

برآوردگر استوار مرزبندی شده تعمیم یافته محتمل در مدل رگرسیون نیمه پارامتری

مهدی روزبه^۱، مرتضی امینی^۲

اگره آمار، دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه سمنان

اگره آمار، دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه تهران

چکیده: در تجزیه و تحلیل مسائل رگرسیونی و به ویژه مدل بندی آماری بسیاری از داده ها مانند داده های اقتصادی، روانشناسی، علوم اجتماعی، علوم پزشکی، مهندسی و غیره با مشکل هم خطی در میان متغیرهای پیشگو و حضور نقاط دورافتاده در مجموعه داده ها مواجه می شویم. در چنین مواقعی برآوردگر کمترین توان های دوم معمولی منجر به برآوردگرهای نادقیق می شود. برای غلبه بر مشکل مشاهده های دورافتاده از روش های استوار استفاده می شود. همچنین برای حل مشکل هم خطی چندگانه استفاده از رگرسیون مرزبندی شده توصیه می شود. از طرف دیگر در شرایطی که واریانس خطاها ناهمگن بوده یا خطاها دارای خودهمبستگی باشند، از روش کمترین توان های دوم تعمیم یافته استفاده می شود. در این مقاله ابتدا یک الگوریتم سریع برای محاسبه برآوردگر کمترین توان های دوم تعمیم یافته پیراسته مرزبندی شده محتمل در مدل رگرسیون نیمه پارامتری پیشنهاد شده و سپس با استفاده از شبیه سازی به روش مونت کارلو و یک داده واقعی، کارایی برآوردگرهای پیشنهادی سنجیده می شود.

واژه های کلیدی: اعتبارسنجی متقابل، برآوردگر کمترین توان های دوم پیراسته، مدل رگرسیون نیمه پارامتری، نقاط دورافتاده، نقطه شکست.

^۱ آدرس الکترونیک مسئول مقاله: مهدی روزبه، mahdi.roozbeh@semnan.ac.ir
^۲ کد موضوع بندی ریاضی (۲۰۱۰): 62J20, 62G35, 62G08

۱ مقدمه

مدل رگرسیون نیمه پارامتری یکی از پرکاربردترین مدل‌های رگرسیون در حوزه یادگیری آماری است که در سال‌های اخیر مورد توجه محققین قرار گرفته است. علاوه بر کاربرد گسترده این مدل‌ها در تحلیل کواریانس، آن‌ها کاربرد زیادی در مسائل اقتصادی دارند که از مهمترین آن‌ها می‌توان به مدل تابع درآمد سرمایه انسانی (ویلیس، ۱۹۸۶) و منحنی دستمزد (بلنچفلور و اسوالد، ۱۹۹۴) اشاره کرد. در هر دو مورد، لگاریتم درآمد شخصی به ویژگی‌های شخصیتی (مانند جنسیت و وضعیت تاهل) و اندازه سرمایه‌های انسانی فرد (مانند تحصیل و تجربه کار در بازار) وابسته است. در نظریه اقتصاد رابطه غیرخطی بین لگاریتم درآمد شخصی و تجربه کار در بازار مورد بررسی قرار می‌گیرد. این در حالی است که رابطه امیدریاضی لگاریتم درآمد شخصی و متغیرهای کیفی جنسیت، وضعیت تاهل و تحصیلات خطی است.

فرض کنید مشاهدات مستقل $(x_1, t_1, y_1), \dots, (x_n, t_n, y_n)$ در مدل

$$y_i = x_i \beta + f(t_i) + \epsilon_i, \quad i = 1, \dots, n, \quad (1)$$

صدق کنند، که در آن y_i ‌ها مشاهدات تک‌بعدی متغیر پاسخ، x_i ‌ها بردارهای p بعدی متغیرهای پیشگو دارای رابطه خطی با متغیر پاسخ، β بردار p بعدی ضرایب نامعلوم، t_i ‌ها مشاهدات متغیر یا متغیرهای پیشگو دارای رابطه غیرخطی با متغیر پاسخ، $f(\cdot)$ تابعی نامعلوم و هموار با دامنه $D \subset \mathcal{R}$ و ϵ_i ‌ها خطاهای تصادفی هستند. در این مدل فرض می‌شود که تابع نامعلوم $f(\cdot)$ هموار بوده و همان‌طور که ملاحظه می‌شود، مدل‌های رگرسیونی نیمه پارامتری انعطاف‌پذیرتر از مدل‌های رگرسیون خطی هستند، زیرا هم دارای قسمت خطی و هم قسمت غیرخطی می‌باشند. در واقع استفاده از این مدل زمانی مفید است که متغیر وابسته y به‌طور خطی با متغیر مستقل x و به‌طور غیرخطی با متغیر مستقل t رابطه داشته باشد. محققان این مدل را به مدل ناپارامتری محض و رگرسیون خطی ترجیح می‌دهند. مدل‌های رگرسیون نیمه پارامتری نخستین بار توسط انگل و همکارانش در سال ۱۹۸۶ در بررسی رابطه میان مصرف ماهیانه برق به‌عنوان متغیر پاسخ و قیمت ماهیانه برق، درآمد ماهیانه و دمای هوا به‌عنوان متغیرهای پیشگو با در نظر گرفتن قیمت ماهیانه برق و درآمد ماهیانه به‌عنوان بخش خطی و دمای هوا به‌عنوان بخش غیرخطی مورد استفاده قرار گرفتند.

هم‌خطی چندگانه یکی از مشکلات متداول در مدل‌های رگرسیون خطی چندگانه است که به دلیل ایجاد تورم در ماتریس کواریانس، باعث کاهش کارایی برآوردگرهای کلاسیک و نهایتاً پیش‌بینی نادرست می‌شود. رگرسیون

تاوانیده^۱ یک راهکار منظم‌سازی است که برآوردگرها را با اضافه کردن یک جمله تاوان به تابع هدف بهبود می‌بخشد. این روش تورم (بزرگی مقدرهای) ماتریس کواریانس برآورد ضرایب را با انقباض برآوردگرها به سمت صفر کاهش می‌دهد. برآوردگر مرزبندی شده^۲ (هول و کنارد، ۱۹۷۰) در واقع برآوردگر رگرسیون تاوانیده با تابع تاوان دوم نرم L_2 ضرایب رگرسیونی است. این روش یکی از پرکاربردترین راهکارها در مقابله با مشکل هم‌خطی چندگانه است. خصوصیات نظری برآوردگر مرزبندی شده به‌طور گسترده‌ای مورد مطالعه قرار گرفته است. به‌عنوان مثال می‌توان به فاریراد (۱۹۷۶)، گلاب و همکاران (۱۹۷۹)، اسپکمن (۱۹۸۸)، آرشی و همکاران (۲۰۱۵) و نوروزی‌راد و آرشی (۱۳۹۶) اشاره کرد. همچنین، این برآوردگر در مدل رگرسیون نیمه‌پارامتری توسط تاباکان و آکدنیز (۲۰۱۰)، آکدنیز و همکاران (۲۰۱۲)، روزبه (۲۰۱۵)، آرشی و ولی‌زاده (۲۰۱۵) و امینی و روزبه (۲۰۱۵)، مورد مطالعه قرار گرفته است. با این حال برآوردگر کلاسیک مرزبندی شده مانند برآوردگر کلاسیک کم‌ترین توان‌های دوم نسبت به مشاهدات دورافتاده حساس است. مشاهدات دورافتاده نقاطی هستند که از مدل اکثریت نقاط تبعیت نمی‌کنند. جانشین‌های استوار روش کلاسیک کم‌ترین توان‌های دوم مانند M -برآوردگرها، S -برآوردگرها، برآوردگر حداقل میانه توان‌های دوم (LMS) و روش کم‌ترین توان‌های دوم پیراسته^۳ (LTS) برای کاهش اثر مشاهدات دورافتاده مورد استفاده قرار گرفته‌اند. برای مروری بر این روش‌ها می‌توان به مارونا و همکاران (۲۰۰۶) مراجعه کرد. روش کم‌ترین توان‌های دوم پیراسته که نخستین بار توسط روسیو (۱۹۸۴) معرفی شد، به دلیل تعریف ساده، نقطه فروریزش بالا (بخش ۴ را ببینید) و استواری در برابر نقاط نافذ (مشاهدات دور افتاده در جهت متغیرهای پیشگو) یکی از متداول‌ترین روش‌های برآورد استوار است. برای یک مجموعه شامل n مشاهده و یک پارامتر صحیح $h \leq n$ ، برآوردگر LTS ابرصفحه حداقل‌کننده مجموع کوچکترین h مربع مانده را به مشاهدات برازش می‌دهد. اگر یکی از متغیرهای پیش‌بین یک تابع دقیق خطی از یک یا چند متغیر پیش‌بین دیگر باشد، مدل رگرسیون دارای هم‌خطی کامل است. هم‌خطی ناقص وقتی اتفاق می‌افتد که یکی از متغیرهای پیش‌بین به‌طور تقریبی یک تابع خطی از یک یا چند متغیر پیش‌بین دیگر باشد. برآوردگرهای استوار کم‌ترین توان‌های دوم پیراسته در مدل رگرسیون توسط روزبه و بابایی کفاکی (۲۰۱۶) و در مدل رگرسیون نیمه‌پارامتری توسط تورکمن و تاباکان (۲۰۱۵) مطالعه شده است. آرشی و نوروزی‌راد (۲۰۱۵) M -برآوردگرهای مرزبندی شده بهبودیافته را در مدل رگرسیون تعمیم

¹Penalized²Ridge³Least Treamed Squares

دادند. در این مقاله، به معرفی حالت تعمیم یافته برآوردگرهای مرزبندی شده استوار محتمل^۴ در مدل رگرسیون نیمه پارامتری برای استفاده در حالت ناهمگنی واریانس خطاها و مدل رگرسیون سری زمانی پرداخته می شود. در بخش ۲ مقدمات نظری و فرضیات اساسی بیان می شوند. بخش ۳ به معرفی برآوردگرهای پارامتری و نیمه پارامتری ضرایب و تابع f اختصاص داده شده است. نقطه فروریزش برآوردگر، برآورد پارامترهای ارب سازی و پهنای باند به ترتیب مطالب بخش های ۴ و ۵ را تشکیل می دهند. در پایان نتایج عددی شامل بررسی شبیه سازی به روش مونت کارلو و تحلیل داده های واقعی در بخش ۶ ارائه شده و در انتها به بحث و نتیجه گیری پرداخته می شود.

۲ فرض های اساسی

شکل ماتریسی مدل رگرسیون نیمه پارامتری به صورت

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{f} + \boldsymbol{\epsilon}, \quad (۲)$$

تعریف می شود، که در آن $\mathbf{y} = (y_1, \dots, y_n)'$ بردار $n \times 1$ از مشاهدات پاسخ، $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$ ماتریس $n \times p$ متغیرهای پیشگوی بخش خطی، $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ بردار $p \times 1$ ضرایب نامعلوم بخش خطی، $\mathbf{f} = (f(t_1), \dots, f(t_n))'$ بردار $n \times 1$ اثر غیرخطی و $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)'$ بردار $n \times 1$ خطای تصادفی هستند. در این مدل فرض ناهمگنی واریانس خطاها یا ناهمبستگی خطاها به صورت

$$E(\boldsymbol{\epsilon}|\mathbf{X}, \mathbf{t}) = \mathbf{o}, \text{Var}(\boldsymbol{\epsilon}|\mathbf{X}, \mathbf{t}) = \sigma^2 \mathbf{V},$$

نوشته می شود، که در آن σ پارامتری مجهول و \mathbf{V} یک ماتریس معین مثبت نامعلوم است. با فرض معلوم بودن $\boldsymbol{\beta}$ ، می توان تابع $f(\cdot)$ را با استفاده از روش هموارسازی تابع رگرسیونی پاسخ $y_i - \mathbf{x}'_i \boldsymbol{\beta}$ بر روی پیشگوی t_i به صورت

$$\hat{f}(t) = \sum_{i=1}^n k_i(t)(y_i - \mathbf{x}'_i \boldsymbol{\beta}), \quad (۳)$$

برآورد کرد، به طوری که k_i توابع وزنی مثبت بوده و به مشاهدات بستگی دارند.

⁴Feasible

در این مقاله از برآوردگر کرنل گاسر-مولر استفاده می‌شود، که در آن وزن‌ها به صورت

$$k_i(t) = \frac{1}{\omega} \int_{s_{i-1}}^{s_i} K\left(\frac{t-s}{\omega}\right) ds, s_i = (t_i + t_{i+1})/2, s_0 = 0, s_1 = 1$$

تعریف می‌شوند (برای مطالعه بیشتر و خصوصیات این برآوردگر می‌توان به یوبانگ (۱۹۹۹)) مراجعه کرد)، به طوری که $K(\cdot)$ تابع هسته و ω پارامتر پهنای باند هستند. همچنین فرض می‌شود که $K(\cdot)$ یک تابع هسته رتبه دوم ($K'' \neq 0$) با خصوصیات

$$\int_{-1}^1 K(u) du = 1, \int_{-1}^1 uK(u) du = 0, \\ M_2 = \int_{-1}^1 u^2 K(u) du \neq 0, \kappa_2 = \int_{-1}^1 K(u)^2 du < \infty,$$

باشد. تحت شرایط بالا بنابر یوبانگ (۱۹۹۹) می‌توان نوشت

$$E(\hat{f}(t) - f(t)) = \frac{\omega^2 f''(t) M_2}{4} + O(\omega^4) + O(n^{-1}) \\ \text{Var}(\hat{f}(t)) = (n\omega)^{-1} \sigma^2 \kappa_2 + O((n\omega)^{-2}).$$

به منظور برآورد بردار پارامتر خطی β به روش مانده‌های جزئی، با جایگذاری مقدار $\hat{f}(\cdot)$ در معادله (۱)، مدل

$$y_i - \sum_{j=1}^n k_j(t_i) y_j = (\mathbf{x}'_i - \sum_{j=1}^n k_j(t_i) \mathbf{x}'_j) \beta + \epsilon_i,$$

حاصل می‌شود. اکنون با در نظر گرفتن $\tilde{y}_i = y_i - \sum_{j=1}^n k_j(t_i) y_j$ و $\tilde{\mathbf{x}}'_i = \mathbf{x}'_i - \sum_{j=1}^n k_j(t_i) \mathbf{x}'_j$ می‌توان مدل اخیر را به صورت $\tilde{y}_i = \tilde{\mathbf{x}}'_i \beta + \epsilon_i$ نوشت، که همان مدل رگرسیون استاندارد است. بنابراین، با تعریف $\tilde{\mathbf{y}} = (\tilde{y}_1, \dots, \tilde{y}_n)'$ و $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}'_1, \dots, \tilde{\mathbf{x}}'_n)'$ که $\tilde{\mathbf{y}} = (\mathbf{I}_n - \mathbf{K})\mathbf{y}$ و $\tilde{\mathbf{X}} = (\mathbf{I}_n - \mathbf{K})\mathbf{X}$ است، برآورد کلاسیک کم‌ترین توان‌های دوم تعمیم‌یافته β در مدل جدید $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta + \epsilon$ به صورت

$$\hat{\beta}_G = (\tilde{\mathbf{X}}' \mathbf{V}^{-1} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \mathbf{V}^{-1} \tilde{\mathbf{y}}, \quad (4)$$

خواهد بود، که در آن $\tilde{\mathbf{y}} = (\tilde{y}_1, \dots, \tilde{y}_n)'$ بردار پاسخ هموار شده و $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}'_1, \dots, \tilde{\mathbf{x}}'_n)'$ ماتریس طرح هموار شده است. جائو و همکاران (۱۹۹۵) ثابت کردند که برآوردگر به دست آمده برای β ، برآوردگری سازگار و دارای توزیع نرمال مجانبی به صورت

$$\sqrt{n}(\hat{\beta}_G - \beta) \xrightarrow{D} N(0, \sigma_\epsilon^2 \Sigma_{\mathbf{x}|t}^{-1}),$$

است. در عمل، برای استنباط راجع به پارامتر β ، باید σ_ϵ^2 و $\Sigma_{x|t}$ برآورد شوند. برآوردگرهای سازگار

$$s_\epsilon^2 = \frac{1}{n} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \hat{\beta}_G)' \mathbf{V}^{-1} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \hat{\beta}_G) \xrightarrow{p} \sigma_\epsilon^2$$

$$\hat{\Sigma}_{x|t} = \frac{1}{n} (\tilde{\mathbf{X}}' \mathbf{V}^{-1} \tilde{\mathbf{X}}) \xrightarrow{p} \Sigma_{x|t}.$$

رهیافتی برای حل این مساله می‌باشند. در مرحله آخر، می‌توان تابع $f(\cdot)$ را با جایگذاری برآوردگر نیمه پارامتری سازگار به دست آمده برای β در معادله (۳)، با استفاده از برآوردگر

$$\hat{f}(t) = \mathbf{k}(t)(\mathbf{y} - \mathbf{x}'_i \hat{\beta}_G), \quad (5)$$

برآورد کرد. بنا بر هاردل و همکاران (۲۰۰۰) برآوردگر $\hat{f}(t)$ یک برآوردگر سازگار برای $f(t)$ است.

۳ برآورد مرزبندی شده کم‌ترین توان‌های دوم تعمیم یافته پیراسته محتمل

با ضرب ماتریس $\mathbf{V}^{-1/2}$ در طرفین معادله $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta + \epsilon$ به منظور ناهمبسته نمودن جملات خطا و ثابت نمودن واریانس آن‌ها، می‌توان مدل جدید

$$\mathbf{V}^{-1/2} \tilde{\mathbf{y}} = \mathbf{V}^{-1/2} \tilde{\mathbf{X}} \beta + \mathbf{V}^{-1/2} \epsilon,$$

را در نظر گرفت. حال با توجه به $\text{Var}(\mathbf{V}^{-1/2} \epsilon) = \sigma^2 \mathbf{I}$ ، می‌توان مانده‌های تبدیل یافته $\mathbf{e}(\beta) = \mathbf{V}^{-1/2} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \beta)$ را با مانده‌های $(\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \beta)$ جایگزین کرد. در این صورت برآوردگر استوار کم‌ترین توان‌های دوم پیراسته تعمیم یافته به صورت

$$\hat{\beta}_{GLTS}(z) = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^h (\mathbf{e}^{\top}(\beta))_{i:n} = \underset{\beta, z \in E_h}{\operatorname{argmin}} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \beta)' \mathbf{V}^{-1/2} \mathbf{Z} \mathbf{V}^{-1/2} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \beta),$$

تعریف می‌شود، به طوری که $(\mathbf{e}^{\top}(\beta))_{1:n} \leq \dots \leq (\mathbf{e}^{\top}(\beta))_{n:n}$ آماره‌های مرتب توان دوم مانده‌ها و \mathbf{Z} ماتریس قطری با عناصر قطر اصلی $\mathbf{z} = (z_1, \dots, z_n)'$ است. اگر مشاهده i ام خوب باشد، z_i مقدار یک و اگر مشاهده i ام پرت باشد، z_i مقدار صفر را اختیار می‌کند. همچنین

$$E_h = \{z; z_i \in \{0, 1\}, i = 1, \dots, n, z' \mathbf{1} = h\}, h \leq n.$$

در حالت کلی $h = \lceil n(1 - \alpha) \rceil$ در نظر گرفته می‌شود که $\alpha \in (0, 1)$ نسبت مشاهدات تاثیرگذار و $\lceil x \rceil$ مقدار صحیح عدد اعشاری x است. در واقع $1 - \alpha$ یک حدس اولیه برای درصد مشاهدات غیر دورافتاده است که توسط برخی از محققان 0.75 در نظر گرفته شده است. با چنین انتخابی برآورد نهایی بر اساس درصد قابل توجهی از مشاهدات به دست آمده و نقطه فروریزش برآوردگر نیز 25% خواهد بود. برخی دیگر از محققان مقدار h را به صورت $h = \lceil n/2 \rceil + \lceil (p+1)/2 \rceil$ پیشنهاد داده‌اند (آفونس و همکاران، ۲۰۱۳).
 در حضور هم‌خطی چندگانه و مشاهدات دورافتاده، می‌توان از برآوردگر کم‌ترین توان‌های دوم پیراسته مرزبندی شده تعمیم‌یافته به صورت

$$\hat{\beta}_{GLTS}(z, \lambda) = \underset{\beta, z \in E_h}{\operatorname{argmin}} Q(z, \beta),$$

استفاده کرد، که در آن $Q(z, \beta) = (\tilde{y} - \tilde{X}\beta)' V^{-1/2} Z V^{-1/2} (\tilde{y} - \tilde{X}\beta) + \lambda \beta' \beta$. با قرار دادن $Z^* = \operatorname{Diag}(z^*)$ که در آن

$$z^* = \underset{E_h}{\operatorname{argmin}} (\tilde{y} - \tilde{X}\beta(z, \lambda))' V^{-1/2} Z V^{-1/2} (\tilde{y} - \tilde{X}\beta(z, \lambda)),$$

$$\beta(z, \lambda) = (\tilde{X}' V^{-1/2} Z V^{-1/2} \tilde{X} + \lambda I_p)^{-1} \tilde{X}' V^{-1/2} Z V^{-1/2} \tilde{y},$$

برآوردگر کم‌ترین توان‌های دوم پیراسته مرزبندی شده تعمیم‌یافته برابر است با

$$\hat{\beta}_{GLTS}(z^*, \lambda) = (\tilde{X}' V^{-1/2} Z^* V^{-1/2} \tilde{X} + \lambda I_p)^{-1} \tilde{X}' V^{-1/2} Z^* V^{-1/2} \tilde{y}. \quad (6)$$

همچنین برآورد استوار تابع $f(\cdot)$ در هر نقطه $t \in D$ با استفاده از رابطه

$$\hat{f}_{GLTS}(t; z^*, \lambda) = \mathbf{k}(t)(\mathbf{y} - \mathbf{X}\hat{\beta}_{GLTS}(z^*, \lambda)), \quad (7)$$

محاسبه می‌شود. با توجه به صورت برآوردگرهای اخیر (در عین کاربردی بودن آن)، تنها در حالتی که ماتریس V معلوم است می‌توان نتایجی دقیق به دست آورد، ولی این حالت مسائل کاربردی و واقعی را پوشش نمی‌دهد. به منظور کاربردی بودن و استفاده از مشاهدات واقعی تنها می‌توان رفتار مجانبی برآوردگرهای اخیر را مورد تجزیه و تحلیل قرار داد. برای مشاهده توجیه و دلیل این مطلب می‌توان به زلنر (۱۹۶۲ و ۱۹۶۳) مراجعه کرد. چون ماتریس کوواریانس خطا یعنی V در عمل معمولاً نامعلوم است، بنابراین استفاده از $\hat{\beta}_{GLTS}(z^*, \lambda)$ غیرممکن بوده و باید ماتریس V بوسیله یک برآوردگر مناسب جایگزین شود. زلنر (۱۹۶۲) ثابت کرد

توان‌های دوم معمولی β یعنی $(\check{X}'\check{X})^{-1}\check{X}'\check{y}$ است. با جایگزینی S در عبارت (۴) بجای V ، برآوردگر کم‌ترین توان‌های دوم تعمیم‌یافته محتمل به صورت

$$\hat{\beta}_{FG} = (\check{X}'S^{-1}\check{X})^{-1}\check{X}'S^{-1}\check{y},$$

به‌دست می‌آید. بنابراین، بر اساس برآوردگر پیشنهادی اخیر توسط زلنر (۱۹۶۲)، می‌توان برآوردگر کم‌ترین توان‌های دوم پیراسته مرزبندی شده تعمیم‌یافته محتمل

$$\hat{\beta}_{FGLTS}(z^*, \lambda) = (\check{X}'S^{-1/2}Z^*S^{-1/2}\check{X} + \lambda I_p)^{-1}\check{X}'S^{-1/2}Z^*S^{-1/2}\check{y}, \quad (۸)$$

را با جایگزینی S در عبارت (۹) بجای V به‌دست آورد. با استفاده از زلنر (۱۹۶۲) می‌توان نتیجه گرفت $\sqrt{n}(\hat{\beta}_{FGLTS}(z^*, \lambda) - \beta) = \sqrt{n}(\hat{\beta}_{GLTS}(z^*, \lambda) - \beta) + O(n^{-1})$ و بنابراین، $\hat{\beta}_{FGLTS}(z^*, \lambda)$ دارای توزیع مجانبی نرمال یکسان هستند. بنابراین، هنگامی که ماتریس کوواریانس جملات خطا نامعلوم است، برآورد استوار محتمل تابع $f(\cdot)$ در هر نقطه $t \in D$ با استفاده از

$$\hat{f}_{FGLTS}(t; z^*, \lambda) = \mathbf{k}(t)(\mathbf{y} - \mathbf{X}\hat{\beta}_{FGLTS}(z^*, \lambda)), \quad (۹)$$

محاسبه می‌شود.

۴ نقطه فروریزش برآوردگر استوار مرزبندی شده محتمل

یک معیار برای میزان استواری یک برآوردگر، نقطه فروریزش نمونه متناهی است.

تعریف ۱. (مارونا و همکاران، ۲۰۰۶) نقطه فروریزش برآوردگر $T = T(\mathbf{Z})$ برای نمونه \mathbf{Z} به حجم n به صورت

$$BP(T; \mathbf{Z}) = \min_m \left\{ \frac{m}{n} : \sup_{\mathbf{Z}^*} \|T(\mathbf{Z}^*)\|^2 = \infty \right\},$$

تعریف می‌شود، که در آن \mathbf{Z}^* نمونه آلوده شده^۵ با جایگزینی $m \leq n$ نقطه از نمونه \mathbf{Z} با مقادیر دلخواه است.

^۵Contaminated

قضیه ۱. نقطه فروریزش برآوردگرهای $\hat{\beta}_{FGLTS}(z^*, \lambda)$ و $\hat{f}_{FGLTS}(t; z^*, \lambda)$ برابر است با

$$BP(\hat{\beta}_{FGLTS}(z^*, \lambda)) = \frac{n-h-1}{n}. \quad (۱۰)$$

برهان. فرض کنید $(\tilde{y}^*, \tilde{X}^*)$ نمونه آلوده شده به مشاهدات پرت (نمونه تخریبی) با جایگزینی $m \leq n-h$ مشاهده باشد. بنابراین تعداد مشاهدات خوب در $(\tilde{y}^*, \tilde{X}^*)$ برابر $n-m \geq h$ بوده و برای نمونه $(\tilde{y}^*, \tilde{X}^*)$ داریم

$$\min_{z \in E_h} Q(z, \circ) = \min_{z \in E_h} \tilde{y}^{*'} S^{-1/2} Z S^{-1/2} \tilde{y}^* \leq \min_{z \in E_h} \tilde{y}' S^{-1/2} Z S^{-1/2} \tilde{y} \leq hM_y^\gamma,$$

به طوری که $M_y = \max_{i=1, \dots, n} |\tilde{y}_i|$. چنانچه $\beta' \beta \geq \frac{hM_y^\gamma + 1}{\lambda}$ ، آنگاه (فرض خلف) رابطه

$$\min_{z \in E_h} Q(z, \beta) \geq \lambda \beta' \beta \geq hM_y^\gamma + 1 > \min_{z \in E_h} Q(z, \circ),$$

برقرار است، که چون $\min_{z \in E_h} Q(z, \hat{\beta}(z^*, k)) \leq \min_{z \in E_h} Q(z, \circ)$ فرض خلف باطل می شود. بنابراین می توان نوشت

$$\hat{\beta}_F(z^*, \lambda)' \hat{\beta}_F(z^*, \lambda) \leq \frac{hM_y^\gamma + 1}{\lambda},$$

و بنابراین $BP(\hat{\beta}_{FGLTS}(z^*, \lambda)) \geq \frac{n-h-1}{n}$. اکنون جای آخرین $m = n-h+1$ مشاهده (\tilde{X}, \tilde{y}) را طوری عوض می کنیم که آخرین m مشاهده (\tilde{X}, \tilde{y}) به (\tilde{X}, \tilde{y}) با $(a, \circ, \dots, \circ)'$ تغییر یابد، که در آن $a > \circ$ ، $M > \circ$ و

$$a^\gamma \geq \max(h-m, \circ) \left(\max_{i=1, \dots, n} |y_i| + M \max_{i=1, \dots, n} \|x_i\| \right)^\gamma + \lambda M^\gamma.$$

با فرض $\beta_M = (M, \circ, \dots, \circ)' \in \mathbb{R}^p$ آخرین m مشاهده $(\tilde{y} - \tilde{X} \beta_M)$ برابرند با

$$(Ma - Ma) = \circ.$$

بنابراین

$$\min_{z \in E_h} Q(z, \beta_M) = \begin{cases} \min_{z \in E_{h-m}} (\tilde{y} - \tilde{X} \beta_M)' S^{-1/2} Z S^{-1/2} (\tilde{y} - \tilde{X} \beta_M) + \lambda M^\gamma, & h > M \\ \lambda M^\gamma, & \text{جاهای دیگر} \end{cases}$$

می توان نوشت

$$\begin{aligned} \min_{z \in E_h} Q(z, \beta_M) &\leq \max(h - m, 0) \left(\max_{i=1, \dots, n} |\check{y}_i| + M \max_{i=1, \dots, n} \|\check{x}_i\| \right)^2 + \lambda M^2 \\ &\leq a^2 - 1. \end{aligned} \quad (11)$$

همچنین برای نمونه آلوده شده و $\beta = (\beta_1, \dots, \beta_p) \in \mathbb{R}^p$ می توان نوشت

$$\min_{z \in E_h} Q(z, \beta) \geq (Ma - a\beta_1)^2.$$

زیرا حداقل یکی از m مشاهده آخر در نمونه آلوده متعلق به مجموعه کوچکترین h مانده است. با در نظر گرفتن

$$\beta_1 \leq (\beta' \beta)^{1/2} \leq M \quad \text{چون } \beta' \beta \leq M^2$$

$$\min_{z \in E_h} Q(z, \beta) \geq a^2 (M - \beta_1)^2 \geq a^2, \quad (12)$$

از نابرابری های (۱۱) و (۱۲)، نابرابری $\hat{\beta}_{FGLTS}(z^*, k)' \hat{\beta}_{FGLTS}(z^*, k) \geq M^2$ نتیجه می شود که به معنی نامتنهایی شدن برآوردگر با میل کردن M به بینهایت و فروریزش آن است و بنابراین برهان تمام است. واضح است که نقطه فروریزش بالاتر به معنی استواری بیشتر برآوردگر در برابر مشاهده های دور افتاده است. از رابطه (۱۰) معلوم می شود که با کوچک شدن h ، نقطه فروریزش بالاتری به دست خواهد آمد. به نظر می آید که اگر $h < n/2$ باشد، حتی نقطه فروریزش بیش از ۵۰٪ به دست می آید. از لحاظ ریاضی اثبات شده است که نقطه فروریزش هر برآوردگر پایای رگرسیونی حداکثر برابر ۵۰٪ است (روسیو و لروی ۲۰۰۳ را ببینید). با این حال دستیابی عملی به نقطه فروریزش بیش از ۵۰٪ نامعقول است، زیرا اگر بیش از ۵۰٪ مشاهدات از الگوی سایر مشاهدات تبعیت نکنند، اکثریت نقاط در واقع این مجموعه خواهند بود. بنابراین توصیه می شود مقادیر $h \geq n/2$ استفاده نشود. معمولاً در عمل مقادیر $\alpha = ۰/۷۵$ یا $\alpha = [(p+1)/2]$ یا $h = \lceil [n/2] \rceil + \lceil [(p+1)/2] \rceil$ پیشنهاد می شوند.

برای محاسبه Z^* و برآوردگرهای معرفی شده روی مجموعه E_h باید تمامی $\binom{n}{h}$ انتخاب از مجموعه $\{1, \dots, n\}$ بررسی شود که به زمان و فضای زیادی نیاز دارد. در این مقاله از روش کمترین توان های دوم پیراسته سریع^۶ که توسط روسیو و ون درینسن (۲۰۰۶) معرفی شده استفاده می شود.

^۶Fast-LTS

۵ تعیین پارامترهای هموارسازی و اریب‌سازی

برای به‌دست آوردن مقادیر بهینه پارامترهای هموارسازی و مرزبندی شده، بنا بر ایوبانک (۱۹۹۹) ابتدا از مقدار بهینه

$$\omega_*^* = n^{-1/5} \left\{ \frac{\kappa_{\gamma} \hat{\sigma}_*^{\gamma}}{M_{\gamma}^{\gamma} J_{\gamma}^{\gamma}} \right\}^{1/5}, \quad (13)$$

برای تعیین پارامتر هموارسازی استفاده می‌شود، که در آن $J_{\gamma} = \int_0^1 \hat{f}''(t) dt$ اندازه همواری برآورد $f(\cdot)$ بوده و

$$\kappa_{\gamma} = \int_{-1}^1 K(u)^{\gamma} du < \infty, \quad M_{\gamma} = \int_{-1}^1 u^{\gamma} K(u) du \neq 0,$$

$$\hat{\sigma}_*^{\gamma} = \frac{\|S^{-1/\gamma}(\mathbf{y} - \mathbf{X}\hat{\beta}_{FG} - \hat{f}_F(t))\|^{\gamma}}{n-p},$$

که در آن $\hat{\beta}_{FG} = (\tilde{\mathbf{X}}'S^{-1}\tilde{\mathbf{X}})^{-1}\tilde{\mathbf{X}}'S^{-1}\tilde{\mathbf{y}}$ برآورد کم‌ترین توان‌های دوم تعمیم‌یافته محتمل β است. اکنون از مقدار بهینه اولیه پارامتر هموارسازی به‌دست آمده در (۱۳) برای تعیین مقدار بهینه اولیه پارامتر مرزبندی شده (λ^*) استفاده می‌شود. این کار براساس معیار اعتبارسنجی تعمیم‌یافته^۷

$$GCV(\omega_*^*, \lambda) = \frac{\frac{1}{n} \|S^{-1/\gamma}(\mathbf{I}_n - \mathbf{L}(\omega_*^*, \lambda))\mathbf{y}\|^{\gamma}}{[1 - \frac{1}{n} \text{tr}(\mathbf{L}(\omega_*^*, \lambda))]^{\gamma}},$$

انجام می‌شود، که در آن

$$\mathbf{L}(\omega_*^*, \lambda) = \mathbf{K} + (\mathbf{I}_n - \mathbf{K})\mathbf{X}(\tilde{\mathbf{X}}'\tilde{\mathbf{X}} + \lambda\mathbf{I}_p)^{-1}\tilde{\mathbf{X}}'(\mathbf{I}_n - \mathbf{K}).$$

ماتریس برازش با استفاده از پهنای باند ω_*^* است. با این حال یافتن مینیمم تابع GCV با استفاده از رابطه اخیر در هر تکرار از الگوریتم کم‌ترین توان‌های دوم پیراسته سریع مدت زمان اجرای الگوریتم را به گونه قابل توجهی افزایش می‌دهد (واضح است که مینیمم‌سازی تابع GCV خود یک فرایند زمان‌بر است). بنابراین برای برآورد پارامترهای پهنای باند و مرزبندی شده در مراحل بعد، از پارامتر اریب‌سازی پیشنهاد شده توسط هورل و کنار (۱۹۷۰) استفاده نموده و با l بار تکرار این الگوریتم برآوردگرهای نهایی بهینه استخراج می‌شوند. برای

⁷Generalized Cross Validation

در مرحله s ام قرار می‌دهیم: $s = 1, \dots, l$

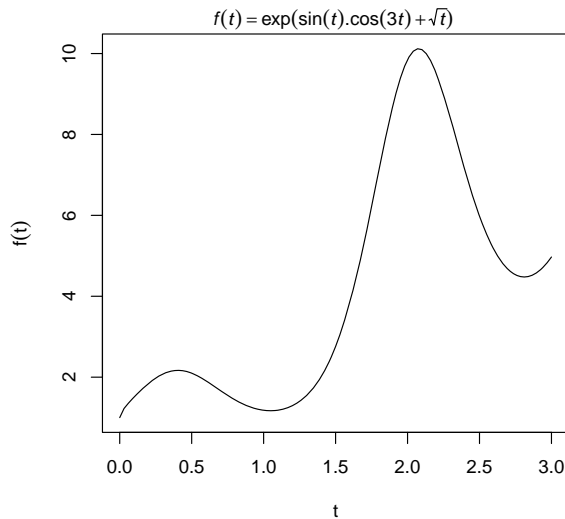
$$\hat{\sigma}_s^2 = \frac{\|S^{-1/2}(\mathbf{y} - \mathbf{X}\hat{\beta}_{FGLTS}(\mathbf{z}_{s-1}^*, \lambda_{s-1}^*) - \hat{f}_{FGLTS}(\mathbf{t}; \mathbf{z}_{s-1}^*, \lambda_{s-1}^*))\|^2}{n - \text{tr}(\mathbf{L}(\omega_{s-1}^*, \lambda_{s-1}^*))},$$

$$\lambda_s^* = \frac{p\hat{\sigma}_s^2}{\|\hat{\beta}_{FGLTS}(\mathbf{z}_{s-1}^*, \lambda_{s-1}^*)\|^2},$$

$$\omega_s^* = n^{-1/5} \left\{ \frac{\kappa_2 \hat{\sigma}_s^2}{M_2^2 J_2^2} \right\}^{1/5}.$$

۶ نتایج عددی

در این بخش برای بررسی کارایی برآوردهای معرفی شده، یک مطالعه شبیه‌سازی به روش مونت کارلو انجام شده و سپس این شیوه برآورد در یک مثال واقعی به کار خواهد رفت. این شبیه‌سازی و تحلیل داده‌های واقعی با استفاده از نرم افزار R نسخه ۲.۳.۳ انجام شده است.



شکل ۱: تابع ناپارامتری در مطالعه شبیه‌سازی

۱.۶ مطالعه شبیه‌سازی به روش مونت کارلو

در این بخش با استفاده از شبیه‌سازی داده‌ها به بررسی و مقایسه برآوردگر پیشنهادی با برآوردگر کلاسیک متداول در مدل نیمه‌پارامتری پرداخته می‌شود. برای رسیدن به ساختار وابستگی، متغیرهای پیشگو بنا بر مدل متغیرهای تصادفی نرمال استاندارد بوده و پارامتر γ تعیین کننده میزان همبستگی خطی بین دو متغیر توضیحی است. متغیر پاسخ از مدل رگرسیون نیمه‌پارامتری

$$y_i = \sum_{j=1}^6 x_{ij} \beta_j + f(t_i) + \epsilon_i, \quad i = 1, \dots, n, \quad (14)$$

با ۵ متغیر پیشگو برای $n = 150$ با $M = 10^3$ تکرار شبیه‌سازی می‌شود، به طوری که پارامترهای قسمت خطی و غیرخطی مدل (این تابع در شکل ۱ رسم شده است) $\beta = (-1, 4, 2, -5, -3)'$ و $\epsilon = (\epsilon_1^T, \epsilon_2^T)^T$ $f(t) = \exp\{\sin(2t) \cos(5t) + \sqrt{t}\}$ ، $t \in [0, 3]$ بوده و خطاها به صورت $\epsilon = (\epsilon_1^T, \epsilon_2^T)^T$ هستند، به طوری که

$$\epsilon_1 (h \times 1) \sim \mathcal{N}_h(0, \sigma^2 V), \quad \sigma^2 = 1/44, \quad v_{ij} = \exp(-9|i-j|),$$

$$\epsilon_2 ((n-h) \times 1) \stackrel{i.i.d.}{\sim} \chi_1^2(15),$$

که در آن نشان دهنده توزیع کای-دو غیر مرکزی با درجه آزادی ν و پارامتر غیرمرکزی δ است. دلیل اصلی انتخاب این ساختار برای خطاها، ایجاد تعدادی داده پرت به منظور ارزیابی برآوردگر پیشنهادی است. به همین منظور برای $h = \lceil [0.67n] \rceil$ خطای اول (مشاهدات خوب) از توزیع نرمال و بقیه (مشاهدات پرت) از توزیع مستقل غیرمرکزی کای-دو استفاده شد. غیرمرکزی بودن باعث منحرف شدن داده‌ها به قسمتی دور از مدل رگرسیون شده که در نتیجه برآورد غیراستوار را به سمت خود منحرف می‌کنند.

نتایج در جدول ۱ برای $n = 150$ و $\gamma = 0/90, 0/95$ با ۳۳ درصد نقاط دورافتاده گزارش شده است. در این جدول مقدار برآورد میانگین توان‌های دوم خطای برآوردگرها به صورت

$$\widehat{\text{mse}}(\hat{\beta}) = \frac{1}{M} \sum_{j=1}^M (\hat{\beta}_j - \beta)' (\hat{\beta}_j - \beta),$$

$$\widehat{\text{mse}}(\hat{f}(t)) = \frac{1}{M} \sum_{j=1}^M (\hat{f}_j(t) - f(t))^2,$$

جدول ۱: مقایسه برآوردگرهای پیشنهادی به روش مونت کارلو

$\gamma = 0/95$		$\gamma = 0/90$		برآوردگر
$\hat{\beta}_{FGLTS}(z^*, \lambda^*)$	$\hat{\beta}_{FG}$	$\hat{\beta}_{FGLTS}(z^*, \lambda^*)$	$\hat{\beta}_{FG}$	
-۱/۰۰۲۵	-۱/۰۱۱۹	-۱/۰۰۱۲	-۱/۰۰۶۴	
۳/۹۴۴۰	۳/۷۳۷۷	۳/۹۷۴۱	۳/۸۵۸۸	
۱/۸۳۲۰	۱/۲۱۳۱	۱/۹۲۲۳	۱/۵۷۶۴	میانگین ضرایب
-۴/۸۵۴۹	-۴/۳۲۰۴	-۴/۹۳۲۹	-۴/۶۳۴۲	
-۲/۹۲۳۶	-۲/۶۴۲۳	-۲/۹۶۴۷	-۲/۸۰۷۵	
۰/۹۹۰۳	۲/۳۵۷۱	۰/۳۶۷۱	۱/۱۷۷۶	$\widehat{mse}(\hat{\beta})$
۰/۳۱۳۹	۰/۴۷۲۹	۰/۲۹۵۹	۰/۵۵۶۲	$\widehat{mse}(\hat{f}(t))$

به دست می آید، که در آن $M = 10^3$ تعداد تکرارهای شبیه سازی است. همان طور که در جدول ۱ دیده می شود، برآوردگرهای قسمت خطی و غیرخطی استوار مرزبندی شده دارای میانگین توان دوم خطای کمتری نسبت به برآوردگرهای غیراستوار هستند. همچنین کمتر بودن اریبی برآوردگرهای قسمت خطی استوار مرزبندی شده به خوبی در این جدول قابل ملاحظه است.

۲.۶ تحلیل داده های واقعی

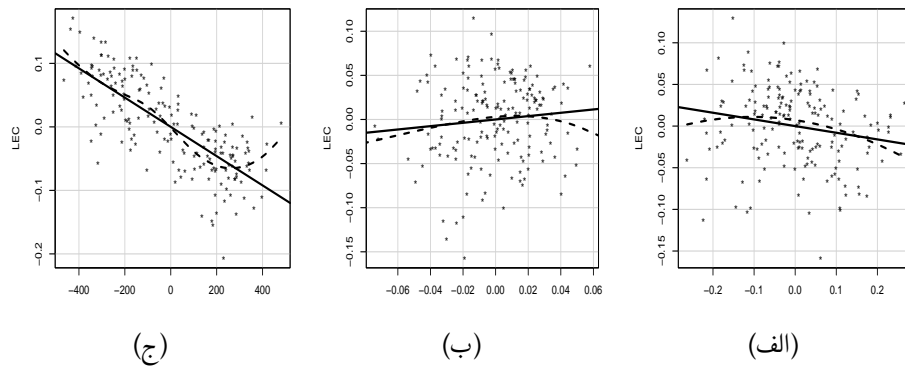
مجموعه داده ها مربوط به مصرف برق است که توسط آکدنیز دوران و همکاران (۲۰۱۲) مطالعه شده است. متغیرها برای $n = 177$ مورد به این ترتیب تعریف می شوند: متغیر پاسخ LEC لگاریتم مصرف ماهیانه برق به ازای هر نفر و متغیرهای پیشگو شامل LI لگاریتم درآمد هر فرد، $LREG$ لگاریتم نسبت قیمت برق به قیمت گاز و $Temp$ جمع شاخص متوسط دما (برای یک ماه خاص) هستند. داده ها مربوط به ۲۰ شهر از شهرهای آلمان است که از سازمان هواشناسی آلمان به دست آمده اند. در اینجا با استفاده از نمودار متغیر افزوده^۸ بخش ناپارامتری مدل را تشخیص دهیم. نمودار متغیر افزوده به طور شهودی اثر هر یک از متغیرهای پیشگو را پس از حذف اثر سایر متغیرهای پیشگو، بر متغیر پاسخ آشکار می کند. با نگاه کردن به نمودار متغیر افزوده در

^۸ Added Variable Plot

شکل ۲ از آنجا که شواهدی مبنی بر این که رابطه جزئی میان متغیر پاسخ مصرف برق و متغیر جمع شاخص میانگین دما ($Temp$) یک رابطه غیرخطی است از نمودار به دست می‌آید، این متغیر بعنوان جزء ناپارامتری مدل در نظر گرفته می‌شود، بنابراین مدل رگرسیون نیمه پارامتری به صورت

$$(LEC)_i = \beta_1(LI)_i + \beta_2(LREG)_i + f(Temp)_i + \epsilon_i, \quad (15)$$

است، که در آن ضریب همبستگی بین متغیرهای پیش‌بین قسمت خطی برابر 0.5697 و نشان‌دهنده وجود هم‌خطی بین آن‌ها است. همچنین شکل ۲ وجود مشاهده‌های دورافتاده زیادی را نشان می‌دهد که می‌توانند بر برآوردهای مرزبندی شده غیر استوار اثر بگذارند. شماره مشاهده‌های دورافتاده که تاثیر مخرب جدی بر برآوردگر کم‌ترین توان‌های دوم می‌گذارند عبارتند از: ۳، ۱۰، ۱۵، ۱۸، ۲۰، ۲۷، ۶۱، ۱۱۱، ۱۲۱، ۱۲۳، ۱۳۳، ۱۳۹، ۱۶۰ و ۱۷۲ که حدود ۸ درصد از کل مشاهدات را تشکیل می‌دهند.

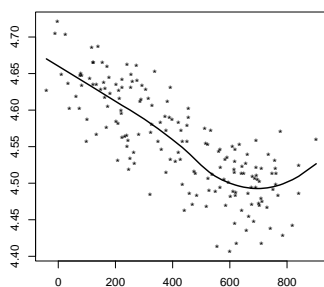


شکل ۲: نمودار افزوده متغیرهای پیشگو در مقابل متغیر پاسخ و برازش کم‌ترین توان‌های دوم (خط ممتد) و برازش ناپارامتری کرنل (خط چین)، الف- لگاریتم درآمد هر فرد، ب- لگاریتم نسبت قیمت برق به گاز، ج- میانگین دما

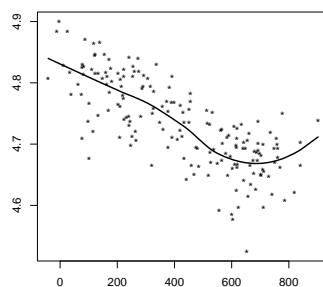
در جدول ۲ مقادیر برآوردگرها و همچنین R^2 و RSS که به ترتیب، مجموع توان‌های دوم مانده‌ها و ضریب تغییرات مدل رگرسیون نیمه پارامتری برازش شده هستند، گزارش شده‌اند. همان‌طور که در این جدول دیده می‌شود، برآوردگرهای استوار مرزبندی شده محتمل دارای مجموع توان دوم خطای کمتری نسبت به برآوردگرهای غیر استوار محتمل بوده و حدود ۴ درصد بیشتر تغییرات متغیر پاسخ را نسبت به برآوردگر غیر استوار بیان می‌کنند. تابع ناپارامتری برازش شده در شکل ۲ رسم شده است.

جدول ۲: مقایسه برآوردگرهای پیشنهادی در داده‌های واقعی

$\hat{\beta}_{FGLTS}(z^*, \lambda^*)$	$\hat{\beta}_{FG}$	برآوردگر
۰/۱۸۸۷	۰/۱۴۳۰	$\hat{\beta}_1$
-۰/۰۷۹۵	-۰/۰۸۱۱	$\hat{\beta}_2$
۰/۲۷۲۱	۰/۳۴۳۳	RSS
۰/۶۷۴۱	۰/۶۳۱۹	R ²



(ب)



(الف)

شکل ۳: الف: برازش تابع ناپارامتری به روش غیراستوار ب: برازش تابع ناپارامتری به روش استوار مرزبندی شده

بحث و نتیجه‌گیری

در این مقاله برآورد مرزبندی شده استوار محتمل ضرایب خطی و توابع غیرخطی در مدل رگرسیون نیمه‌پارامتری در حضور داده‌های پرت و هم‌خطی چندگانه به کمک الگوریتم سریع کم‌ترین توان‌های دوم مرزبندی شده پیراسته با فرض همبسته بودن جملات خطا، ارائه شد. نقطه فروریزش برآوردگرهای معرفی شده به دست آمد و از لحاظ نظری نشان داده شد که این نقطه فروریزش می‌تواند بیش از ۵۰ درصد باشد. با این حال از دیدگاه عملی دستیابی به نقطه فروریزش بیش از ۵۰ درصد امکان‌پذیر نیست.

نتایج عددی مطالعه شبیه‌سازی نشان می‌دهد که استفاده از برآوردگر استوار مرزبندی شده تعمیم‌یافته محتمل در مدل رگرسیون نیمه‌پارامتری، دارای تاثیر به سزایی در بهبود برآوردگرهای ضرایب خطی و توابع غیرخطی مدل در برازش مدل به داده‌های همبسته و دارای مشاهدات دورافتاده و دارای هم‌خطی چندگانه است.

تقدیر و تشکر

نویسندگان مقاله ضمن تشکر از اعضای محترم هیئت تحریریه مجله، از پیشنهادهای و نظرات ارزشمند داوران و ویراستار محترم مقاله که موجب ارتقاء سطح آن گردید کمال تشکر و قدردانی را دارند.

مراجع

نوروزی‌راد، م. و آرشی، م. (۱۳۹۶)، مطالعه رفتار حدی برآوردگرهای انقباضی در مدل رگرسیون تاوانیده با نرم مستطیلی، مجله علوم آماری ایران، ۱۱، ۱۴۹-۱۷۴.

Akdeniz Duran, E., Hardle, W. K. and Osipenko, M. (2012), Difference-based Ridge and Liu Type Estimators in Semiparametric Regression Models, *Journal of Multivariate Analysis*, **105**, 164-175.

Arashi, M., Janfada, M. and Norouzirad, M. (2015), Singular Ridge Regression with Stochastic Constraints, *Communication in Statistics - Theory and Methods*, **44**, 1281-1292.

Arashi, M. and Norouzirad, M. (2015), Improved Ridge M-estimators, 46th Annual Iranian Mathematics Conference, Yazd University, Yazd, Iran.

Arashi, M. and Valizadeh, T. (2015), Performance of Kibria's Methods in Partial Linear Ridge Regression Model, *Statistical Papers*, **56**, 231-246.

Alfons, A. Croux, C. and Gelper, S. (2013), Sparse Least Trimmed Squares Regression for Analyzing High-dimensional Large Data Sets, *The Annals of Applied Statistics*, **7**, 226-248.

Amini, M. and Roozbeh, M. (2015), Optimal Partial Ridge Estimation in Restricted Semiparametric Regression Models, *Journal of Multivariate Analysis*, **136**, 26-40.

- Blanchfower, D. G. and Oswald A. J. (1994), *The Wage Curve*, MIT Press Cambridge.
- Engle, R. F., Granger, C. W. J., Rice, J. and Weiss, A. (1986), Semiparametric Estimates of the Relation Between Weather and Electricity Sales, *Journal of the American Statistical Association*, **81**, 310-320.
- Eubank, R.L. (1999), *Nonparametric Regression and Spline Smoothing*, Second edition, Marcel Dekker, New York.
- Farebrother, R. (1976), Further Results on the Mean Square Error of Ridge Regression, *Journal of the Royal Statistical Society Ser. B*, **38**, 248-250.
- Golub, G., Heath, M. and Wahba, G. (1979), Generalized Cross Validation as a Method for Choosing a Good Ridge Parameter, *Technometrics*, **21**, 215-223.
- Hardle, W., Liang, H. and Gao, J. (2000), *Partially Linear Models*. Physika Verlag, Heidelberg.
- Hoerl, A. E. and Kennard, R. W. (1970), Ridge Regression: Biased Estimation for Non-orthogonal Problems, *Technometrics*, **12**, 55-67.
- Gao, J. T., Hong, S. Y. and Liang, H. (1995), Convergence Rates of a Class of Estimates in Partly Linear Models, *Acta Mathematica Sinica*, **38**, 658-669.
- Maronna, R. A., Martin, D. R. and Yohai, V. J. (2006), *Robust Statistics: Theory and Methods*, John Wiley, New York.
- Roosbeh, M. (2016), Robust Ridge Estimator in Restricted Semiparametric Regression Models, *Journal of Multivariate Analysis*, **147**, 127-144.
- Roosbeh, M. and Babaie-Kafaki, S. (2016), Extended Least Trimmed Squares Estimator in Semiparametric Regression Models with Correlated Errors, *Journal of Statistical Computation and Simulation*, **186**, 357-372.

- Rousseeuw, P. J. (1984), Least Median of Squares Regression, *Journal of the American Statistical Association*, **79**, 871-880.
- Rousseeuw, P. J. and Leroy, A. M. (2003), *Robust Regression and Outlier Detection*, 2nd ed. Wiley, Hoboken.
- Rousseeuw, P. J. and van Driessen, K. (2006), Computing LTS Regression for Large Data Sets, *Data Mining and Knowledge Discovery* **12**, 29-45.
- Speckman, P. (1988), Kernel Smoothing in Partial Linear Models, *Journal of the Royal Statistical Society Ser. B.*, **50**, 413-436.
- Tabakan, G. and Akdeniz, F. (2010), Difference-based Ridge Estimator of Parameters in Partial Linear Model, *Statistical Papers*, **51**, 357-368.
- Turkmen, A. S. and Tabakan, G. (2015), Outlier Resistant Estimation in Difference-based Semiparametric Partially Linear Models, *Communications in Statistics - Simulation and Computation*, **44**, 417-432.
- Willis, R. J. (1986), Wage Determinants: A Survey and Reinterpretation of Human Capital Earnings Functions in: Ashenfelter, O. and Layard, R. *The Handbook of Labor Economics*, Vol.1 *North Holland-Elsevier Science Publishers Amsterdam*, **1**, 525-602.
- Zellner, A. (1962), An Efficient Method of Estimating Seemingly Unrelated Regression and Tests for Aggregation Bias, *Journal of the American Statistical Association*, **58**, 977-992.
- Zellner, A. (1963), Estimators for Seemingly Unrelated Regressions: Some Finite Sample Results, *Journal of the American Statistical Association*, **58**, 977-992.

Feasible Generalized Ridge Robust Estimator in Semiparametric Regression Models

Mahdi Roozbeh¹, Morteza Amini²

¹Department of Statistics, Semnan University, Semnan, Iran.

²Department of Statistics, University of Tehran, Tehran Iran.

Abstract: In many fields such as econometrics, psychology, social sciences, medical sciences, engineering, etc., we face with multicollinearity among the explanatory variables and the existence of outliers in data. In such situations, the ordinary least-squares estimator leads to an inaccurate estimate. The robust methods are used to handle the outliers. Also, to overcome multicollinearity ridge estimators are suggested. On the other hand, when the error terms are heteroscedastic or correlated, the generalized least squares method is used. In this paper, a fast algorithm for computation of the feasible generalized least trimmed squares ridge estimator in a semiparametric regression model is proposed and then, the performance of the proposed estimators is examined through a Monte Carlo simulation study and a real data set.

Keywords: Breakdown point, Generalized cross validation, Least trimmed squares estimator, Outliers, Semiparametric regression model.

Mathematics Subject Classification (2010): 62G08, 62G35, 62J20.