

مدل‌های نیمه پارامتری استوار تنک در داده‌های با بعد بالا

مهدی روزبه^۱، منیره معنوی^۲

تاریخ دریافت: ۱۴۰۰/۰۴/۲۷

تاریخ پذیرش: ۱۴۰۰/۱۲/۲۸

چکیده:

تحلیل و مدل‌بندی داده‌های با بعد بالا یکی از چالش‌برانگیزترین مسائل روز دنیا است. تفسیر این داده‌ها کاری ساده نیست و نیازمند استفاده از روش‌های مدرن است. روش‌های تاوانیده یکی از رایج‌ترین راه‌های تحلیل داده‌های با بعد بالاست. همچنین مدل‌های رگرسیونی و تحلیل آن‌ها به شدت تحت تأثیر مشاهدات پرت قرار می‌گیرند. روش کمترین توان‌های دوم پیراسته یکی از بهترین روش‌های استوار برای از بین بردن تأثیر تخریبی این نقاط است. مدل‌های نیمه پارامتری که ترکیبی از هر دو نوع مدل‌های پارامتری و ناپارامتری هستند، بسیار انعطاف‌پذیرند. این مدل‌ها زمانی مفید هستند که هم بخش پارامتری و هم بخش ناپارامتری در مدل وجود دارد. هدف اصلی این مقاله تحلیل مدل‌های نیمه پارامتری در داده‌های با بعد بالا با حضور نقاط پرت با استفاده از روش لاسو تنک استوار است. در انتها، کارایی برآوردگر پیشنهادی با استفاده از تحلیل داده‌هایی واقعی در مورد تولید ویتامین B2 سنجیده می‌شود.

واژه‌های کلیدی: داده‌های با بعد بالا، روش لاسو، روش کمترین توان‌های دوم پیراسته، روش کمترین توان‌های دوم پیراسته تنک، مدل‌های نیمه پارامتری.

۱ مقدمه

دارد که بار محاسباتی برخی از روش‌ها می‌تواند به صورت نمایی با بعد افزایش یابد. دومین مورد که واسرمن آن را نفرین آماری ابعاد نامید به این صورت است که اگر داده‌ها دارای بعد d باشند، به نمونه‌ای با اندازه n مورد نیاز است که به صورت نمایی با d رشد کند.

در گذشته، کار با داده‌های کلاسیک که دارای حداکثر چند ده متغیر توضیحی بوده‌اند، بسیار ساده بود و روش‌های کلاسیک نظیر کمترین توان‌های دوم نتایج قابل قبولی ارائه می‌دادند؛ اما این روش‌ها در حضور داده‌های بعد بالا برآوردی ارائه نمی‌دهند. چراکه برای محاسبه برآوردگر کمترین توان‌های دوم یعنی $\hat{\beta}_{OLS} = (X^T X)^{-1} X^T y$ نیاز به محاسبه وارون ماتریس $X^T X$ است. ولی این ماتریس در این وضعیت رتبه کامل نبوده و در نتیجه وارون‌پذیر نخواهد بود.

بنابراین استفاده از روش‌های کلاسیک در این شرایط سودمند نخواهد بود و مستلزم به‌کارگیری روش‌های دیگری خواهیم بود. در تحلیل رگرسیونی داده‌های بعد بالا تعداد زیاد متغیرهای توضیحی، محقق را با چالشی جدی روبه‌رو می‌کند و همیشه جدالی بین دقت

امروزه با گسترش روزافزون علم، دانش و فناوری، روش‌های نوین و دقیقی برای اندازه‌گیری، جمع‌آوری و ثبت اطلاعات ابداع شده و این امر باعث ظهور و گسترش داده‌های بعد بالا شده است. از اواخر سال ۱۹۹۰ کار با این مجموعه داده‌ها یعنی داده‌هایی که در آن تعداد متغیرهای توضیحی (p) بسیار بیشتر از تعداد مشاهدات (n) است، شروع شد [۱۱]. افزایش بیش از پیش ابعاد داده‌ها به یک مسئله اساسی در مبحث داده‌کاوی تبدیل شده است. تجزیه، تحلیل و تفسیر این داده‌ها امری دشوار و بسیار قابل تأمل است که واسرمن [۲۶] در کتاب خود به این امر اشاره نموده و از این پدیده تحت عنوان نفرین ابعاد یاد کرده است. اصطلاحی که معمولاً به بلمن نسبت داده می‌شود. تقریباً این بدان معنی است که با افزایش بعد مشاهدات، برآورد پارامترها با سرعت زیادی، بسیار دشوار می‌شود. حداقل دو نوع از این نفرین وجود دارد. اولین مورد، نفرین محاسباتی ابعاد است. این امر به این واقعیت اشاره

^۱ هیئت علمی گروه آمار، دانشگاه سمنان، سمنان، ایران (نویسنده مسئول: mahdi.roozbeh@semnan.ac.ir)

^۲ دانش‌آموخته کارشناسی ارشد آمار، دانشگاه سمنان، سمنان، ایران.

^۳ Curse of dimensionality

^۴ Overfitting

و ناپارامتری^۱ تقسیم کرد. مدل‌های رگرسیونی پارامتری صرفاً به منظور بررسی ارتباطات خطی میان متغیر پاسخ و متغیرهای توضیحی به کار گرفته می‌شوند. بدین سبب کار کردن با آن‌ها بسیار ساده و لذت‌بخش است و تحلیل نتایج آن پیچیدگی عجیبی ندارد. مدل‌های رگرسیونی ناپارامتری زمانی که تمام متغیرهای توضیحی ارتباط غیرخطی با متغیر پاسخ دارند، مورد استفاده قرار می‌گیرند. مدل‌های ناپارامتری ممکن است، نماینده بهتری برای مجموعه داده‌های واقعی باشند، اما متأسفانه تجزیه و تحلیل آن‌ها کار ساده‌ای نیست. مدل‌های نیمه پارامتری از دو قسمت پارامتری و ناپارامتری تشکیل شده‌اند. مدل‌های رگرسیونی نیمه پارامتری انعطاف پذیرتر از مدل‌های رگرسیونی ناپارامتری و پارامتری هستند زیرا هم ویژگی‌های مدل رگرسیونی پارامتری را دارا هستند و هم ویژگی‌های مدل رگرسیونی ناپارامتری [۳]. به عبارت دیگر مدل رگرسیونی نیمه پارامتری انعطاف پذیری یک مدل رگرسیونی ناپارامتری و قدرت توضیحی یک مدل رگرسیونی پارامتری را حفظ می‌کند. انگل و همکارانش [۱۲] برای نخستین بار این مدل‌ها را به منظور مدل‌سازی مصرف ماهیانه برق خانوارها ابداع کردند و به طور رسمی آن را در آمار به کار بردند. متغیرهای مورد بررسی آن‌ها مصرف برق ماهیانه خانوارها، درآمد ماهیانه خانوارها، قیمت ماهیانه برق و دمای هوا بود. متغیر مصرف ماهیانه برق خانوارها متغیر پاسخ و سایر متغیرها، توضیحی هستند. در بررسی‌های اولیه انگل و همکارانش دریافتند که رابطه بین درجه حرارت و مصرف ماهیانه برق بسیار غیرخطی است؛ زیرا مصرف ماهیانه برق هم در درجه حرارت پایین و هم در درجه حرارت بالا افزایش می‌یابد که این ارتباط قابل چشم‌پوشی نبوده چراکه این ارتباط بسیار مهم بوده و فروش ماهیانه برق در شرایط عادی هوا بسیار متفاوت از شرایط غیرعادی بوده است. از سوی دیگر ارتباط میان مصرف ماهیانه برق و درآمد خانوارها و همچنین ارتباط میان مصرف ماهیانه برق و قیمت ماهیانه برق نیز خطی بوده است. به همین سبب برای مدل‌سازی این مورد نه استفاده از مدل‌های پارامتری و نه به کارگیری مدل‌های ناپارامتری راهگشا نبود. انگل و همکارانش در این شرایط وادار به استفاده از روش دیگری بودند و این امر باعث ظهور مدل‌های نیمه پارامتری شد. این مدل‌ها به سرعت جایگاه خود را در حیطه‌های مختلف علوم یافتند به طوری که کاربردهایشان در علوم اقتصادی و تحلیل بقا زیانزد است. مدل‌های نیمه پارامتری پیشرفت

سرعت و هزینه وجود دارد. از سوی دیگر، وجود متغیرهایی در مدل که با متغیر پاسخ ارتباطی ندارند، منجر به بیش‌برازش^۴ می‌شود. چنین مدلی، مدلی بسیار پیچیده برای داده‌ها است. در این شرایط، مدل با تغییرات جهشی سعی در پوشش داده‌های حاصل از نمونه و حتی مقدارهای خطاها می‌کند. در حالی که این مدل باید منعکس‌کننده رفتار جامعه باشد. به کارگیری تمامی متغیرهای توضیحی زمانی طولانی صرف کرده و هزینه‌های محاسباتی بسیاری را بر محقق تحمیل می‌کند. در این شرایط این ابهام مطرح می‌شود که شاید حضور تمامی متغیرها برای برازش مدل الزامی نباشد و بدین سان اغلب محققان با در نظر گرفتن فرض تنکی^۵ سعی در کاهش بعد دارند و از زیرمجموعه‌ای مناسب از متغیرها برای برازش مدل کمک می‌گیرند و از این رو علاقه‌مند به شناسایی این زیرمجموعه مناسب هستند. روش‌های متنوعی برای رویارویی با داده‌های با بعد بالا وجود دارد. روش کمترین توان‌های دوم تاوانیده^۶ یکی از روش‌های بسیار مفید برای تحلیل داده‌های با بعد بالا است.

نقاط پرت یکی دیگر از معضلات مدل‌های رگرسیونی است. بارنت و لوئیس [۹] نقطه پرت را مشاهده‌ای (مجموعه‌ای از مشاهدات) تعریف کردند که به نظر می‌رسد دارای باقی‌مانده‌ای متفاوت با باقی‌مانده‌ی سایر نقاط داده‌ها باشد. این نقاط، تأثیرات بسیار بدی روی عملکرد مدل‌های رگرسیونی دارند. روش‌های استوار^۷، روش‌هایی هستند که برای از بین بردن این اثرات استفاده می‌شوند. روش کمترین توان‌های دوم پیراسته^۸ یکی از مشهورترین روش‌های استوار است که در این مقاله از آن استفاده می‌شود. این روش توسط روسو و لروی [۲۲] معرفی شد.

مدل‌های پارامتری با تمام مزایا و محاسن خود، نمی‌توانند نماینده خوبی از آنچه در دنیای واقعی رخ می‌دهد، باشند. در دنیای واقعی به ندرت می‌توان مجموعه داده‌هایی که در آن ارتباط خطی بین متغیر توضیحی و متغیر پاسخ وجود داشته باشد، یافت. در عمل مجموعه داده‌های دنیای واقعی بسیار پیچیده و دارای ارتباطات غیرخطی‌اند که در مواردی حتی با به کارگیری تبدیل‌های قوی نظیر تبدیل باکس-کاکس، لگاریتمی و ... نیز ارتباط مابین آن‌ها خطی نمی‌شود. در این شرایط ناگزیر به استفاده از روش‌های دیگری خواهیم بود.

مدل‌های رگرسیونی را می‌توان به سه دسته پارامتری، نیمه پارامتری^۹

⁵Sparsity

⁶Penalized least squares method

⁷Robust methods

⁸Least trimmed squares method

⁹Semiparametric

¹⁰Nonparametric

است. در مدل (۱) متغیر پاسخ y به طور خطی با متغیرهای توضیحی x_1, \dots, x_p در ارتباط است و با متغیر توضیحی t ارتباط غیرخطی دارد. واضح است که به دلیل دو قسمتی بودن این مدل‌ها (پارامتری و ناپارامتری) مسئله برآورد پارامترها کمی پیچیده خواهد شد. در میان کتب و مقالات آماری روش‌های متفاوتی به منظور برآورد پارامترهای این مدل‌ها توسط محققان مختلف ارائه شده است. یکی از روش‌های جالب و پرکاربرد روش پیشنهادی فن و همکارانش [۱۴] است. او و همکارانش طی دو مرحله پارامترهای مدل را برآورد نمودند. به این ترتیب که در مرحله اول قسمت پارامتری مدل و در مرحله دوم قسمت ناپارامتری مدل برآورد می‌شود. قدم اول به منظور برآورد پارامترها به روش فن [۱۴] معلوم در نظر گرفتن بردار پارامتر β است. با این فرض می‌توان نوشت

$$y - X\beta = f(t) + \varepsilon \quad (2)$$

اکنون با بازنویسی رابطه (۲) خواهیم داشت

$$\hat{y} = f(t) + \varepsilon \quad (3)$$

به طوری که در آن $\hat{y} = y - X\beta$ است. بدیهی است که مدل (۳) یک مدل ناپارامتری است. کافی است که تابع f برآورد شود؛ اما برآورد قسمت ناپارامتری خود امری بسیار مهم و حائز اهمیت است. روش‌های گوناگونی به منظور ارائه قسمت ناپارامتری و تجزیه و تحلیل مدل‌های ناپارامتری وجود دارد. در این مقاله از روش کرنل برای برآورد قسمت ناپارامتری استفاده شده است. تابع کرنل یک تابع تک مدی است که بیشترین وزن را به نقطه مورد نظر (نقطه‌ای که قرار است در آن تابع برآورد شود)، در اینجا t می‌دهد. غالباً تابع کرنل یک تابع مقارن است که به فاصله‌های مساوی از هر دو سمت وزن یکسان اختصاص می‌دهد، البته این امر یک قاعده کلی نیست و $K(\cdot)$ می‌تواند وزن‌های غیر یکسان نیز توزیع نماید. کرنل یک برآوردگر اریب با واریانس متناهی است. کرنل $K(\cdot) \geq 0$ در حقیقت تابع همواری است که سه ویژگی

$$\int K(t)d(t) = 1, \int tK(t)d(t) = 0, 0 < \sigma_t^2 \equiv \int t^2 K(t)d(t) < \infty$$

را دارا باشد.

با ظهور کرنل، کاربرد فراوان و محبوبیت آن سبب پیدایش توابع مختلف شد. از محبوب‌ترین کرنل‌ها می‌توان به کرنل یکنواخت

بسیار سریعی داشته و مورد توجه بسیاری از محققان و کارشناسان آمار قرار گرفته است. اسپکمن [۲۳] برآورد باقی‌مانده‌ی جزئی پارامتر قسمت خطی و تابع ناپارامتری، واریانس و اریبی مجانبی برآوردگرها را به دست آورد. بونیا [۱۰] روش‌های انتخاب متغیرهای کمکی سازگار را در مدل‌های نیمه‌پارامتری بر اساس روش کمترین توان‌های دوم تاوانیده پیشنهاد نمود. او نشان داد که برآوردگر انتخاب شده از قسمت خطی به طور مجانبی نرمال است. آکدیز و همکاران [۴] نیز با استفاده از این مدل‌ها اقدام به مدل‌سازی مصرف ماهیانه برق در حضور متغیرهای توضیحی درآمد، دمای هوا و نسبت نرخ برق به گاز کردند. روزبه و آرشی [۲۱]، آرشی و همکاران [۷]، امینی و روزبه [۶]، روزبه [۱۹، ۲۰]، یوزباشی و همکاران [۲۵] و روزبه و چاچی [۳] در این حیطه تحقیقات ارزشمندی انجام دادند.

برای مشاهده کامل مدل‌های نیمه‌پارامتری و کاربردهایش کتاب هاردل [۱۶] منبع بسیار خوبی است که به خواننده پیشنهاد می‌شود.

در بخش ۳ روش کمترین توان‌های دوم تاوانیده معرفی می‌شود. روش کمترین توان‌های دوم پیراسته در بخش ۴ بیان می‌شود. در بخش ۵ مدل نیمه‌پارامتری کمترین توان‌های دوم پیراسته روی داده‌های مربوط به ریپولادین به کار گرفته می‌شود.

۲ مدل‌های نیمه پارامتری

مدل رگرسیونی نیمه‌پارامتری به صورت

$$y = X\beta + f(t) + \varepsilon, \quad (1)$$

تعریف می‌شود که در آن

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix},$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}, f(t) = \begin{bmatrix} f(t_1) \\ f(t_2) \\ \vdots \\ f(t_n) \end{bmatrix}, \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix},$$

که متغیر پاسخ y ، ماتریس طرح X ، بردار پارامتر مجهول β ، $f(\cdot)$ تابع حقیقی مقدار مجهول از مقادیر متغیر t ، و بردار خطا ε با شرایط

$$E(\varepsilon) = 0, \quad E(\varepsilon\varepsilon^T) = \sigma^2 \mathbf{I}_{(p+1) \times (p+1)}$$

می‌شد. همان‌طور که در بخش قبل مطرح شد روش محبوب کمترین توان‌های دوم در مواجهه با داده‌های با بعد بالا قادر به محاسبه برآورد ضرایب نبود و بدین‌سان ایده تاوانیدن این برآوردگر میان محققان آماری شکل گرفت. در واقع برآوردگرهای تاوانیده حاصل از به‌کارگیری روش‌های کمترین توان‌های دوم یا ماکسیمم درستنمایی تاوانیده در مدل‌بندی رگرسیونی است. در حقیقت این برآوردگرها از روش بهینه‌سازی یک تابع درجه دوم (البته نه همیشه) نسبت به یک برآوردگر تاوانیده به دست می‌آید. ایده اصلی در این برآوردگرها، این است که ضرایب برآورد شده متغیرهای توضیحی کم‌اهمیت در یک مدل تنک، به سمت صفر منقبض شود. در اکثر روش‌های تاوانیده، این متغیرها از مدل خارج می‌شوند و همین امر موجب ساده‌تر شدن تفسیرپذیری مدل خواهد شد [۱].

مسئله کمینه‌سازی در روش کمترین توان‌های دوم معمولی در مدل فوق به صورت

$$\min_{\beta} \{(\tilde{y} - \tilde{X}\beta)^T(\tilde{y} - \tilde{X}\beta)\}, \quad (5)$$

است. با محاسبه مشتق تابع هدف عبارت (۵) نسبت به بردار پارامتر β و برابر صفر قرار دادن آن و حل معادله حاصل، برآوردگر کمترین توان‌های دوم به‌سادگی یافت می‌شود. به عبارت دیگر در این روش خطا برحسب β کمینه می‌شود. واضح است که اگر مقادیر واقعی β بزرگ باشد، برآوردگر حاصل، فارغ از احتساب بزرگی β ، فاصله زیادی از مقدار واقعی β خواهد داشت و به همین سبب خطای برآورد زیاد شده و دقت آن کاهش می‌یابد [۱]. یک روش بسیار سودمند و تقریباً ساده برای مقابله با این مشکل، تاوانیدن مقادیر بزرگ β است بدین‌سان به‌جای کمینه کردن تابع هدف در روش کمترین توان‌های دوم، مسئله کمینه‌سازی تاوانیده یعنی

$$\min_{\beta} \{(\tilde{y} - \tilde{X}\beta)^T(\tilde{y} - \tilde{X}\beta)\}, \quad \text{s.t. } p(\beta) \leq t, \quad (6)$$

جایگزین می‌شود. عدد ثابت $t \geq 0$ پارامتر تنظیم‌کننده^{۲۰} نامیده می‌شود که اگر برابر صفر باشد، مدل تنها شامل عرض از مبدأ است و اگر بی‌نهایت باشد، مدل کامل (با حضور تمامی متغیرها) خواهد بود. $p(\beta)$

(مستطیلی)^{۱۱}، کرنل مثلثی^{۱۲}، کرنل درجه چهارم^{۱۳}، کرنل آپاشینکو^{۱۴}، کرنل سه وزنی^{۱۵}، کرنل سه مکعبی^{۱۶}، کرنل کسینوسی^{۱۷}، کرنل لاپلاسیان گوسی^{۱۸} و کرنل نرمال (گاوسی)^{۱۹} اشاره کرد.

همان‌طور که ذکر شد هدف یافتن برآورد $f(\cdot)$ است. با فرض هموار بودن (دارا بودن مشتق اول و دوم) $f(\cdot)$ می‌توان نوشت

$$\hat{f}(t) = \sum_{i=1}^n k_i(t)(y_i - \mathbf{x}_i^T \beta)$$

که در آن $k_i(\cdot)$ تابع کرنل و \mathbf{x}_i^T سطر i ام ماتریس \mathbf{X} است. در حقیقت در این روش برای برآورد قسمت ناپارامتری از میانگین وزنی استفاده می‌شود. اکنون به منظور برآورد بردار ضرایب رگرسیونی β ، \hat{f} در مدل (۱) جایگذاری می‌شود بدین ترتیب به‌سادگی می‌توان بردار پارامترها را برآورد نمود. به عبارت دیگر می‌بایست مدل را (۱) به صورت

$$\mathbf{y} - \hat{f}(t) = \mathbf{X}\beta + \varepsilon \quad (4)$$

بازنویسی نمود. اکنون مدل (۴) یک مدل پارامتری است که در درون خود اثر ناپارامتری را پنهان نموده است. با ساده‌سازی مدل (۴) خواهیم داشت:

$$y_i - \sum_{j=1}^n k_j(t_i)y_j = (\mathbf{x}_i^T - \sum_{j=1}^n k_j(t_i)\mathbf{x}_j^T)\beta + \varepsilon_i \\ \Rightarrow \tilde{y}_i = \tilde{\mathbf{x}}_i^T \beta + \varepsilon_i, i = 1, \dots, n.$$

که در آن $\tilde{\mathbf{x}}_i^T = \mathbf{x}_i^T - \sum_{j=1}^n k_j(t_i)\mathbf{x}_j^T$ و $\tilde{y}_i = y_i - \sum_{j=1}^n k_j(t_i)y_j$ است. با در نظر گرفتن فرم ماتریسی مدل به صورت $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta + \varepsilon$ برآورد ضرایب رگرسیونی با استفاده از روش کمترین توان‌های دوم معمولی به صورت زیر است.

$$\hat{\beta} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}.$$

۳ روش کمترین توان‌های دوم تاوانیده

با ظهور داده‌های با بعد بالا و ناتوانی روش‌های کلاسیک برای تحلیل و تفسیر مدل‌های رگرسیونی نیاز به معرفی روش‌های جدیدی احساس

¹¹Uniform kernel or boxcar kernel or rectangular kernel

¹²Triangular kernel

¹³Quartic kernel or biweight kernel

¹⁴Epanechnikov kernel

¹⁵Triweight kernel

¹⁶Tricube kernel

¹⁷Cosine kernel

¹⁸Laplacian of Gaussian kernel

¹⁹Gaussian kernel or normal kernel

²⁰Tuning parameter

در ادامه به معرفی چندین روش تاوانیده پرداخته می‌شود.

۱.۳ رگرسیون لاسو

رگرسیون لاسو^{۲۴} توسط رابرت تیبیشیرانی در سال ۱۹۹۶ ابداع شد. ایده اولیه لاسو برگرفته از پیشنهاد جالب لئو برایمن بود [۲۴]. روش گروتی برخی از متغیرها را حذف می‌کند و مابقی آن‌ها را منقبض می‌کند؛ بنابراین این روش مدل‌های قابل تفسیر ارائه می‌دهد. این روش نسبتاً پایدار است و همچنین مقیاس پایاست [۲۴]. وابستگی گروتی به علامت و مقدار برآورد کمترین توان‌های دوم نقص بزرگی است که بی‌اعتمادی محققان به این روش را به ارمغان آورده است. بدین جهت می‌توان گفت زمانی که با داده با بعد بالا مواجه باشیم و یا زمانی که متغیرها بسیار همبسته‌اند، برآوردهای کمترین توان‌های دوم ضعیف عمل کرده و به تبع آن گروتی نیز عملکرد ضعیفی داشته و مورد اطمینان نیست. این در حالی است که لاسو برخلاف گروتی از استفاده مستقیم از برآوردهای کمترین توان‌های دوم اجتناب می‌کند. مسئله کمینه‌سازی لاسو به صورت

$$\min_{\beta} \left\{ \sum_{i=1}^n (\tilde{y}_i - \sum_{j=1}^p \tilde{x}_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\},$$

است.

این محدودیت به‌طور طبیعی تمایل به تولید ضرایبی که صفر هستند، دارد و از این رو مدل‌های قابل تفسیر ارائه می‌دهد. روش لاسو برآوردهای پایدار، اریب و پیوسته تولید می‌کند [۲۴].

دقیقاً روی یکی از گوشه‌های لوزی باشد، (همانند شکل ۱) تنها یک متغیر انتخاب شده و دیگری صفر می‌شود. همان‌طور که ذکر شد، در روش لاسو به سبب نوع ناحیه تاوان (لوزی و گوشه داشتن لوزی) برخی از متغیرها دقیقاً برابر صفر می‌شوند. از منظر دیگر می‌توان استدلال کرد که چون تابع تاوان روش لاسو در نقطه صفر مشتق‌ناپذیر است، مقادیر بزرگ‌تری از λ منجر به حذف متغیرهای پیشگوی کم‌تأثیر در مدل می‌گردد و در نتیجه عمل برآورد پارامتر و انتخاب متغیر به‌صورت هم‌زمان انجام می‌گیرد [۲].

از معایب روش لاسو می‌توان به عملکرد ضعیف آن در حضور متغیرهای

نیز تابعی از بردار پارامتر است که تابع تاوان^{۲۱} نامیده شده و با توجه به نوع آن روش‌های تاوانیده متفاوتی ایجاد می‌شود. برای یافتن کمینه (بیشینه) یک تابع چند متغیره که با یک یا چند محدودیت مواجه است، روش لاگرانژ به کار گرفته می‌شود و با استفاده از روش لاگرانژ چند متغیره مسئله کمینه‌سازی (۶) به صورت زیر بازنویسی می‌شود [۱۳]

$$\min_{\beta} \left\{ (\tilde{y} - \tilde{X}\beta)^T (\tilde{y} - \tilde{X}\beta) + \lambda p(\beta) \right\}. \quad (7)$$

در عبارت (۷) پارامتر نامنفی λ ، پارامتر منظم‌سازی^{۲۲} یا پارامتر تاوان نامیده می‌شود. اگر λ بی‌نهایت باشد، تنها عرض از مبدأ در مدل حضور خواهد داشت و اگر برابر صفر باشد، قسمت تاوان از مدل حذف شده و مدل کمترین توان‌های دوم حاصل خواهد شد.

پارامترهای λ و t نقش یکسانی در مدل ایفا می‌کنند و کنترل بزرگی یا کوچکی ناحیه محدودیت (تاوان) به عهده این دو پارامتر است. جالب‌توجه است که رابطه این دو پارامتر از لحاظ بزرگی معکوس است. لازم به ذکر است که در دو مسئله کمینه‌سازی (۶) و (۷) از تابع زیان توان دوم خطا^{۲۳} استفاده شده است که بیانگر روش کمترین توان‌های دوم تاوانیده است؛ اما به‌طور کلی مسئله کمینه‌سازی تاوانیده به صورت زیر است:

$$\min_{\beta} \left\{ L(\tilde{y}, \tilde{X}\beta) + \lambda p(\beta) \right\},$$

که در آن $L(\tilde{y}, \tilde{X}\beta)$ تابع زیان نامنفی برای نیکویی برازش، $p(\beta)$ تابع تاوان نامنفی و λ پارامتر منظم‌سازی است که تعادل بین نیکویی برازش و پیچیدگی مدل (حضور متغیرهای زیاد) را برقرار می‌کند [۱۷].

لاسو به معنی عملگر انتخاب و کمترین قدرمطلق انقباضی^{۲۵} است. لاسو برخی ضرایب را منقبض و بقیه را صفر می‌کند.

برای یافتن درک بصری از روش لاسو با فرض وجود تنها دو متغیر توضیحی در مدل رگرسیونی، ناحیه محدودیت به صورت $|\beta_1| + |\beta_2| \leq t$ یک لوزی است. در شکل ۱ لوزی زرد رنگ ناحیه تاوان است. محل برخورد کانتورهای RSS (مجموع توان‌های دوم باقی‌مانده‌ها) روش کمترین توان‌های دوم با لوزی، مختصات برآوردگر لاسو را نمایش می‌دهد. با عمود کردن از این نقطه بر محور افقی، $\hat{\beta}_1$ و با عمود کردن از آن بر محور عمودی، $\hat{\beta}_2$ حاصل می‌شود. اگر محل برخورد

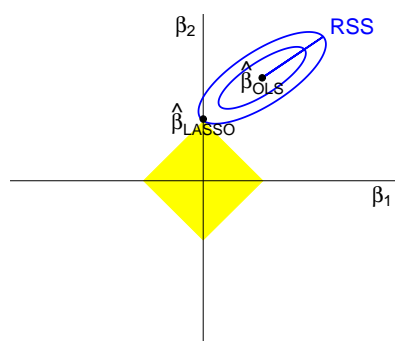
²¹Penalty function

²²Regularization parameter

²³Square error loss function

²⁴Lasso

²⁵Least Absolute Shrinkage and Selection Operator (LASSO)



شکل ۱: ناحیه محدودیت لاسو.

توضیحی همبسته اشاره نمود. در این شرایط لاسو بین متغیرهای همبسته خصوصاً متغیرهایی که به صورت گروهی همبسته‌اند تنها یکی از آن‌ها را برمی‌گزیند و متغیر انتخاب شده لزوماً بهترین متغیر نیست. همچنین، لاسو برای یک مدل رگرسیونی خطی با p متغیر توضیحی و n مشاهده، حداکثر n متغیر را انتخاب می‌کند، بنابراین اگر متغیرهای بیشتری (بیشتر از n) در مدل معنی‌دار باشند، توسط لاسو انتخاب نمی‌شوند.

در روش کمترین توان‌های دوم پیراسته تعداد h تا از کوچک‌ترین توان‌های دوم باقی‌مانده‌های ($\hat{\varepsilon}_i = e_i = y_i - \hat{y}_i$) مسئله کمینه‌سازی مدل رگرسیونی خطی چندگانه به صورت

$$\min \left\{ \sum_{i=1}^h (e^{\downarrow})_{i:n} \right\}, \quad (8)$$

در نظر گرفته می‌شود، به طوری که $(e^{\downarrow})_{1:n} \leq (e^{\downarrow})_{2:n} \leq \dots \leq (e^{\downarrow})_{n:n}$ توان دوم باقی‌مانده‌های مرتب‌شده هستند. لازم به ذکر است که ابتدا باقی‌مانده‌ها به توان دو می‌رسند، سپس مرتب می‌شوند. برخی h را پارامتر پیراسته^{۲۶} می‌نامند؛ اما برای تعیین این‌که مشاهده‌ی i ام یک مشاهده‌ی خوب است یا نه لازم است تابع نشانگر z_i را به صورت زیر تعریف نماییم.

$$z_i = \begin{cases} 1 & \text{مشاهده } i \text{ ام دورافتاده نباشد} \\ 0 & \text{مشاهده } i \text{ ام دورافتاده باشد} \end{cases}$$

و در نهایت ماتریس Z که یک ماتریس قطری است به صورت زیر حاصل

$$Z = \begin{bmatrix} z_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & z_n \end{bmatrix}.$$

بنابراین می‌توان مسئله کمینه‌سازی (۸) را به صورت

$$\begin{aligned} & \min_{\beta} \left\{ (\tilde{y} - \tilde{X}\beta)^T Z (\tilde{y} - \tilde{X}\beta) \right\}; \\ & e^T z = h, \text{ به طوری که} \\ & z_i \in \{0, 1\}, \quad i = 1, \dots, n, \end{aligned} \quad (9)$$

بازنویسی کرد که در آن $e = (1, \dots, 1)_{n \times 1}^T$ و $z = (z_1, \dots, z_n)^T$ است. با حل مسئله بهینه‌سازی (۹) برآوردگر کمترین توان‌های دوم پیراسته حاصل می‌شود. اگر $h = n$ باشد، کمترین توان‌های دوم پیراسته معادل کمترین توان‌های دوم معمولی خواهد بود.

یک نکته بحث‌برانگیز در این روش انتخاب پارامتر پیراسته است؛ زیرا h می‌تواند از $0, \dots, n$ باشد اگر این پارامتر به درستی انتخاب نشود، یا اطلاعات از دست می‌رود و یا نتایج حاصل تحریف می‌شود. روسو و لروی [۲۲] بایان قضیه‌ای در کتاب خود نشان دادند که با انتخاب $\left[\left[\frac{p+1}{4} \right] \right] + \left[\left[\frac{n}{4} \right] \right]$ برای h ، نقطه‌ی فروریزش (معیار استواری یک برآوردگر) این روش به ماکزیم مقدار ممکن خود (۵۰ درصد) خواهد رسید؛ یعنی حتی اگر نیمی از مشاهدات پرت باشند، این روش همچنان استوار خواهد ماند؛ اما این امر همیشه و در هر شرایطی صادق نیست، به خصوص در داده با بعد بالا ضعف این قضیه به وضوح آشکار می‌شود.

²⁶Trimmed parameter

۱.۴ نقطه فروریزش برآوردگر کمترین توان‌های دوم پیراسته

با استفاده از قضیه زیر که توسط آلفونز و همکاران [۵] مطرح شد، نقطه شکست یا نقطه فروریزش^{۲۸} برآوردگر کمترین توان‌های دوم پیراسته حاصل می‌شود.

قضیه ۱.۴. فرض کنید $\rho(x)$ یک تابع زیان متقارن و محدب باشد، به طوری که $\rho(0) = 0$ و به ازای $x \neq 0$ ، $\rho(x) > 0$. همچنین $\rho(x) = (\rho(x_1), \dots, \rho(x_n))^T$ زیرمجموعه‌ای با اندازه $h \leq n$ و برآوردگر رگرسیونی

$$\hat{\beta} = \min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^h (\rho(\tilde{y}_i - \tilde{X}_i \beta))_{i:n} + h\lambda \sum_{j=1}^p |\beta_j| \right\},$$

در نظر بگیرید، به طوری که

$$(\rho(\tilde{y}_i - \tilde{X}_i \beta))_{1:n} \leq \dots \leq (\rho(\tilde{y}_i - \tilde{X}_i \beta))_{n:n}$$

تابع زیان مرتب‌شده هستند، در این صورت نقطه فروریزش برآوردگر رگرسیونی $\hat{\beta}$ به صورت

$$\varepsilon^*(\hat{\beta}, Z^*) = \frac{n - h + 1}{n},$$

است که در آن Z^* به عنوان تمامی نمونه‌هایی که دارای نقاط پرت و تأثیرگذار هستند، در نظر گرفته می‌شوند. این نمونه‌ها جایگزین هر زیرمجموعه‌ی m تایی از مشاهدات، در داده‌ها با مقادیر دلخواه می‌شوند.

با استفاده از قضیه ۱.۴ و با فرض این که $\rho(x) = x^2$ باشد، نقطه فروریزش برآوردگر کمترین توان‌های دوم پیراسته تنک برابر $\frac{n - h + 1}{n}$ است. مقادیر کوچک h نقطه فروریزش بالاتر دارد و با فرض این که h به اندازه کافی کوچک باشد، نقطه فروریزش بزرگ‌تر از ۵۰٪ است [۵].

۵ نتایج کاربردی با داده ریبولوین

در این بخش به منظور بررسی و تحلیل مدل‌های نیمه پارامتری پیراسته در ابعاد بالا از داده‌های ریبولوین استفاده می‌شود. ریبولوین که به عنوان ویتامین B2 نیز شناخته می‌شود، یکی از ویتامین‌های B است که همه محلول در آب هستند. از مهم‌ترین ویتامین‌های لازم برای حیات موجود زنده است. دلیل انتخاب این نام، رنگ زرد این ویتامین است که ناشی از وجود حلقه فلاوینی موجود در ساختمان آن است. ریبولوین

²⁷Sparse least trimmed squares

²⁸Breakdown point

یکی دیگر از معایب این روش یافتن h مشاهده عادی نیز امری بسیار دشوار و زمان‌بر است، زیرا باید $\binom{n}{h}$ مدل مختلف تشکیل شود و باهم مقایسه شوند. همین امر موجب پیچیدگی روش کمترین توان‌های دوم پیراسته و عدم استقبال محققان از آن تا قبل از تدوین بسته‌های نرم‌افزاری بود.

۴ روش کمترین توان‌های دوم پیراسته تنک

همان‌طور که ذکر شد روش لاسو تفسیر مدل‌های با بعد بالا را ساده می‌کند اما در حضور نقاط پرت مقاوم نیست و عملکرد ضعیفی دارد. به طور مثال نقطه فروریزش روش لاسو $\frac{1}{n}$ است؛ یعنی تنها حضور یک مشاهده پرت می‌تواند کلیه نتایج را زیر سؤال برده و بی‌اعتماد کند. علاوه بر این روش کمترین توان‌های دوم پیراسته در حالتی که $h < p$ باشد، نمی‌تواند محاسبه شود [۵].

روش‌های تاوانیده استوار، به عنوان روش‌های کمترین توان‌های دوم پیراسته تنک^{۲۷} نیز مشهور است. مسئله بهینه‌سازی این روش به صورت

$$\min_{\beta} \left\{ (\tilde{y} - \tilde{X}\beta)^T Z (\tilde{y} - \tilde{X}\beta) + h\lambda p(\beta) \right\};$$

$$e^T z = h, \text{ به طوری که}$$

$$z_i \in \{0, 1\}, \quad i = 1, \dots, n, \quad (10)$$

خواهد بود به طوری که در آن $Z = (1, \dots, 1)_{n \times 1}^T$ ، $e = (1, \dots, 1)_{n \times 1}^T$ یک ماتریس قطری با درایه‌های روی قطر اصلی $z = (z_1, \dots, z_n)_{n \times 1}^T$ ، h پارامتر پیراسته و λ پارامتر تاوان است. انتخاب پارامتر پیراسته در ابعاد بالا امری به‌غایت دشوار است.

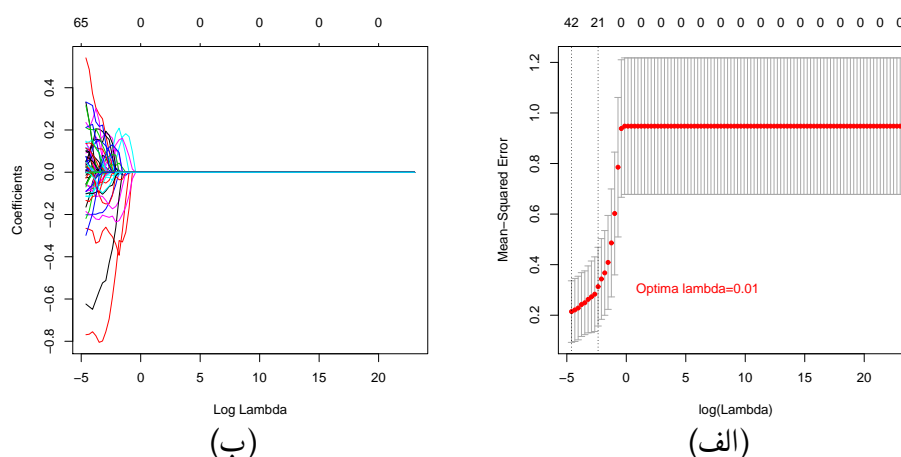
روش کمترین توان‌های دوم پیراسته مزایایی به شرح زیر داراست:

- عملکرد پیش‌بینی را از طریق کاهش واریانس بهبود می‌بخشد،
- به سبب انتخاب هم‌زمان، مدل تفسیرپذیری بیشتری را تضمین می‌کند،
- در حضور مجموعه داده با بعد بالا از مشکلات محاسباتی روش‌های کلاسیک جلوگیری می‌نماید.

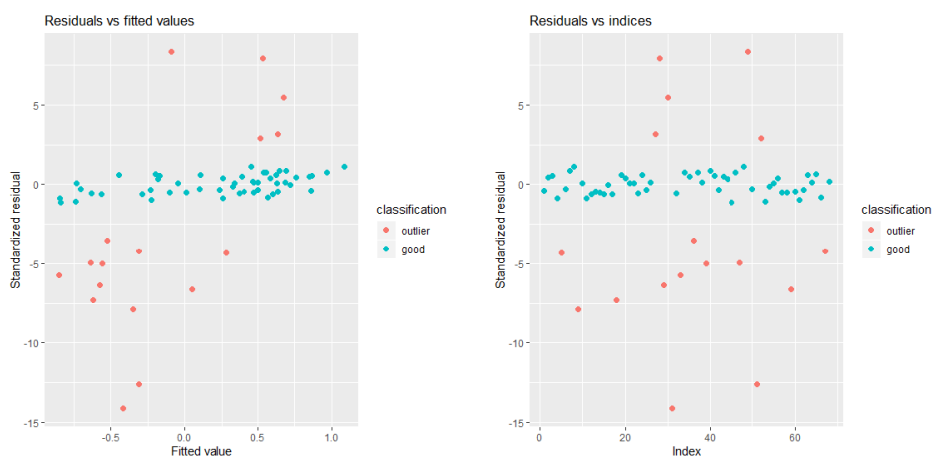
تا ۱/۳ میلی‌گرم برای بزرگسالان و برای زنان باردار یا شیرده ۱/۶ است. شایع‌ترین علت کمبود ریوفلاوین در بدن رژیم غذایی نامناسب است. کمبود ریوفلاوین هم‌چنین می‌تواند در افرادی که دچار نقص در فعالیت‌های کبدی هستند ایجاد شود، چراکه مانع از استفاده مناسب از ویتامین‌ها می‌شود.

این داده‌ها در بسته نرم‌افزاری "hdi" نرم‌افزار R موجود است. در این مجموعه داده متغیر پاسخ لگاریتم نرخ تولید ریوفلاوین است. این مجموعه داده دارای ۴۰۸۸ متغیر توضیحی بوده که هرکدام نشان‌دهنده لگاریتم سطح ژن‌ها است. با استفاده از روش تاوانیده لاسو استوار که در بسته نرم‌افزاری "robustHD" نرم‌افزار R موجود است، تعداد ۳۸ ژن به‌عنوان متغیرهای توضیحی مؤثر بر تغییرات و پیش‌بینی متغیر پاسخ بر اساس مقدار بهینه پارامتر تاوان $\lambda = 0.0040871$ شناسایی شدند. در شکل ۲ نمودارهای اعتبارسنجی متقابل 10^{-1} -سطحی و نمودار برآورد ضرایب به ازای مقادیر مختلف پارامتر تاوان نشان داده شده است. همچنین نمودار نقاط پرت شناسایی‌شده در شکل ۳ نمایش داده شده است.

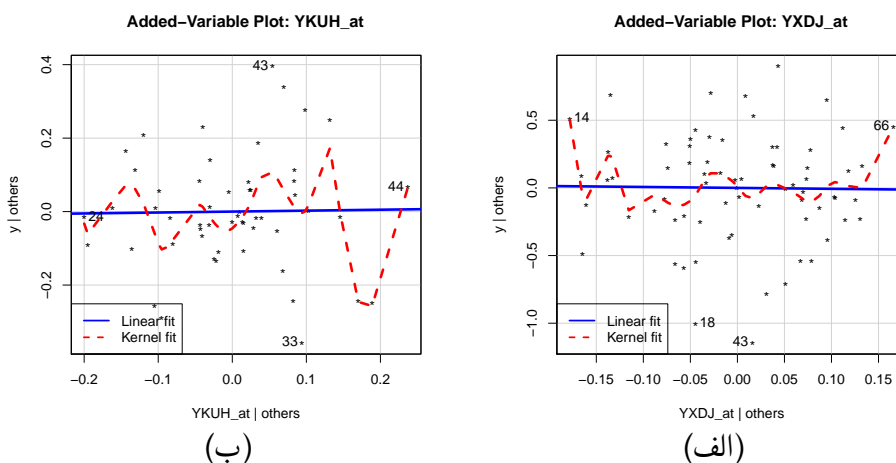
به‌طور طبیعی در برخی غذاها وجود دارد، به بعضی از محصولات غذایی اضافه می‌شود و به‌عنوان یک مکمل رژیم غذایی در دسترس است. این ویتامین یکی از اجزای اساسی دو کوآنزیم اصلی، فلاوین مونونوکلوئوتید (FMN) و فلاوین آدنین دینوکلوئوتید (FAD) است. این کوآنزیم‌ها نقش عمده‌ای در تولید انرژی، عملکرد سلولی، رشد و نمو و متابولیسم چربی‌ها، داروها و استروئیدها دارند. علاوه بر این، ریوفلاوین به حفظ سطح طبیعی هموسیستین، یک اسیدآمینو در خون کمک می‌کند. بیشترین درصد ریوفلاوین پس از مصرف مواد گوشتی، سبزیجات، لبنیات و مقداری کمتر در مغزها و تخمه‌ها، حبوبات و سبوس غلات موجود توسط بدن دریافت می‌شود. بیشتر ریوفلاوین توسط روده کوچک جذب و مقدار کمی در کبد، قلب و کلیه‌ها ذخیره‌شده و مابقی از طریق ادرار از بدن دفع خواهد شد. وضعیت ریوفلاوین در افراد سالم به‌طور معمول اندازه‌گیری نمی‌شود؛ اما در صورت بروز مشکل و تشخیص کمبود (وفور) این ویتامین از طریق بررسی گلوبول‌های قرمز خون و یا بررسی میزان ریوفلاوین دفع شده از طریق ادرار اندازه‌گیری خواهد شد. میزان میانگین ریوفلاوین موردنیاز بدن روزانه بین ۱/۱



شکل ۲: نمودارهای مربوط به پارامتر تاوان لاسو، الف: اعتبارسنجی متقابل 10^{-1} -سطحی، ب: برآورد ضرایب به ازای مقادیر مختلف پارامتر تاوان به روش لاسو.



شکل ۳: نمودار شناسایی نقاط پرت.



شکل ۴: نمودار افزوده متغیرهای پیشگو دارای رابطه غیرخطی با متغیر پاسخ و برازش کمترین توان‌های دوم (خط ممتد) و برازش ناپارامتری کرنل (خط چین)، الف: روش لاسو غیر استوار، ب: روش لاسو استوار

مجدداً به‌طور شهودی در شکل ۴ بررسی و نتیجه به‌دست‌آمده در بالا تأیید شد. نمودار متغیر افزوده به‌طور شهودی اثر هر یک از متغیرهای پیشگو را پس از حذف اثر سایر متغیرهای پیشگو، بر متغیر پاسخ آشکار می‌کند.

عنصر ناپارامتری مطابق [۸] با محاسبه

$$s_i^2 = \frac{1}{n - p_1 - 1} (\tilde{y} - \bar{X}[-i]\hat{\beta})^T (\tilde{y} - \bar{X}[-i]\hat{\beta}),$$

$$i = 1, \dots, 4088, \quad p_1 = 38,$$

در جدول ۱ مقادیر برآوردگرها، انحراف از معیار و همچنین RSS، R^2 و CV که به ترتیب، مجموع توان‌های دوم مانده‌ها، ضریب تعیین و اعتبارسنجی متقابل مدل رگرسیون نیمه‌پارامتری برازش شده هستند، گزارش شده‌اند. همان‌طور که در این جدول دیده می‌شود، برآوردگرهای استوار لاسو دارای مجموع توان دوم خطا و اعتبارسنجی متقابل کمتری نسبت به برآوردگرهای غیر استوار لاسو بوده و حدود ۳۵ درصد بیشتر تغییرات متغیر پاسخ را نسبت به برآوردگر غیر استوار بیان می‌کنند.

شناسایی می‌شود که در آن \tilde{y} متغیر پاسخ (لگاریتم تولید ریوفلاوین)، $\bar{X}[-i]$ ماتریس طرح و $\bar{X}[-i]$ ماتریس طرح فراهم‌شده با حذف i امین ستون است. بدیهی است عنصر ناپارامتری، متغیر توضیحی است که با حذف آن میزان s_i^2 کمینه خواهد شد. با استفاده از روش بالا به ترتیب متغیرهای توضیحی “YKUH_at” و “YXDJ_at” با استفاده از رویکردهای لاسو غیراستوار و لاسو استوار به‌عنوان عضو ناپارامتری شناسایی شدند. همچنین، با استفاده از نمودار متغیر افزوده^{۲۹} نیز رابطه غیرخطی اجزای ناپارامتری مدل‌های لاسو غیراستوار و لاسو استوار

²⁹Added Variable Plot

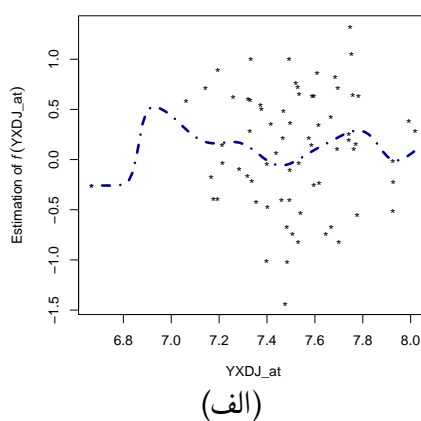
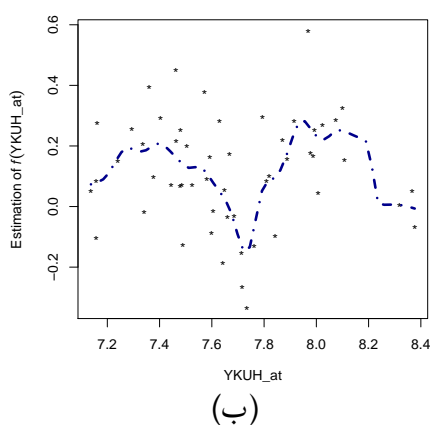
داده‌های با بعد بالا به کمک الگوریتم سریع کم‌ترین توان‌های دوم پیراسته تنک، معرفی شد. نقطه فروریزش برآوردگر معرفی شده به دست آمد و از لحاظ نظری نشان داده شد که این نقطه فروریزش می‌تواند بیش از ۵۰ درصد باشد. با این حال از دیدگاه عملی دستیابی به نقطه فروریزش بیش از ۵۰ درصد امکان‌پذیر نیست.

نتایج عددی مطالعه داده‌های واقعی نشان می‌دهد که استفاده از برآوردگر استوار لاسو تنک در مدل رگرسیون نیمه پارامتری، دارای تأثیر به‌سزایی در بهبود برآوردگرهای ضرایب خطی و توابع غیرخطی مدل در برازش مدل به داده‌های با بعد بالا و دارای مشاهدات دورافتاده است.

پس از برآورد ضرایب متغیرهای مؤثر شناسایی شده به روش تاوانیده لاسو غیراستوار و استوار، توابع ناپارامتری با استفاده از روش کرنل برازش گردید. مقدار بهینه پارامتر پهنای باند به روش کرنل به ترتیب برای متغیرهای توضیحی ناپارامتری "YXDJat" و "YKUHat" برابر ۰/۱۲۴۳ و ۰/۹۴۵ با رویکردهای لاسو غیراستوار و لاسو استوار تعیین شد و ۳۷ متغیر توضیحی دیگر به‌عنوان بخش خطی مدل در نظر گرفته شدند. توابع ناپارامتری برازش شده در شکل ۵ رسم شده است.

۶ بحث و نتیجه‌گیری

در این مقاله برآورد تاوانیده استوار لاسو برای ضرایب خطی و توابع غیرخطی در مدل رگرسیون نیمه پارامتری در حضور داده‌های پرت در



شکل ۵: الف: برازش تابع ناپارامتری به روش غیراستوار لاسو ب: برازش تابع ناپارامتری به روش استوار لاسو

جدول ۱: ارزیابی برآوردهای پیشنهادی برای متغیرهای توضیحی مؤثر در داده‌های واقعی.

لاسواستوار		لاسو غیراستوار		روش
انحراف استاندارد	برآورد	انحراف استاندارد	برآورد	
۰/۰۲۲۰	-۰/۰۶۵۲۶	۰/۰۳۳۹۸	-۰/۰۵۶۲۵	عرض از مبدأ
۰/۰۲۴۱۵	۰/۰۳۰۱۴۶	۰/۰۱۶۰۴۳	-۰/۰۲۲۷۰۳	ALD_at
۰/۰۷۳۲۷	-۰/۰۴۵۸۸۷	۰/۰۴۶۷۱۵	-۰/۰۹۱۷۲	ALDX_at
۰/۰۳۴۶۳	۰/۰۴۸۰۵۹	۰/۰۳۷۰۵۳	۱/۱۱۳۷۹	CSBA_at
۰/۰۱۴۰۱	-۰/۰۳۴۴۱۶	۰/۰۸۶۹۴	-۰/۰۶۸۲۵۶	CTAA_at
۰/۰۱۳۹۶	-۰/۰۵۱۳۷	۰/۰۱۲۵۳۲	-۰/۰۱۷۵۳۷	DEAD_at
۰/۰۱۶۹۹	-۰/۰۵۶۰۵	۰/۰۱۵۰۱۹	۰/۰۷۹۰۱	IOLJ_at
۰/۰۵۸۶۷	۰/۰۴۵۸۸۸	۰/۰۳۵۶۴۵	۰/۰۴۱۲۰۱	LACA_at
۰/۰۳۱۹۳	-۰/۰۱۳۵۱۴	۰/۰۳۷۰۳۷	۰/۰۶۹۶۳۴	MMGE_at
۰/۰۷۳۰۹	۰/۰۲۴۰۸	۰/۰۴۵۸۰۷	-۰/۰۴۲۵۳۱	MREBH_at
۰/۰۴۹۶۵	-۰/۰۲۶۱۳۶	۰/۰۳۸۲۳۶	۱/۵۲۳۲۹	PCKA_at
۰/۰۳۶۲۰	۰/۰۳۶۴۳۰	۰/۰۳۵۴۷۳	-۰/۰۳۲۲۷۲	PKSI_at
۰/۰۷۶۲۱	-۰/۰۵۰۱۲	۰/۰۶۰۶۳۱	۰/۰۵۰۹۵۱	PKSM_at
۰/۰۱۵۳۰۳	۱/۲۶۹۰۰	۰/۰۶۲۴۴۵	۱/۶۷۶۶۱	PTA_at
۰/۰۶۰۲۹	-۰/۰۳۵۱۴۵	۰/۰۴۴۲۹۶	-۰/۰۳۶۶۱۵	RIBT_at
۰/۰۱۴۰۶۸	-۰/۰۶۶۷۰۱	۰/۰۳۴۱۱۲	-۱/۰۲۴۱۸۹	SPOIISA_at
۰/۰۴۸۵۱	۰/۰۸۶۷۲	۰/۰۶۹۳۵۸	-۰/۰۷۲۴۴۷	SPOISB_at
۰/۰۱۰۴۲	-۰/۰۵۹۵۴	۰/۰۸۹۰۳	۰/۰۵۴۰۶۸	XKDS_at
۰/۰۲۳۱۹	۰/۰۵۱۲۸۹	۰/۰۲۳۱۲۵	-۰/۰۱۲۵۰۷	YCEI_at
۰/۰۵۴۰۰	۰/۰۷۲۰۹	۰/۰۳۵۱۷۷	-۰/۰۷۲۲۷	YCGR_at
۰/۰۶۷۹۴	-۰/۰۹۹۰۸۷	۰/۰۴۳۵۳۸	-۰/۰۳۶۴۱۴	YDBM_at
۰/۰۸۴۰۰	-۰/۰۱۰۰۲۲	۰/۰۴۵۱۵۷	-۰/۰۳۲۷۰۰	YEBC_at
۰/۰۱۰۵۴۱	۱/۱۴۶۷۹	۰/۰۳۷۲۲۶	۰/۰۹۲۸۲۹	YFHE_r_at
۰/۰۵۹۵۱	-۱/۰۶۴۵۸	۰/۰۳۷۶۶۶	۰/۰۲۱۷۱	YHCL_at
۰/۰۵۱۳۳	۰/۰۶۵۲۲۲	۰/۰۶۴۵۹۵	-۰/۰۲۷۱۵۰	YJBJ_at
۰/۰۴۲۱۳	-۰/۰۲۷۵۱۷	۰/۰۴۹۵۵۳	۱/۸۴۶۰۶	YJCJ_at
۰/۰۳۱۹۲	-۰/۰۳۵۶۳۰	۰/۰۳۱۶۶۵	-۱/۰۳۶۷۹۰	YKUH_at
۰/۰۱۶۶۱۰	-۰/۰۴۴۷۱۰	۰/۰۲۲۵۶۳	۰/۰۲۶۲۷۰	YORB_i_at
۰/۰۱۴۷۹	-۰/۰۴۸۷۷۲	۰/۰۶۲۰۷۵	-۲/۰۲۲۲۷۴	YQET_at
۰/۰۱۷۶۵	-۰/۰۲۷۹۱۲	۰/۰۹۶۶۹	-۰/۰۱۰۹۶۱	YQHJ_at
۰/۰۲۳۴۶	۰/۰۵۹۳۶۴	۰/۰۱۹۳۵۲	-۰/۰۶۴۴۳۵	YUNL_at
۰/۰۲۱۰۹	۰/۰۴۹۰۱۱	۰/۰۲۲۳۲۳	۰/۰۴۴۰۲۵	YWFM_at
۰/۰۷۸۸۲	-۱/۰۰۰۰۳۵	۰/۰۱۷۷۶۶	-۰/۰۲۸۱۲۶	YWPB_at
۰/۰۰۵۴۲	۰/۰۲۷۳۸۹	۰/۰۵۶۲۹	۰/۰۶۲۷۹	YXDJ_at
۰/۰۵۵۶۲	-۰/۰۸۴۵۲۷	۰/۰۳۶۱۳۶	-۱/۰۸۴۹۷۹	YXER_at
۰/۰۲۹۵۶	-۰/۰۲۳۱۸	۰/۰۱۴۹۷۱	-۰/۰۱۱۳۸۴	YXJL_at
۰/۰۷۸۲۱	۰/۰۱۲۰۶۷	۰/۰۴۱۲۵۲	۰/۰۶۳۰۱۴	YXLD_at
۰/۰۱۴۰۹	۰/۰۴۳۸۳۱	۰/۰۲۰۹۰۸	۰/۰۸۷۹۸	YXLE_at
	۲/۰۰۱۳	۲/۰۲۵۱۳		RSS
	۰/۰۸۹۲۸	۰/۰۵۴۰۵		R ²
	۰/۰۱۶۱۸	۰/۰۲۷۸۰		CV

مراجع

- [۱] آرشی، م.، صادقی، ح؛ و طباطبایی م. (۱۳۹۸)، استنباط آماری، انتشارات دانشگاه صنعتی شاهرود، شاهرود.
- [۲] جذن، س؛ و امینی، م. (۱۳۹۶)، برآورد استوار نسبت به مشاهده‌های دورافتاده در رگرسیون خطی در حضور هم‌خطی چندگانه، اندیشه آماری، ۴۴، ۹۳-۱۱۰.
- [۳] روزبه، م؛ و چاچی، ج. (۱۳۹۵)، برآورد تفاضلی استوار مدل‌های خطی جزئی، مجله علوم آماری، ۱۰(۱)، ۹۵-۱۱۲.
- [4] Akdeniz, D.E., Hardle, W.K., and Osipenko, M. (2012), Difference Based Ridge and Liu Type Estimators in Semiparametric Regression Models, *Journal of Multivariate Analysis*. **105(1)**, 164-175.
- [5] Alfons, A., Croux, C., and Gelper, S. (2013), Sparse Least Trimmed Square Regression for Analyzing High-Dimensional Large Data Set, *The Annals of Applied Statistics*. **7(1)**, 226-248.
- [6] Amini, M., and Roozbeh, M. (2015), Optimal Partial Ridge Estimation in Restricted Semiparametric Regression Models, *Journal of Multivariate Analysis*. **136(2)**, 26-40.
- [7] Arashi, M., Asar, Y., and Yuzbasi B. (2021), SLASSO: a scaled LASSO for multicollinear situations, *Journal of Statistical Computation and Simulation*. **91(15)**, 3170-3183.
- [8] Arashi, M., and Roozbeh, M. (2019) , Some Improved Estimation Strategies in High-dimensional Semiparametric Regression Models with Application to Riboflavin Production Data, *Journal of Statistical Papers*, **60**, 667-686.
- [9] Barnett, V., and Lewis, T. (1994), *Outliers in Statistical Data*. JohnWiley ad Sons, New York.
- [10] Bunea, F. (2004), Consistent Covariate Selection and Post Model Selection Inference in Semiparametric Regression, *The Annals of Mathematical Statistics*. **32(3)**, 898-927.
- [11] Efron, B., and Hastie, T. (2017), *Computer Age Statistical Inference*. Cambridge University Press, Cambridge.
- [12] Engle, R. F., Granger, C.W.J., Rice, J., and Weiss, A. (1986), Semiparametric Estimation of the Relation Between Weather and Electricity Sale, *Journal of the American statistical Association*. **81(394)**, 310-320.
- [13] Everitt, B., and Hothorn, T. (2011), *An Introduction to Applied Multivariate Analysis with R*. Springer, New York Dordrecht, Heidelberg, London.
- [14] Fan, J., Hardle, W., and Mammen, E. (1998), Direct Estimation of Low Dimensional Components in Additive Models, *The Annals of Statistics*. **26(3)**, 943-971.
- [15] Gannaz, I. (2007), Robust Estimation and Wavelet Thresholding in Partially Linear Models, *Statistics and Computing*. **17(4)**, 293-310.
- [16] Hardle, W.k., Laing, H., and Gao, J. (2000), *Partially Linear Models*. PhysikaVerlag, Heidelberg.
- [17] Li, B., and Yu, Q. (2009), Robust and Sparse Bridge Regression, *Statistics and its interface*, **2(4)**, 481-491.
- [18] Nadaraya, E.A. (1964), On Estimating Regression, *Theory of Probability and its Applications*. **9(1)**, 141-142.
- [19] Roozbeh, M. (2015), Shrinkage Ridge Estimators in Semiparametric Regression Models, *Journal of Multivariate Analysis*. **136(2)**, 56-74.

- [20] Roozbeh, M. (2018), Optimal QR-based estimation in partially linear regression models with correlated errors using GCV criterion, 117, *Computational Statistics & Data Analysis*. **117**, 45-61.
- [21] Roozbeh, M., and Arashi, M. (2013), Feasible Ridge Estimator in Partially Linear Models, *Journal of Multivariate Analysis*. **116(4)**, 35-44.
- [22] Rousseeuw, P.J., and Leroy, A.M. (1987), *Robust Regression and Outlier Detection*. John Wiley, New York.
- [23] Speakman, P. (1988), Kernel Smoothing in Partial Linear Models, *Journal of the Royal Statistical Society: Series B (Methodological)*. **50(1)**, 413-436.
- [24] Tibshirani, R. (1996), Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society. Series B (Methodological)*. **58(1)**, 267-288.
- [25] Yuzbasi B., Arashi, M., and Akdeniz, F. (2021), Penalized regression via the restricted bridge estimator, *Soft Computing*. **25(13)**, 8401-8416
- [26] Wasserman, L. (2006), *All of Nonparametric statistics*. Springer Science and Business Media, New York.
- [27] Watson, G. S. (1964), Smooth Regression Analysis, *Sankhyā: The Indian Journal of Statistics, Series A*. **26(4)**, 359-372.

Sparse robust semiparametric models in high-dimensional data

Mahdi Roozbeh¹, Monireh Manavi²

Abstract:

Analysis and modeling the high-dimensional data is one of the most challenging problems faced by the world nowadays. Interpretation of such data is not easy and needs to be applied to modern methods. The penalized methods are one of the most popular ways to analyze the high-dimensional data. Also, the regression models and their analysis are affected by the outliers seriously. The least trimmed squares method is one of the best robust approaches to solve the corruptive influence of the outliers. Semiparametric models, which are a combination of both parametric and nonparametric models, are very flexible models. They are useful when the model contains both parametric and nonparametric parts. The main purpose of this paper is to analyze semiparametric models in high-dimensional data with the presence of outliers using the robust sparse Lasso approach. Finally, the performance of the proposed estimator is examined using a real data analysis about production of vitamin B2.

Keywords: High-dimensional data, Lasso method, Least trimmed squares method, Semiparametric model, Sparse least trimmed squares method.

¹Faculty of mathematics, Semnan university, Semnan, Iran.

²Master's degree graduate, statistics and Computer science, Semnan university, Semnan, Iran.