

مروری بر مشکلات متداول رگرسیون و راه‌حل‌های آنها

زهرا جعفریان مورکانی^۱، حیدرعلی مردانی فرد^۲

تاریخ دریافت: ۱۴۰۱/۰۷/۲۰

تاریخ پذیرش: ۱۴۰۱/۱۲/۰۲

چکیده:

مدل رگرسیون خطی معمولی به صورت $Y = X\beta + \varepsilon$ است و برآورد پارامتر β عبارت است از: $\hat{\beta} = (X'X)^{-1}X'Y$. با این حال در هنگام استفاده از این برآوردگر به صورت عملی، ممکن است مشکلات خاصی مانند مشکل انتخاب متغیر، هم‌خطی، مدل با ابعاد بالا، کاهش بعد و وجود خطای اندازه‌گیری به وجود آید که استفاده از برآوردگر $\hat{\beta}$ را مشکل می‌سازد. در اغلب این مشکلات، مسئله اصلی عدم معکوس‌پذیری ماتریس $X'X$ است. برای رفع آن‌ها راه‌حل‌های متعددی ارائه شده است. در این مقاله ضمن مروری بر این مشکلات، مجموعه‌ای از راه‌حل‌های معمول و متداول و همچنین چند روش خاص و پیشرفته (که کمتر مورد اقبال همگان است ولی با این حال توانایی بالقوه‌ای در رفع هوشمند این مشکلات دارند) را بررسی می‌کنیم.

واژه‌های کلیدی: هم‌خطی، کاهش بعد، خطای اندازه‌گیری، رگرسیون ستیخی، شبه‌برآورد، رگرسیون معکوس قطعه‌قطعه شده

۱ مقدمه

است. به‌طور معمول برآورد β با روش کمترین مربعات (که در صورت نرمال و مستقل بودن خطاها معادل با برآورد حداکثر درست‌نمایی است) عبارت است از:

$$\hat{\beta} = (X'X)^{-1}X'y, \quad (۳)$$

با این حال، در استفاده از برآورد فوق لازم است به برخی پیش‌فرض‌های مهم دقت کرد. هنگام استفاده از داده‌های همبسته یا داده‌های با ابعاد بالا در رگرسیون خطی با معکوس ناپذیری ماتریس $X'X$ مواجه می‌شویم که محاسبه $\hat{\beta}$ را غیرممکن و یا آن را بسیار ناپایدار می‌کند به‌گونه‌ای که با اندک تغییر جزئی در داده‌ها مقدار برآوردگر به‌طور قابل‌ملاحظه‌ای تغییر می‌کند. اخیراً، لی و بین [۸] پیشنهاد کردند که خطاهایی به پیش‌بینی‌کننده‌ها با ساختار کوواریانس شناخته شده افزوده شود و برآورد حاصل را «شبه برآورد» نامیدند. واگر و همکاران [۱۲] نشان دادند که در حالتی خاص، داده‌های دارای خطا، ارتباط تنگاتنگی با برآوردگر معروف ستیخی دارد. یک روش جدید انتخاب متغیر بر اساس شبه‌فاصله اطمینان می‌باشد که برای پیش‌بینی‌کننده‌های همبسته و مشکل "p بزرگ، n کوچک" پیشنهاد می‌شود. در این مقاله، روش شبه برآورد

یکی از متداول‌ترین و پرکاربردترین روش‌های آماری در تحلیل داده‌ها رگرسیون است که در حالت معمول آن، اثر یک یا چند متغیر روی یک متغیر دیگر به صورت یک رابطه عموماً خطی سنجیده می‌شود. در مدل رگرسیون خطی، رابطه بین متغیر وابسته Y و متغیرهای مستقل X_0, X_1, \dots, X_p به صورت زیر در نظر گرفته می‌شود:

$$y = X'\beta + \varepsilon, \quad (۱)$$

که در آن بردار β بردار ضرایب رگرسیونی، $X = (X_0, X_1, \dots, X_p)'$ بردار متغیرهای مستقل با $X_0 \equiv 1$ بوده و ε خطای مدل است. وقتی متغیرها را روی یک نمونه n تایی اندازه‌گیری کرده باشیم مدل به صورت ماتریسی زیر نوشته می‌شود:

$$y = X\beta + \varepsilon, \quad (۲)$$

که در آن X ماتریس داده‌های متغیرهای مستقل است و بردار y مقادیر متغیر وابسته (پاسخ) و بردار ε هم بردار خطاهای تصادفی

^۱ دانشجوی کارشناسی ارشد رشته آمار، گروه ریاضی، دانشگاه یاسوج

^۲ دانشگاه یاسوج، گروه ریاضی (نویسنده مسئول): h_mardanifard@yu.ac.ir

نرمال می‌توان به صورت زیر از روی مجموع مربعات خطا محاسبه کرد [۵]:

$$AIC = \frac{1}{n}(SSE + 2d\hat{\sigma}^2)$$

$$BIC = \frac{1}{n}(SSE + d \ln(n)\hat{\sigma}^2)$$

طبیعی است که هر چه مقدار این آماره‌ها کمتر باشد مناسب مدل بیشتر است.

۳.۲ رگرسیون تاوانیده

روش حداقل مربعات معمولی هنگامی که تعداد متغیرها بیشتر از نمونه باشد عملکرد نامیدکننده‌ای دارد. رگرسیون تاوانیده یا روش‌های انقباضی، به این صورت است که مدل رگرسیونی جریمه می‌شود و اعمال این تاوان به متغیرها اجازه می‌دهد تا ضریبی نزدیک به صفر یا مساوی صفر داشته باشند. توجه کنید که انقباض نیاز به انتخاب یک پارامتر تنظیم‌کننده^۶ (مثلاً λ) دارد که تعیین‌کننده‌ی میزان انقباض است و با روش اعتبارسنجی متقابل تعیین می‌شود [۴]. در رگرسیون تاوانیده با تابع تاوان $P_\lambda(\beta)$ برآورد پارامترها از روی بهینه‌سازی تابع زیر به دست می‌آید:

$$Q_\lambda(\beta) := SSE(\beta) + P_\lambda(\beta),$$

که در آن $SSE(\beta) = \sum_i (y_i - \beta_0 - \sum_j \beta_j x_{ij})^2$.

۱.۳.۲ رگرسیون ستیغی

رگرسیون ستیغی^۷ برای رسیدن به پیش‌بینی بهتر در برخورد با هم‌خطی، توسط هورل و کنارد [۴]، پیشنهاد شد. در این روش مقدار تاوان به صورت $P_\lambda(\beta) = \sum \beta_j^2$ است و در نتیجه برآورد ستیغی با مینیم کردن عبارت زیر به دست می‌آید:

$$SSE(\beta) + \lambda \sum_{j=1}^p \beta_j^2.$$

در مسائل بهینه‌سازی، این موضوع معادل است با مینیم کردن $SSE(\beta)$ با قید $\sum \beta_j^2 = c$. در سمت چپ شکل ۱ نمودار سطح مقطع و ناحیه قید بهینه‌سازی فوق رسم شده است.

پیشنهادی و انتخاب متغیر بر اساس آن‌هم مورد مطالعه قرار می‌گیرد.

۲ ابزارها و مفاهیم مفید رگرسیون

۱.۲ ضریب تعیین و آماره C_p مالو

ضریب تعیین (R^2) کسری از کل تغییرات متغیر پاسخ است که توسط متغیرهای مستقل تبیین می‌شود. مقدار این ضریب همواره بین ۰ و ۱ است و اگر به یک نزدیک باشد، نشان‌دهنده‌ی برازش بهتر خط رگرسیون است و اگر به صفر نزدیک‌تر باشد، خط رگرسیون به‌خوبی به داده‌ها برازش داده نشده است. معادله ضریب تعیین به صورت زیر است:

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

اگر در مقایسه مدل‌ها تعداد پارامترها برابر نباشد از ضریب تعیین تعدیل‌شده استفاده می‌شود که در آن اثر تعداد متغیرها تعدیل شده است:

$$R_{adj}^2 = 1 - \frac{n-1}{n-p} \times \frac{SSE}{SST} = 1 - \frac{n-1}{n-p} (1 - R^2)$$

علاوه بر این‌ها، آماره C_p مالو^۳ برای اولین بار توسط مالو [۹] به صورت زیر تعریف شد:

$$C_p = \frac{MSE_p}{MSE} - (n - 2p)$$

که در آن $MSE = \frac{SSE}{n-p}$ و SSE_p از یک مدل با $k = p - 1$ متغیر به وجود آمده است. در اینجا k و p نشان‌دهنده تعداد متغیرهای مستقل و تعداد پارامترها در مدل می‌باشد. امید ریاضی این آماره عبارت است از: $E(C_p) = p$ و در استفاده از این آماره متغیری برای ورود به مدل مناسب‌تر است که C_p کمتری دارد [۲].

۲.۲ معیار اطلاع آکائیکه

معیارهای اطلاع مانند معیار اطلاع آکائیکه^۴ (به اختصار AIC) و معیار اطلاع بیزی^۵ (به اختصار BIC) از روی مقدار بهینه تابع درست‌نمایی تعریف می‌شود که در مدل‌های رگرسیونی با توزیع خطای

^۳Mallows' C_p

^۴Akaike information criterion

^۵Bayesian information criterion

^۶Tuning parameter

^۷Ridge regression

ویژه ماتریس کواریانس است [۳]. این روش برای داده‌ها با حجم نمونه زیاد مورد استفاده قرار می‌گیرد. اولین جهت مؤلفه اصلی داده‌ها همان جهتی است که بیشترین تغییرات مشاهدات در آن وجود دارند و باعث تعریف خطی می‌شود که تا حد امکان به داده‌ها نزدیک می‌باشد. مدل رگرسیونی (۱) را در نظر بگیرید که در آن بردار تصادفی X دارای امید ریاضی μ و ماتریس واریانس-کواریانس Σ است. هدف این است که بجای بردار X بتوان از بردار

$$X^* = (X_1^*, \dots, X_q^*)'$$

استفاده کرد به طوری که $q \ll p$ باشد در حالی که X^* نسبت به X اطلاعات زیادی از دست ندهد. برای یافتن چنین X^* راه‌های زیادی وجود دارد که یکی از این روش‌ها تجزیه طیفی ماتریس واریانس-کواریانس است. تجزیه طیفی به این صورت عمل می‌کند که باید Σ به صورت $P\Lambda P'$ نمایش داده شود. اگر مقادیر λ_j ها از بزرگ به کوچک مرتب شوند، عملاً X_j^* ها به ترتیب بیشترین تا کمترین واریانس مرتب می‌شوند. به این ترتیب می‌توان از q تا X_j^* ابتدایی، یعنی q تا بردار ویژه متناسب با q تا بزرگ‌ترین مقادیر ویژه استفاده کرد.

۵.۲ رگرسیون معکوس قطعه‌قطعه شده

برآورد ماتریس واریانس-کواریانس با استفاده از قطعه‌بندی دامنه‌ی متغیر پاسخ، برای اولین بار توسط لی^{۱۰} [۷] مطرح شد. این روش، مشهورترین روش کاهش بعد در رگرسیون است و برای کاهش بعد متغیر مستقل X بدون انجام هیچ‌گونه فرآیند برازش مدل پارامتری یا ناپارامتری ارائه می‌شود که در عوض رگرسیون متغیر پاسخ تک متغیره y در برابر X چند متغیره اجرا می‌گردد. هدف رگرسیون معکوس قطعه‌قطعه شده^{۱۱} (به اختصار SIR)، کشف یک تصویر از متغیر کمکی p -بعدی X ، روی یک زیرفضای خطی k بعدی است که حاوی اکثر اطلاعات در خصوص پاسخ y است. هر پارامتر β که در زیرفضای k بعدی وجود داشته باشد، به عنوان «جهت» کاهش بعد مؤثر^{۱۲} (به اختصار e.d.r.) بیان و کوچک‌ترین e.d.r. زیرفضای مرکزی نامیده می‌شود که دارای کمترین بعد و بیشترین واریانس است. لی [۶] تشریح کرد هنگامی که y تغییر می‌کند،

نتیجه این بهینه‌سازی، برآورد زیر را به دست می‌دهد:

$$\hat{\beta}^{(R)} = (X'X + \lambda I)^{-1} X'y.$$

۲.۳.۲ رگرسیون لاسو

رگرسیون لاسو^۸ برای اولین بار توسط تیشیرانی [۱۱] پیشنهاد شد که انتخاب متغیر و برآورد پارامتر به‌طور هم‌زمان توسط آن، انجام می‌گیرد. در این روش تابع تاوان به صورت مجموع قدر مطلق ضرایب تعریف می‌شود و در نتیجه باید عبارت زیر برای برآورد ضرایب مینیمم شود:

$$SSE(\beta) + \lambda \sum_{j=1}^p |\beta_j|.$$

مشابه رگرسیون ستیغی، در اینجا هم مینیمم سازی فوق معادل با مینیمم کردن تابع $SSE(\beta)$ با قید $\sum |\beta_j| = c$ است. نمودار سطح مقطع تابع هدف و همچنین ناحیه قید این بهینه‌سازی در سمت راست شکل ۱ رسم شده است.

۳.۳.۲ مقایسه رگرسیون ستیغی و لاسو

در هر دو روش ستیغی و لاسو، محل اولین برخورد سطح مقطع‌ها با ناحیه محدودیت، جواب رگرسیون تاوانیده می‌باشد و ناحیه محدودیت مربوط به رگرسیون ستیغی دایره‌ای شکل و نقاط روی دایره نسبت به هم ارجحیتی ندارند؛ بنابراین، احتمال برخورد بیضی و ناحیه روی دایره در یک نقطه مشخص (مانند نقاط روی محورها) صفر است و هرچند که ممکن است برآورد ضرایب رگرسیون ستیغی به صفر خیلی نزدیک باشند ولی دقیقاً صفر نمی‌شود. در طرف مقابل، ناحیه محدودیت مربوط به لاسو دارای تیزی و گوشه می‌باشد، بنابراین ناحیه شرطی و سطح مقطع‌ها با احتمال مثبت در محورها و نقاط گوشه برخورد داشته و ممکن است در ابعاد بالاتر، تعداد بیشتری از ضرایب صفر شوند.

۴.۲ رگرسیون مؤلفه‌های اصلی

روش تحلیل مؤلفه‌های اصلی برای اولین بار در سال ۱۹۰۱ توسط کارل پیرسون^۹ دانشمند انگلیسی بیان شد که شامل تجزیه مقادیر

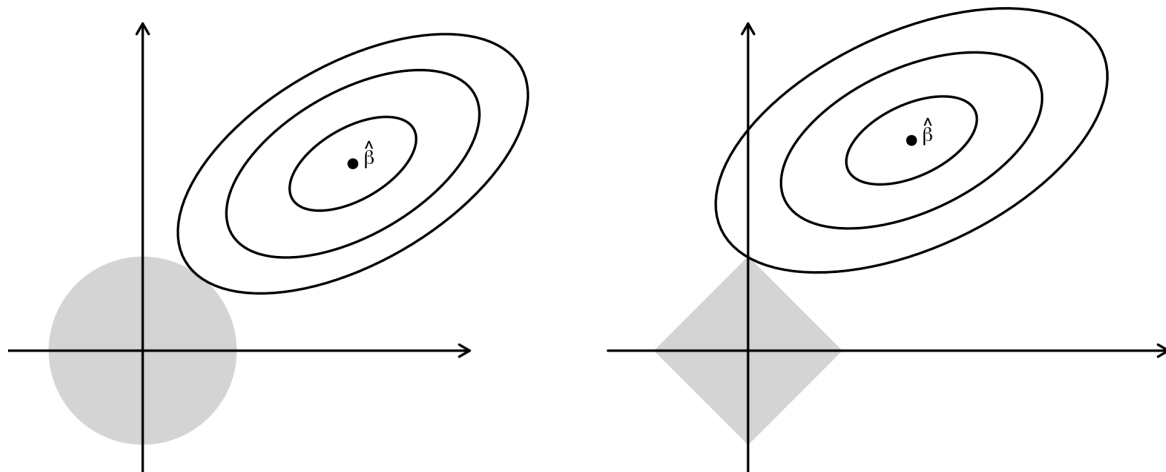
^۸LASSO

^۹Karl Pearson

^{۱۰}Li

^{۱۱}Sliced Inverse Regression

^{۱۲}effective dimension reduction



شکل ۱: مقایسه نمودار سطح مقطع‌های رگرسیون ستیغی (چپ) و لاسو (راست)

و سپس مقادیر ویژه و بردارهای ویژه \hat{V} محاسبه می‌شود. در واقع \hat{V} برآوردی از $cov(\hat{E}(Z|y))$ می‌باشد.

۵. قرار می‌دهیم:

$$\hat{\beta}_k = \hat{\eta}_k \hat{\Sigma}_x^{-1/2}, \quad k = 1, \dots, K$$

که در آن $\hat{\eta}_k$ ها بردارهای ویژه متناظر با K مقدار ویژه بزرگ ماتریس \hat{V} هستند.

رابطه $E(\mathbf{X}|y)$ موجب ترسیم یک منحنی می‌شود که در R^p ، منحنی رگرسیون معکوس نام می‌گیرد. در صورتی که متغیرهای کمکی \mathbf{X} طوری استاندارد شده باشند که میانگین آن بردار صفر و ماتریس واریانس-کواریانس آن ماتریس همانی باشد، آنگاه این زیرفضا با فضای e.d.r. منطبق و سبب شناسایی جهت‌های اصلی تغییرات می‌شود. الگوریتم رگرسیون معکوس قطعه‌قطعه شده به وسیله‌ی لی [۶] به صورت زیر مطرح شد:

۱. ابتدا X به صورت زیر استاندارد می‌شود:

$$\tilde{x}_i = \hat{\Sigma}_x^{-1/2} (\mathbf{x}_i - \bar{\mathbf{x}}), \quad i = 1, \dots, n,$$

که در آن $\bar{\mathbf{x}}$ و $\hat{\Sigma}_x$ به ترتیب ماتریس واریانس-کواریانس نمونه و میانگین نمونه X هستند.

۲. محدوده y ، به H قطعه، یعنی I_1, \dots, I_H به صورت صعودی تقسیم می‌گردد و نسبتی از y_i ها که در قطعه h ام قرار می‌گیرد \hat{P}_h نامیده می‌شود؛ یعنی، $\hat{P}_h = n^{-1} \sum \delta_h(y_i)$ ، می‌باشد که در آن $\delta_h(y_i)$ مقادیر ۰ یا ۱ را بسته به اینکه آیا y_i در قطعه h قرار گرفته است یا نه، به خود می‌گیرد.

۳. در هر قطعه، میانگین نمونه \tilde{x}_i ها که با \hat{m}_h نشان داده می‌شود محاسبه می‌شود:

$$\hat{m}_h = \frac{1}{n} \hat{P}_h \sum \tilde{x}_i = \frac{1}{n_h} \sum \tilde{x}_i,$$

که در آن n_h تعداد مشاهدات در هر قطعه می‌باشد. در واقع، \hat{m}_h برآوردی از $E(Z|y)$ در هر قطعه هستند.

۴. ماتریس کواریانس وزنی $\hat{V} = \sum_{h=1}^H \hat{P}_h \hat{m}_h \hat{m}_h'$ تشکیل شده

۶.۲ شبه برآورد

ایده شبه برآورد توسط لی و بین [۸] برای مدل‌های رگرسیونی با خطای اندازه‌گیری ارائه شد. هنگامی که اندازه نمونه کمتر از تعداد متغیرهای پیش‌بین است، به علت معکوس پذیر نبودن ماتریس $(X'X)$ ، روش رگرسیون معمولی درست عمل نمی‌کند و در این حالت، یکی از روش‌های مناسب استفاده از شبه برآورد است [۱۰]. در مدل‌های معمول رگرسیونی با خطای اندازه‌گیری فرض می‌شود که بردار $1 \times q$ بُعدی پیش‌بینی کننده W با پیش‌بین‌های واقعی X از طریق یک معادله خطی مانند زیر مرتبط است:

$$W = \gamma + \Gamma X + \delta, \quad (۴)$$

که در آن γ و Γ به ترتیب، بردار غیر تصادفی $1 \times q$ بُعدی و ماتریس غیر تصادفی $p \times q$ بُعدی بوده و δ یک بردار خطای تصادفی با توزیع نرمال چند متغیره مستقل از X و Y می‌باشد. شبه برآورد هنگامی مطرح می‌شود که متغیر مستقل، دارای خطای اندازه‌گیری باشد و برآوردگر به دست آمده به عنوان «شبه برآورد» نام می‌گیرد که

۲.۲ شبه فاصله اطمینان

از آنجا که می توان برای یک Σ_δ چندین شبه نمونه ایجاد کرد، می توان چندین شبه برآورد به دست آورد و توزیع نمونه برداری از شبه برآورد مستقیماً تقریب زده و منجر به یک روش انتخاب متغیر جدید برای شبه برآوردها می شود. فرض کنید تولید شبه مشاهدات N بار انجام شود و برای هر $1, 2, \dots, p$ قرار دهید:

$$\bar{\beta}_j = N^{-1} \sum_{k=1}^N \hat{\beta}_{kj},$$

$$S.E.(\bar{\beta}_j) = \frac{\sum_{k=1}^N (\hat{\beta}_{kj} - \bar{\beta}_j)^2}{N-1},$$

حال مشابه تئوری بوت استرپ^{۱۳}، می توان برای β_j فاصله اطمینانی به صورت زیر ساخت:

$$\beta_j \in \left(\bar{\beta}_j \pm Z_{\alpha} \times S.E.(\bar{\beta}_j) \right)$$

این فاصله اطمینان را اصطلاحاً یک شبه فاصله اطمینان^{۱۴} (به اختصار PCI) می نامیم زیرا با استفاده از شبه مشاهدات به دست می آید. حال هر کدام از β_j هایی که فاصله اطمینان آن ها شامل عدد صفر می شود در مدل تأثیرگذار نیستند، بنابراین صفر در نظر گرفته می شوند. این روش انتخاب متغیر اصطلاحاً S-Pseudo نام دارد. بر اساس PCI، با این روش می توان یک شبه برآورد پراکنده شده^{۱۵} به دست آورد.

۱.۷.۲ سطح مقطع ها

در هندسه به شکلی که با بریدن یک شیء سه بعدی توسط شیء دوبعدی به وجود می آید، اصطلاحاً سطح مقطع گفته می شود و دانستن سطح مقطع های تابع معادل با دانستن تابع است زیرا ارتباط یک به یک بین آن ها برقرار می باشد. اگر تابع ما، یک فرم درجه دوم با ماتریسی مانند Σ باشد، هر Σ سطح مقطع های خاص خود را دارد. در واقع در فرم های درجه دوم اگر ماتریس Σ معین مثبت و متقارن باشد، سطح مقطع های $x' \Sigma x$ به فرم بیضی هستند اما ممکن است مورب، افقی یا قائم باشند که دانستن اینکه بیضی ها کجا و به چه شکل هستند معادل با دانستن ماتریس Σ است.

جهت های اصلی و فرعی بیضی هایی که سطح مقطع های $x' \Sigma x$ هستند، همان بردارهای ویژه Σ بوده و اندازه قطرهای بیضی هم

ساده ترین حالت معادله (۴) با مقادیر $p = q$ و $\Gamma = I_p$ است؛ یعنی $W = X + \delta$

معمولاً، یک نمونه کمکی وجود دارد که اطلاعاتی در مورد رابطه بین پیش بینی کننده اصلی X و پیش بینی کننده جایگزین W مانند $\Sigma_{wx} = cov(W, X)$ ارائه می کند. با استفاده از این برآورد کوواریانس، می توان از پیش بینی کننده جایگزین W استفاده کرد که تا حد امکان با پیش بینی کننده واقعی X هماهنگ شود. بدین ترتیب، این روش در بعضی مواقع ناچار می باشد به متغیرهای پیش بین خطایی بیفزاید. مثالی از متغیرهای با خطای اندازه گیری، داده های فشارخون است که معمولاً خود داده ها، واقعی نیستند و تقریبی هستند.

۱.۶.۲ شبه برآورد و انتخاب متغیر

برآورد پارامترها در روش شبه برآورد به دو صورت زیر است:

$$\hat{\beta}^{pseudo} = (X'X + n\Sigma_\delta)^{-1} X'y \quad (5)$$

$$\tilde{\beta}^{pseudo} = (W'W)^{-1} W'y \quad (6)$$

در روش شبه برآورد، هنگامی که از رابطه (۵) استفاده می شود، نگرانی به وجود نمی آید زیرا برای هر دو حالت $n < p$ و $n > p$ درست عمل می کند اما رابطه دوم یعنی (۶)، زمانی که $n < p$ است، جواب مطلوبی ندارد و برای استفاده از این معادله وقتی $Z = (X, y)$ کل نمونه های مشاهده شده باشد، یک نمونه جایگزین ترکیبی به صورت $Z^* = (W^*, y^*)$ به دست می آید که y^* تکرار بردار y به اندازه m بار می باشد و ماتریس W^* ، m بار با استفاده از معادله (۴) تولید می شود. تعداد تکرارها یعنی m هم طوری پیدا می شود که اندازه نمونه جدید ($m \times n$) بزرگ تر یا مساوی تعداد متغیرها باشد. به طور دقیق تر:

$$y^* = (y', \dots, y')',$$

$$W^* = (W_1', \dots, W_m')',$$

$$W_i = X + \delta_i.$$

می توان ثابت نمود که برآورد به دست آمده با استفاده از شبه مشاهدات (یعنی $\tilde{\beta}^{pseudo}$) وقتی تعداد تکرارها خیلی زیاد شود (وقتی $m \rightarrow \infty$) به برآورد $\hat{\beta}^{pseudo}$ میل می کند.

¹³Bootstrap

¹⁴Pseudo Confidence Interval

¹⁵Sparsed pseudo-estimation

تعدیل شده، آماره C_p مالو، معیار اطلاع آکائیکه و بیزی برای این کار می‌تواند استفاده شود. همچنین از روش‌های شبه برآورد و رگرسیون تاوانیده لاسو هم می‌توان برای انتخاب متغیرهای مهم استفاده نمود.

۲.۳ هم خطی

به وجود ارتباط خطی بین متغیرهای پیش‌بین در مدل رگرسیون خطی چندگانه اصطلاحاً هم خطی گفته می‌شود. مشکل هم خطی باعث کاهش دقت برآورد ضرایب رگرسیون شده و خطای استاندارد $\hat{\beta}_j$ را زیاد می‌کند به نحوی که عمده ضرایب متغیرهای واقعاً مهم هم معنادار نمی‌شوند. برای مواجهه با هم خطی می‌توان از روش‌های شبه برآورد، رگرسیون مؤلفه‌های اصلی و رگرسیون تاوانیده‌ی ستیگی استفاده نمود. این روش‌ها، یا با ترکیب متغیرهای پیش‌بین همبسته، تعداد کمتری متغیر پیش‌بین ناهمبسته می‌سازند و یا اساساً با انقباض پارامترهای مدل جلوی ناپایداری و زیاد شدن خطای استاندارد آن‌ها را می‌گیرند.

۳.۳ مدل با ابعاد بالا

اخیراً، در سال ۱۹۹۰ استفاده از داده‌هایی که در آن تعداد متغیرهای پیش‌بین (p) بسیار بیشتر از تعداد مشاهدات (n) می‌باشد، توسط افرون و هستی [۳]، آغاز شد. افزایش زیاد ابعاد داده‌ها یک مشکل اساسی در مبحث داده‌کاوی هم هست. اصطلاح *مدل با ابعاد بالا*^{۱۷} مشخصاً زمانی به کار می‌رود که در یک نمونه، تعداد پارامترهای نامعلوم، بیشتر از اندازه نمونه باشد ($n \ll p$). نخستین بار ایده روش برآورد کننده ستیگی آستانه^{۱۸} توسط شاو و دنگ [۱۰] ارائه شد. در این روش، پیشنهاد می‌گردد که برآوردگر رگرسیون ستیگی، زمانی که بردار ضرایب دارای مقادیر کوچک است، آستانه‌گذاری شود. مدل خطی $y = X\beta + \varepsilon$ را در نظر بگیرید. در برخی از مواقع، نیازی به برآورد خود بردار β نیست و هدف پیش‌بینی یا برآورد متغیر پاسخ است؛ بنابراین، تصویر β روی $R(X)$ در نظر گرفته می‌شود که به صورت $\theta = X'(X'X)^{-1}X\beta = QQ'\beta$ می‌باشد. توجه کنید که همواره $\theta \in R(X)$ بوده و $\theta = \beta$ است اگر و فقط اگر $\beta \in R(X)$ باشد، علاوه بر این همواره داریم: $X\theta = X\beta$ و مدل اول را می‌توان به صورت $y = X\theta + \varepsilon$ هم نوشت که در این حالت برآورد θ برای

متناسب با λ_i یا همان مقادیر ویژه می‌باشند و به این خاطر است که مقادیر ویژه و بردارهای ویژه برای شناخت Σ بکار گرفته می‌شود.

۸.۲ جهت‌های اصلی هسین

این روش بر اساس یافتن جهت‌های اصلی ماتریس هسین^{۱۶} (به اختصار phd) است که نخستین بار توسط لی [۷] بیان شده است. ماتریس هسین با $H(x)$ نمایش داده می‌شود. فرض کنید ماتریس A به صورت $A = \bar{H}(x)\Sigma_x$ تعریف شده باشد. به کمک روش تجزیه طیفی این ماتریس به صورت زیر درمی‌آید:

$$A = \sum \lambda_j e_j e_j' = P\Lambda P'$$

که در آن Λ ماتریسی قطری متشکل از مقادیر ویژه‌ی A روی قطر اصلی آن بوده و ستون‌های ماتریس P متشکل از بردارهای ویژه A است:

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p),$$

$$P = \begin{pmatrix} e_1 & e_2 & \dots & e_p \end{pmatrix}.$$

کاهش بعد در تجزیه طیفی به این صورت است که λ_i ها مرتب می‌شوند و از جایی به بعد حذف می‌گردند. به بیان دیگر، بجای Λ از

$$\Lambda^* = \text{diag}(\lambda_1, \dots, \lambda_q, 0, \dots, 0)$$

استفاده می‌شود که بعضی از مقادیر ویژه صفر شده‌اند. در واقع، $PA\Lambda^*P'$ به وجود می‌آید. دقت کنید که رابطه مقادیر و بردار ویژه برای ماتریس A به صورت $Ae_i = \lambda_i e_i$ است اما از آنجاکه هر بردار در راستای قطرهای بیضی دارای خاصیت $Ae_i = \lambda_i e_i$ می‌باشد، باید محدودیت $b'\Sigma b = 1$ اعمال شود که ماتریس هسین آن به صورت $\bar{H}(x) = E(H(x))$ است.

۳ مشکلات معمول در رگرسیون

۱.۳ انتخاب متغیر

هدف از انتخاب متغیر در رگرسیون، یافتن بهترین مجموعه پیش‌بینی کننده‌ها از میان بسیاری از متغیرها برای گنجاندن در یک مدل است. رگرسیون گام‌به‌گام با بهره‌گیری از معیارهای ضریب تبیین

¹⁶Principal Hessian Directions

¹⁷High-dimensional model

¹⁸Threshold ridge method

$E(y|X) = f(X) = h(\beta'_1 X, \dots, \beta'_k X) = h(W_1, \dots, W_k)$
هدف یافتن β_j ها است به طوری که کمترین خطای ممکن را داشته باشیم. از روش های زیر برای کاهش بعد می توان استفاده کرد:

- رگرسیون مؤلفه اصلی (PCR).
- رگرسیون معکوس قطعه قطعه شده (SIR).
- جهت های اصلی هسین (PHD).

۵.۳ وجود خطای اندازه گیری

برای نخستین بار ادکاک [۱]، مدل های رگرسیونی با خطای اندازه گیری را مطرح کرد. مدلی که در آن متغیر پیش بین یا متغیر پاسخ یا هر دو با خطا اندازه گیری شوند را مدل رگرسیونی با خطای اندازه گیری می گویند و اگر متغیر پیش بین X دارای خطای اندازه گیری باشد، در عمل به جای مشاهده x_i متغیر تصادفی $w_i = x_i + \delta_i$ مشاهده می شود که در آن x_i متغیر مشاهده نشده، w_i متغیر مشاهده شده و δ_i متغیری تصادفی است که بیان کننده خطای اندازه گیری است. در مدل با خطای اندازه گیری با توجه به رابطه $var(W) = var(X) + var(\delta) > var(X)$ واریانس متغیر پیش بین قابل مشاهده W از واریانس متغیر پیش بین مشاهده نشده X بزرگ تر است. در مواجهه با مدل با خطای اندازه گیری می توان از روش های شبه برآورد و رگرسیون های تاوانیده استفاده کرد.

۴ بحث و نتیجه گیری

مدل های رگرسیون خطی ممکن است با چالش های متعددی در کاربردهای عملی روبرو شوند، از جمله انتخاب متغیر، همبستگی، بعد بالا و خطای اندازه گیری. این مسائل می توانند استفاده از برآوردگر سنتی را به چالش بکشانند. در این مقاله، یک مجموعه جامع از روش های متداول و پیشرفته برای حل این مشکلات ارائه شده است. استفاده از این روش ها می تواند در رفع این چالش ها موثر باشد.

استنباط در مورد پارامترهای $X\theta = X\beta$ و پیش بینی y کافی است. درحالی که β ممکن است حاوی تعداد زیادی مؤلفه صفر باشد، θ می تواند هیچ مؤلفه صفری نداشته باشد، اگرچه بسیاری از اجزای θ ممکن است نزدیک به صفر باشند. برآوردگر رگرسیون ستیغی آستانه ای به این صورت تاوانیده می شود که اگر قدر مطلق هر یک از مؤلفه های $\hat{\theta}$ از یک حد آستانه کمتر شود آنگاه آن ضریب برابر صفر قرار داده می شود و در غیر این صورت مقدار آن ضریب به $\hat{\theta}$ منتقل می شود:

$$\tilde{\theta}_j = \hat{\theta}_j I_{|\hat{\theta}_j| > a_n}, \quad j = 1, \dots, p.$$

تعیین مقدار آستانه به صورت زیر انجام می شود:

$$a_n = Cn^{-\alpha}, \quad 0 < \alpha \leq \frac{1}{p}, \quad C > 0,$$

و مقادیر دلخواه α و C به n بستگی ندارند. از این آستانه می توان به عنوان یک روش انتخاب هم متغیر استفاده نمود. روش های شبه برآورد، روش های کاهش ابعاد از جمله SIR و PHD، رگرسیون تاوانیده از جمله رگرسیون ستیغی و لاسو، رگرسیون مؤلفه های اصلی (PCA) می توانند در مدل های با ابعاد بالا استفاده شوند. دقت کنید که روش های کلاسیک از جمله رگرسیون بهترین زیرمجموعه با استفاده از معیارهای مجموع توان های دوم خطا، C_p مالو و AIC و BIC و... در این گونه مدل ها کارایی ندارند زیرا برآورد پارامترها در این روش ها وابسته به معکوس ماتریس $X'X$ است که در این مدل ها معکوس پذیر نیست.

۴.۳ کاهش بعد

هدف از کاهش بعد، برآورد جهت بردارهای ستون X یا فضای ستونی X است که فضای کاهش ابعاد مجموعه ای از پیش بینی کننده های مهم از میان تمام ترکیبات خطی X به دست می دهد. فرض کنید $y = g(X, \varepsilon)$ باشد که X تشکیل یک فضای p بُعدی می دهد. اکنون، هدف ساختن زیر فضای k بُعدی از طریق ترکیبات خطی X است که $y = g(\beta'_1 X, \dots, \beta'_k X, \varepsilon)$ تقریباً همان فضای قبلی را تولید کند. اکنون، برای کاهش بعد فرض کنید

مراجع

- [1] Adcock, R.J. (1877). Note on the method of least squares. *Analyst*, **4**, 183-184.
- [2] Chowdhury, M.Z.I. and Turin, T.C. (2020). Variable selection strategies and its importance in clinical prediction modelling. *Family medicine and community health*, **8(1)**, e000262.
- [3] Hastie, T., Tibshirani, R., and Friedman, J. (2001). *The elements of statistical learning*. Springer series in statistics. New York, NY, USA.
- [4] Horel, A.E. (1962). Applications of ridge analysis to regression problems. *Chem. Eng. Progress.*, **58**, 54-59.
- [5] Kuha, J. (2004). AIC and BIC: Comparisons of Assumptions and Performance. *Sociological Methods & Research*, **33(2)**, 188-229.
- [6] Li, K.C. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association*, **86(414)**, 316-327.
- [7] Li, K.C. (1992). On principal Hessian directions for data visualization and dimension reduction: Another application of Stein's lemma. *Journal of the American Statistical Association*, **87(420)**, 1025-1039.
- [8] Li, B. and Yin, X. (2007). On surrogate dimension reduction for measurement error regression: an invariance law. *The Annals of Statistics*, **35(5)**, 2143-2172.
- [9] Mallows, C.L. (1973). Some Comments on Cp. *Technometrics*, **15(4)**, 661-675.
- [10] Shao, J. and Deng, X. (2012). Estimation in high-dimensional linear models with deterministic design matrices. *The Annals of Statistics*, **40(2)**, 812-831.
- [11] Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, **58(1)**, 267-288.
- [12] Wager, S., Wang, S., and Liang, P.S. (2013). Dropout training as adaptive regularization. *Proceedings of the 26th International Conference on Neural Information Processing Systems*, **1**, 351-359.

On common linear regression problems and their solutions

ZahraJafarian Moorakani¹ and Heydar Ali Mardani-Fard²✉

Abstract:

The traditional linear regression model is represented as $Y = X\beta + \varepsilon$, with the estimated parameter β calculated as $\hat{\beta} = (X'X)^{-1}X'Y$. However, when implementing this estimator in practical applications, several issues may arise, such as variable selection, collinearity, high dimensionality, dimension reduction, and measurement error, which can make it challenging to use the above estimator. The primary problem in most of these cases is the singularity of the matrix $X'X$. A variety of solutions have been proposed to address these problems. In this article, we review these issues and present a comprehensive set of common solutions, as well as some advanced and less commonly used methods, that have the potential to address these problems in an intelligent manner.

Keywords: Colliniarity, dimension reduction, measurement error, ridge regression, psuedo estimation, SIR regression.

¹Yasouj University, Yasouj, IRAN

²✉Yasouj University, Yasouj, IRAN (corresponding author: h_mardanifard@yu.ac.ir)