

مدل رگرسیون لاسو بیزی با خطای نامتقارن در بعد بالا

زهرا خادم بشیری، علی شادرخ، مسعود یارمحمدی

گروه آمار، دانشگاه پیام نور، تهران

تاریخ دریافت: ۱۳۹۸/۱۰/۱۷ تاریخ پذیرش و انتشار: ۱۳۹۹/۰۶/۰۲

چکیده: یکی از بحث‌های چالشی در مدل‌های رگرسیونی انتخاب مدل بهینه است، بدین شکل که چگونه می‌توان متغیرهای توضیحی مهم و متغیرهای قابل اغماض را مشخص کرده و رابطه بین متغیر پاسخ و متغیرهای توضیحی را به‌طور ساده‌تر بیان نمود. با توجه به محدودیت‌های مربوط به انتخاب متغیر به روش کلاسیک نظیر انتخاب گام به گام، می‌توان از روش‌های رگرسیون تاوانیده استفاده کرد. یکی از مدل‌های رگرسیون تاوانیده، مدل رگرسیونی لاسو است که در آن فرض می‌شود خطاها از توزیع نرمال پیروی می‌کنند. در این مقاله، مدل رگرسیون لاسو بیزی با خطایی با توزیع نامتقارن و وجود متغیرهای توضیحی از بعد بالا معرفی می‌شود. سپس با شبیه‌سازی و تحلیل داده‌های واقعی، عملکرد مدل پیشنهادی مورد بحث و بررسی قرار می‌گیرد.

واژه‌های کلیدی: رگرسیون تاوانیده، رگرسیون لاسو بیزی، توزیع لاپلاس نامتقارن، توزیع آلفا چوله‌نرمال.

۱ مقدمه

مدل رگرسیون خطی چندگانه

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

را در نظر بگیرید، که در آن X یک ماتریس $n \times (p + 1)$ ، بردار پارامترهای $p + 1$ بعدی و ε متغیر تصادفی خطا با $p + 1$ بعد بوده که معمولاً از توزیع نرمال با میانگین صفر و واریانس σ^2 پیروی می‌کند.

در صورتی که تعداد متغیرهای توضیحی در مدل رگرسیونی بیشتر از تعداد مشاهدات باشد، احتمال وجود همبستگی میان متغیرهای توضیحی زیاد شده و در نتیجه ممکن است ماتریس $X'X$ وارون پذیر نبوده و برآورد یکتایی برای پارامتر β به دست نیاید. جستجو برای یک مدل رگرسیونی بهینه را انتخاب متغیر یا انتخاب زیر مجموعه می‌نامند. روش‌های کلاسیک زیادی برای انتخاب متغیر وجود دارند که برخی از آن‌ها مانند رگرسیون گام به گام بر اساس دنباله‌ای از آزمون‌های فرض یا برآوردهایی از نوع میانگین توان دوم خطا یا دیگر ملاکهای زیان هستند. این روش‌ها بیشتر بر موضوع انتخاب متغیر متمرکز شده و به برآورد ضرایب رگرسیونی نمی‌پردازند. همچنین در این روش‌ها، تغییرات کوچک در داده‌ها موجب می‌شود یک متغیر به جای دیگری انتخاب شده و نتایج انتخاب‌ها متفاوت شوند، به عبارت دیگر این روش‌ها نسبت به تغییرات کوچک داده‌ها پایدار نیستند. به عنوان یک روش جایگزین برای انتخاب مدل بهینه، می‌توان از روش‌های انقباضی نظیر رگرسیون ستیغی^۱ (هورل و کنار، ۱۹۷۰) و رگرسیون لاسو^۲ (LASSO) (تیشیرانی، ۱۹۹۶) استفاده کرد. در این روش‌ها، برآوردهای ضرایب رگرسیونی محدود شده و به سمت صفر کاهش می‌یابند. این کاهش می‌تواند به طور معنی‌داری واریانس ضرایب رگرسیونی را کاهش دهد.

در روش رگرسیون ستیغی برآورد ضرایب رگرسیونی به صفر میل کرده و از آنجایی که مقادیر آن دقیقاً صفر نشده، از مدل رگرسیونی (۱) به طور کامل حذف نمی‌گردند. بنابراین روش رگرسیون ستیغی را نمی‌توان برای هدف انتخاب متغیر استفاده کرد. روش رگرسیون لاسو به عنوان روشی برای انتخاب متغیر و برآورد پارامتر به طور همزمان، مورد توجه قرار گرفته است. از مزایای اصلی و مهم رگرسیون لاسو می‌توان افزایش دقت پیش‌بینی و بهبود تفسیر مدل رگرسیونی را نام برد. علی‌رغم مزایای این روش، برآوردگر لاسو محدودیت‌هایی نیز داشته و در مواردی که گروهی از متغیرهای توضیحی همبستگی بالایی دارند، این برآوردگر قادر به انتخاب یکی از متغیرهای همبسته بوده و از انتخاب دیگر متغیرهای توضیحی صرف‌نظر می‌کند. برای حل این مسئله، زو و هستی (۲۰۰۵) روش رگرسیون شبکه ارتجاعی^۳ را ارائه دادند که در آن ترکیبی از تابع تاوان رگرسیون ستیغی و رگرسیون لاسو استفاده شده است. مسئله دیگر در روش رگرسیون لاسو، تخمین خطاهای استاندارد برای برآورد پارامترهای ضرایب رگرسیونی است، زیرا الگوریتم‌هایی از قبیل الگوریتم لارس (افرون و همکاران، ۲۰۰۴) فقط برآوردهای نقطه‌ای پارامترهای رگرسیونی را فراهم کرده و استفاده از بوت استرپ برای محاسبه برآورد خطاهای استاندارد و در نهایت محاسبه بازه اطمینان برای پارامترهای مدل رگرسیونی از لحاظ محاسباتی مشکل است (مالیک و نینجان، ۲۰۱۳). روش بوت

¹Ridge Regression

²Least Absolute Shrinkage and Selection Operator

³Elastic net regression

استرپ را **نایت و فو (۲۰۰۰)** برای برآورد خطای استاندارد مورد استفاده قرار داده و تأکید کردند، وقتی برآورد برخی از ضرایب رگرسیونی صفر می‌شوند، برآوردهای بوت استرپ خطای استاندارد به‌طور مجانبی اریب هستند. در این میان روش‌های بیزی علاوه بر اینکه برآوردهای استاندارد برای ضرایب رگرسیونی را فراهم می‌کنند، در حالاتی که تعداد مشاهدات از پارامترها کمتر هستند، نیز عملکرد خوبی دارند. رگرسیون لاسوی بیزی (تیبشیرانی، ۱۹۹۶)، یک روش انتخاب مدل است که در آن خطاها دارای توزیع نرمال فرض شده و چگالی پیشینی برای ضرایب مدل رگرسیونی نیز توزیع لاپلاس متقارن در نظر گرفته شده است. با توجه به توصیه تیبشیرانی، **پارک و کسلا (۲۰۰۸)** استفاده از تابع چگالی لاپلاس متقارن را پیشنهاد دادند. در عمل فرض نرمال بودن مانده‌ها در مدل‌های خطی می‌تواند محدود شود. مسئله انحراف از فرض نرمال بودن وقتی رخ می‌دهد که نمونه شامل داده‌های دور افتاده باشد یا توزیع خطاها نامتقارن باشند. یک روش مواجهه با داده‌های غیر نرمال، بررسی روش‌های تبدیل‌کننده مناسب و نرمال کردن آن‌ها است (مارچنکو و گنتون، ۲۰۱۰).

روش دیگر در برخورد با داده‌های غیرنرمال، جستجوی روش‌های انعطاف‌پذیرتر برای مدل‌سازی داده‌های غیرنرمال است. به عبارت دیگر، توزیع‌هایی با پارامترهای اضافی برای تنظیم این عدم تقارن استفاده شوند. برای مثال می‌توان توزیع خطاها با دم‌های سنگین‌تر را در نظر گرفت. توزیع لاپلاس، دم‌های سنگین‌تری نسبت به نرمال داشته و می‌تواند انتخاب مناسبی برای توزیع خطاها باشد. مزیت مهم دیگر توزیع لاپلاس نسبت به دیگر توزیع‌های دم‌سنگین نظیر توزیع تی - استیودنت آن است که همه گشتاورهای مرکزی آن متناهی است. همچنین برای اصلاح کردن عدم تقارن، می‌توان توزیع‌های چوله‌تی، چوله‌نرمال را در نظر گرفت که نسخه‌های چوله‌ای از توزیع‌های تی - استیودنت و نرمال هستند (مارچنکو و گنتون، ۲۰۱۰). به طور کلی، خانواده توزیع‌های چوله‌بیضوی ارائه شده توسط **ساهو و دی (۲۰۰۳)**، عدم تقارن در کلاسی از توزیع‌های متقارن بیضوی را لحاظ می‌کنند. ساده‌ترین نمایش از خانواده چوله بیضوی همان‌طور که توسط **آزالینی (۱۹۸۵)** معرفی شده، توزیع چوله‌نرمال است که در مقایسه با توزیع نرمال، علاوه بر پارامترهای مقیاس و مکان، پارامتر شکل برای تنظیم عدم تقارن توزیع را نیز دارد که برای مطالعه رفتار نامتقارن مجموعه داده‌های تجربی در حوزه‌های مختلف استفاده می‌شود. **الورا (۲۰۱۰)** خانواده‌ای از توزیع‌ها را معرفی کرد که هر دو حالت تک مدی بودن و دو مدی بودن شکل را حمایت می‌کنند. این خانواده جدید توزیع را آلفا چوله‌نرمال^۱ نامیده و با نماد $ASN(\alpha)$ نشان می‌دهند، که در آن α نشان‌دهنده پارامتر عدم تقارن و $ASN(0)$ بیانگر توزیع نرمال متقارن است.

¹Alpha Skew Normal Distribution

هدف مقاله حاضر، بیان تحلیل بیزی متفاوتی برای داده‌های نامتقارن و با خطای آلفا چوله‌نرمال است. در این راستا در بخش ۲ به توسیع روش لاسوی بیزی در مواردی که توزیع خطاهای مدل چوله است، پرداخته و سپس مدل برای مسئله رگرسیونی با بعد بالا ارایه می‌شود. در بخش ۳ پارامتر تنظیم‌کننده مدل رگرسیونی تاوانیده را برآورد کرده و در بخش ۴ با استفاده از روش‌های شبیه‌سازی، دو مثال در شرایط مختلف از نظر همبستگی میان متغیرهای توضیحی، تعداد پارامترها و تعداد مشاهدات مدل مورد بررسی قرار می‌گیرند. در ادامه با تحلیل داده‌های واقعی به ارزیابی روش‌های رگرسیونی پرداخته می‌شود.

۲ مدل رگرسیون لاسوی بیزی نامتقارن

در مدل رگرسیون خطی چندگانه (۱) اغلب اوقات برآوردگر کمترین توان‌های دوم خطا انتخاب مناسبی است، اما در مواردی که متغیرهای توضیحی همبستگی بالا داشته باشند، این برآوردگر کارایی لازم را ندارد (آرست و همکاران، ۱۳۹۸). در حضور همبستگی بین متغیرهای توضیحی، مقادیر واریانس‌ها و کواریانس‌ها برای برآوردهای کمترین توان‌های دوم خطای ضرایب رگرسیونی، بزرگ خواهد شد. در واقع اعضای قطر ماتریس $(X'X)^{-1}$ عبارتند از $\frac{1}{1-R_j^2}$ ، $j = 1, \dots, p$ ، که در آن R_j^2 ضریب تعیین چندگانه از رگرسیون متغیر توضیحی j ام نسبت به $p-1$ متغیر توضیحی باقیمانده است که در صورت وجود همخطی شدید میان متغیرهای توضیحی، R_j^2 نزدیک به واحد خواهد بود و چون واریانس ضریب رگرسیونی j ام برابر با $V(\hat{\beta}_j) = (1 - R_j^2)^{-1} \sigma^2$ است، همخطی شدید موجب می‌شود که واریانس برآورد کمترین توان‌های دوم ضریب رگرسیونی β_j بسیار زیاد شود (نتر و همکاران، ۱۹۹۶).

روش‌های رگرسیون تاوانیده برای بهبود برآوردهای کمترین توان‌های دوم خطا ارایه شده‌اند. در این روش‌ها با افزایش کمی اریبی، واریانس مدل کاهش می‌یابد. یکی از روش‌های تاوانیده که اخیراً توجه زیادی را به خود جلب کرده، مدل رگرسیون لاسو است. دلیل محبوبیت این روش آن است که می‌تواند به طور همزمان برآورد پارامترهای رگرسیونی و انتخاب متغیرهای مهم را انجام داده و منجر به مدل‌های رگرسیونی دقیق با تفسیرپذیری بالا شود. به عبارت دیگر در این روش، شانس آمدن متغیرهایی که توان پیش‌بینی پایینی دارند، در حضور دیگر متغیرهای توضیحی در مدل، کم است. بنابراین این روش را می‌توان برای هدف انتخاب متغیر استفاده کرد. دو الگوریتم بسیار مهم برای محاسبه برآورد لاسو، الگوریتم‌های لارس و مختصات مسیر است (افرون و همکاران، ۲۰۰۴). با وجود مزایای روش رگرسیون لاسو، در حالت $p \gg n$ ، این روش با مشکل مواجه است. زیرا روش لاسو نمی‌تواند بیش از اندازه نمونه n ، متغیرهای توضیحی را انتخاب کند و حداکثر قادر به انتخاب n متغیر است (زو و هستی، ۲۰۰۵). برآورد

لاسو برای پارامترهای مدل رگرسیون خطی چندگانه را می‌توان به صورت برآورد بیزی تحت چگالی پیشینی لاپلاس متقارن برای پارامترهای مدل رگرسیونی نیز تفسیر کرد. **پارک و کسلا (۲۰۰۸)** در روش لاسو بیزی، تابع چگالی لاپلاس را به صورت آمیخته‌ای از چگالی‌های نرمال بیان کرده و نمونه‌گیر گیبس را برای تولید نمونه‌هایی از توزیع پسینی اتخاذ نمودند.

۲.۱ مدل رگرسیون لاسوی بیزی

در مدل رگرسیون خطی چندگانه (۱)، توزیع خطاها نرمال با میانگین صفر و واریانس σ^2 است و مجموع توان دوم مانده‌ها نیز به صورت

$$RSS = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2,$$

خواهد بود. هدف روش کمترین توان‌های دوم این است که مقادیری از β_0, \dots, β_p را بیابد که کمیت RSS را مینیمم کند. در مدل رگرسیون لاسو با در نظر گرفتن محدودیت $\lambda \sum_{j=1}^p |\beta_j|$ ، برآوردهای ضرایب رگرسیونی محاسبه می‌شوند. در واقع عبارت $RSS + \lambda \sum_{j=1}^p |\beta_j|$ نسبت به β مینیمم می‌شود. λ پارامتر تعدیل‌کننده غیر منفی است و برای کنترل سطح انقباض به کار می‌رود. یعنی اگر مقدارش برابر صفر باشد، همه متغیرها در مدل حضور خواهند داشت و اگر مقدار آن افزایش یابد، تعداد متغیرهای مستقل در مدل کاهش می‌یابد. مقدار این پارامتر معمولاً توسط روش اعتبارسنجی متقابل تعیین می‌شود. در مدل رگرسیون لاسو بیزی برای پارامترهای β_j توزیع پیشینی لاپلاس به صورت $\pi(\beta) = \prod_{j=1}^p \frac{\lambda}{\gamma} e^{-\lambda|\beta_j|}$ ، در نظر گرفته می‌شود. بنابراین توزیع پسینی به شرط \bar{y} برای پارامترها به صورت

$$\pi(\beta, \sigma^2 | \bar{y}) \propto \pi(\sigma^2) (\sigma^2)^{-(n-1)/2} \exp \left\{ \frac{-1}{2\sigma^2} (\bar{y} - \mathbf{X}\beta)^T (\bar{y} - \mathbf{X}\beta) - \lambda \sum_{j=1}^p |\beta_j| \right\},$$

است، که در آن $\bar{y} = (y - \bar{y})$. برای هر مقدار ثابت σ^2 ، در صورت وجود نما برای توزیع پسین، برآورد بیز همان نمای پسین خواهد بود که به λ و انتخاب چگالی پیشینی برای σ^2 بستگی دارد (**پارک و کسلا، ۲۰۰۸**).

۲.۲ مدل رگرسیون لاسو بیزی با توزیع خطای آلفا چوله‌نرمال

در روش رگرسیون لاسو بیزی، فرض می‌شود خطاها دارای توزیع نرمال بوده و توزیع مشاهدات به صورت $y|\beta, \sigma^2 \sim N(X\beta, \sigma^2 I_n)$ در نظر گرفته می‌شود. همچنین متناسب با توزیع مشاهدات، چگالی پیشینی پارامترهای β لاپلاس متقارن است. اما در عمل توزیع خطاها همیشه نرمال نیست. بنابراین در مواردی که توزیع خطاها غیرنرمال است، رگرسیون لاسو بیزی کارایی لازم را برای انتخاب بهترین مدل ندارد. در این بخش فرض می‌شود مشاهدات \mathbf{Y} دارای توزیع آلفا چوله‌نرمال الیورا (۲۰۱۰) به صورت

$$f_Y(\mathbf{Y}|\beta, \sigma^2, k) = \prod_{i=1}^n \left(\frac{(1 - ky_i)^2 + 1}{2 + k^2} \right) \phi(\mathbf{Y}),$$

$$\phi(\mathbf{Y}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{ -\frac{1}{2\sigma^2} (\mathbf{Y} - X\beta)^T (\mathbf{Y} - X\beta) \right\},$$

است، به طوری که $k \in R$ همان پارامتر α در توزیع آلفا چوله‌نرمال بوده و میزان چولگی توزیع را تعیین می‌کند. با در نظر گرفتن تابع چگالی پیشینی لاپلاس نامتقارن

$$\pi(\beta|\sigma^2) \propto \prod_{j=1}^p \frac{\lambda}{\sigma^2} \frac{k}{1 + k^2} \begin{cases} \exp(-\frac{k\lambda}{\sigma^2} \beta_j) & \beta_j < 0 \\ \exp(\frac{\lambda}{k\sigma^2} \beta_j) & \beta_j \geq 0, \end{cases}$$

برای بردار ضرایب رگرسیونی β چگالی گامای معکوس با پارامتر γ, a برای پارامتر σ^2 به صورت

$$\pi(\sigma^2) = \frac{\gamma^a}{\Gamma(a)} \left(\frac{1}{\sigma^2} \right)^{a+1} \exp\left(-\frac{\gamma}{\sigma^2}\right),$$

و چگالی پسینی توأم β, σ^2 عبارت است از

$$\pi(\beta, \sigma^2 | \mathbf{Y}) \propto \pi(\sigma^2) (\sigma^2)^{-\frac{n-1}{2}} \exp\left\{ -\frac{1}{2\sigma^2} (\mathbf{Y} - X\beta)^T (\mathbf{Y} - X\beta) \right\} \pi(\beta | \sigma^2). \quad (2)$$

قضیه ۱. (کوتز و همکاران، ۲۰۰۱) اگر $Z \sim N(0, 1)$ ، $X \sim \text{Asymmetric Laplace}(k, 1)$ ، $W \sim \text{Exp}(1)$ ، همچنین Z و W مستقل باشند، آنگاه $X \stackrel{d}{=} kW + \sqrt{W}Z$.

برهان: متغیر تصادفی W با تابع چگالی e^{-w} را در نظر بگیرید. تابع مشخصه عبارت $kW + W^{\frac{1}{2}}Z$

را به شرط W می‌توان به صورت

$$E[e^{it(kW+\sqrt{W}Z)}] = \int_0^{+\infty} e^{itkw} E(e^{it\sqrt{w}Z}) e^{-w} dw,$$

به دست آورد. توجه داشته باشید که $E(e^{it\sqrt{w}Z}) = \phi_z(t\sqrt{w}) = e^{-\frac{1}{2}t^2 w}$ و $\phi_z(s) = e^{-\frac{1}{2}s^2}$ تابع مشخصه نرمال استاندارد است. بنابراین

$$E[e^{it(kW+\sqrt{W}Z)}] = \int_0^{+\infty} e^{itkw} e^{-\frac{1}{2}t^2 w} e^{-w} dw = \int_0^{+\infty} e^{-w(1+\frac{1}{2}t^2-itk)} dw ,$$

که با حل این انتگرال به سادگی تابع مشخصه تابع چگالی لاپلاس نامتقارن به صورت

$$\psi(t) = \frac{1}{1 + \frac{1}{2}t^2 - itk},$$

به دست می‌آید. چون بین توزیع‌های احتمال و توابع مشخصه رابطه یک به یک وجود دارد یا به عبارت دیگر هر دو متغیر تصادفی X_1, X_2 دارای توزیع احتمال یکسان هستند اگر و فقط اگر تابع مشخصه‌های این دو یکسان باشند، می‌توان نتیجه گرفت $X \stackrel{d}{=} kW + \sqrt{W}Z$. با توجه به قضیه ۱ رابطه

$$f_X(x) = \int_0^{+\infty} f_Z\left(\frac{x-kw}{\sqrt{w}}\right) \frac{1}{\sqrt{w}} f_W(w) dw = \int_0^{+\infty} \frac{1}{\sqrt{2\pi w}} e^{-\frac{1}{2}\left(\frac{x-kw}{\sqrt{w}}\right)^2} e^{-w} dw,$$

برای تابع چگالی لاپلاس نامتقارن برقرار است. بنابراین با فرض $\tau_j^2 = w$ در قضیه ۱، چگالی پیشینی لاپلاس نامتقارن مرتبط با پارامتر ضرایب مدل رگرسیونی (۱) را به صورت

$$\pi(\beta, \tau^2 | \lambda, \sigma^2) \propto \prod_{j=1}^p \pi(\beta | \tau_j^2, \sigma^2) \pi(\tau_j^2 | \lambda),$$

می‌توان نوشت. فرض می‌شود τ_j^2 دارای توزیع نمایی با پارامتر $\frac{\lambda}{2}$ است. بنابراین

$$\beta | \tau_1^2, \dots, \tau_p^2, \sigma^2 \sim N_p(k\tau, \sigma^2 \mathbf{D}_\tau), \tau = (\tau_1^2, \dots, \tau_p^2), \mathbf{D}_\tau = \text{diag}(\tau). \quad (۳)$$

با در نظر گرفتن چگالی (۳) در چگالی پسینی توأم (۲)، توزیع توأم به دست آمده، به منظور استنباط در مورد پارامترهای مدل رگرسیونی بیزی استفاده می‌شود. ابتدا برای تولید نمونه‌هایی از تابع چگالی پسینی توأم با استفاده از الگوریتم گیبس، باید توابع چگالی پسینی حاشیه‌ای برای پارامترهای β, τ, σ^2 محاسبه شوند. برای یافتن چگالی پسینی حاشیه‌ای پارامتر β ، ابتدا عباراتی از چگالی پسینی توأم (۲) که شامل β هستند در نظر گرفته می‌شود:

$$\begin{aligned} & \exp\left(-\frac{1}{2\sigma^2}(\mathbf{Y} - X\beta)^T(\mathbf{Y} - X\beta)\pi(\beta|\sigma^2)\right) \\ &= \prod_{j=1}^p \frac{1}{\sqrt{2\pi\sigma^2\tau_j^2}} \exp\left(-\frac{1}{2\sigma^2\tau_j^2}(\beta_j - k\tau_j)^2\right) \\ &\sim \exp\left(-\frac{1}{2\sigma^2}(\mathbf{Y} - X\beta)^T(\mathbf{Y} - X\beta)\right) \exp\left(-\frac{1}{2\sigma^2}(\beta - k\tau)^T \mathbf{D}_\tau^{-1}(\beta - k\tau)\right) \\ &= \exp\left[-\frac{1}{2\sigma^2}\left[(\beta' \mathbf{D}_\tau^{-1} \beta + (k\tau)' \mathbf{D}_\tau^{-1} (k\tau) - (k\tau)' \mathbf{D}_\tau^{-1} \beta - \beta' \mathbf{D}_\tau^{-1} k\tau)\right]\right] \\ &= \exp\left[-\frac{1}{2\sigma^2}\left[y'y + \beta' \mathbf{A} \beta - 2\beta' (\mathbf{X}' \mathbf{Y} + \mathbf{D}_\tau^{-1} k\tau)\right]\right] \\ &\sim \exp\left[\frac{-1}{2\sigma^2}\left[(\beta - \mathbf{A}^{-1}(\mathbf{X}' \mathbf{Y} + \mathbf{D}_\tau^{-1} k\tau))\right]^T \mathbf{A} (\beta - \mathbf{A}^{-1}(\mathbf{X}' \mathbf{Y} + \mathbf{D}_\tau^{-1} k\tau))\right. \\ &\quad \left.+ y'(\mathbf{I} - (\mathbf{X}' + (k\tau)' \mathbf{D}_\tau^{-1}) \mathbf{A}^{-1} (\mathbf{X} + \mathbf{D}_\tau^{-1} k\tau)) \mathbf{Y}\right], \end{aligned}$$

بنابراین تابع چگالی پسینی حاشیه‌ای برای بردار پارامتر ضرایب β ، چگالی نرمال با پارامتر میانگین $\mathbf{A}^{-1}(\mathbf{X}' \mathbf{Y} + \mathbf{D}_\tau^{-1} k\tau)$ و واریانس $\sigma^2 \mathbf{A}^{-1}$ است. برای یافتن تابع چگالی پسینی حاشیه‌ای σ^2 عبارات شامل σ^2 از تابع چگالی پسینی توأم (۲) جدا می‌شوند، بنابراین

$$\begin{aligned} & \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{Y} - X\beta)^T(\mathbf{Y} - X\beta)\right\} \\ & \times \prod_{j=1}^p \exp\left(-\frac{1}{\sqrt{2\pi\sigma^2\tau_j^2}}\right) \exp\left(-\frac{1}{2\sigma^2}(\beta - k\tau)^T \mathbf{D}_\tau^{-1}(\beta - k\tau)\right) \left(\frac{1}{\sigma^2}\right)^{a+1} \exp\left(\frac{-\gamma}{\sigma^2}\right) \\ & \propto (\sigma^2)^{-\frac{n}{2} - \frac{p}{2} - a - 1} \exp\left[-\frac{1}{\sigma^2} \left(\frac{(\mathbf{Y} - X\beta)^T(\mathbf{Y} - X\beta)}{2} \right. \right. \\ & \quad \left. \left. + \frac{\tau\sigma^2(\beta - k\tau)^T \mathbf{D}_\tau^{-1}(\beta - k\tau)}{2} + \gamma\right)\right]. \end{aligned}$$

در نتیجه

$$\sigma^2 | \beta, \tau \sim \mathbf{IG} \left(\frac{\mathbf{n}}{2} + \frac{\mathbf{p}}{2} + \mathbf{a}, \frac{(\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta) + (\beta - \mathbf{k}\tau)^T \mathbf{D}_\tau^{-1} (\beta - \mathbf{k}\tau)}{2} + \gamma \right).$$

برای یافتن تابع چگالی پسینی حاشیه‌ای τ_j^2 ، عبارات شامل τ_j^2 از تابع چگالی پسینی (۲) جدا شده و داریم:

$$\begin{aligned} & \prod_{j=1}^p (\pi \sigma^2 \tau_j^2)^{-\frac{1}{2}} \exp \left(-\frac{1}{2\tau_j^2} (\beta_j - k\tau_j^2)^2 \right) \frac{\lambda^2}{\gamma} \exp \left(-\frac{\lambda^2}{\gamma} \tau_j^2 \right) \\ & \sim \prod_{j=1}^p (\sigma^2 \tau_j^2)^{-\frac{1}{2}} \exp \left(-\frac{1}{2\sigma^2 \tau_j^2} (\beta_j - k\tau_j^2)^2 \right) \exp \left(-\frac{\lambda^2}{\gamma} \tau_j^2 \right). \end{aligned}$$

برای تولید نمونه از τ_j^2 از الگوریتم متروپلیس هستینگ استفاده می‌شود.

۳ برآورد پارامتر تنظیم λ

برای برآورد پارامتر تنظیم مدل رگرسیونی تاوانیده می‌توان از روش تعیین ابر پیشین برای پارامتر λ استفاده کرد. در این روش می‌توان به جای λ ، λ^2 در نظر گرفته و کلاسی از چگالی‌های پیشینی گاما به صورت

$$\pi(\lambda^2) \propto \frac{\delta^r}{\Gamma(r)} (\lambda^2)^{r-1} e^{-\delta \lambda^2}, r > 0, \delta > 0,$$

برای λ^2 تعیین نمود. این چگالی پیشینی را نیز وارد الگوریتم گیبس کرده و طبق روال قبل عباراتی که شامل λ^2 هستند را جدا کرده و به صورت

$$(\lambda^2)^{p+r-1} \exp \left(-\lambda^2 \left(\frac{1}{\gamma} \sum_{j=1}^p \tau_j^2 + \delta \right) \right).$$

نوشت. با توجه به فرم این تابع، می‌توان گفت تابع چگالی توزیع حاشیه‌ای پسینی برای λ^2 توزیع گاما با پارامتر شکل $p+r$ و پارامتر مقیاس $\delta > 0$ ، $\frac{\tau_j^2}{\gamma} + \delta$ است. لازم به ذکر است که تابع چگالی پسینی حاشیه‌ای پارامترهای σ^2, τ, β بدون تغییر باقی می‌مانند.

۴ مطالعه شبیه‌سازی

در یک مطالعه شبیه‌سازی چهار مدل رگرسیونی بیزی، شامل رگرسیون ستیغی بیزی هیانگ (۱۹۷۵)، رگرسیون شبکه ارتجاعی بیزی لی و لین (۲۰۱۰)، رگرسیون لاسوی بیزی پارک و کسلا (۲۰۰۸) و رگرسیون لاسوی بیزی نامتقارن برای انتخاب مدل بهینه رگرسیونی، با هم مقایسه می‌شوند. طبق فرمولبندی رگرسیون ستیغی بیزی، مدل بیز سلسله مراتبی به صورت

$$\begin{aligned} y|X, \beta, \sigma^2 &\sim N_n(X\beta, \sigma^2 I_n), \\ \beta|\tau^2 &\sim \prod_{j=1}^p N(\cdot, \tau_j^2), \\ \tau_j^2 &\sim Inv - \chi^2(v, s^2), \end{aligned}$$

و در مدل رگرسیون شبکه ارتجاعی بیزی، مدل سلسله مراتبی به صورت

$$\begin{aligned} y|\beta, \sigma^2 &\sim N(X\beta, \sigma^2 I_n), \\ \beta|\tau, \sigma^2 &\sim \prod_{j=1}^p N\left(\cdot, \left(\frac{\lambda_j}{\sigma^2} \frac{\tau_j}{\tau_j - 1}\right)^{-1}\right), \\ \tau|\sigma^2 &\sim \prod_{j=1}^p TG\left(\frac{1}{\nu}, \frac{\lambda_j \sigma^2}{\lambda_j^2}, (1, \infty)\right), \\ \pi(\sigma^2) &\sim \frac{1}{\sigma^2}, \end{aligned}$$

تعریف می‌شود که TG بیانگر توزیع گامای بریده شده^۱ است. در ادامه شبیه‌سازی با دو مثال متفاوت در شرایط مختلف از نظر همبستگی میان متغیرهای توضیحی، تعداد پارامترها و تعداد مشاهدات مدل بررسی شده است. برای حل مسئله انتخاب مدل بیزی بهینه، ملاک‌هایی از قبیل DIC بردلی و همکاران (۲۰۰۲)، BIC اصلاح شده (mBIC) بوگدان و همکاران (۲۰۰۴) و BIC تعمیم یافته (EBIC) چن و چن (۲۰۰۸) پیشنهاد شده‌اند.

فرض کنید برای مشاهدات \mathbf{Y} ، تابع چگالی و $\theta = (\theta_1, \dots, \theta_p)$ بردار پارامترهای مدل

¹Truncated Gamma distribution

است. ملاک EBIC به صورت

$$EBIC = -2 \log L(\theta|\mathbf{Y}) + p \log n + 2\gamma \log p,$$

تعریف شود، که در آن $L(\theta|\mathbf{Y})$ تابع درستنمایی و $\gamma \in [0, 1]$ پارامتر تنظیم‌کننده است. ملاک BIC حالت خاصی از EBIC با $\gamma = 0$ و ملاک mBIC نیز معادل با EBIC با $\gamma = 1$ است. ملاک DIC نیز تعمیمی از ملاک AIC برای مسائل انتخاب مدل بی‌زی است و به صورت

$$DIC = \overline{D(\theta)} + P_D,$$

تعریف می‌شود، که در آن $\overline{D(\theta)} = -2 \log L(\theta|\mathbf{Y})$ را انحراف نامیده و تابعی از θ است. عبارت اول، به صورت امید انحراف تحت تابع چگالی پسینی پارامتر به صورت

$$\overline{D(\theta)} = E_{\theta|\mathbf{Y}}[D(\theta)] = E_{\theta|\mathbf{Y}}[-2 \log L(\theta|\mathbf{Y})],$$

تعریف می‌شود. مؤلفه دوم، تعداد پارامترهای مؤثر یا P_D را اندازه می‌گیرد که به صورت اختلاف بین میانگین پسین انحراف و انحراف محاسبه شده در میانگین پسین پارامترها ($\bar{\theta}$) تعریف می‌شود. یعنی

$$P_D = \overline{D(\theta)} - D(\bar{\theta}) = E_{\theta|\mathbf{Y}}[D(\theta)] - D(\mathbf{E}_{\theta|\mathbf{Y}}[\theta]) = \mathbf{E}_{\theta|\mathbf{Y}}[-2 \ln L(\theta|\mathbf{Y})] + 2 \ln L(\bar{\theta}|\mathbf{Y}).$$

با دوباره‌آرایی عبارت P_D داریم: $\overline{D} = D(\bar{\theta}) + P_D$ ، بنابراین $DIC = D(\bar{\theta}) + 2P_D$ و محاسبه آن با استفاده از الگوریتم‌های زنجیره مارکف مونت کارلو از قبیل گیبس اغلب بدیهی و ساده است. برآورد \overline{D} با میانگین گرفتن از مقادیر شبیه‌سازی شده از $D(\theta)$ محاسبه می‌شود. مقدار پارامترهای مؤثر P_D را نیز می‌توان با محاسبه $D(\theta)$ در میانگین نمونه‌ای مقادیر شبیه‌سازی شده پارامترها و کم کردن آن از برآورد $D(\bar{\theta})$ به دست آورد.

مثال ۰۱. مقادیر مشاهدات X, \mathbf{Y} با در نظر گرفتن فرضیات $n = 100, p = 10$ به این صورت تولید می‌گردد که مشاهدات X را از توزیع نرمال چند متغیره استاندارد نمونه‌گیری کرده و بردار ضرایب نیز به صورت $(0, \dots, 0, 32, -28, 4, 5, -3, 3)$ β در نظر گرفته می‌شود. بردار n بعدی \mathbf{e} از مقادیر خطا با نمونه‌گیری از تابع چگالی آلفا چوله‌نرمال با میزان چولگی $k = 0.8$ و سپس طبق مدل (۱) مشاهدات

Y به دست می‌آیند. سپس مشاهدات Y را مرکزی کرده به گونه‌ای که میانگین صفر داشته باشند و همچنین مقادیر متغیر کمکی X نیز استاندارد می‌شوند. با استفاده از شبیه‌سازی و انجام الگوریتم گیبس مقادیر ملاکهای DIC و EBIC با $\gamma = 1$ برای چهار مدل مختلف رگرسیون ریج بیزی، رگرسیون لاسوی بیزی، رگرسیون شبکه ارتجاعی بیزی و رگرسیون لاسوی بیزی نامتقارن در جدول ۱ آمده است. طبق هر دو ملاک، مینیمم مقادیر DIC و EBIC، مدل با بهترین پیش‌بینی را فراهم می‌کند. معمولاً اگر اختلاف مقادیر هر کدام از این ملاکها در دو مدل متفاوت، بیشتر از 10 باشد، مدل با مقادیر DIC و EBIC کمتر به عنوان مدل بهینه انتخاب می‌شود. بنابراین طبق نتایج جدول ۱ می‌توان گفت مدل لاسوی بیزی نامتقارن نسبت به روش‌های دیگر در این داده‌های شبیه‌سازی، عملکرد بهتری در انتخاب متغیر دارد.

جدول ۱. مقادیر DIC و EBIC به ازای اندازه نمونه 100 و تعداد پارامترها 10

EBIC	DIC	مدل رگرسیون
۲۳۵/۳۴	۹۴/۲	ستیغی بیزی
۱۱۸/۶	۸۴/۲۳	لاسوی بیزی
۴۳۲/۸۷	۳۶۸/۲	شبکه ارتجاعی بیزی
۱۰۴/۲	۶۲/۲	لاسوی بیزی نامتقارن

مثال ۲. در این مثال مانند مثال ۱ عمل می‌شود، با این تفاوت که فرض می‌شود $n = 100$ ، $p = 400$ ، یعنی $n \gg p$ است. همچنین میان متغیرهای توضیحی، همبستگی به صورت $\text{Cov}(X)_{ij} = \rho^{|i-j|}$ لحاظ شده است. مقادیر DIC و EBIC با $\gamma = 1$ حاصل از شبیه‌سازی برای چهار مدل مختلف را در جدول ۲ مشاهده می‌کنید. واضح است که ملاک DIC و EBIC برای مدل رگرسیون لاسوی بیزی نامتقارن به طور قابل ملاحظه‌ای کمتر از سایر مدل‌ها است و در نتیجه بهتر از سه مدل دیگر عمل می‌کند.

جدول ۲. مقادیر DIC و EBIC به ازای اندازه نمونه 100 و تعداد پارامترها 400

EBIC	DIC	مدل رگرسیون
۲۷۹۰/۶	-۱۹۹/۳۱	ستیغی بیزی
۱۹۷۴/۵	۲۱۹/۳۱	لاسوی بیزی
۴۸۹۶/۸	-۳۹۶	شبکه ارتجاعی بیزی
۱۸۴۵/۲	۷۵/۹۸	لاسوی بیزی نامتقارن

۵ تحلیل داده واقعی

از آنجائی که وجود ریوفلاوین یا ویتامین B۲ در بدن انسان، باعث آزادسازی انرژی از کربوهیدرات‌ها، پروتئین و چربی می‌شود و این امر در رشد و ترمیم پوست، مو، ناخن و مفاصل ضروری و باعث حفاظت سیستم ایمنی بدن در مقابل بیماری‌ها است. اندازه‌گیری میزان این ویتامین امری ضروری است. در این اندازه‌گیری متغیر پاسخ، لگاریتم میزان ریوفلاوین است.

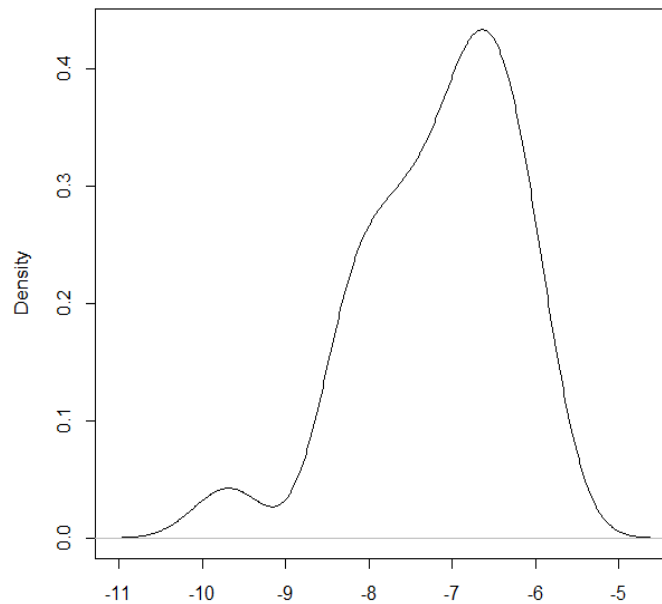
داده‌ها شامل $p = 4088$ ژن متفاوت برای $n = 71$ نمونه است که با توجه به اندازه متغیر و اندازه نمونه، مجموعه داده‌ای با ابعاد بالا محسوب می‌شود (بولمن و همکاران، ۲۰۱۴). مجموعه داده‌های ریوفلاوین از طریق <https://rdrr.io/cran/hdi/man/riboflavin.html> قابل دسترس است و نمودار توزیع آن‌ها در شکل ۱ ارائه شده است. با توجه به نمودار، مقداری چولگی در داده‌ها مشاهده می‌شود. همچنین با توجه به مقدار احتمال ($P\text{-value} = 0.005$) به دست آمده از آزمون شاپیرو و ویلک (۱۹۶۵) در سطح خطای ۵ درصد، فرض صفر که نرمال بودن داده‌ها است، رد شده و توزیع داده‌های ریوفلاوین از توزیع نرمال متفاوت است. پس از اجرای روش‌های رگرسیون لاسو بیزی، رگرسیون لاسو بیزی نامتقارن، رگرسیون ستیغی بیزی و رگرسیون شبکه ارتجاعی بیزی برای داده‌های ریوفلاوین، نتایج در جدول ۳ ارائه شده است. با توجه به مقادیر DIC، EBIC، AIC، BIC و CV در چهار روش رگرسیونی بیزی، نتیجه می‌شود که روش رگرسیون لاسو بیزی نامتقارن بهتر از بقیه روش‌ها به انتخاب مدل بهینه می‌پردازد.

جدول ۳. مقادیر ملاک‌های DIC، EBIC، AIC، BIC و CV برای داده‌های ریوفلاوین

BIC	AIC	CV	EBIC	DIC	مدل رگرسیون
۶۵۹/۲	۶۵۰/۱	۵۶۶/۲	۱۷۸۹۲/۹	۶۲۵/۱۹	ستیغی بیزی
۳۶۵/۱	۳۶۲/۴	۵۶۳/۴	۱۷۸۵۰/۶	۳۵۴/۷	لاسوی بیزی
۵۸۰/۳	۵۷۵/۲	۵۶۵/۱	۱۷۸۷۲/۲	۵۶۵/۱	شبکه ارتجاعی بیزی
۲۹۶/۴	۲۸۹/۱	۵۶۲/۳	۱۷۶۴۵/۴	۲۶۲/۲	لاسوی بیزی نامتقارن

بحث و نتیجه‌گیری

تعمیمی از مدل رگرسیون لاسو بیزی با توزیع خطاهای آلفا چوله‌نرمال وقتی $n \gg p$ ارائه شد. در مدل پیشنهادی از تابع چگالی پیشینی لاپلاس نامتقارن برای پارامترهای ضرایب رگرسیونی استفاده شده و



شکل ۱. نمودار توزیع داده‌های ریپوفلاوین

پارامترهای مدل نیز با الگوریتم‌های MCMC برآورد شده است. با توجه به نتایج شبیه‌سازی و تحلیل داده‌های واقعی، می‌توان نتیجه گرفت که بر اساس ملاک‌های DIC، EBIC، AIC، BIC و CV مدل پیشنهادی رگرسیون لاسوی بیزی نامتقارن در مقایسه با مدل‌های رگرسیون لاسوی بیزی، رگرسیون ستیغی بیزی و رگرسیون شبکه ارتجاعی بیزی برآزش بهتری برای داده‌ها فراهم می‌کند.

تقدیر و تشکر

نویسندگان مقاله از سر دبیر، داوران و ویراستار محترم مجله در ارزیابی این مقاله قدردانی و تشکر می‌کنند.

مراجع

آرست، م.، آرشی، م. و ربیعی، م. (۱۳۹۸)، مطالعه رفتار برآوردگر انقباضی تحت یک قید خطی در مدل رگرسیون تاوانیده، مجله علوم آماری، ۱۳، ۱-۱۴.

- Azzalini, A. (1985), A Class of Distributions Which Includes the Normal Ones, *Scandinavian Journal of Statistics*, **12**, 171-178.
- Bogdan, M., Doerge, R. W. and Ghosh, J. (2004), Modifying the Schwarz Bayesian Information Criterion to Locate Multiple Interacting Quantitative Trait Loci, *Genetics*, **167**, 989-999.
- Buhlmann, P., Kalisch, M. and Meier, L. (2014), High-Dimensional Statistics with a View Towards Applications in Biology, *Annual Review of Statistics and Its Applications*, **1**, 255-278.
- Chen, J., and Chen, Z. (2008), Extended Bayesian Information Criteria for Model Selection with Large Model Space, *Biometrika*, **95**, 759-771.
- Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R. (2004), Least Angle Regression, *Annals of Statistics*, **32**, 407-499.
- Elal-Olivero, D. (2010), Alpha-Skew-Normal Distribution, *Proyecciones Journal of Mathematics*, **29**, 224 - 240.
- Hoerl, A. E., and Kennard, R. W. (1970), Ridge Regression: Biased Estimation for Nonorthogonal Problems, *Technometrics*, **12**, 55-67.
- Hsiang, T. C. (1975), A Bayesian View on Ridge Regression, *The Statistician*, **24**, 267-268.
- Kotz, S., Koubowski, T. and Podgorski, K. (2001), *The Laplace Distribution and Generalizations*, Birkhauser Basel, New York.
- Knight, K., and Fu, W. (2000), Asymptotics for LASSO-type Estimators, *Annals of Statistics*, **28**, 1356-1378.

- Li, Q., and Lin, N. (2010), The Bayesian Elastic Net, *Bayesian Annals of Statistics*, **5**, 151-170.
- Mallick, H., and Nengjun, Y. (2013), Bayesian Methods for High Dimensional Linear Models, *Journal of Biometrics & Biostatistics*, **1**, 1-13.
- Marchenko, V., and Genton, M. (2010), A Suite of Commands for Fitting the Skew-normal and Skew-t Models, *Stata Journal*, **10**, 507-539.
- Neter, J., Kutner, H., Wasserman W. and Nachtsheim, J. (1996), *Applied Linear Regression Models*, McGraw-Hill College.
- Park, T., and Casella, G. (2008), The Bayesian LASSO, *Journal of the American Statistical Association*, **103**, 681-687.
- Sahu, S. K., and Dey, M. D. (2003), A New Class of Multivariate Skew Distributions with Applications to Bayesian Regression Models, *Canadian Journal of Statistics*, **31**, 129-150.
- Spiegelhalter, D., Bradley, P., and Van der Linde, A. (2002), Bayesian Measures of Model Complexity and Fit (with Discussion), *Journal of the Royal Statistical Society, Series B*, **64**, 583-639.
- Shapiro, S. S. and Wilk, M. B. (1965), An Analysis of Variance Test, *Biometrika*, **52**, 591-611.
- Tibshirani, R. (1996), Regression Shrinkage and Selection via the LASSO, *Journal of the Royal Statistical Society, Series B*, **58**, 267-288.
- Zou, H., and Hastie, T. (2005), Regularization and Variable Selection via the Elastic Net, *Journal of the Royal Statistical Society* , **67**, 301-320.

Journal of Statistical Sciences, Spring and Summer, 2021
Vol. 15, No. 1, pp 81-96
DOI: 10.29252/jss.15.1.149

Bayesian LASSO Regression with Asymmetric Error in High Dimensional

Khadem Bashiri, Z., Shadrokh, A., Yarmohammadi, M.
Department of Statistics, Payame Noor University, Tehran, Iran.

Abstract: One of the most critical discussions in regression models is the selection of the optimal model, by identifying critical explanatory variables and negligible variables and more easily express the relationship between the response variable and explanatory variables. Given the limitations of selecting variables in classical methods, such as stepwise selection, it is possible to use penalized regression methods. One of the penalized regression models is the Lasso regression model, in which it is assumed that errors follow a normal distribution. In this paper, we introduce the Bayesian Lasso regression model with an asymmetric distribution error and the high dimensional setting. Then, using the simulation studies and real data analysis, the performance of the proposed model's performance is discussed.

Keywords: Penalized regression, Bayesian LASSO regression, Asymmetric Laplace distribution, Alpha-skew normal distribution.

Mathematics Subject Classification (2010): 62F15, 62G35, 65C10.