

## مدل‌های رگرسیونی برای تحلیل داده‌های چوله دو مدی

حسن میرزاوند<sup>۱</sup>، مجید جعفری خالدی<sup>۲</sup>

تاریخ دریافت: ۱۴۰۰/۱۰/۲۰

تاریخ پذیرش: ۱۴۰۱/۰۵/۲۴

### چکیده:

برای استنباط آماری در مورد پارامترهای مدل رگرسیونی نیاز به فرض توزیع مشخصی بر روی عبارت خطای تصادفی می‌باشد. یک فرض اساسی در مدل رگرسیون خطی این است که عبارت خطای تصادفی از یک توزیع نرمال پیروی کند. با این حال، در پژوهش‌های آماری گاهی با داده‌هایی مواجه می‌شویم که توزیع آن‌ها چولگی و دو مدی را ارائه می‌دهند، و دیگر نمی‌توان از فرض توزیع نرمال برای تحلیل آن‌ها استفاده کرد. یک رویکرد مرسوم برای حل این مسئله به‌کارگیری آمیخته‌ای از مدل‌های چوله نرمال است. اما در این‌گونه مدل‌ها تعداد پارامترها به نحو فزاینده‌ای افزایش می‌یابد که این خود برآزش مدل‌ها به داده‌ها را دشوار می‌نماید. بعلاوه مدل‌های آمیخته خود درگیر مسائلی مانند شناسا ناپذیری هستند. در این حالت یک راه‌حل مناسب استفاده از توزیع‌های منعطفی است، که بتوانند چولگی و دو مدی بودن داده‌ها را در مدل بندی لحاظ کنند. تاکنون روش‌های مختلفی ارائه شده که بر مبنای توسعه توزیع چوله نرمال، توزیع‌های دو مدی نامتقارن ایجاد شده‌اند. در این مقاله از این روش‌ها برای ساخت و معرفی مدل رگرسیونی منعطف نسبت به مدل‌های رگرسیون مبتنی بر توزیع چوله نرمال و آمیخته‌ای از دو توزیع چوله نرمال استفاده شده و با به‌کارگیری مثال شبیه‌سازی عملکرد آن‌ها مورد بررسی قرار می‌گیرد. سپس نحوه کاربست آن‌ها در یک مثال کاربردی مربوط به مجموعه داده‌های اسب‌دوانی نشان داده می‌شود.

واژه‌های کلیدی: چولگی، توزیع‌های دو مدی، تقارن، توزیع‌های آمیخته، رگرسیون.

### ۱ مقدمه

$N(x_i^T \beta, \sigma^2)$  که تابع چگالی آن به صورت

$$f_{Y_i}(y_i | x_i, \theta) = \phi\left(\frac{y_i - x_i^T \beta}{\sigma}\right); \quad y \in \mathbb{R} \quad (2)$$

که در آن  $\phi$  تابع چگالی استاندارد و  $\theta = (\beta^T, \sigma)$  بردار پارامترهای مدل است.

این در حالی است که در بسیاری از کاربردها مشاهدات انحراف زیادی از توزیع نرمال دارند و مدل‌سازی بر پایه فرض نرمال به برآوردهای غیرمنطقی از پارامترهای مدل منجر می‌شود. این خصوصیات را می‌توان با تبدیل داده‌ها مرتفع نمود، که نرمال بودن تقریبی مشاهدات را نتیجه می‌دهد. با این حال برخی از مشکلات این روش به شرح زیر هستند:

• پارامترها ممکن است تفسیرپذیری را در مقیاس تبدیل شده از

نظریه سنتی تحلیل رگرسیون بر پایه نرمال بودن توزیع مشاهدات در سطوح مختلف متغیر کمکی نباشد است، در این صورت فرض می‌شود که عبارت خطای تصادفی مدل از توزیع نرمال پیروی می‌کند.

تعریف ۱.۰۱. یک مدل رگرسیون خطی با باقیمانده‌های نرمال که با نماد  $N-LR$  نمایش داده می‌شود، به صورت

$$Y_i = x_i^T \beta + \epsilon_i, \quad i = 1, \dots, n \quad (1)$$

تعریف می‌شود، که در آن  $Y_i$  متغیر پاسخ،  $x_i = (1, x_{i1}, \dots, x_{ip})^T$  بردار مقادیر متغیر مستقل با بُعد  $(p+1) \times 1$ ،  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$ ، در بردار ضرایب رگرسیونی و خطای تصادفی  $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$  در این صورت با توجه به خواص توزیع نرمال، متغیر پاسخ  $Y_i | x_i \stackrel{iid}{\sim}$

<sup>۱</sup> گروه آمار، دانشکده علوم ریاضی، دانشگاه تربیت مدرس

<sup>۲</sup> عضو هیئت‌علمی گروه آمار، دانشکده علوم ریاضی، دانشگاه تربیت مدرس (نویسنده مسئول: jafari-m@modares.ac.ir)

دست بدهند.

• گاهی وقت‌ها پیدا کردن تبدیل مناسب زمان‌بر است و بعضاً ممکن است تبدیلی که داده‌ها را نرمال نماید، یافت نشود.

از این رو از دیدگاه عملی، نیاز به ارائه یک مدل نظری مناسب است، که از تبدیل داده‌ها جلوگیری کند. جایگزین کردن توزیع نرمال با توزیع‌های مناسب‌تری که قابلیت مدل‌بندی داده‌های نامتقارن را نیز داشته باشند، راهکاری است که در سال‌های اخیر مورد توجه بسیاری از محققین قرار گرفته است. در این راستا، آزالینی [۴] یک نتیجه اساسی را برای توسعه مدل‌های چوله-مقارن ارائه می‌کند.

لم ۲۰۱. فرض کنید  $f_0$  یک تابع چگالی مقارن حول صفر باشد و  $G$  یک تابع توزیع به‌طوری‌که  $G'$  تابع چگالی آن مقارن حول صفر است. در این صورت:

$$f_Z(z | \alpha) = 2f_0(z)G(\alpha z) \quad -\infty < z < \infty \quad (3)$$

برای هر  $\alpha \in \mathbb{R}$  یک تابع چگالی است. از نماد  $Z \sim Sf_0(\alpha)$  برای نشان دادن این واقعیت که  $Z$  دارای چگالی (۳) است، استفاده خواهد شد.

با استفاده از لم ۲۰۱ می‌توان طیف گسترده‌ای از توزیع‌های نامتقارن تولید کرد. که از جمله مشهورترین آن‌ها توزیع چوله‌نرمال است.

تعریف ۳۰۱. اگر  $Z \sim SN(\alpha)$  باشد، تابع چگالی آن به صورت

$$f(z | \alpha) = 2\phi(z)\Phi(\alpha z) \quad z \in \mathbb{R} \quad (4)$$

می‌باشد، که در آن  $\phi$  و  $\Phi$  به ترتیب تابع چگالی و تابع توزیع نرمال استاندارد می‌باشند.  $\alpha$  پارامتر کنترل چولگی است، به ازای  $\alpha = 0$  توزیع مقارن و به چگالی نرمال استاندارد تبدیل می‌شود. لازم به ذکر است تابع چگالی (۴) به ازای مقادیر مختلف پارامتر  $\alpha$  تغییر می‌کند، که نشان‌دهنده شناساپذیری بودن این خانواده از توزیع‌ها است. به بیان دقیق‌تر، به ازای  $\alpha_1, \alpha_2 \in \mathbb{R}$  به طوری‌که  $\alpha_1 \neq \alpha_2$  آنگاه  $f(z | \alpha_1) \neq f(z | \alpha_2)$  می‌باشد. توزیع چوله‌نرمال یک متغیره نخستین بار توسط آزالینی [۴] ارائه شد. هنز [۱۰] یک روش ساخت متغیرهای تصادفی با توزیع چوله‌نرمال معرفی کرد، سپس گشتاورهای آن را به دست آورد. این توزیع علی‌رغم ویژگی‌های خوبی که دارد، در تخمین پارامتر شکل  $\alpha$  با مشکلاتی مواجه است. به‌طور خاص، برای حجم‌های نمونه متوسط، برآوردگر ماکسیمم درست‌نمایی با احتمال مثبت، نامتناهی است ([۱۵] را مشاهده کنید). [۱۶] از یک تابع امتیاز اصلاح‌شده به‌عنوان معادله درست‌نمایی برای برآورد پارامتر شکل استفاده

می‌کند و نشان می‌دهد که برآوردگر ماکسیمم درست‌نمایی اصلاح‌شده همیشه متناهی است.

تعریف ۴۰۱. اگر  $Z \sim SN(\alpha)$  باشد، آنگاه تابع مولد گشتاور آن به صورت

$$M_Z(t) = 2 \exp\left\{\frac{t^2}{\gamma}\right\} \Phi\left(\frac{\alpha}{\sqrt{1+\alpha^2}}t\right), \quad (5)$$

است.

با استفاده از تابع مولد گشتاور به‌سادگی گشتاورهای اول و دوم به صورت

$$E(Z) = \frac{\partial}{\partial t} M_Z(t)|_{t=0} \Rightarrow E(Z) = \sqrt{\frac{2}{\pi}}\lambda,$$

$$E(Z^2) = \frac{\partial^2}{\partial t^2} M_Z(t)|_{t=0} = 1 \Rightarrow Var(Z) = 1 - \frac{2}{\pi}\lambda^2,$$

به دست می‌آیند، که در آن  $\lambda = \frac{\alpha}{\sqrt{1+\alpha^2}}$  است.  $E(Z)$  تابعی صعودی از  $\lambda$  است و  $Var(Z)$  تابعی نزولی از  $|\lambda|$  است. آزالینی [۴] ضرایب چولگی و کشیدگی متغیر تصادفی چوله‌نرمال را به ترتیب به صورت

$$\sqrt{\beta_1} = \frac{1}{\gamma}(\pi - \pi)(\frac{E^2(Z)}{Var(Z)})^{\frac{1}{2}},$$

$$\beta_2 = 2(\pi - 3)(\frac{E^2(Z)}{Var(Z)})^2,$$

محاسبه نمود. بر این اساس آزالینی [۴] نشان داد که مقدار ضریب چولگی و ضریب کشیدگی در بازه  $(-0.995, 0.995)$  و  $(-0.869, 0.869)$  تغییر می‌کند. اضافه کردن پارامترهای مکان و مقیاس به توزیع چوله‌نرمال، انعطاف‌پذیری آن را بیشتر می‌کند و با تغییر مقدار پارامترها می‌توان کنترل بیشتری روی توزیع داشت. تبدیل  $X = \mu + \sigma Z$ ، با  $\mu \in \mathbb{R}$  و  $\sigma > 0$  را در نظر می‌گیریم. متغیر تصادفی  $X$  دارای توزیع چوله‌نرمال با تابع چگالی

$$f_X(x) = \frac{2}{\sigma} \phi\left(\frac{x-\mu}{\sigma}\right) \Phi\left(\alpha \frac{x-\mu}{\sigma}\right), \quad (6)$$

در ادامه مدل رگرسیون مبتنی بر این توزیع  $SN$  معرفی می‌شود.

تعریف ۵۰۱. یک مدل رگرسیون خطی با باقیمانده‌های چوله‌نرمال که با نماد  $SN-LR$  نمایش داده می‌شود، به صورت

$$Y_i = x_i^T \beta + \epsilon_i, \quad i = 1, \dots, n \quad (7)$$

تعریف می‌شود، که در آن  $Y_i$  متغیر پاسخ،  $x_i = (1, x_{i1}, \dots, x_{ip})^T$  بردار مقادیر متغیر مستقل با بعد  $(p+1) \times 1$ ،  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  بردار ضرایب رگرسیونی و خطای تصادفی  $\epsilon_i \stackrel{iid}{\sim} SN(0, \sigma, \alpha)$ .

با استفاده از مطالعات شبیه‌سازی و همچنین تحلیل مجموعه داده‌های اسب‌دوانی (*horses*) نشان داده می‌شود. همچنین در بخش ۶، به بحث و نتیجه‌گیری پرداخته خواهد شد.

این صورت با توجه به خواص توزیع  $SN$ ، متغیر پاسخ  $Y_i|x_i \stackrel{iid}{\sim}$  تابع چگالی آن به صورت

$$f_{Y_i}(y_i|x_i, \theta) = \frac{2}{\sigma} \phi\left(\frac{y_i - x_i^\top \beta}{\sigma}\right) \Phi\left(\alpha \frac{y_i - x_i^\top \beta}{\sigma}\right); \quad y \in \mathbb{R} \quad (۸)$$

که در آن  $\phi$  تابع چگالی استاندارد و  $\theta = (\beta^\top, \sigma, \alpha)$  بردار پارامترهای مدل است. امید ریاضی و واریانس  $Y_i$  به صورت

$$E(Y_i|x_i) = x_i^\top \beta + c\alpha$$

$$Var(Y_i|x_i) = \sigma + (1 - c^2)\alpha^2$$

که در آن  $c = \sqrt{2/\pi}$  می‌باشد.

## ۲ مدل رگرسیونی مبتنی بر آمیخته‌ای از دو توزیع چوله‌نرمال

یکی از راهکارهایی که به منظور مدل‌سازی چولگی و دو مدی می‌توان استفاده نمود به کارگیری مدل‌های آمیخته متناهی است، که در ادامه به معرفی آن‌ها پرداخته شده است. توزیع‌های آمیخته متناهی، یعنی توزیع‌هایی که به صورت مجموع موزون از چند توزیع با مؤلفه‌های ساده‌تر نوشته می‌شود، بسیار مورد توجه قرار دارند. معمولاً توزیع مؤلفه‌ها با توجه به نوع داده‌های مورد بررسی و نیز به دلیل جذابیت‌های محاسباتی به کلاس‌های خاصی از تابع چگالی‌های پارامتری محدود می‌شوند. به‌طور کلی مدل آمیخته متناهی بر اساس آمیخته موزونی از تابع چگالی‌های دلخواه به صورت

$$f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta}) = \sum_{k=1}^K \eta_k f_k(\mathbf{y}|\boldsymbol{\theta}_k), \quad (۹)$$

تعریف می‌شود که در آن  $\mathbf{y} = (y_1, \dots, y_n)'$  بردار مشاهدات،  $K$  تعداد مؤلفه‌ها و  $\eta_k$ ‌ها وزن‌های آمیزنده<sup>۳</sup> مؤلفه‌ها هستند، به طوری که

$$\forall k: \eta_k \geq 0, \quad \sum_k \eta_k = 1.$$

در رابطه (۹)،  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_K)'$  بردار وزن‌های آمیزنده و  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1', \dots, \boldsymbol{\theta}_K')$  بردار پارامترهای نامعلوم تابع چگالی‌ها را نمایش می‌دهند. همچنین  $f_k(\mathbf{y}|\boldsymbol{\theta}_k)$  تابع چگالی مؤلفه  $k$ ام با پارامترهای  $\boldsymbol{\theta}_k$  است. در رابطه (۹) تعداد مؤلفه‌ها ثابت در نظر گرفته شده است، اما در عمل تعداد مؤلفه‌ها نامعلوم است و بایستی به همراه وزن‌های آمیزنده و پارامترهای موجود در تابع چگالی مؤلفه‌ها، از روی داده‌های موجود تعیین و برآورد شوند. برای درک بهتر از مدل آمیخته متناهی، آمیخته‌ای از توزیع‌های نرمال در شکل ۲ نمایش داده شده است.

این مدل رگرسیونی توسط افراد مختلفی مورد مطالعه قرار گرفته است، که از جمله می‌توان به مطالعات آرلانو-واله و همکاران [۳]، لاجوس و همکاران [۱۱] و کانچو و همکاران [۶] اشاره کرد، که از توزیع چوله‌نرمال برای برازش مدل رگرسیونی استفاده کرده‌اند. اهمیت توزیع چوله‌نرمال و مبنا قرار گرفتن آن در بسیاری از مباحث استنباط آماری و مدل‌سازی به این دلیل است که این توزیع علی‌رغم داشتن قابلیت مدل‌بندی مشاهدات نامتقارن، شباهت زیادی به توزیع نرمال داشته و آن را نه به‌عنوان یک توزیع حدی، بلکه به‌عنوان یک عضو در بردارد. توجه کنید که مدل رگرسیون مبتنی بر توزیع چوله‌نرمال تک مدی است، اما در عمل داده‌هایی وجود دارند که نیازمند برازش با توزیع‌های نامتقارن دومی هستند، در این صورت این مدل از کارایی لازم برخوردار نیست. به‌عنوان مثال داده‌های حداقل دمای کشور در بهمن‌ماه سال ۱۳۹۶ جمع‌آوری شده در ۳۸۱ ایستگاه هواشناسی را در نظر بگیرید. شکل ۱ موقعیت جغرافیایی ایستگاه‌ها و هیستوگرام حداقل دمای هوا را به نمایش می‌گذارد. هیستوگرام نشان می‌دهد توزیع داده‌ها چوله به راست بوده و دو مدی است. لازم به ذکر است که می‌توان از متغیر ارتفاع ایستگاه‌ها از سطح دریا به‌عنوان متغیر کمکی در تحلیل رگرسیونی این داده‌ها استفاده نمود. اکنون سؤالی که مطرح می‌شود این است که چگونه می‌توان چولگی و دومی موجود در داده‌ها را مدل‌بندی نمود؟ ادامه این مقاله به صورت زیر سازمان‌دهی شده است. در بخش ۲، مدل رگرسیون مبتنی بر آمیخته‌ای از دو توزیع چوله‌نرمال مطالعه می‌گردد. در ادامه در بخش ۳، توزیع‌های نامتقارن دو مدی معرفی می‌گردد، سپس برخی از خواص و ویژگی‌های آن‌ها مورد مطالعه قرار می‌گیرند. در بخش ۴، مدل رگرسیون مبتنی بر توسعه‌های دو مدی توزیع چوله‌نرمال معرفی می‌شود. در بخش ۵، نحوه کاربست مدل‌های رگرسیونی توسعه‌یافته،

<sup>3</sup> Mixing Weights

$Mix.SN(\circ, \circ, \sigma_1, \sigma_2, \alpha_1, \alpha_2, \eta_1, \eta_2)$  نمایش داده می‌شود، در این صورت متغیر پاسخ دارای توزیع آمیخته دو مؤلفه‌ای از توزیع چوله‌نرمال به صورت

$$Y_i|x_i \sim \sum_{j=1}^2 \eta_j f_{SN}(y_i|\theta_j(x_i)) \quad i = 1 \dots n \quad (11)$$

است، که با نماد  $Mix.SN - LR$  نمایش داده می‌شود. در رابطه (۱۱) بردار  $(\mu_j(x_i), \sigma_j^2, \alpha_j)^T$  و  $\theta_j(x_i) = (\mu_j(x_i), \sigma_j^2, \alpha_j)^T$  بردار ضرایب  $\beta_j = (\beta_{j0}, \beta_{j1}, \dots, \beta_{jp})^T$  و  $(1, x_{i1}, \dots, x_{ip})^T$  رگرسیونی زیر جامعه  $Z$  و  $\eta = (\eta_1, \eta_2)$  بردار ضرایب آمیختگی را نشان می‌دهد.

اما همان‌طور که بیان شد به دلیل بُعد زیاد پارامترها و مشکلات شناسانپذیری مبتلابه این‌گونه مدل‌ها استفاده از آن‌ها چالش‌برانگیز است. مک‌لاچان و پیل [۱۴] و مارین و همکاران [۱۳] مشکلات جدی را در مورد برآورد چنین مدل‌هایی ارائه می‌دهند. حال سؤالی که در اینجا مطرح می‌شود این است که چگونه مدل رگرسیون منعطفی معرفی کنیم که هم چولگی و دو مدی بودن را کنترل کند و هم از مشکلات مبتلابه مدل‌های آمیخته رنج نبرد؟

یک‌راه حل دیگر برای تحلیل رگرسیونی زمانی که چولگی و دو مدی به‌طور هم‌زمان وجود دارد، استفاده از توسعه‌های دو مدی توزیع چوله‌نرمال است که در بخش ۳ معرفی می‌شوند.

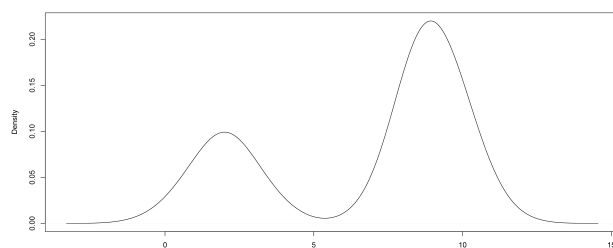
### ۳ معرفی توزیع‌های نامتقارن دو مدی

در این بخش، توسعه‌هایی از توزیع چوله‌نرمال به حالت دو مدی، بررسی می‌شود. از جمله خانواده جدیدی از توزیع‌های چوله‌نرمال انعطاف‌پذیر بانام آلفا - چوله‌نرمال که توسط الال-آلیویرو [۷] معرفی شد، که داده‌های تک مدی و دو مدی را پوشش می‌دهد. به صورت زیر معرفی کردند.

تعریف ۱۰۳. اگر متغیر تصادفی  $Z$  دارای تابع چگالی

$$f(z; \alpha) = \frac{1 + (1 - \alpha z)^2}{2 + \alpha^2} \phi(z) \quad z \in \mathbb{R} \quad (12)$$

که در آن  $\alpha \in \mathbb{R}$ ، در این صورت  $Z$  را متغیر تصادفی آلفا-چوله‌نرمال با پارامتر  $\alpha$  نامیده و به صورت  $Z \sim ASN(\alpha)$  نشان می‌دهند. در مدل  $ASN(\alpha)$  چون تنها از یک پارامتر برای کنترل چولگی و مد توزیع استفاده می‌شود، چندان کارآمد به نظر نمی‌رسد.



شکل ۲. نمودار تابع چگالی آمیخته از دو توزیع نرمال با وزن‌های ۰.۳ و ۰.۷.

یکی از پرکاربردترین مدل‌های آمیخته، برای زمانی که گروه‌های متفاوت از داده‌ها دارای چولگی باشند، مدل آمیخته‌ای از چوله‌نرمال‌ها است. که جایگزین مناسب‌تری نسبت به مدل آمیخته‌ای از نرمال‌ها است. آرلانو-واله و همکاران [۲] یک مدل چوله-مقارن معرفی می‌کنند، که آمیخته‌ای از توزیع‌های چوله‌نرمال است.

تعریف ۱۰۲. مدل آمیخته متناهی چوله‌نرمال دو مؤلفه‌ای، آمیخته‌ای از ۲ توزیع چوله‌نرمال یک متغیره به صورت (۱۰)

$$f(y|\theta_1, \theta_2, \eta) = \eta f_{SN}(y; \mu_1, \sigma_1^2, \alpha_1) + (1 - \eta) f_{SN}(y; \mu_2, \sigma_2^2, \alpha_2)$$

است که در آن

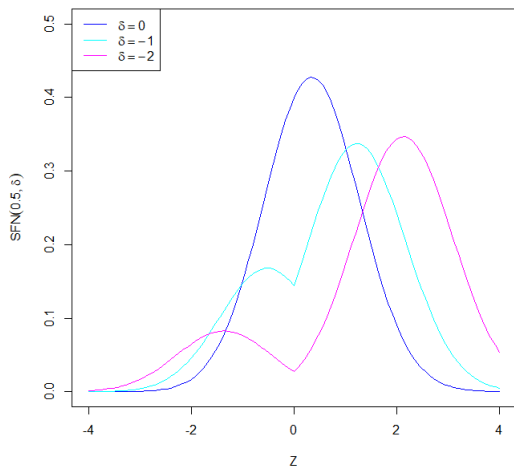
$$f_{SN}(y; \mu_k, \sigma_k^2, \alpha_k) = \frac{1}{\sigma_k} \phi\left(\frac{y - \mu_k}{\sigma_k}\right) \Phi\left(\alpha_k \left(\frac{y - \mu_k}{\sigma_k}\right)\right) \quad k = 1, 2$$

تابع چگالی گروه  $k$ ام با پارامترهای  $\theta_k = (\mu_k, \sigma_k, \alpha_k)$  که در آن پارامتر مکان،  $\sigma_k$  پارامتر مقیاس و  $\alpha_k$  پارامتر چولگی مربوط به گروه  $k$ ام هستند، این توزیع با نماد  $Y \sim Mix.SN(\mu_1, \mu_2, \sigma_1, \sigma_2, \alpha_1, \alpha_2, \eta)$  نمایش داده می‌شود. برای مطالعه بیشتر پیرامون آمیخته‌ای از توزیع‌های چوله‌نرمال و چوله‌تی به لین و همکاران [۱۲] و بهرامی [۱] مراجعه کنید.

از طرفی، هنگامی که جامعه ساختاری ناهمگن دارد و از زیرجوامعی چوله تشکیل شده است، (به‌عنوان مثال، داده‌های درآمد خانوارهای روستایی ایران تحلیل شده توسط بهرامی [۱])، می‌توان آمیخته‌ای متناهی از توزیع‌های چوله‌نرمال را به منظور مدل‌سازی خطای تصادفی استفاده نمود. در این حالت، می‌توان از مدل رگرسیونی مبتنی بر آمیخته‌ای از چوله‌نرمال‌ها برای خطای تصادفی به‌عنوان جایگزین مدل سنتی استفاده کرد.

تعریف ۲۰۲. فرض کنید، خطای تصادفی مدل دارای توزیع آمیخته‌ای از دو توزیع چوله‌نرمال باشد که با نماد  $\epsilon_i \stackrel{iid}{\sim}$

با نماد  $Z \sim SFN(\alpha, \delta)$  نمایش داده می‌شود. لازم به ذکر است تابع چگالی (۱۴) به ازای مقادیر مختلف پارامترهای  $\alpha$  و  $\delta$  تغییر می‌کند، که نشان‌دهنده شناساپذیر بودن این توزیع است (گومز و همکاران [۹]).



شکل ۱۰.۳. نمایش توزیع  $SFN(\alpha, \delta)$  به ازای مقادیر مختلف پارامترهای  $\alpha$  و  $\delta$ .

شکل ۱۰.۳ نمایش گرافیکی توزیع  $SFN(\alpha, \delta)$ ، برای مقادیر مختلف پارامترها را نشان می‌دهد. با توجه به این شکل، پارامتر  $\alpha$  پارامتر چولگی و  $\delta$  پارامتر دو مدی ساز به نحوی است که به ازای مقادیر کمتر از صفر دو مدی بودن را برای توزیع به ارمغان می‌آورد. خواص زیر مستقیماً از تعریف ۳.۳ به دست می‌آیند.

گزاره ۴.۳. فرض کنید متغیر تصادفی  $Z$  دارای توزیع  $SFN(\alpha, \delta)$  باشد، در این صورت

• اگر  $\alpha = 0$  و  $\delta = 0$  آنگاه  $Z \sim N(0, 1)$

• اگر  $\delta = 0$  آنگاه  $Z \sim SN(\alpha)$

• اگر  $\alpha = 0$  آنگاه  $f(z|\alpha = 0, \delta) = c_\delta \phi(|z| + \delta)/2$

• اگر  $\alpha \rightarrow +\infty$  آنگاه  $f(z|\alpha, \delta) = c_\delta \phi(z + \delta)I(z \geq 0)$

• اگر  $\alpha \rightarrow -\infty$  آنگاه  $f(z|\alpha, \delta) = c_\delta \phi(z - \delta)I(z < 0)$

گزاره ۵.۳. فرض کنید  $Z \sim SFN(\alpha, \delta)$  است. اگر  $\delta < 0$ ، در این صورت  $Z$  یک متغیر تصادفی با تابع چگالی دو مدی می‌باشد.

نتیجه ۶.۳. فرض کنید  $Z \sim SFN(\alpha, \delta)$  است. اگر  $\alpha = 0$  و  $\delta < 0$ ، آنگاه  $Z$  یک متغیر تصادفی با تابع چگالی دو مدی متقارن است.

تعریف ۷.۳. فرض کنید  $Z \sim SFN(\alpha, \delta)$  با  $\alpha, \delta \in \mathbb{R}$ . خانواده توزیع‌های چوله‌نرمال-انعطاف‌پذیر با پارامترهای مکان و مقیاس

رده توزیع‌های معرفی‌شده توسط الال-الیویرو [۷]، منجر به تعمیم و بسط سایر توزیع‌ها توسط محققان دیگری شد، که در ادامه به آن‌ها پرداخته می‌شود. در ساخت و معرفی توسعه‌های دو مدی توزیع چوله‌نرمال، فرع ۲.۳ نقش اساسی ایفا می‌کند.

نتیجه ۲.۳. (آزالینی و کاپتانو [۵]) با فرض آنکه  $f_0, G', W$  چنان باشد که درلم ۲.۱ بیان شد و  $X \sim G', Y \sim f_0$  متغیرهای تصادفی مستقل باشند، آنگاه متغیر تصادفی  $Z$ ، که به صورت

$$Z = \begin{cases} Y & X \leq W(Y) \\ -Y & X > W(Y) \end{cases}$$

تعریف می‌شود، دارای تابع چگالی (۳) است.

• اگر در فرع ۲.۳، قرار دهیم  $W(Y) = \alpha Y$ ، که در آن  $\alpha$  یک عدد حقیقی است و  $X \sim N(0, 1)$  و همچنین  $Y$  دارای توزیع نرمال بریده‌شده  $Y \sim c_\delta \phi(y + \delta)I(y \geq 0)$  باشد، در این صورت متغیر تصادفی  $Z$  دارای توزیع چوله‌نرمال دو مدی معرفی‌شده توسط گومز و همکاران [۹] است.

• اگر در فرع ۲.۳، قرار دهیم  $W(Y) = \alpha Y$ ، که در آن  $\alpha$  یک عدد حقیقی است و  $X \sim N(0, 1)$  و  $Y \sim BN(\delta)$  باشد، در این صورت متغیر تصادفی  $Z$  دارای توزیع چوله‌نرمال دو مدی معرفی‌شده توسط الال-الیویرو و همکاران [۸] است.

در ادامه به تفصیل این دو توسعه توزیع چوله‌نرمال مورد بحث و بررسی قرار می‌گیرد.

## ۱۰.۳ توزیع چوله‌نرمال انعطاف‌پذیر ( $SFN$ )

تعریف ۳.۳. (گومز و همکاران [۹]) متغیر تصادفی  $Z$  دارای توزیع چوله‌نرمال-انعطاف‌پذیر با پارامترهای  $\alpha$  و  $\delta$  نامیده می‌شود، هرگاه تابع چگالی آن به صورت

$$f(z|\alpha, \delta) = c_\delta \phi(|z| + \delta)\Phi(\alpha z); \quad z \in \mathbb{R} \quad (13)$$

یا

$$f(z|\alpha, \delta) = \begin{cases} c_\delta \phi(z - \delta)\Phi(\alpha z) & z \leq 0 \\ c_\delta \phi(z + \delta)\Phi(\alpha z) & z > 0 \end{cases} \quad (14)$$

باشد، که در آن  $\phi$  و  $\Phi$  به ترتیب نشان‌دهنده تابع چگالی و تابع توزیع نرمال استاندارد،  $\alpha, \delta \in \mathbb{R}$  و  $c_\delta^{-1} = (1 - \Phi(\delta))$  است. این توزیع

تعریف ۱۰.۳. اگر  $X$  یک متغیر تصادفی با تابع چگالی

$$f(x|\delta) = \left(\frac{1+\delta x^\lambda}{1+\delta}\right)\phi(x); \quad x \in \mathbb{R}, \delta \geq 0 \quad (17)$$

باشد، در این صورت  $X$  دارای توزیع نرمال دو مدی با پارامتر شکل  $\delta$  است، که با نماد  $X \sim BN(\delta)$  نشان داده می‌شود.

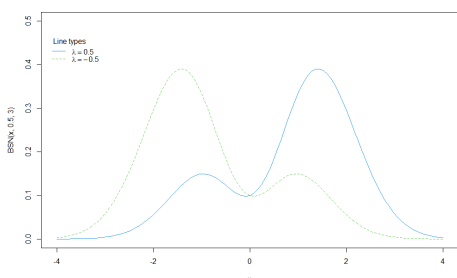
گزاره ۱۱.۳. فرض کنید  $X$  و  $Y$  متغیرهای تصادفی مستقل به‌گونه‌ای که  $X \sim BN$  و  $Y \sim N(0, 1)$ ، اگر  $W = \sqrt{\frac{\delta}{1+\delta}}X + \sqrt{\frac{1}{1+\delta}}Y$  آنگاه  $W \sim BN(\delta)$  است.

تعریف ۱۲.۳. اگر متغیر تصادفی  $X$  دارای تابع چگالی

$$f(x|\alpha, \delta) = 2\left(\frac{1+\delta x^\lambda}{1+\delta}\right)\phi(x)\Phi(\alpha x); \quad x \in \mathbb{R}, \alpha \in \mathbb{R}, \delta \geq 0 \quad (18)$$

آنگاه  $X$  دارای توزیع چوله نرمال دو مدی با پارامترهای  $\alpha$  و  $\delta$ ، که با نماد  $X \sim BSN(\alpha, \delta)$  نشان داده می‌شود. شکل ۲.۳ تابع چگالی معادله (۱۸) را به ازای مقادیر مختلف پارامترهای  $\alpha$  و  $\delta$  نشان می‌دهد. با توجه به این شکل  $\alpha$  پارامتر چولگی و  $\delta$  پارامتر کنترل‌کننده مد توزیع است، که به ازای مقادیر  $\delta = 0$  توزیع تک مدی است.

لازم به ذکر است تابع چگالی (۱۸) به ازای مقادیر مختلف پارامترهای  $\alpha$  و  $\delta$  تغییر می‌کند، که نشان‌دهنده شناساپذیر بودن این توزیع است (گومز و همکاران [۹]).



شکل ۲.۳. تابع چگالی توزیع  $BSN(\alpha, \delta)$  برای  $\delta = 3$  و مقادیر مختلف  $\alpha$ .

گزاره ۱۳.۳. اگر  $X \sim BN(\alpha, \delta)$ ، آنگاه خواص زیر برقرار می‌باشند.

- اگر  $\alpha = 0$  آنگاه  $X \sim BN(\delta)$
- اگر  $\alpha = 0$  و  $\delta = 0$  آنگاه  $X \sim N(0, 1)$
- اگر  $\delta = 0$  آنگاه  $X \sim SN(\alpha)$
- اگر  $\delta \rightarrow \infty$  آنگاه  $X \sim BN(\alpha)$
- خاصیت ناوردایی، اگر  $W \sim BSN(\alpha, \delta)$  و  $Z \sim BN(\delta)$  آنگاه  $|X| \stackrel{d}{=} |Z|$  هم توزیع می‌باشند.

به صورت توزیع تبدیل  $X = \mu + \sigma Z$  تعریف می‌شود، که در آن  $\mu \in \mathbb{R}$  و  $\sigma > 0$ . تابع چگالی  $X$  به صورت

$$f(x|\mu, \sigma, \alpha, \delta) = \frac{c_\delta}{\sigma} \phi\left(\frac{|x-\mu|}{\sigma} + \delta\right) \Phi\left(\alpha \frac{x-\mu}{\sigma}\right) \quad (15)$$

است. این توزیع به صورت  $X \sim SFN(\theta)$  نمایش داده می‌شود که در آن  $\theta = (\mu, \sigma, \alpha, \delta)$ .

گشتاور اول تا چهارم برای  $Z \sim SFN(\alpha, \delta)$  عبارت‌اند از:

$$\begin{aligned} E(Z) &= -c_\delta[\delta\Phi(-\delta) - 2\lambda\phi(\lambda\delta)\Phi(-\delta\lambda/\alpha)] + 2\delta\mu_0(-\delta, \alpha, \delta) \\ E(Z^\lambda) &= c_\delta[(1+\delta^\lambda)\Phi(-\delta) - \delta\phi(\delta)] \\ E(Z^3) &= -c_\delta[\delta(3+\delta)\Phi(-\delta) - 2\lambda\phi(\lambda\delta)\{[2+(3+\lambda^4-3\lambda^2)\delta^\lambda + \lambda^2\alpha^\lambda]\Phi(-\delta\lambda/\alpha) + [(\lambda^2\delta/\alpha)(2+1/\alpha^\lambda) - 3\lambda\delta/\alpha]\phi(\delta\lambda/\alpha)\}] \\ &\quad + 2\delta(3+\delta^\lambda)\mu_0(-\delta, \alpha, \delta) \\ E(Z^4) &= c_\delta[(3+6\delta^2+\delta^4)\Phi(-\delta) - \delta(5+\delta^2)\phi(\delta)] \end{aligned}$$

که در آن  $\lambda = \frac{\alpha}{\sqrt{1+\alpha^2}}$  است و  $\mu_0(-\delta, \alpha, \delta)$  باید به صورت عددی محاسبه شود (گومز و همکاران [۹]).

## ۲.۳ توزیع چوله نرمال دو مدی ( $BSN$ )

در این زیر بخش به بررسی توسعه دیگری از توزیع چوله نرمال پرداخته می‌شود که بر اساس فرع ۲.۳ ساخته شده است. این توزیع توسط الال-الیویرو و همکاران [۸] معرفی شده است. ابتدا توزیع نرمال دو مدی معرفی می‌شود.

تعریف ۸.۳. اگر  $X$  یک متغیر تصادفی با تابع چگالی

$$f(x) = x^\lambda \phi(x); \quad x \in \mathbb{R} \quad (16)$$

باشد، که در آن  $\phi$  تابع چگالی نرمال استاندارد، در این صورت  $X$  دارای توزیع نرمال دو مدی نامیده می‌شود، که با نماد  $X \sim BN$  نشان داده می‌شود (الال-الیویرو [۷]).

تذکر ۹.۳. اگر  $X \sim BN$  با تابع چگالی  $f(x)$ ، آنگاه  $f(x)$  دو مدی است. این موضوع را می‌توان با توجه به اینکه  $f'(x) = x\phi(x)(2-x^\lambda)$  نتیجه گرفت، که کمترین مقدار تابع  $f(x)$  در  $x = 0$  و بیشترین مقدار در  $x = \sqrt{2}$  و  $x = -\sqrt{2}$  است. توجه داشته باشید که  $f(\sqrt{2}) = f(-\sqrt{2})$  است، این واقعیت با اضافه کردن یک پارامتر اضافی که ارتفاع را در مدها کنترل می‌کند، منجر به شکل‌گیری یک مدل انعطاف‌پذیر می‌شود که در ادامه مطرح خواهد شد.

معرفی می‌شود، که در آن  $Y_i$  متغیر پاسخ،  $x_i = (1, x_{i1}, \dots, x_{ip})^T$  بردار مقادیر متغیر مستقل با بعد  $(p+1) \times 1$ ،  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  بردار ضرایب رگرسیونی و خطای تصادفی  $\epsilon_i \stackrel{iid}{\sim} SFN(0, \sigma, \alpha, \delta)$  در این صورت با توجه به خواص توزیع  $SFN$ ، متغیر پاسخ  $Y_i|x_i \stackrel{iid}{\sim} SFN(x_i^T \beta, \sigma, \alpha, \delta)$  که تابع چگالی آن به صورت

$$f_{Y_i}(y_i|x_i, \theta) = \frac{c_\delta}{\sigma} \phi\left(\frac{|y_i - x_i^T \beta|}{\sigma} + \delta\right) \Phi\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right) \quad (21)$$

که در آن  $\phi$  و  $\Phi$  به ترتیب تابع چگالی و تابع توزیع نرمال استاندارد و  $\theta = (\beta^T, \sigma, \alpha, \delta)$  بردار پارامترهای مدل است.

لگاریتم تابع درستنمایی با توجه به نمونه مشاهده شده  $y_1, \dots, y_n$  به صورت

$$l(\theta) = \sum_{i=1}^n l_i(\theta) = C - \frac{n}{2} \left( \log\left(\frac{\sigma}{c_\delta}\right) + \delta^2 \right) - \frac{\delta}{\sigma} \sum_{i=1}^n |y_i - x_i^T \beta| - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \sum_{i=1}^n \log \Phi\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right) \quad (22)$$

است. با مشتق گرفتن از  $l(\theta)$  نسبت به بردار پارامترهای مدل خواهیم داشت:

$$\frac{\partial l(\theta)}{\partial \beta_j} = \sum_{i=1}^n x_{ij} \left( \frac{y_i - x_i^T \beta}{\sigma^2} \right) + \frac{\delta}{\sigma} \sum_{i=1}^n x_{ij} \text{sign}(y_i - x_i^T \beta) - \frac{\alpha}{\sigma} \sum_{i=1}^n x_{ij} \eta\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right) = 0$$

$$\frac{\partial l(\theta)}{\partial \sigma} = -\frac{n}{2\sigma} + \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - x_i^T \beta)^2 - \frac{\alpha}{\sigma^2} \sum_{i=1}^n \eta\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right) (y_i - x_i^T \beta) + \frac{\delta}{\sigma^2} \sum_{i=1}^n |y_i - x_i^T \beta| = 0$$

$$\frac{\partial l(\theta)}{\partial \alpha} = \frac{1}{\sigma} \sum_{i=1}^n \eta\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right) (y_i - x_i^T \beta) = 0$$

$$\frac{\partial l(\theta)}{\partial \delta} = (c_\delta \phi(\delta) - \delta) - \frac{1}{\sigma} \sum_{i=1}^n |y_i - x_i^T \beta| = 0$$

که در آن  $\eta(a) = \phi(a)/\Phi(a)$ . بنابراین برآوردگر  $ML$  از حل معادلات بالا به دست می‌آید، که می‌بایست روش‌های عددی مورد استفاده قرار گیرد.

تعریف ۱۴.۳. فرض کنید  $Z \sim BSN(\alpha, \delta)$  با  $\alpha, \delta \in \mathbb{R}$ . خانواده توزیع‌های چوله‌نرمال دو مدی با پارامترهای مکان و مقیاس به صورت  $X = \mu + \sigma Z$  تعریف می‌شود، که در آن  $\mu \in \mathbb{R}$  و  $\sigma > 0$ . تابع چگالی  $X$  به صورت

$$f_X(x; \mu, \sigma, \alpha, \delta) = \frac{\sigma^2 + \delta(x - \mu)^2}{\sigma^2(1 + \delta)} \phi\left(\frac{x - \mu}{\sigma}\right) \times \Phi\left(\alpha \frac{x - \mu}{\sigma}\right), \quad x \in \mathbb{R} \quad (19)$$

است. این توزیع به صورت  $X \sim BSN(\theta)$  نمایش داده می‌شود که در آن  $\theta = (\mu, \sigma, \alpha, \delta)$ .

تعریف ۱۵.۳. فرض کنید  $Z \sim BSN(\alpha, \delta)$  آنگاه گشتاورهای اول تا چهارم برای این توزیع عبارت‌اند از:

$$E(X) = \frac{b\lambda}{(1 + \delta)(1 + \alpha^2)} [(1 + 2\alpha)\delta^2 + 3\alpha + 1]$$

$$E(X^2) = \frac{1 + 3\alpha}{1 + \alpha}$$

$$E(X^3) = \frac{b\lambda}{(1 + \delta)(1 + \alpha^2)^2} [(2 + 8\delta)\alpha^3 + (\delta + 2\delta^2)\alpha^2 + 15\delta + 3]$$

$$E(X^4) = \frac{3 + 15\delta}{1 + \delta}$$

که در آن  $\lambda = \sqrt{\frac{2}{\pi}}\alpha$  و  $b = 1/\sqrt{(1 + \alpha^2)}$  است.

در این بخش توسعه‌های توزیع چوله‌نرمال به حالت دو مدی مورد مطالعه قرار گرفت. در ادامه نحوه به‌کارگیری این توزیع‌ها در ساخت مدل‌های رگرسیونی منعطف ارائه می‌شود.

## ۴ مدل‌های رگرسیونی منعطف برای داده‌های نامتقارن دومی

در این بخش مدل‌های رگرسیون مبتنی بر توزیع‌های  $SFN$  و  $BSN$  به‌عنوان یک فعالیت پژوهشی جدید معرفی می‌شوند.

### ۱.۴ مدل رگرسیون مبتنی بر توزیع $SFN$

تعریف ۱.۴. یک مدل رگرسیون خطی با باقیمانده‌های  $SFN$  که با نماد  $SFN - LR$  نمایش داده می‌شود، به صورت

$$Y_i = x_i^T \beta + \epsilon_i, \quad i = 1, \dots, n \quad (20)$$

## ۲.۴ مدل رگرسیون مبتنی بر توزیع BSN

تعریف ۲.۴. یک مدل رگرسیون خطی با باقیمانده‌های BSN که با نماد  $BSN - LR$  نمایش داده می‌شود، به صورت

$$Y_i = x_i^T \beta + \epsilon_i, \quad i = 1, \dots, n \quad (23)$$

معرفی می‌شود، که در آن  $Y_i$  متغیر پاسخ،  $x_i = (1, x_{i1}, \dots, x_{ip})^T$  بردار مقادیر متغیر مستقل با بعد  $(p+1) \times 1$ ،  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  بردار ضرایب رگرسیونی و خطای تصادفی  $\epsilon_i \stackrel{iid}{\sim} BSN(0, \sigma, \alpha, \delta)$ .

در این صورت با توجه به خواص توزیع BSN، متغیر پاسخ  $Y_i | x_i \stackrel{iid}{\sim} BSN(x_i^T \beta, \sigma, \alpha, \delta)$  که چگالی آن به صورت

$$f_{Y_i}(y_i | x_i, \theta) = \frac{\sigma^\nu + \delta(y_i - x_i^T \beta)^\nu}{\sigma^\nu(1 + \delta)} \phi\left(\frac{y_i - x_i^T \beta}{\sigma}\right) \times \Phi\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right); \quad y \in \mathbb{R} \quad (24)$$

که در آن  $\phi$  و  $\Phi$  به ترتیب تابع چگالی و تابع توزیع نرمال استاندارد و  $\theta = (\beta^T, \sigma, \delta, \alpha)$  بردار پارامترهای مدل است.

لگاریتم تابع درستنمایی با توجه به نمونه مشاهده شده  $y_1, \dots, y_n$  به صورت

$$l(\theta) = \sum_{i=1}^n l_i(\theta) = -\nu n \log(\sigma) - n \log(1 + \delta) - \frac{1}{\nu \sigma^\nu} \sum_{i=1}^n (y_i - x_i^T \beta)^\nu + \sum_{i=1}^n \log(\sigma^\nu + \delta(y_i - x_i^T \beta)^\nu) + \sum_{i=1}^n \log \Phi\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right)$$

است. با مشتق گرفتن از  $l(\theta)$  نسبت به بردار پارامترهای مدل خواهیم داشت:

$$\begin{aligned} \frac{\partial l(\theta)}{\partial \beta_j} &= \sum_{i=1}^n x_{ij} \left( \frac{y_i - x_i^T \beta}{\sigma^\nu} \right) - \nu \delta x_{ij} \sum_{i=1}^n \frac{(y_i - x_i^T \beta)}{\sigma^\nu + \delta(y_i - x_i^T \beta)^\nu} - \frac{1}{\sigma^\nu} \sum_{i=1}^n x_{ij} \eta\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right) = 0 \\ \frac{\partial l(\theta)}{\partial \sigma} &= -\frac{\nu n}{\sigma} + \frac{1}{\sigma^\nu} \sum_{i=1}^n (y_i - x_i^T \beta)^\nu + \nu \sum_{i=1}^n \frac{1}{\sigma^\nu + \delta(y_i - x_i^T \beta)^\nu} - \frac{1}{\sigma^\nu} \sum_{i=1}^n (y_i - x_i^T \beta) \eta\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right) = 0 \\ \frac{\partial l(\theta)}{\partial \delta} &= -\frac{n}{1 + \delta} + \sum_{i=1}^n \frac{(y_i - x_i^T \beta)^\nu}{\sigma^\nu + \delta(y_i - x_i^T \beta)^\nu} = 0 \\ \frac{\partial l(\theta)}{\partial \alpha} &= \frac{1}{\sigma} \sum_{i=1}^n (y_i - x_i^T \beta) \eta\left(\alpha \frac{y_i - x_i^T \beta}{\sigma}\right) = 0 \end{aligned}$$

که در آن  $\eta(a) = \phi(a)/\Phi(a)$ . بنابراین برآوردگر  $ML$  از حل معادلات بالا به روش عددی به دست می‌آید.

در این راستا از الگوریتم نیلدر-مید در تابع  $optim$  موجود در نرم افزار  $R$  استفاده می‌شود. این الگوریتم یکی از شناخته شده ترین روش های بدون مشتق است، که فقط از مقادیر تابع برای جستجوی مینیمم استفاده می‌کند. این رویکرد شامل ساخت یک سیمپلکس از  $n + 1$  نقطه و حرکت این سیمپلکس در جهت صحیح است. برای آشنایی بیشتر با این الگوریتم به این سایت مراجعه کنید.

## ۵ کاربرد و ارزیابی مدل‌های رگرسیونی منعطف

در این بخش با استفاده از شبیه سازی و مثال کاربردی مربوط به مجموعه داده های اسب دوانی (*horses*)، نحوه کاربست مدل های رگرسیونی معرفی شده در این مقاله مورد ارزیابی قرار می‌گیرد.

### ۱.۵ مطالعه شبیه سازی

در این زیر بخش، شبیه سازی مختصری از مدل های رگرسیونی  $BSN - LR$  و  $SFN - LR$  برای بررسی رفتار برآوردگرها در نمونه های متناهی ارائه می‌شود. برای محاسبه برآوردگرهای  $MLE$  از الگوریتم  $Nelder - Mead$  در تابع  $optim$  موجود در نرم افزار  $R$  استفاده می‌شود. شبیه سازی از مدل رگرسیونی

$$Y_i = 2 + 5X_{i1} + 3X_{i2} + \epsilon_i \quad i = 1, \dots, n$$

صورت می‌گیرد، که در آن ابتدا خطای تصادفی دارای توزیع  $\epsilon_i \sim BSN(0, 4, 0.5, -3)$  و بار دیگر  $\epsilon_i \sim SFN(0, 4, 0.5, -3)$  است. حجم نمونه ها ۱۰۰، ۵۰۰ و ۱۰۰۰ می‌باشد. همان طور که از جداول ۱ و ۲ مشاهده می‌شود، بر اساس ویژگی های مجانبی برآوردگرهای  $MLE$  با افزایش حجم نمونه مقدار اریبی برآوردها کمتر می‌شود، مقدار خطای برآوردگرها و همچنین  $AIC$  مدل نیز کاهش می‌یابد. در ادامه این زیر بخش، با استفاده از شبیه سازی، مدل های  $SFN - LR$  و  $BSN - LR$  با  $Mix.SN - LR$  مورد مقایسه قرار می‌گیرند. برای محاسبه برآوردگرهای  $MLE$  مدل های  $SFN - LR$  و  $BSN - LR$  از تابع  $optim$  و همچنین برای برازش مدل رگرسیون آمیخته ای از دو توزیع چوله نرمال از کتابخانه  $FMsmsnReg$  موجود در نرم افزار  $R$  استفاده می‌شود. ابتدا، خطای تصادفی مدل با استفاده از آمیخته ای از

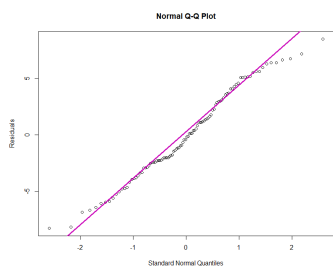


اسب‌های مسابقه (Starters) و نسبت بردها در شروع قبلی (Ratio) هستند. بنابراین، مدل رگرسیونی

$$Position_i = \beta_0 + \beta_1 Starters_i + \beta_2 Ratio_i + \epsilon_i$$

در نظر گرفته می‌شود. ابتدا، مدل رگرسیون خطی با دو متغیر مستقل به روش OLS با دستور lm به داده‌ها برازش داده می‌شود، در ادامه برای بررسی فرض نرمال بودن، باقیمانده‌های مدل محاسبه می‌شود. حال برای بررسی این فرض به صورت شهودی از نمودار چندک-چندک (Q-Q plot)، و در ادامه برای تحلیل دقیق‌تر از آزمون کولموگروف اسمیرنوف بر روی باقیمانده‌های مدل استفاده می‌شود.

نمودار چندک-چندک یا Q-Q plot به منظور مقایسه دو توزیع به کار گرفته می‌شود. از چنین نمودارهایی حتی می‌توان مطابقت توزیع داده‌ها را با یک توزیع مشخص، مورد بررسی قرار داد. نمودار چندک-چندک بر اساس چندک‌های دو توزیع ترسیم می‌شود. در محور افقی چندک‌های توزیع اول و در محور عمودی نیز چندک‌های متناظر برای توزیع دوم مشخص می‌شود. اگر این مقادیر را در یک صفحه مختصات دکارتی ترسیم کنیم، یک نمودار چندک چندک یا Q-Q plot ایجاد کرده‌ایم.



شکل ۲۰۵. نمودار چندک-چندک مربوط به باقیمانده‌های مدل OLS.

جدول ۰۴. نتایج آزمون کولموگروف-اسمیرنوف بر روی باقیمانده‌های

مدل OLS.

آزمون	p-value
کولموگروف-اسمیرنوف	$3.9 \times 10^{-12}$

با توجه به شکل ۲۰۵ مشاهده می‌شود، که باقیمانده‌ها به طور کامل روی خط قرار نمی‌گیرند، و از طرفی پی-مقدار مربوط به آزمون کولموگروف-اسمیرنوف از سطح معنی‌داری ۰/۰۵ کمتر است (با توجه به جدول ۲۰۵)، پس می‌توان نتیجه گرفت که فرض نرمال بودن باقیمانده‌های مدل برقرار نمی‌باشد. در ادامه هیستوگرام باقیمانده‌ها رسم می‌شود و توزیع‌های معرفی شده در این مقاله، بر روی آن برازش شده است. با توجه به شکل ۲۰۵، هیستوگرام باقیمانده‌ها تک مدی و متقارن نمی‌باشد، بنابراین می‌توان از توزیع‌های SFN، Mix.SN

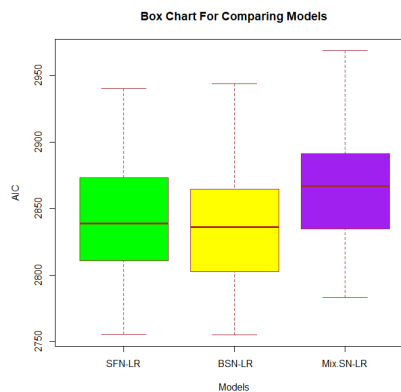
دو توزیع چوله‌نرمال  $\epsilon \sim Mix.SN(0, 0, 1, 1, 0.7, 0.7, 0.3)$  نمونه‌ای تصادفی به حجم ۱۰۰۰ با تکرار ۵۰ تولید می‌شود، سپس مدل رگرسیونی

$$Y_i = 4 + 2X_{1i} - 2X_{2i} + \epsilon_i \quad i = 1, \dots, n$$

شبه‌سازی می‌شود، در ادامه مدل‌های رگرسیونی معرفی شده یعنی SFN-LR و BSN-LR با Mix.SN-LR به داده‌ها برازش می‌شود و مقادیر AIC، Bias( $\hat{\beta}$ ) و  $MSE(\hat{\beta})$  به صورت

$$Bias(\hat{\beta}) = (\hat{\beta}) - \beta_{True}, \quad MSE(\hat{\beta}) = \frac{1}{50} \sum_{i=1}^{50} (\hat{\beta} - \beta_{True})^2$$

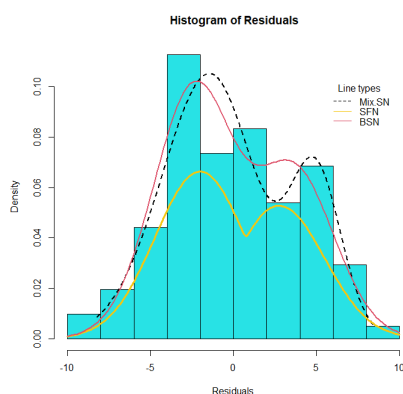
محاسبه می‌شود. نتایج در جدول ۳ ارائه شده است. در انتها نمودار جعبه‌ای برای مقادیر AIC تحت مدل‌های رگرسیونی مختلف رسم می‌شود (شکل ۱۰۵ را ملاحظه کنید). با توجه به این شکل و جدول ۳ مشاهده می‌شود باینکه مجموعه داده‌ها از مدل Mix.SN-LR تولید شده است، عملکرد مدل BSN-LR نسبت به سایر مدل‌ها مناسب‌تر است. این امر می‌تواند به دلیل انعطاف‌پذیری بیشتر و تعداد پارامترهای کمتر مدل‌های چوله‌نرمال دو مدی نسبت به مدل‌های آمیخته‌ای از دو توزیع چوله‌نرمال باشد.



شکل ۱۰۵. نمودار جعبه‌ای مقادیر AIC تحت مدل‌های مختلف برای ۵۰ تکرار شبه‌سازی.

## ۲۰۵ مثال کاربردی

در این زیر بخش به عنوان یک مثال انگیزشی، مجموعه داده‌ای را که قبلاً توسط فوریز (۱۹۹۸) به وسیله یک مدل رگرسیون نرمال تجزیه و تحلیل شده است، در نظر گرفته می‌شود. مجموعه داده‌های اسب‌دوانی (horses)، شامل نتایج مربوط به هر اسب در ۸ مسابقه متوالی می‌باشد. داده‌ها شامل ۱۴ متغیر که هر کدام ۱۰۲ بار مشاهده شده است. در اینجا، متغیر پاسخ موقعیت پایان (Position) و متغیرهای کمکی شامل تعداد



شکل ۳.۵. هیستوگرام باقیمانده‌ها همراه با برازش توابع چگالی توزیع‌های معرفی شده.

و *BSN* برای مدل‌سازی این داده‌ها استفاده کرد. در جدول ۴ خلاصه برازش مدل‌های رگرسیونی به داده‌ها و برآورد پارامترهای مدل آورده شده است. به‌عنوان یک شاخص سنجش و انتخاب مدل مناسب، از معیار ارزیابی *AIC* استفاده شده است. شاخص *AIC* توسط دانشمند ژاپنی آمار، «هیروتاگا آکایکه» (*Hirotsugu Akaike*) در سال‌های ۱۹۷۰ برای تشخیص مدل مناسب از بین مدل‌های موجود، معرفی شد و امروز به‌عنوان یک ابزار مهم در انتخاب مدل‌ها بر اساس تابع درستنمایی به کار گرفته می‌شود.

فرض کنید در یک مدل آماری  $k$  تعداد پارامترهای مدل باشد. اگر  $\hat{L}$  را ماکسیمم تابع درستنمایی تحت مدل مفروض در نظر بگیریم، معیار ارزیابی *AIC* توسط رابطه

$$AIC = 2k - 2 \log \hat{L}$$

## ۶ بحث و نتیجه‌گیری

در این مقاله، به بررسی رویکردهای مختلف در تحلیل رگرسیونی برای حالتی که داده‌ها به‌طور هم‌زمان چولگی و دو مدی بودن را نشان می‌دهند، پرداخته شد. به‌ویژه مدل‌های رگرسیونی مبتنی بر توزیع‌های چوله‌نرمال دو مدی معرفی گردید، که این‌گونه مدل‌ها می‌توانند به‌عنوان جایگزینی برای مدل‌های آمیخته به کار روند. با استفاده از مطالعات شبیه‌سازی و مثال کاربردی، عملکرد مناسب‌تر این‌گونه مدل‌ها نشان داده شد. با این‌حال در بسیاری از مسائل کاربردی ممکن است با داده‌هایی که ساختار همبستگی فضایی دارند مواجه شویم، در نتیجه فرض استقلال نیز برای این داده‌ها برقرار نمی‌باشد. در این حالت توسعه مدل‌های پیشنهادی به حالت فضایی به‌عنوان کار آینده مورد توجه نویسنده مقاله قرار دارد.

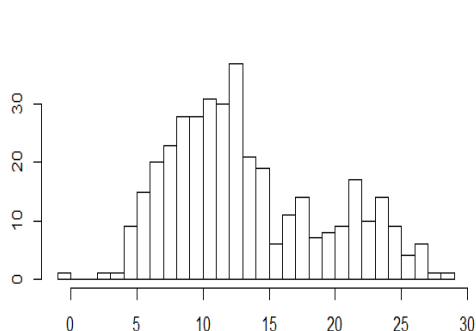
قابل‌محاسبه است. بنابراین مناسب‌ترین مدل برحسب معیار اطلاع آکایکه، دارای کمترین مقدار *AIC* است. با توجه به اینکه مقدار *AIC* برای مدل رگرسیونی مبتنی بر توزیع *BSN* نسبت به مدل‌های دیگر کمتر است، در نتیجه این مدل برای برازش به داده‌ها عملکرد مناسب‌تری دارد.

جدول ۱: نتایج شبیه‌سازی از مدل رگرسیونی *BSN - LR*.

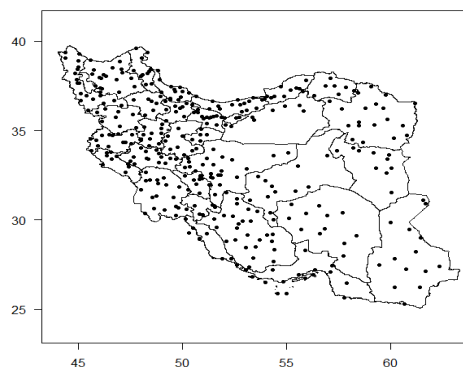
$n = 1000$				$n = 500$				$n = 100$				
فاصله اطمینان ۹۵٪		خطای برآورد	برآورد <i>MLE</i>	فاصله اطمینان ۹۵٪		خطای برآورد	برآورد <i>MLE</i>	فاصله اطمینان ۹۵٪		خطای برآورد	برآورد <i>MLE</i>	پارامترها
بالا	پایین			بالا	پایین			بالا	پایین			
۲.۵۵۹۹	۱.۸۰۹۹	۰.۱۶۵۸	۲.۲۳۳۹	۲.۸۵۷۴	۱.۶۳۵۱	۰.۳۱۱۸	۲.۲۴۶۲	۳.۰۵۵۹	۰.۳۲۸۳	۰.۷۶۹۸	۱.۵۲۷۱	$\beta_1$
۵.۲۲۵۶	۲.۷۷۶۸	۰.۱۱۲۵	۵.۰۰۱۲	۵.۳۷۳۴	۴.۶۷۱۶	۰.۱۷۹۰	۵.۰۲۲۵	۵.۴۶۶۴	۳.۸۵۸۴	۰.۴۱۰۲	۴.۶۶۲۴	$\beta_2$
۳.۲۸۲۸	۲.۸۳۷۵	۰.۱۱۴۱	۳.۰۶۱۱	۳.۴۴۶۱	۲.۶۳۷۰	۰.۱۸۰۸	۲.۸۹۱۵	۳.۷۲۴۱	۲.۱۶۶۲	۰.۴۰۲۵	۲.۸۵۵۲	$\beta_3$
۴.۰۶۸۷	۳.۷۴۶۹	۰.۰۸۲۰	۳.۹۰۷۸	۴.۴۶۱۹	۳.۸۰۰۷	۰.۱۳۳۱	۴.۱۸۱۳	۴.۶۷۹۲	۳.۳۵۸۷	۰.۳۳۶۸	۴.۰۱۹۰	$\sigma$
۳.۵۴۲۷	۲.۲۵۳۵	۰.۱۷۱۴	۳.۰۱۲۳	۳.۸۷۹۸	۱.۵۵۰۲	۰.۵۹۴۳	۲.۷۱۵۰	۳.۱۱۲۳	۰.۱۶۷۱	۰.۵۵۱۲	۱.۶۴۹۷	$\delta$
-۰.۴۲۰۷	-۰.۵۸۰۴	۰.۰۳۰۷	-۰.۵۰۰۵	-۰.۴۴۲۱	-۰.۷۱۲۷	۰.۰۶۹۰	-۰.۵۷۷۴	-۰.۱۰۰۱	-۰.۷۲۶۵	۰.۱۵۹۸	-۰.۲۱۲۳	$\alpha$
-۰.۵۹۲۵۸۳۹				-۰.۲۹۹۰۱۷۴				-۰.۵۹۹۳۸۷۲				<i>AIC</i>

## مراجع

[۱] بهرامی، م. (۱۳۹۴)، مدل‌های آمیخته متناهی مبتنی بر توزیع‌های چوله متقارن، پایان نامه دوره ارشد، دانشگاه تربیت مدرس.



(ب)



(الف)

شکل ۱: (الف) موقعیت مکانی ۲۸۱ ایستگاه هواشناسی در ایران و (ب) هیستوگرام متغیر حداقل دمای هوا.

جدول ۲: نتایج شبیه‌سازی از مدل رگرسیونی  $SFN - LR$ .

$n = 1000$				$n = 500$				$n = 100$				
فاصله اطمینان ۹۵%		فاصله اطمینان ۹۵%		فاصله اطمینان ۹۵%		فاصله اطمینان ۹۵%		فاصله اطمینان ۹۵%		فاصله اطمینان ۹۵%		
بالا	پایین	خطای برآورد	برآورد MLE	بالا	پایین	خطای برآورد	برآورد MLE	بالا	پایین	خطای برآورد	برآورد MLE	
۲,۳۷۸۱	۱,۵۱۷۴	-۰,۲۱۹۵	۱,۹۴۷۸	۲,۵۳۶۶	۱,۳۹۶۰	-۰,۲۹۰۹	۱,۸۶۶۳	۲,۷۳۱۴	-۰,۱۸۱۶	-۰,۷۳۳۱	۱,۷۷۵۰	
۵,۳۵۰۷	۲,۸۵۶۱	-۰,۱۲۶۱	۵,۱۰۲۴	۵,۳۰۳۸	۳,۶۴۹۴	-۰,۱۶۶۶	۴,۹۷۶۱	۶,۱۱۴۶	۴,۵۷۹۵	-۰,۳۹۱۶	۵,۳۳۷۰	
۳,۳۱۴۰	۲,۸۱۸۷	-۰,۱۱۸۱	۳,۰۸۰۳	۳,۲۴۰۱	۲,۷۱۹۳	-۰,۱۷۸۷	۳,۰۶۹۷	۳,۲۰۴۶	۱,۵۰۲۴	-۰,۳۳۴۲	۲,۳۵۳۵	
۴,۲۴۸۶	۲,۸۶۳۰	-۰,۰۹۸۳	۴,۰۵۵۸	۴,۱۵۳۲	۳,۶۳۳۸	-۰,۱۲۲۵	۳,۸۹۳۵	۴,۸۷۱۱	۳,۵۸۸۷	-۰,۳۲۷۱	۴,۲۲۹۹	
-۰,۵۵۴۰	-۰,۲۵۲۰	-۰,۰۲۶۰	-۰,۵۰۳۰	-۰,۵۳۳۸	-۰,۳۹۹۶	-۰,۰۳۳۹	-۰,۶۶۶۲	-۰,۷۴۱۳	-۰,۳۷۷۹	-۰,۰۹۲۷	-۰,۵۵۹۶	
-۲,۸۲۲۰	-۳,۱۵۸۷	-۰,۰۸۰۷	-۲,۰۰۰۴	-۲,۸۸۰۶	-۲,۳۳۳۰	-۰,۱۱۵۴	-۲,۱۰۶۸	-۲,۴۰۱۷	-۲,۳۹۸۱	-۰,۱۵۳۱	-۲,۸۹۹۹	
-۶۱۳۸-۸۳				-۳۰۴۱۸۸۹				-۶۰۶۰۱۹۷				AIC

جدول ۳: خلاصه برازش مدل‌های رگرسیونی به داده‌های شبیه‌سازی شده.

Mir.SN - LR			SFN - LR			BSN - LR			مقدار واقعی پارامتر	پارامترها
MSE( $\hat{\beta}$ )	Bias( $\hat{\beta}$ )	MLE	MSE( $\hat{\beta}$ )	Bias( $\hat{\beta}$ )	MLE	MSE( $\hat{\beta}$ )	Bias( $\hat{\beta}$ )	MLE		
۰,۰۲۵۰	-۰,۱۵۶۲	۴,۱۵۶۲	۰,۳۷۹۵	-۰,۲۶۲۲	۳,۷۳۷۷	۱,۷۱۳۸	-۱,۲۹۴۸	۲,۷۰۵۱	۴	$\beta_1$
۰,۰۰۱۰	-۰,۰۰۰۲	۱,۹۹۹۸	۰,۰۰۱۳	۰,۰۰۱۷	۲,۰۰۱۷	۰,۰۰۱۰	۰,۰۰۲۹	۲,۰۰۲۹	۲	$\beta_2$
۰,۰۰۰۷	۰,۰۰۲۲	-۱,۹۹۷۷	۰,۰۰۱۴	۰,۰۰۵۵	-۱,۹۹۴۴	۰,۰۰۰۷	۰,۰۰۲۳	-۱,۹۹۷۶	-۲	$\beta_3$
۲۸۶۷,۳۶۹			۲۸۴۳,۸۸۱			۲۸۳۷,۴۵۹			AIC	

[2] Arellano-Valle, R. B., Gómez, H. W., and Quintana, F. A. (2004). A New Class of Skew-Normal Distributions. *Communications in Statistics-Theory and Methods*, **33**(7), 1465-1480.

[3] Arellano-Valle, R. B., Castro, L. M., Genton, M. G., and Gómez, H. W. (2008). Bayesian Inference for Shape Mixtures of Skewed Distributions, with Application to Regression Analysis. *Bayesian Analysis*, **3**(3), 513-539.

[4] Azzalini, A. (1985). A Class of Distributions Which Includes the Normal Ones. *Scandinavian journal of statistics*, **12**, 171-178.

[5] Azzalini, A., and Capitanio, A. (1999). Statistical Applications of the Multivariate Skew Normal Distribution. *Journal of the Royal Statistical Society, Series B*, **61**(3), 579-602.

جدول ۴: نتایج برازش مدل‌های رگرسیونی در نظر گرفته شده به داده‌ها.

<i>BSN - LR</i>		<i>SFN - LR</i>		<i>Mix.SN - LR</i>		پارامترها
خطای برآورد	MLE برآورد	خطای برآورد	MLE برآورد	خطای برآورد	MLE	
۲٫۰۷	۵٫۶۰	۱٫۳۱	۰٫۹۶	۰٫۰۰۹	۲٫۰۴	$\beta_0$
۰٫۸۱	۰٫۳۰	۰٫۰۹	۰٫۴۸	۰٫۱۵	۰٫۵۰۴	$\beta_1$
۳٫۲۲	-۶٫۱۰	۲٫۹۱	-۰٫۵۲	۰٫۲۳	-۸٫۴۳	$\beta_2$
-	-	-	-	۳٫۹۳	۳٫۷۰	$\mu_1$
-	-	-	-	۲٫۱۷	-۲٫۶۳	$\mu_2$
۰٫۲۳	۲٫۸۲	۰٫۳۸	۲٫۸۲	۱٫۸۰	۶٫۵۲	$\sigma_1$
-	-	-	-	۳٫۳۴	۹٫۵۱	$\sigma_2$
۰٫۱۹	-۰٫۲۶	۰٫۱۲	-۰٫۰۲	۶٫۰۷	۰٫۷۴	$\alpha_1$
-	-	-	-	۱٫۱۲	-۱٫۳۲	$\alpha_2$
۰٫۷۲	۱٫۴۳	۰٫۳۰	-۰٫۸۵	-	-	$\delta$
-	-	-	-	۱٫۲۰	۰٫۴۱	$\eta$
۵۵۰٫۷۰۷۹		۵۵۴٫۴۰۱۶		۵۷۸٫۲۰۵		<i>AIC</i>

- [6] Cancho, V. G., Dey, D. K., Lachos, V. H., and Andrade, M. G. (2011). Bayesian Nonlinear Regression Models with Scale Mixtures of Skew-Normal Distributions: Estimation and case Influence Diagnostics. *Computational Statistics and Data Analysis*, **55**(1), 588-602.
- [7] Elal-Olivero, D. (2010), Alpha-Skew-Normal Distribution. *Proyecciones Journal of Mathematics*. **29**(3), 224–240., 11, 12.
- [8] Elal-Olivero, D., Olivares-Pacheco, J. F., Venegas, O., Bolfarine, H., and Gomez, H. W. (2020). On Properties of the Bimodal Skew-Normal Distribution and an Application. *Mathematics*, **8**(5), 703.
- [9] Gómez, H. W., Elal-Olivero, D., Salinas, H. S., and Bolfarine, H. (2011). Bimodal Extension Based on the Skew-Normal Distribution with Application to Pollen Data. *Environmetrics*, **22**(1), 50-62.
- [10] Henze, N. (1986). A Probabilistic Representation of the Skew-Normal Distribution, *Scandinavian Journal of Statistics*. **13**, 271-275.
- [11] Lachos, V., Bolfarine, H., Arellano-Valle, R. and Montenegro, L. (2007). Likelihood Based Inference for Multivariate Skew-Normal Regression Models. *Communications in Statistics-Theory and Methods*, **36**, 1769–1786.
- [12] Lin, T. I., Lee, J. C., and Yen, S. Y. (2007). Finite Mixture Modelling Using the Skew-Normal Distribution. *Statistica Sinica*, **17**: 81–92.
- [13] Marin, J. M., Mengersen, K., and Robert, C. P. (2005). Bayesian Modelling and Inference on Mixtures of Distributions. *Handbook of Statistics*, **25**, 459-507.
- [14] McLachan, G. and Peel, D. (2000). *Mixture Models: Inference and Applications to Clustering*, Marcel Dekker, N.Y.
- [15] Pewsey, A. (2003), The Characteristic Functions of the Skew-Normal and Wrapped Skew-Normal Distributions, In 27 Congreso Nacional de Estadística e Investigación Operativa, pp. 4383-4386.
- [16] Sartori, N. 2006. Bias Prevention of Maximum Likelihood Estimates for Scalar Skew-Normal and Skew-t Distributions. *J. Statist. Plann. Inference*, **136**, 4259–4275.

## پیوست:

تولید نمونه از توزیع *BSN*:

```

1 BSN=function(n,alpha,lambda){
2 x1=rchisq(n,3)
3 x2=rnorm(n)
4 u=c();w=c();t=c();s=c()
5 for (i in 1:n) {
6 u[i]=sample(c(-1,1),size = 1,prob = c(1/2,1/2))
7 w[i]=sqrt(x1[i])*u[i]
8 t[i]=sqrt(alpha/(1+lambda))*w[i]+sqrt(1/(1+lambda))*x2[i]
9 s[i]=sample(c(t[i],-t[i]),size = 1,prob = c(pnorm(alpha*t[i]),1-pnorm(alpha*t[i])))
10 }
11 return(s)
12 }

```

تولید نمونه از توزیع *SFN*:

```

1 library(TruncatedNormal)
2 SFN=function(n,alpha,delta){
3 y1=c();w1=c();s1=c();z=c()
4 for(i in 1:n){
5 y1[i]=rtnorm(1,-delta,lb=0,ub=Inf,sd=1)
6 w1[i]=pnorm(alpha*y1[i])
7 s1[i]=sample(c(1,-1),size =1,replace = TRUE,prob = c(w1[i],1-w1[i]))
8 z[i]=s1[i]*y1[i]
9 }
10 return(z)
11 }

```

برازش مدل رگرسیون به روش *OLS* و بررسی نرمال بودن باقیمانده‌های مدل:

```

1 library(FMsmnReg)
2 library(e1071)
3 data("horses")
4 head(horses)
5 str(horses)
6 attach(horses)
7 model.ols=lm(Position~Starters+Ratio)
8 skewness(model.ols$res)

```

```

9 kurtosis(model.ols$res)
10 sd(model.ols$res)
11 mean(model.ols$res)
12 #####
13 hist(model.ols$res,breaks = 10,col=13,main= "Residuals of ols method",xlab = "Ordinary
    residuals",freq = FALSE)
14 ks.test(model.ols$res,"pnorm")
15 qqnorm(model.ols$res,ylab="Residuals",xlab = "Standard Normal Quantiles")
16 qqline(model.ols$res,col=6,lwd=3)

```

برازش مدل‌های رگرسیونی معرفی شده در بخش ۴:

```

1 ##### BSN-LR
2 regression.BSN=function(y,x1,x2,beta){
3 b0=beta[1]
4 b1=beta[2]
5 b2=beta[3]
6 alpha=beta[4]
7 lambda=beta[5]
8 sigma=beta[6]
9 lhod=2*(((sigma^2)+lambda*(y-b0-b1*x1-b2*x2)^2)/(sigma^3*(1+lambda)))*dnorm((y-b0-b1*x1-b2
    *x2)/sigma)*pnorm(alpha*(y-b0-b1*x1-b2*x2)/sigma)
10 return(-sum(log(lhod)))
11 }
12 MLE.BSN=optim(c(3,2,-1,1,-1,2),regression.BSN,y=Position,x1=Starters,x2=Ratio,hessian =
    TRUE)
13 MLE.BS=MLE.BSN$par
14 ObsInfo=MLE.BSN$hess #Observed Fisher information matrix
15 Vhat=solve(ObsInfo) #Inverse of observed Fisher information
16 Std.Errors=sqrt(diag(Vhat))
17 #Obtain the MLEs,estimated std errors, and approx Wald 95% Cis
18 Wald.table.BSN=cbind(MLE.BS,Std.Errors,LowerBound=MLE.BS-qnorm(0.975)*Std.Errors,
19 UpperBound=MLE.BS+qnorm(0.975)*Std.Errors)
20 Wald.table.BSN
21 row.names(Wald.table.BSN )=c("beta0","beta1","beta2","alpha","delta","sigma")
22 Wald.table.BSN
23 K=6
24 AIC=2*K-2* MLE.BSN$value

```

```

25 natijeh.BSN=list(Wald.table.BSN,AIC=AIC)
26 ##### SFN-LR
27 regression.SFN=function(beta,y,x1,x2){
28 b0=beta[1]
29 b1=beta[2]; b2=beta[3]; alpha=beta[4];delt=beta[5];sigma=beta[6];cdeta=1/(1-pnorm(delt))
30 lhd= (cdeta/sigma)*dnorm((abs(y-b0-b1*x1-b2*x2)/sigma)+delt)*pnorm(alpha*(y-b0-b1*x1-b2*x2
    )/sigma)
31 return(-sum(log(lhd)))
32 }
33 MLE.SFN=optim(c(3,2,-1,1,-1,2),regression.SFN,y=Position,x1=Starters,x2=Ratio,hessian =
    TRUE)
34 MLE.SF=MLE.SFN$par
35 ObsInfo=MLE.SFN$hess #Observed Fisher information matrix
36 Vhat=solve(ObsInfo) #Inverse of observed Fisher information
37 Std.Errors=sqrt(diag(Vhat))
38 #Obtain the MLEs,estimated std errors, and approx Wald 95% Cis
39 Wald.table.SFN=cbind(MLE.SF,Std.Errors,LowerBound=MLE.SF-qnorm(0.975)*Std.Errors,
40 UpperBound=MLE.SF+qnorm(0.975)*Std.Errors)
41 Wald.table.SFN
42 row.names(Wald.table.SFN)=c("beta0","beta1","beta2","alpha","delta","sigma")
43 Wald.table.SFN
44 K=6
45 AIC=2*K-2*MLE.SFN$value
46 natijeh.SFN=list(Wald.table.SFN,AIC=AIC)
47 ##### Mix.sn.LR
48 data("horses")
49 head(horses)
50 attach(horses)
51 FMsmnReg()
52 x=cbind(1,Starters,Ratio)
53 parCN <- FMsmnReg( Position,x, g=2, get.init = TRUE, criteria = TRUE, group = FALSE,shape
    =c(0,0),
54 family = "Skew.normal", error = 10^-4,obs.prob= FALSE)

```

## Regression Models for Analyzing of Bimodal and Skewed Data

Hassan Mirzavand<sup>1</sup>, and Majid Jafari Khaledi<sup>2</sup>

### Abstract:

For statistical inference about the parameters of the regression model, it is necessary to assume a specific distribution for the random error term. A basic assumption in the linear regression model is that the random error term follows a normal distribution. However, in statistical research, sometimes the distribution of the data display both skewness and bimodality, and in such situations, it is inappropriate to use the normal distribution for statistical analysis. A conventional approach to overcome this problem is to use a mixture of normal models. But in such models, the number of parameters increases substantially, which makes it difficult to fit these models to the data. In addition, the mixed models suffer from the non-identifiability issues. In this case, a suitable solution is to use flexible distributions which can simultaneously handle the skewness and bimodality of the data in the modeling structure. So far, various methods have been proposed, which were created based on the development of the skew-normal distribution. In this article, these asymmetric bimodal distributions are used to build and introduce a flexible regression model compared to the regression models based on the normal distribution as well as a mixture of two normal distributions. Their performance is evaluated using a simulation example. Then, the usefulness of the method is demonstrated through a practical example related to the horse data set.

**Keywords:** Skewness, Bimodal distributions, Symmetric, Mixture of distributions, Regression .

---

<sup>1</sup>Tarbiat Modares University, Tehran

<sup>2</sup>Tarbiat Modares University, Tehran.