamos pagal formulę $Y(t) = a + bt + \xi(t)$, kur $a = 1$, $b = 1$.

2 eksperimentas. Ši eksperimentą aprašysime detaliau.

Tarkime, kad $\xi(t)$ yra stacionarus baltas triukšmas, bet prielaida apie liekanų normalųjį pasiskirstymą yra netenkinama, t.y. liekanos $\xi(t)$ yra pasiskirstę ne pagal Gauso, o pagal simetrinį Pareto dėsnį:

$$F_{\xi}(y) = \mathbf{P}\{\xi < y\} = \begin{cases} 1 - \frac{1}{2} \left(\frac{c_0}{c_0 + y} \right)^{\alpha}, & \text{kai } y > 0; \\ \frac{1}{2} \left(\frac{c_0}{c_0 - y} \right)^{\alpha}, & \text{kai } y < 0. \end{cases}$$

Čia $\alpha > 0$ ir c_0 yra konstantos.

1 lentelė.
Parametro įvertinimo paklaidos modulio vidurkiai

Įvertinimo metodo numeris	$\bar{\Delta}$, kai $\alpha = 1$	$\bar{\Delta}$, kai $\alpha = 2$	$\bar{\Delta}$, kai $\alpha = 3$	$\bar{\Delta}$, kai $\alpha = 4$
1.	25.49342	1.33624	0.552668	0.339779
2.	21.62688	1.358871	0.567551	0.344082
3.	5.566214	0.620278	0.359143	0.257478
4.	6.541385	0.610019	0.361619	0.25801
5.	26.07149	1.331635	0.554122	0.340496
6.	7.225453	0.865068	0.42345	0.286761
7.	0.817534	0.385431	0.251884	0.187065
8.	0.708726	0.370205	0.249438	0.199592
9.	0.811271	0.369532	0.244961	0.18019
10.	0.914457	0.379111	0.253962	0.207839

Atsitiktinės sekos buvo generuojamos su parametrais $a = 1$, $b = 1$, $c_0 = 1$, $\alpha = 1, 2, 3, 4$.

3 eksperimentas. Šio eksperimento metu buvo gauti parametro θ įverčiai, kai imitys buvo užterštos išsiskiriančiais stebėjimais. Atsitiktinės imtys buvo generuojamos pagal formulę $Y(t) = a + bt + \xi(t) + \eta(t)$, kur $a = 1$, $b = 1$, $\xi(t) \sim N(0, 1)$, $\eta(t) = \begin{cases} \gamma_k, & \text{jei } t = \tau_k \\ 0, & \text{jei } t \notin \{\tau_1, \tau_2, \dots\} \end{cases}$, $\gamma_1, \gamma_2, \dots \sim N(0, 49)$, laiko tarpai $\tau_1, \tau_2 - \tau_1, \tau_3 - \tau_2, \dots$ yra pasiskirstę pagal Puasono dėsnį su parametrais $\lambda = 5, 10, 15$.

4 eksperimentas. Šio eksperimento metu buvo gauti parametro θ įverčiai, kai yra netenkinama modelio (1) trečioji prielaida, t.y. kai liekanų dispersija kinta laike. Atsitiktinės imtys buvo generuojamos pagal formulę: $Y(t) = a + bt + \sigma(t)\varepsilon(t)$, kur $a = 1$, $b = 1$, $\varepsilon(t) \sim N(0, 0.25)$, $\sigma^2(t) = \exp(1 + 0, 02t)$.

5 eksperimentas. Šio eksperimento metu buvo gauti parametro θ įverčiai, kai yra netenkinama modelio (1) ketvirtoji prielaida, t.y. kai liekanos $\xi(t)$ yra tarpusavyje koreliuotos. Atsitiktinės imtys buvo generuojamos pagal formulę: $Y(t) = a + bt + \xi(t)$, kur $a = 1$, $b = 1$, $\xi(t)$ yra $ARMA(1, 1)$ procesas: $\xi(t) + \alpha\xi(t-1) = \tilde{\varepsilon}(t) + \beta\tilde{\varepsilon}(t-1)$, $\alpha = 0.95$, $\beta = 0.5$.

4. Išvados ir pasiūlymai

Atliktų eksperimentų rezultatai trumpai yra pateiliami šiose išvadose.

1. Pirmojo eksperimento rezultatai rodo, kad tuo atveju kai visos modelio (1) prielaidos yra tenkinamos, tiksliausi įverčiai yra gauti mažiausių kvadratų metodu, nedaug skiriasi ir kitų šio metodo modifikacijų tikslumas. Netiksliausias šiuo atveju yra šeštasis metodas, kurio parametro įvertinimo modulio vidurkis yra apie 10% didesnis negu mažiausių kvadratų metodo. Medianinio

įvertinimo metodas bei įvairios jo modifikacijos šiame eksperimente pasirodė esą mažiau tikslūs už mažiausių kvadratų metodą. Jų visų tikslumas yra žymiai mažesnis.

2. Antrojo eksperimento rezultatai rodo, kad tuo atveju, kai yra netenkinama modelio (1) penktoji prielaida apie liekanų normalųjį pasiskirstymą, t.y. kai liekanos $\xi(t)$ yra pasiskirstę ne pagal Gauso, o pagal simetrinį Pareto dėsnį, geriausi rezultatai yra pasiekti taikant medianinį įvertinimo metodą bei įvairias jo modifikacijas. Tiksliausias šiuo atveju yra medianinis įvertinimo metodas su svoriais dėl išsiskiriančių stebėjimų (9 metodas). Mažiausių kvadratų metodu ir jo modifikacijų pagalba gauti rezultatai yra mažiau tikslūs.

3. Trečiojo eksperimento rezultatai rodo, kad tuo atveju, kai liekanos yra pasiskirstę pagal Gauso dėsnį, tačiau imtis yra užteršta išsiskiriančiais stebėjimais, tiksliausi įverčiai yra medianinio įvertinimo metodo bei įvairių jo modifikacijų. Čia tiksliausias yra medianinis įvertinimo metodas su svoriais dėl išsiskiriančių stebėjimų (9 metodas). Netiksliausias pasirodė esąs penktasis metodas, kurio pavyzdžiui parametro įvertinimo modulio vidurkis, kai imtis yra užteršta apie 10% išsiskiriančių stebėjimų, yra apie 100% didesnis negu aštuntojo metodo.

4. Iš ketvirtojo eksperimento rezultatų matyti, kad tuo atveju kai yra netenkinama modelio (1) trečioji prielaida, t.y. liekanų dispersija kinta laike, tiksliausi rezultatai yra gauti mažiausių kvadratų metodo su kintančia dispersija (2 metodas). Nedaug skiriasi ir kitų šio metodo modifikacijų tikslumas. Medianinio įvertinimo metodų grupės gauti rezultatai yra mažiau tikslūs. Netiksliausias pasirodė esąs mažiausių kvadratų metodas, kurio parametro įvertinimo modulio vidurkis yra apie 50% didesnis negu antrojo metodo.

5. Iš penktojo eksperimento rezultatų matyti, kad kai yra netenkinama modelio (1) ketvirtoji prielaida, t.y. kai liekanos $\xi(t)$ yra tarpusavyje koreliuotos ($\xi(t)$ yra *ARMA* procesas), tiksliausi įverčiai yra šeštojo metodo, nors visų mažiausių kvadratų metodo modifikacijų tikslumas skiriasi nedaug. Medianinio įvertinimo metodo įverčiai šiuo atveju yra mažiau tikslūs.

6. Apskaičiuojant parametro θ medianinį įvertį praktikoje galima naudoti rekurentinę procedūrą:

$$\hat{\theta}^{(k+1)} = \arg \min_{\theta} \sum_{t=1}^n (Y(t) - \theta^{T(k)} x(t))^2 w^{(k)}(t),$$

$$\text{kur } w^{(0)}(t) = 1, \quad w^{(k)}(t) = 1/|Y(t) - \theta^{T(k)} x(t)|, \quad k = 1, 2, \dots$$

Ši procedūra realizuojama panaudojant standartinę apibendrinto mažiausių kvadratų metodo procedūrą, kuri yra visuose didesniuose statistiniuose pakeituose. Jos pagalba gaunamas įvertis konverguoja į medianinį parametro įvertį, kai iteracijų skaičius k auga.

Literatūra

- [1] R. Salėtis, *Apie tiesinės regresijos modelių identifikavimo metodus*, MII, Vilnius (1997).
- [2] П. Хьюбер, *Робастность в статистике*, Мир, Москва (1984).
- [3] С.А. Айвзян, В.С. Мхитарян, *Прикладная статистика и основы эконометрики*, Юнити, Москва (1998).
- [4] Ю. Тюрин, Г. Симонова, *Знаковый анализ линейных моделей*, Обзорение прикладной и промышленной математики, 215–278 (1993).

On identification methods of linear regression models

R. Salėtis

The aim of this paper is to discuss the efficiency of the least squares and least absolute values estimation methods, in such situations, then the sample is with outliers, the variance of the residuals changes in the time and the residuals are correlated. A range of different modifications of these methods in order to improve accuracy of the estimation is offered. All the discussed procedures are fulfilled in a PC by means of SAS (*Statistical Analysis System*).