

تحلیل بیزی رگرسیون چوله‌نرمال آمیخته

افشین فلاح، رامین کاظمی، حسن خسروی

گروه آمار، دانشگاه بین‌المللی امام خمینی

تاریخ دریافت: ۱۳۹۵/۸/۱۳ تاریخ آخرین بازنگری: ۱۳۹۶/۱۱/۲۱

چکیده: تحلیل رگرسیونی به‌طور سنتی با فرض همگن بودن جامعه و نرمال بودن توزیع متغیر پاسخ صورت می‌پذیرد. این در حالی است که در بسیاری از کاربردها، به‌دلیل ناهمگنی مشاهدات، وجود نقاط دور افتاده، چولگی یا ترکیبی از آن‌ها، مشاهدات ساختاری ناهمگن با زیرجوامعی چوله-متقارن را نشان می‌دهند. در چنین حالاتی، می‌توان آمیخته‌ای متناهی از توزیع‌های چوله-متقارن را برای مدل‌بندی جامعه مورد استفاده قرار داد. در این مقاله رهیافت بیزی تحلیل رگرسیونی تحت فرض ناهمگن بودن جامعه و چوله-متقارن بودن توزیع زیرجوامع، با استفاده از آمیخته‌ای متناهی از توزیع‌های چوله‌نرمال مورد توجه قرار گرفته است. به منظور ارزیابی رهیافت پیشنهادی و مقایسه آن با مدل فراوانی‌گرا، از یک مطالعه شبیه‌سازی و یک مثال کاربردی استفاده شده است.

واژه‌های کلیدی: تحلیل رگرسیون بیزی، توزیع آمیخته متناهی، چولگی، الگوریتم EM، الگوریتم گیبز.

۱ مقدمه

نظریه سنتی تحلیل رگرسیونی بر پایه همگن بودن جامعه و نرمال بودن توزیع مشاهدات بنا شده 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. 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 زمینه، می‌توان به مطالعات

آدرس الکترونیکی نویسنده مسئول مقاله: افشین فلاح، a.fallah@sci.ikiu.ac.ir

کد موضوع‌بندی ریاضی (۲۰۱۰): 62J02، 62F15

آرلانواله و همکاران (۲۰۰۸)، لاجوس و همکاران (۲۰۱۰)، کانچو و همکاران (۲۰۱۰a) 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. از طرفی، هنگامی که جامعه ساختاری ناهمگن دارد و از زیرجوامعی چوله-متقارن تشکیل شده است، می‌توان آمیخته‌ای متناهی از توزیع‌های چوله-متقارن را مدنظر قرار داد. توزیع‌های آمیخته متناهی، به دلیل انعطاف‌پذیری قابل توجه آن‌ها، در ادبیات آماری به دفعات در زمینه‌های مختلف مورد استفاده قرار گرفته‌اند و اغلب قادرند توصیف مناسبی از جوامع پیچیده و ناهمگن فراهم سازند. تحلیل رگرسیونی تحت توزیع‌های آمیخته متناهی، اولین بار توسط کوانت (۱۹۷۲) و کوانت و رامسی (۱۹۷۸) مورد مطالعه قرار گرفت. با معرفی الگوریتم EM توسط دمپستر و همکاران (۱۹۷۷)، محققانی بسیاری از جمله اِتکین و ویلسون (۱۹۸۰)، دسارو و کورن (۱۹۸۸)، تورنر (۲۰۰۰)، لاجوس و همکاران (۲۰۰۱، ۲۰۱۰)، کانچو و همکاران (۲۰۱۰b) و لیو و همکاران (۲۰۱۱) تحلیل رگرسیونی چوله‌نرمال را از دیدگاه فراوانی‌گرایانه مورد مطالعه قرار دادند. با توجه به دشواری‌های ذاتی برآورد پارامترها در توزیع‌های آمیخته متناهی، و این که خواص بهینه برآوردگرهای ماکسیمم درست‌نمایی تنها در حالت مجانبی تحقق می‌یابند، رهیافت فراوانی‌گرایانه تحلیل رگرسیونی تحت توزیع‌های آمیخته متناهی از کارایی لازم برخوردار نیست. نشان داده شده است تابع درست‌نمایی مبتنی بر توزیع چوله‌نرمال، به خصوص در حالت آمیخته بودن توزیع، نسبت به پارامترهای چولگی کران‌دار نیست. افزون بر این، در زمینه کاربست الگوریتم EM به‌عنوان عمومی‌ترین راه برای محاسبه برآوردگرهای ماکسیمم درست‌نمایی در توزیع‌های آمیخته متناهی، ملاحظات و دشواری‌هایی وجود دارد، که از آن جمله می‌توان به عدم همگرایی الگوریتم EM در برآورد پارامترهای چولگی اشاره کرد. این مشکلات توسط برناردی و همکاران (۲۰۱۲) مورد توجه قرار گرفته است. با توجه به مشکلات رهیافت فراوانی‌گرا، در این مقاله رهیافت بیزی تحلیل رگرسیونی آمیخته متناهی چوله‌نرمال مورد مطالعه قرار گرفته است. برای این منظور، با در نظر گرفتن توزیع‌های پیشین مناسب، توزیع‌های پسین شرطی کامل پارامترهای مدل محاسبه و یک الگوریتم گیبز برای نمونه‌گیری از توزیع پسین توسعه داده شده است.

در بخش ۲، رهیافت فراوانی‌گرایانه تحلیل رگرسیونی آمیخته متناهی چوله‌نرمال، به صورت اجمالی مورد بحث قرار می‌گیرد. برای این منظور، توزیع چوله‌نرمال و برخی از ویژگی‌های شاخص آن با تاکید بر مشکلات مربوط به برآورد پارامترهای این توزیع در رهیافت فراوانی‌گرا به اختصار شرح داده می‌شود. در

بخش ۳، برای تحلیل رگرسیونی تحت مدل چوله‌نرمال آمیخته متناهی یک مدل بیزی پیشنهاد شده است. برای ارزیابی مدل بیزی پیشنهادی و مقایسه آن با همتای فراوانی‌گرایانه رقیب، یک مطالعه شبیه‌سازی در بخش ۴ اجرا شده است. در بخش ۵، یک مثال کاربردی در زمینه هزینه‌های خالص خانوارهای شهری در ایران استفاده شده است. بحث و نتیجه‌گیری در بخش ۶ آورده شده است.

۲ تحلیل فراوانی‌گرایانه رگرسیون چوله‌نرمال آمیخته

متغیر تصادفی Z دارای توزیع چوله‌نرمال با پارامتر چولگی α ، پارامتر مکان ξ و پارامتر مقیاس ω است $(SN(\xi, \omega, \alpha))$ ، هرگاه تابع چگالی آن به صورت

$$f_{SN}(y|\xi, \omega, \alpha) = \frac{2}{\omega} \phi\left(\frac{y-\xi}{\omega}\right) \Phi\left(\alpha \frac{y-\xi}{\omega}\right), \quad \omega > 0, \xi, \alpha \in \mathbb{R}, \quad (1)$$

باشد، که در آن $\phi(\cdot)$ و $\Phi(\cdot)$ به ترتیب نشان دهنده توابع چگالی و توزیع تجمعی متغیر تصادفی نرمال استاندارد هستند (آزالی، ۱۹۸۵). این توزیع در واقع آمیخته‌ای از دو توزیع نرمال و نیم‌نرمال است. فرض کنید متغیرهای تصادفی مستقل T و U به ترتیب دارای توزیع‌های نرمال $N(0, \omega^2)$ و نرمال بریده‌شده $TN(0, \omega^2) \mathbb{I}_t(0, +\infty)$ باشند. در این صورت متغیر تصادفی

$$Y = \xi + \frac{\alpha}{\sqrt{1+\alpha^2}} T + \frac{1}{\sqrt{1+\alpha^2}} U, \quad (2)$$

دارای توزیع (۱) خواهد بود و

$$T|Y=y \sim N(\delta(\alpha)(y-\xi), \omega^2(1-\delta^2(\alpha))) \mathbb{I}_t(0, +\infty). \quad (3)$$

به‌علاوه، به کمک روابط مربوط به گشتاورهای توزیع نرمال بریده‌شده (به‌عنوان مثال، بار و همکاران، ۱۹۹۹ را ببینید)، می‌توان نشان داد که دو گشتاور اول توزیع (۳) که در اجرای الگوریتم EM مورد نیاز هستند، به ترتیب عبارتند از

$$\mathbb{E}(T|Y=y) = \delta(\alpha)(y-\xi) + \frac{\phi\left(\frac{\delta(\alpha)(y-\xi)}{\omega\sqrt{1-\delta^2(\alpha)}}\right)}{\Phi\left(\frac{\delta(\alpha)(y-\xi)}{\omega\sqrt{1-\delta^2(\alpha)}}\right)} \omega\sqrt{1-\delta^2(\alpha)}, \quad (4)$$

$$\begin{aligned} \mathbb{E}(T^{\gamma} | Y = y) &= \delta^{\gamma}(\alpha)(y - \xi)^{\gamma} + \omega^{\gamma}(1 - \delta^{\gamma}(\alpha)) \\ &+ \frac{\phi\left(\frac{\delta(\alpha)(y - \xi)}{\omega\sqrt{1 - \delta^{\gamma}(\alpha)}}\right)}{\Phi\left(\frac{\delta(\alpha)(y - \xi)}{\omega\sqrt{1 - \delta^{\gamma}(\alpha)}}\right)} \delta(\alpha)(y - \xi) \omega^{\gamma}(1 - \delta^{\gamma}(\alpha)). \end{aligned} \quad (5)$$

در بسیاری از کاربردهای عملی، توزیع‌های مشاهدات علاوه بر چوله بودن، چندنمایی نیز هستند. این وضعیت زمانی پیش می‌آید که جامعه ناهمگن بوده و از زیرجوامعی با ویژگی‌های مختلف تشکیل شده است. در این حالت، می‌توان از مدل رگرسیونی چوله‌نرمال با متغیرهای پاسخ آمیخته به‌عنوان جایگزین مدل سنتی استفاده کرد. چنانچه متغیر پاسخ دارای توزیع آمیخته g مولفه‌ای باشد، مدل رگرسیونی را می‌توان به‌صورت

$$Y_i | \mathbf{x}_i \sim \sum_{j=1}^g \eta_j f_{SN}(y_i | \boldsymbol{\theta}_j(\mathbf{x}_i)), \quad i = 1, \dots, n, \quad (6)$$

نوشت، که در آن $\mathbf{x}'_i = (x_{i1}, \dots, x_{ip})$ ، $\xi_j(\mathbf{x}_i) = h(\mathbf{x}_i, \boldsymbol{\beta}_j)$ ، $\boldsymbol{\theta}_j(\mathbf{x}_i) = (\xi_j(\mathbf{x}_i), \omega_j^{\gamma}, \alpha_j)$ بردار متغیرهای تبیینی، $\boldsymbol{\beta}_j = (\beta_{j1}, \dots, \beta_{jp})$ بردار ضرایب رگرسیونی زیرجامعه j ام و $\boldsymbol{\eta} = (\eta_1, \dots, \eta_p)$ بردار ضرایب آمیختگی را نشان می‌دهد. در این صورت، مجموعه پارامترهای مدل را می‌توان به‌صورت $\boldsymbol{\Psi} = (\eta_1, \dots, \eta_g, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_g)$ نمایش داد. در رهیافت فراوانی‌گرا، برآورد پارامترها در توزیع‌های آمیخته متناهی معمولاً از طریق الگوریتم EM صورت می‌پذیرد. برای این منظور، مشاهداتی را که از یک توزیع آمیخته متناهی به‌دست می‌آیند، از آن جهت که معلوم نیست هر مشاهده به کدام مؤلفه از توزیع آمیخته منتسب است، مشاهدات ناقص می‌نامند. سپس، مسأله را به کمک مجموعه‌ای از متغیرهای نشان‌گر در قالب داده‌های کامل بازنویسی نموده و مراحل E (امیدگیری) و M (ماکسیم‌سازی) به صورت متوالی تکرار می‌شوند. بر این اساس، بردار نشان‌گر متناظر با مشاهده i ام، $\mathbf{Z}_i = (Z_{i1}, \dots, Z_{ig})$ ، دارای توزیع چندجمله‌ای $M(1, \eta_1, \dots, \eta_g)$ است. بعلاوه، چون هر زیرجامعه دارای یک توزیع چوله‌نرمال است، روابط

$$[T_i | Z_{ij} = 1] \sim TN(0, \omega_j^{\gamma}), \quad [T_i | Y_i = y_i, Z_{ij} = 1] \sim TN(\mu_{T_{ij}}, \sigma_{T_j}^{\gamma}), \quad (7)$$

برقرارند، که در آن‌ها $\mu_{T_{ij}} = \delta(\alpha_j)(y_i - \mathbf{x}'_i \boldsymbol{\beta}_j)$ و $\sigma_{T_j} = \omega_j \sqrt{1 - \delta^{\gamma}(\alpha_j)}$ بر این اساس،

توابع درستنمایی و لگ درستنمایی داده‌های کامل را می‌توان به ترتیب به صورت

$$\begin{aligned}
 L_c(\Psi, \mathbf{T}|\mathbf{y}) &= \prod_{i=1}^n \prod_{j=1}^g \left\{ \left[\frac{\eta_j}{\pi \omega_j^2 \sqrt{1 - \delta^2(\alpha_j)}} \right]^{z_{ij}} \right\} \\
 &\times \prod_{j=1}^g \left\{ \exp \left\{ \frac{-1}{\sqrt{\omega_j^2(1 - \delta^2(\alpha_j))}} \sum_{i=1}^n z_{ij} [(y_i - \mathbf{x}_i' \beta_j)^2 \right. \right. \\
 &\quad \left. \left. - \sqrt{t_i} \delta(\alpha_j) (y_i - \mathbf{x}_i' \beta_j) + t_i^2] \right\} \right\}, \quad (8) \\
 \ell_c(\Psi, \mathbf{T}|\mathbf{y}) &= \sum_{i=1}^n \sum_{j=1}^g z_{ij} \left\{ \log(\eta_j) - \log(\pi \omega_j^2) - \frac{1}{2} \log(1 - \delta^2(\alpha_j)) \right\} \\
 &- \sum_{g=1}^g \frac{1}{\sqrt{\omega_j^2(1 - \delta^2(\alpha_j))}} \left[\sum_{i=1}^n z_{ij} (y_i - \mathbf{x}_i' \beta_j)^2 \right. \\
 &\quad \left. - \sqrt{2} \delta(\alpha_j) \sum_{i=1}^n z_{ij} t_i (y_i - \mathbf{x}_i' \beta_j) + \sum_{i=1}^n z_{ij} t_i^2 \right],
 \end{aligned}$$

نوشت، که در آن $\mathbf{y} = (y_1, \dots, y_n)'$ از این رو، برآورد ماکسیم درستنمایی پارامترها در تکرار $(k+1)$ ام الگوریتم به صورت

$$\begin{aligned}
 \hat{\eta}_j^{(k+1)} &= \frac{1}{n} \sum_{i=1}^n \hat{z}_{ij}^{(k)}, \\
 \hat{\beta}_j^{(k+1)} &= (\mathbf{X}' \mathbf{Z}_j^{*(k)} \mathbf{X})^{-1} (\mathbf{X}' \mathbf{Z}_j^{*(k)} \mathbf{y} - \delta(\hat{\alpha}_j^{(k)}) \mathbf{X}' \hat{\mathbf{S}}_{\sqrt{j}}^{(k)}), \\
 \hat{\omega}_j^{2(k+1)} &= \frac{(\mathbf{y} - \mathbf{X}' \hat{\beta}_j^{(k+1)})' \mathbf{Z}_j^{*(k)} (\mathbf{y} - \mathbf{X}' \hat{\beta}_j^{(k+1)})}{2(1 - \delta^2(\hat{\alpha}_j^{(k)})) \sum_{i=1}^n \hat{z}_{ij}^{(k)}} \\
 &\quad + \frac{-\sqrt{2} \delta(\hat{\alpha}_j^{(k)}) (\mathbf{y} - \mathbf{X}' \hat{\beta}_j^{(k+1)})' \hat{\mathbf{S}}_{\sqrt{j}}^{(k)} + \mathbf{J}' \mathbf{S}_{\sqrt{j}}^{(k)}}{2(1 - \delta^2(\hat{\alpha}_j^{(k)})) \sum_{i=1}^n \hat{z}_{ij}^{(k)}}, \\
 \hat{\alpha}^{(k+1)} &= \arg \max_{\alpha_j} \sum_{i=1}^n \log \left(\sum_{j=1}^g \hat{\eta}_j^{(k+1)} f_{SN}(y_i | \hat{\beta}_j^{(k+1)}, \hat{\omega}_j^{2(k+1)}, \alpha_j) \right),
 \end{aligned}$$

به دست می‌آیند، که در آن‌ها $\hat{\mathbf{S}}_{\sqrt{j}}^{(k)} = (\hat{s}_{1\sqrt{j}}^{(k)}, \hat{s}_{2\sqrt{j}}^{(k)}, \dots, \hat{s}_{n\sqrt{j}}^{(k)})'$ ، $\mathbf{Z}_j^{*(k)} = \text{diag}(\hat{z}_{1j}^{(k)}, \hat{z}_{2j}^{(k)}, \dots, \hat{z}_{nj}^{(k)})$

$$\mathbf{J} = (1, \dots, 1)', \mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]', \mathbf{S}_{\Psi_j}^{(k)} = (\hat{s}_{\Psi_{1j}}^{(k)}, \hat{s}_{\Psi_{2j}}^{(k)}, \dots, \hat{s}_{\Psi_{nj}}^{(k)})'$$

$$\hat{z}_{ij}^{(k)} = \mathbb{E}_{\hat{\Psi}^{(k)}}(Z_{ij} | \mathbf{y}) = \frac{\hat{\eta}_j^{(k)} f_{SN}(y_i | \mathbf{x}_i' \hat{\beta}_j^{(k)}, \hat{\omega}_j^{(k)}, \hat{\eta}_j^{(k)})}{\sum_{j=1}^g \hat{\eta}_j^{(k)} f_{SN}(y_i | \mathbf{x}_i' \hat{\beta}_j^{(k)}, \hat{\omega}_j^{(k)}, \hat{\eta}_j^{(k)})},$$

$$\hat{s}_{\Psi_{1j}}^{(k)} = \mathbb{E}_{\hat{\Psi}^{(k)}}(Z_{ij} T_i | \mathbf{y}) = \hat{z}_{ij}^{(k)} [\hat{\mu}_{T_{ij}}^{(k)} + \frac{\phi(\hat{\alpha}_j^{(k)} (\frac{y_i - \mathbf{x}_i' \hat{\beta}_j^{(k)}}{\hat{\omega}_j^{(k)}}))}{\Phi(\hat{\alpha}_j^{(k)} (\frac{y_i - \mathbf{x}_i' \hat{\beta}_j^{(k)}}{\hat{\omega}_j^{(k)}}))} \hat{\sigma}_{T_i}^{(k)}],$$

$$\hat{s}_{\Psi_{2j}}^{(k)} = \mathbb{E}_{\hat{\Psi}^{(k)}}(Z_{ij} T_i^2 | \mathbf{y}) = \hat{z}_{ij}^{(k)} [\hat{\mu}_{T_{ij}}^{(k)2} + \hat{\sigma}_{T_i}^{(k)2} + \frac{\phi(\hat{\alpha}_j^{(k)} (\frac{y_i - \mathbf{x}_i' \hat{\beta}_j^{(k)}}{\hat{\omega}_j^{(k)}}))}{\Phi(\hat{\alpha}_j^{(k)} (\frac{y_i - \mathbf{x}_i' \hat{\beta}_j^{(k)}}{\hat{\omega}_j^{(k)}}))} \hat{\mu}_{T_{ij}}^{(k)} \hat{\sigma}_{T_i}^{(k)}].$$

ملاحظه می‌شود که محاسبه برآورد پارامترهای مدل رگرسیونی چوله‌نرمال آمیخته متناهی با پیچیدگی‌ها و دشواری‌های زیادی همراه است. برخی از این مشکلات حتی در حالت غیرآمیخته بودن توزیع متغیر پاسخ نیز وجود دارند که توسط بسیاری از محققان مورد بررسی قرار گرفته‌اند. پیوسی (۲۰۰۰) نشان داد که برآورد گشتاوری پارامتر چولگی همواره وجود ندارد. چیوگنا (۲۰۰۵)، فیگوئراس و همکاران (۲۰۰۷) و سارتوری (۲۰۰۶) نشان دادند که برآورد پارامتر چولگی حتی از طریق روش‌های عددی نیز به سادگی ممکن نیست. هم‌چنین، نشان داده شده است که اگر همه مشاهدات مثبت باشند، برآورد ماکسیمم درست‌نمایی پارامتر چولگی به دلیل کراندار نبودن تابع درست‌نمایی وجود ندارد (برای جزئیات بیشتر: نتون، ۲۰۰۴ و دی، ۲۰۱۰ را ببینید). این مشکلات در حالت آمیخته بودن توزیع بیشتر نمود پیدا می‌کند، به طوری که کران‌دار نبودن تابع درست‌نمایی نسبت به پارامترهای چولگی مربوط به مولفه‌های تشکیل‌دهنده یک توزیع آمیخته موجب می‌شود در همگرایی الگوریتم EM مشکلاتی بروز نماید. افزون بر این مشکلات، چون خواص بهینگی برآوردگرهای ماکسیمم درست‌نمایی تنها در حالات بزرگ نمونه‌ای بروز می‌نمایند، رهیافت بیزی به طور خاص برای نمونه‌های کوچک و متوسط به دلیل کارکرد مطلوب‌تر جذابیت‌های بیشتری دارد. در بخش بعد، نشان داده می‌شود که رهیافت بیزی مسأله بسیار سراسر بوده و مشکلات نظری و دشواری‌های عملی رهیافت فراوانی‌گرا را ندارد.

۳ تحلیل بیزی رگرسیون چوله‌نرمال آمیخته

در این بخش تحلیل رگرسیونی چوله‌نرمال آمیخته متناهی از دیدگاه بیزی مورد بررسی قرار گرفته و یک مدل بیزی برای آن توسعه داده شده است. لازمه تحلیل بیزی آن است که پارامترهای مدل متغیرهایی تصادفی تلقی شده و برای آن‌ها توزیع‌های پیشین مناسب لحاظ شود. چنانچه ریچاردسون و گرین (۱۹۹۷) اشاره نموده‌اند، در تحلیل بیزی مدل‌های آمیخته متناهی نمی‌توان تمام توزیع‌های پیشین را ناآگاهی بخش در نظر گرفت. زیرا همیشه این امکان وجود دارد که به یک یا چند زیرجامعه هیچ مشاهده‌ای تعلق نگیرد و از این رو مشاهدات درباره آن مولفه‌ها آگاهی بخش نباشند، که این امر موجب ناسره شدن توزیع پسین خواهد شد. در این مقاله، به پیروی از گلن و همکاران (۲۰۰۸) برای ضرایب رگرسیونی و پارامتر واریانس مربوط به هر مؤلفه از توزیع آمیخته، به ترتیب توزیع‌های پیشین نرمال و وارون گاما در نظر گرفته شده‌اند. برای ضرایب آمیختگی، $\eta = (\eta_1, \dots, \eta_g)$ ، نیز از توزیع پیشین دیریکله استفاده شده است. این انتخاب یکی از انتخاب‌های مرسوم در ادبیات توزیع‌های آمیخته محسوب می‌شود. بر این اساس، برای پارامترهای مدل توزیع‌های پیشین

$$\begin{aligned} \beta_j &\sim N_p(\mu_j, \Sigma_j), & j = 1, \dots, g, \\ \omega_j^{-1} | b &\sim \Gamma(a, b), & j = 1, \dots, g; \quad a, b \in \mathcal{R}^+, \\ b &\sim \Gamma(v_1, v_2), & v_1, v_2 \in \mathcal{R}^+, \\ \eta &\sim \text{Dirichlet}(h_1, \dots, h_g), & h_1, \dots, h_g \in \mathcal{R}^+ \\ \delta(\alpha_j) &\sim U(-1, 1), & j = 1, \dots, g. \end{aligned} \quad (9)$$

در نظر گرفته شده است، که در آن‌ها $\mu_j, \Sigma_j, v_1, v_2, a$ و (h_1, \dots, h_g) ، ابرپارامترهای مدل بیزی هستند و توسط تحلیل‌گر تعیین یا برآورد می‌شوند. با در نظر گرفتن یک توزیع پیشین برای ابرپارامتر b ، در توزیع پیشین وارون گاما، تلاش شده است استواری بیشتری به استنباط‌های بیزی حاصل القاء شود. هم‌چنین، به منظور پرهیز از پیش‌داوری درباره سهم مولفه‌های توزیع آمیخته از نمونه تصادفی، معمولاً ابرپارامترهای توزیع پیشین دیریکله به صورت $h_1, \dots, h_g = h$ برابر در نظر گرفته می‌شوند. این توزیع‌های پیشین انتخاب‌های بی‌طرفانه‌ای محسوب می‌شوند که برخی از مزایای مدل‌سازی بیزی برپایه توزیع‌های پیشین مزدوج را نیز به همراه دارند. این انتخاب‌ها تحلیل بیزی را در حالتی مد نظر قرار می‌دهند

که اطلاعات پیشین قابل توجه و قوی درباره پارامترهای مدل در دسترس تحلیل‌گر قرار ندارد. اگر اطلاعات پیشین قابل توجهی در دسترس باشند، می‌توان توزیع‌های پیشین آگاهی‌بخش‌تری را مد نظر قرار داد. برای ملاحظه جزئیات بیشتر درباره پیشین‌های انتخاب شده، به‌عنوان نمونه، ریچاردسون و گرین (۱۹۹۷) و برناردی و همکاران (۲۰۱۲) را ملاحظه نمائید.

برای تعیین ابرپارامترهای توزیع پیشین نرمال، داده‌ها را بر اساس میزان شباهت آن‌ها به g گروه تقسیم کرده و از برآورد ضرایب رگرسیونی و ماتریس کواریانس متناظر با آن‌ها، برای تعیین ابر پارامترهای هر گروه استفاده می‌شود. با توجه به تابع درستنمایی داده‌های کامل، رابطه (۸)، توزیع پسین داده‌های کامل پس از انجام محاسبات جبری، به صورت

$$\begin{aligned} \pi(\Psi, \mathbf{T}|\mathbf{y}) &= L_c(\Psi, \mathbf{T}|\mathbf{y}) \times \pi(\Psi) \propto \prod_{j=1}^g \prod_{i=1}^n \left\{ \frac{\eta_j}{\pi \omega_j^\nu \sqrt{1 - \delta^\nu(\alpha_j)}} \right. \\ &\times \exp\left\{ -\frac{1}{2\omega_j^\nu (1 - \delta^\nu(\alpha_j))} [(y_i - \mathbf{x}'_i \boldsymbol{\beta}_j)^\nu - 2\delta(\alpha_j) t_i (y_i - \mathbf{x}'_i \boldsymbol{\beta}_j) + t_i^\nu] \right\}^{Z_{ij}} \\ &\times \prod_{j=1}^g \left\{ |\boldsymbol{\Sigma}_j|^{-p/\nu} \exp\left\{ -\frac{1}{\nu} (\boldsymbol{\beta}_j - \boldsymbol{\mu}_j)' \boldsymbol{\Sigma}_j^{-1} (\boldsymbol{\beta}_j - \boldsymbol{\mu}_j) \right\} b^a (\omega_j^{-\nu})^{a-1} e^{-b\omega_j^{-\nu}} \eta_j^h \right\} \\ &\times b^{v_1-1} e^{-v_2 b}, \end{aligned}$$

به دست می‌آید. با توجه به پیچیدگی توزیع پسین (۱۰) برآوردگرهای بیزی پارامترهای مدل، فاقد صورت بسته هستند. از این رو، لازم است از روش‌های مونت‌کارلو زنجیر مارکوفی برای نمونه‌گیری از توزیع پسین و تقریب کمیت‌های پسینی مورد علاقه استفاده نمود. برای این منظور، در ادامه یک الگوریتم گیبز برای نمونه‌گیری از توزیع پسین توسعه داده شده است. لازمه استفاده از الگوریتم گیبز شناخت توزیع‌های پسین شرطی کامل پارامترها است. این توزیع‌ها پس از انجام محاسبات جبری به صورت

$$\begin{aligned} \mathbf{Z}_i^{(k+1)} | \text{others} &\sim M(1, \eta_1^*, \dots, \eta_g^*), \\ T_i^{(k+1)} | \text{others} &\sim TN(\delta(\alpha_j^{(k)})(y_i - \mathbf{x}'_i \boldsymbol{\beta}_j^{(k)}), \omega_j^{\nu(k)} (1 - \delta^\nu(\alpha_j^{(k)}))), \\ b^{(k+1)} | \text{others} &\sim \Gamma(v_1 + ga, v_2 + \sum_{j=1}^g \omega_j^{-\nu(k)}), \\ \boldsymbol{\eta}^{(k+1)} | \text{others} &\sim \text{Dirichlet}(h + n_1^{(k+1)}, \dots, h + n_g^{(k+1)}), \\ \boldsymbol{\beta}_j^{(k+1)} | \text{others} &\sim N(\mathbf{M}_j^{(k+1)}, \mathbf{S}_j^{(k+1)}), \end{aligned}$$

$$\omega_j^{-\nu(k+1)} | others \sim \Gamma(a + n_j^{(k+1)}, b^{(k+1)} + \gamma),$$

$$\pi(\boldsymbol{\delta} | others) \propto \prod_{i=1}^n \prod_{j=1}^g \{ (1 - \delta^\nu(\alpha_j))^{-\frac{1}{\nu}} \exp\left\{ -\frac{1}{\nu \omega_j^{\nu(k+1)} (1 - \delta^\nu(\alpha_j))} \left[(y_i - \mathbf{x}'_i \boldsymbol{\beta}_j^{(k+1)})^\nu - \nu \delta(\alpha_j) t_i^{(k+1)} (y_i - \mathbf{x}'_i \boldsymbol{\beta}_j^{(k+1)}) + t_i^{\nu(k+1)} \right] \right\} Z_{ij}^{(k+1)}, \quad (10)$$

به دست می‌آیند، که در آن‌ها

$$\eta_j^* = \frac{\eta_j^{(k)} f_{SN}(y_i | \mathbf{x}'_i \boldsymbol{\beta}_j^{(k)}, \omega_j^{\nu(k)}, \alpha_j^{(k)})}{\sum_{m=1}^g \eta_m^{(k)} f_{SN}(y_i | \mathbf{x}'_i \boldsymbol{\beta}_m^{(k)}, \omega_m^{\nu(k)}, \alpha_m^{(k)})},$$

$$n_j^{(k+1)} = \sum_{i=1}^n Z_{ij}^{(k+1)},$$

$$\mathbf{Z}_j^{*(k+1)} = \text{diag}(Z_{1j}^{(k+1)}, \dots, Z_{nj}^{(k+1)}),$$

$$\mathbf{S}_j^{(k+1)} = \left(\frac{\mathbf{X}' \mathbf{Z}_j^{*(k+1)} \mathbf{X}}{\omega_j^{\nu(k)} (1 - \delta^\nu(\alpha_j^{(k)}))} + \boldsymbol{\Sigma}_j^{-1} \right)^{-1},$$

$$\mathbf{M}_j^{(k+1)} = \mathbf{S}_j^{(k+1)} \left(\frac{\mathbf{X}' \mathbf{Z}_j^{*(k+1)} \mathbf{y} - \delta(\alpha_j^{(k)}) \mathbf{X}' \mathbf{Z}_j^{*(k+1)} \boldsymbol{\Gamma}^{(k+1)}}{\omega_j^{\nu(k)} (1 - \delta^\nu(\alpha_j^{(k)}))} + \boldsymbol{\Sigma}_j^{-1} \boldsymbol{\mu}_j \right),$$

$$\gamma = \frac{1}{\nu (1 - \delta^\nu(\alpha_j^{(k)}))} \left[(\mathbf{y} - \mathbf{X} \boldsymbol{\beta}_j^{(k+1)})' \mathbf{Z}_j^{*(k+1)} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}_j^{(k+1)}) - \nu \delta(\alpha_j^{(k)}) (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}_j^{(k+1)})' \mathbf{Z}_j^{*(k+1)} \boldsymbol{\Gamma}^{(k+1)} + \boldsymbol{\Gamma}'^{(k+1)} \mathbf{Z}_j^{*(k+1)} \boldsymbol{\Gamma}^{(k+1)} \right].$$

به دلیل عدم اطلاع از ثابت چگالی‌ساز در توزیع پسین شرطی کامل پارامتر $\boldsymbol{\delta} = (\delta(\alpha_1), \dots, \delta(\alpha_g))$ در رابطه (۱۰)، برای نمونه‌گیری از این توزیع از الگوریتم متروپولیس-هاستینگز استفاده شده است. برای این منظور، به پیروی از لین و همکاران (۲۰۰۷)، ابتدا $\delta(\alpha_j)$ را به صورت

$$\delta^*(\alpha_j) = \log \left\{ \frac{1 + \delta(\alpha_j)}{1 - \delta(\alpha_j)} \right\} = \nu \tanh^{-1}(\delta(\alpha_j)),$$

تبدیل کرده و سپس الگوریتم متروپولیس- هاستینگز برای نمونه‌گیری از توزیع

$$g(\delta^*) = \pi(\delta(\delta^*)) \prod_{j=1}^g J_{\delta^*(\alpha_j)},$$

به کار گرفته می‌شود، که در آن $\delta^* = (\delta^*(\alpha_1), \dots, \delta^*(\alpha_g))$ و $J_{\delta^*(\alpha_j)} = \frac{2e^{\delta^*(\alpha_j)}}{(1 + e^{\delta^*(\alpha_j)})^2}$ زاکوبی تبدیل است. برای این منظور، به پیروی از گلن و همکاران (۱۹۹۶)، یک توزیع نرمال $-g$ -متغیره با میانگین $\delta^{*(k)}$ و ماتریس کواریانس $c^2 \Sigma_{\delta^*}^{(k)}$ به‌عنوان توزیع پیشنهادی لحاظ شده است که در آن $c \approx \frac{2/4}{\sqrt{g}}$ و مقدار $\Sigma_{\delta^*}^{(k)}$ توسط وارون ماتریس اطلاع فیشر برآورد می‌شود. پس از استخراج نمونه‌های $\delta^*(\alpha_j)$ ، $j = 1, \dots, g$ ، می‌توان مقادیر $\delta(\alpha_j)$ و در نهایت α_j را به‌دست آورد. این تبدیل پارامتر $\delta(\alpha_j)$ با دامنه تغییرات $(-1, 1)$ را به پارامتر $\delta^*(\alpha_j)$ با دامنه تغییرات \mathbb{R} تبدیل می‌کند و این امکان را فراهم می‌سازد که بتوان از توزیع پیشنهادی گلن و همکاران (۱۹۹۶) که یک توزیع پیشنهادی کارا است، استفاده نمود. در واقع از این طریق می‌توان به یک الگوریتم متروپولیس-هاستینگز با نرخ پذیرش مطلوب برای نمونه‌گیری از توزیع پسین شرطی کامل مورد نظر دست یافت.

ذکر این نکته ضروری است که یکی از مشکلات جدی در تحلیل بیزی توزیع‌های آمیخته، وضعیتی است که از آن تحت عنوان جابجایی‌پذیری اندیس‌ها یاد می‌شود. در چنین وضعیتی، توزیع پسین توام نسبت به جایگشت پارامترها ناوردا بوده و از این رو زنجیرهای مارکوف شبیه‌سازی‌شده از توزیع پسین توام پارامترها، شناسایی‌پذیر نیستند. بر این اساس، استنباط‌های پسینی برپایه توزیع‌های پسین حاشیه‌ای و از راه محاسبه کمیت‌های پسینی مانند میانگین پسین، گمراه‌کننده خواهد بود. روش‌های مختلفی برای برخورد با مشکل جابجایی‌پذیری اندیس‌ها، پیشنهاد شده است که خواننده علاقه‌مند می‌تواند برای ملاحظه برخی از این روش‌ها، استفنس (۲۰۰۰) و جاسارا و همکاران (۲۰۰۵) را ملاحظه نماید. در این مقاله، برای مواجهه با این مشکل از بسته کتابخانه‌ای *label.switching* در نرم‌افزار R که توسط پایاستامولیس (۲۰۱۵) توسعه داده شده است، بهره برده‌ایم. این بسته نرم‌افزاری مجموعه‌ای از جایگشت‌ها را برای مرتب‌سازی مجدد زنجیرهای مارکوف شبیه‌سازی‌شده از توزیع پسین، فراهم می‌سازد. هرگونه استنباط پسینی بر پایه این زنجیرهای مجدداً مرتب شده، صورت می‌پذیرد.

۴ مطالعه شبیه‌سازی

در این بخش برای ارزیابی رهیافت بیزی پیشنهادی یک مدل رگرسیونی آمیخته دو مؤلفه‌ای به صورت

$$Y_i | \mathbf{x}_i \sim \eta f_{SN}(\beta_{11}x_{i1} + \beta_{12}x_{i2}, \omega_1^2, \alpha_1) + (1 - \eta) f_{SN}(\beta_{21}x_{i1} + \beta_{22}x_{i2}, \omega_2^2, \alpha_2), \quad (11)$$

$i = 1, 2, \dots, n$ در نظر گرفته شده است. با فرض $(\beta_{11}, \beta_{12}) = (9, 15)'$ ، $(\beta_{21}, \beta_{22}) = (-1, -6)'$ ، $\omega_1^2 = 0.5$ ، $\omega_2^2 = 0.5$ و $\eta = 0.5$ ، مقادیر متغیرهای پاسخ از توزیع (۱۱) شبیه‌سازی شده‌اند. به منظور لحاظ نمودن وضعیت‌های مختلف ممکن در ارزیابی رهیافت بیزی پیشنهادی، برای مؤلفه اول-دوم توزیع آمیخته (۱۱)، چهار وضعیت مختلف چوله به راست-چوله به راست، چوله به راست-مقارن، چوله به چپ، و مقارن-مقارن، به ترتیب برای مؤلفه‌های اول و دوم توزیع آمیخته، در نظر گرفته شده است.

برای ارزیابی تأثیر اندازه نمونه بر کارایی برآوردگرهای بیزی، مطالعه شبیه‌سازی به ازای مقادیر مختلف اندازه نمونه 10^1 ، 20^1 ، 40^1 و 80^1 تکرار شده است. چون رهیافت پیشنهادی بیزی است، از در نظر گرفتن اندازه نمونه‌های بزرگتر پرهیز شده است، چرا که مزیت اصلی رهیافت بیزی عمدتاً مربوط به اندازه‌های کوچک نمونه‌ای است و برای نمونه‌های بزرگ به دلیل غلبه درستنمایی بر پیشین، نتایج دو رهیافت بیزی و فراوانی‌گرا مشابه خواهند بود. به منظور لحاظ نمودن عدم قطعیت و تغییرپذیری حاکم بر فرایند تولید نمونه‌های تصادفی، تمام محاسبات ۱۰۰۰ بار تکرار و مقادیر جذر میانگین توان‌های دوم خطای برآوردگرهای ضرایب رگرسیونی، برای وضعیت‌های چهارگانه یادشده محاسبه و در جداول ۱ الی ۴ ارائه شده‌اند.

برای اطمینان از همگرایی زنجیرهای مارکوف تولید شده به توزیع مانای متناظر خود، از معیار پیشنهادی گلמן و روبین (۱۹۹۲) استفاده شده است. برای این منظور تابع *gelman.diag* در بسته کتابخانه‌ای *Coda* (پلومر و همکاران، ۲۰۰۶) از نرم‌افزار *R* (گروه مرکزی *R*، ۲۰۱۳) مورد استفاده قرار گرفته است. لازم به ذکر است که زمان اجرای الگوریتم‌ها بسته به تعداد مؤلفه‌های توزیع آمیخته، مقادیر مختلف ضرایب آمیختگی و اندازه نمونه‌های مختلف، متفاوت است. به عنوان مثال، در یک توزیع آمیخته دو مؤلفه‌ای، به ازای ضریب آمیختگی 0.5 و اندازه نمونه 40^1 ، زمان همگرایی زنجیر مارکوف تولید شده از توزیع پسین در یک بار اجرای الگوریتم گیبز، اندکی کمتر از 20^1 دقیقه بوده است.

با توجه به مقادیر جذر میانگین توان‌های دوم خطای مربوط به ضرایب رگرسیونی، ملاحظه می‌شود

جدول ۱. مقادیر جذر میانگین توان دوم خطای برآوردگرهای ضرایب رگرسیونی برای مدل‌های رقیب در وضعیت متقارن-متقارن $(\alpha_1 = 0, \alpha_2 = 0)$.

n	چوله‌نرمال آمیخته			نرمال آمیخته		
	۲۰	۴۰	۸۰	۲۰	۴۰	۸۰
β_{11}	۰/۰۰۰۰۱	$< 10^{-12}$	$< 10^{-12}$	$< 10^{-12}$	$< 10^{-12}$	$< 10^{-12}$
β_{12}	۰/۰۰۰۰۱	$< 10^{-12}$	$< 10^{-12}$	$< 10^{-12}$	$< 10^{-12}$	$< 10^{-12}$
β_{21}	۰/۰۰۰۰۳	$< 10^{-12}$	۰/۰۰۰۰۱	۰/۰۰۰۰۷	$< 10^{-12}$	$< 10^{-12}$
β_{22}	۰/۰۰۰۰۴	$< 10^{-12}$	$< 10^{-12}$	۰/۰۰۰۰۵	$< 10^{-12}$	$< 10^{-12}$
ω_1^2	۰/۰۰۰۲۳	۰/۰۰۰۸۹	۰/۰۰۰۳۶	۰/۰۰۰۹۲	۰/۰۰۰۰۳	۰/۰۰۰۰۸
ω_2^2	۰/۰۰۰۲۰	۰/۰۰۰۹۳	۰/۰۰۰۴۱	۰/۰۰۰۱۳	۰/۰۰۰۲۹	۰/۰۰۰۰۶
α_1	۰/۰۰۰۰۱	$< 10^{-12}$	۰/۰۰۰۰۴	-	-	-
α_2	۰/۰۰۰۰۴	۰/۰۰۰۰۱	۰/۰۰۰۰۱	-	-	-
η	۰/۰۰۰۱۵	$< 10^{-12}$	$< 10^{-12}$	۰/۰۰۰۲۰	$< 10^{-12}$	$< 10^{-12}$

جدول ۲. مقادیر جذر میانگین توان دوم خطای برآوردگرهای ضرایب رگرسیونی برای مدل‌های رقیب در وضعیت چوله به راست-متقارن $(\alpha_1 = 4, \alpha_2 = 0)$.

n	چوله‌نرمال آمیخته			نرمال آمیخته		
	۲۰	۴۰	۸۰	۲۰	۴۰	۸۰
β_{11}	۰/۰۱۵۳	۰/۰۰۰۶۴	۰/۰۰۰۱۹	۰/۰۷۵۴	۰/۰۵۳۲	۰/۰۸۱۰
β_{12}	۰/۰۲۵۷	۰/۰۰۰۸۳	۰/۰۰۰۲۲	۰/۱۰۸۷	۰/۰۸۵۳	۰/۰۹۰۴
β_{21}	۰/۰۰۰۰۱	۰/۰۰۰۰۱	$< 10^{-12}$	۰/۰۰۰۰۴	۰/۰۰۰۰۱	$< 10^{-12}$
β_{22}	۰/۰۰۰۰۲	$< 10^{-12}$	$< 10^{-12}$	۰/۰۰۰۰۲	۰/۰۰۰۰۱	$< 10^{-12}$
ω_1^2	۰/۰۱۲۱	۰/۰۱۰۳	۰/۰۰۰۳۷	۰/۰۷۵۲	۰/۰۶۶۴	۰/۰۶۱۴
ω_2^2	۰/۰۰۰۲۲	۰/۰۰۰۶۰	۰/۰۰۰۳۶	۰/۰۱۴۹	۰/۰۰۰۲۳	۰/۰۰۰۰۸
α_1	۹/۰۹۷۲	۵/۳۶۶۷	۱/۴۸۳۷	-	-	-
α_2	۰/۰۰۰۰۱	۰/۰۰۰۰۹	$< 10^{-12}$	-	-	-
η	۰/۰۰۰۱۴	$< 10^{-12}$	$< 10^{-12}$	۰/۰۰۰۲۰	$< 10^{-12}$	$< 10^{-12}$

که وقتی جامعه ناهمگن بوده و از زیرجوامعی با ویژگی‌های متفاوت تشکیل شده است، توزیع‌های آمیخته متناهی به خوبی ناهمگنی و تمایز بین زیرجوامع را مدنظر قرار می‌دهند. بعلاوه، هنگامی که توزیع‌های مربوط به هر زیرجامعه دارای ساختار نامتقارن هستند، امکان مدل‌بندی این عدم تقارن به کمک توزیع‌های آمیخته چوله‌نرمال به خوبی فراهم آورده می‌شود. کاهش مقادیر جذر میانگین توان‌های دوم خطا ضرایب رگرسیونی با افزایش اندازه نمونه حاکی از سازگاری برآوردگرهای پیشنهادی است.

۵ مثال کاربردی

در این بخش نحوه کاربست رهیافت بیزی پیشنهادی را در قالب یک مثال کاربردی شرح می‌دهیم. برای این منظور، بخشی از داده‌های هزینه و درآمد خانوارها، مربوط به سرشماری نفوس و مسکن سال ۱۳۸۸، مورد استفاده قرار گرفته است. جدول ۵ مقادیر درآمد ۳۶ خانوار ایرانی (به هزار ریال) به همراه هزینه و بعد خانوارها را نشان می‌دهد. هدف بررسی تأثیر متغیرهای تبیینی بعد و درآمد، بر متغیر پاسخ هزینه‌های خالص خانوارها است. بافت‌نگار هزینه خالص خانوارها که در شکل ۱ رسم شده است، به‌خوبی ساختار ناهمگن، دونمایی و نامتقارن مشاهدات تحت بررسی را نشان می‌دهد.

با این وجود نمی‌توان در مورد تعداد مولفه‌های توزیع آمیخته تنها بر پایه شکل توزیع قضاوت نمود، چرا که در بسیاری حالات نماهای مختلف یک توزیع آمیخته متناهی دارای هم‌پوشانی بوده و به‌راحتی قابل تمایز نیستند. در اغلب کاربردها تعداد مولفه‌های یک توزیع آمیخته متناهی مشخص نیست و می‌بایست به‌عنوان یک پارامتر بر اساس مشاهدات موجود برآورد شود. در طی دهه‌های گذشته روش‌های مختلفی برای این منظور پیشنهاد شده است. یکی از عمومی‌ترین روش‌ها در این زمینه، روش مبتنی بر معیارهای انتخاب مدل است. در این روش‌ها مدل‌هایی با تعداد مولفه‌های مختلف به مشاهدات برازش و سپس تعداد مولفه‌ها به‌صورت $\hat{g} = \arg \min_{g \in G} C(\hat{\Psi}_g, g)$ برآورد می‌شود، که در آن G مجموعه مقادیر ممکن برای تعداد مولفه‌ها، $\hat{\Psi}$ برآوردگر پارامترهای مدل با فرض آنکه مدل g مولفه دارد و $C(\cdot)$ یک معیار انتخاب مدل مانند معیار اطلاع آکائیک (AIC)، شوارتز (BIC) و هانان-کوئین (HQC) یا

۳۱۰ تحلیل بیزی رگرسیون چوله‌نرمال آمیخته

جدول ۳. مقادیر جذر میانگین توان دوم خطای برآوردگرهای ضرایب رگرسیونی برای مدل‌های رقیب در وضعیت چوله به راست-چوله به چپ ($\alpha_1 = 4, \alpha_2 = -4$).

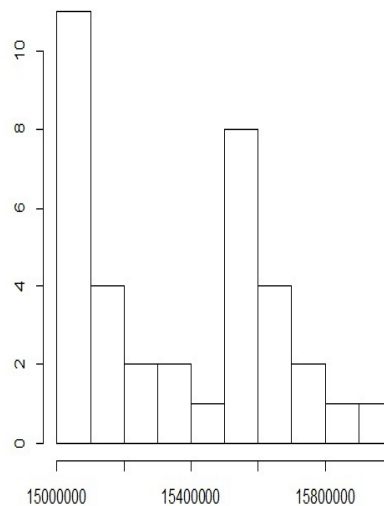
n	چوله‌نرمال آمیخته			نرمال آمیخته		
	۲۰	۴۰	۸۰	۲۰	۴۰	۸۰
β_{11}	۰/۰۲۵۹	۰/۰۰۵۸	۰/۰۰۱۰	۰/۰۸۸۷	۰/۰۵۶۰	۰/۰۷۵۶
β_{12}	۰/۰۴۱۲	۰/۰۰۵۶	۰/۰۰۱۴	۰/۱۲۶۳	۰/۰۵۳۴	۰/۰۶۴۵
β_{21}	۰/۰۸۱۶	۰/۰۰۵۶	۰/۰۰۱۵	۰/۳۸۴۶	۰/۰۵۵۵	۷۶۴/۰
β_{22}	۰/۰۰۷۷	۰/۰۰۵۲	۰/۰۰۱۴	۰/۰۱۹۲	۰/۰۵۲۶	۰/۰۶۵۹
ω_1^2	۰/۰۲۶۳	۰/۰۰۱۴	۰/۰۰۳۸	۰/۰۸۶۶	۰/۰۶۳۵	۰/۰۵۸۹
ω_2^2	۰/۰۲۱۱	۰/۰۰۳۳	۰/۰۰۴۸	۰/۰۷۹۴	۰/۰۶۲۶	۰/۰۶۱۳
α_1	۱۰/۰۱۳۵	۵/۳۰۷۹	۱/۲۶۱۶	-	-	-
α_2	۹/۹۷۲۷	۵/۰۲۵۶	۱/۱۷۶۹	-	-	-
η	۰/۰۰۱۳	$< 10^{-12}$	$< 10^{-12}$	۰/۰۰۱۸	$< 10^{-12}$	۰/۰۰۲۹

جدول ۴. مقادیر جذر میانگین توان دوم خطای برآوردگرهای ضرایب رگرسیونی برای مدل‌های رقیب در وضعیت چوله به راست-چوله به راست ($\alpha_1 = 4, \alpha_2 = 4$).

n	چوله‌نرمال آمیخته			نرمال آمیخته		
	۲۰	۴۰	۸۰	۲۰	۴۰	۸۰
β_{11}	۰/۰۶۰۸	۰/۰۱۱۷	۰/۰۰۲۵	۰/۱۴۶۷	۰/۰۸۷۳	۰/۰۷۲۴
β_{12}	۰/۰۲۲۰	۰/۰۱۱۰	۰/۰۰۱۹	۰/۰۵۲۸	۰/۰۸۱۰	۰/۰۵۹۲
β_{21}	۰/۰۲۲۵	۰/۰۰۹۹	۰/۰۰۲۱	۰/۰۴۸۰	۰/۰۸۱۶	۰/۰۷۶۱
β_{22}	۰/۰۴۵۸	۰/۰۱۰۰	۰/۰۰۲۲	۰/۱۱۲۳	۰/۰۷۹۶	۰/۰۵۹۹
ω_1^2	۰/۰۳۳۵	۰/۰۱۶۵	۰/۰۰۶۱	۰/۰۹۶۰	۰/۰۶۸۲	۰/۰۶۲۶
ω_2^2	۰/۰۳۰۳	۰/۰۱۶۹	۰/۰۰۷۷	۰/۰۹۴۳	۰/۰۷۳۳	۰/۰۶۶۸
α_1	۱۱/۴۶۲۲	۶/۰۸۶۲	۱/۷۹۵۹	-	-	-
α_2	۱۱/۵۷۴۰	۵/۹۸۹۳	۱/۸۷۴۲	-	-	-
η	۰/۰۰۱۵	$< 10^{-12}$	$< 10^{-12}$	۰/۰۰۲۱	$< 10^{-12}$	$< 10^{-12}$

جدول ۵. مقادیر درآمد ۳۶ خانوار ایرانی (به ۱۰۰۰ ریال) به همراه هزینه و بعد خانوارها.

درآمد	بعد خانوار	هزینه	ردیف	درآمد	بعد خانوار	هزینه	ردیف
۲۰۰۷۰	۳	۱۵۰۰۴۶/۴	۱۹	۱۹۰۶۰	۲	۱۵۰۱۴	۱
۲۴۱۲۰	۴	۱۵۶۱۰/۴۸	۲۰	۹۸۵۰	۳	۱۵۰۶۲/۴	۲
۱۵۵۸۰	۲	۱۵۱۲۳	۲۱	۳۵۰۴۰	۴	۱۵۲۹۷/۷۲	۳
۱۵۳۲۹/۹۷	۲	۱۵۰۵۳/۹۷	۲۲	۸۹۵۲	۲	۱۵۵۴۳/۵۲	۴
۱۳۸۵۰۰	۱	۱۵۱۰۸	۲۳	۵۰۵۰۰	۱	۱۵۰۷۲	۵
۱۷۴۳۰	۳	۱۵۴۴۶/۸	۲۴	۳۳۶۰۰	۲	۱۵۰۲۰	۶
۳۵۹۵/۲	۶	۱۵۰۳۱/۲	۲۵	۲۷۶۰۰	۱	۱۵۵۰۸/۸	۷
۱۶۸۰۰	۲	۱۵۰۳۱/۶	۲۶	۳۰۳۰۰	۱	۱۵۵۴۹/۲	۸
۴۳۹۶۰	۷	۱۵۰۴۶/۶	۲۷	۹۴۳۰	۱	۱۵۰۳۸	۹
۴۳۵۰۰	۴	۱۵۵۸۲/۱۲	۲۸	۶۱۷۶۰	۴	۱۵۰۵۴	۱۰
۱۵۰۴۶	۲	۱۵۶۷۸	۲۹	۱۳۶۵۰۰	۲	۱۵۱۴۲/۲	۱۱
۸۷۱۰	۱	۱۵۳۳۵/۶۶	۳۰	۱۹۶۰۰	۲	۱۵۱۵۵/۳۸	۱۲
۱۹۱۱۰	۳	۱۵۵۶۳/۶۸	۳۱	۳۴۲۲۰	۴	۱۵۳۹۷/۷	۱۳
۲۴۷۲۰	۴	۱۵۵۵۴	۳۲	۱۶۸۰۰	۱	۱۵۸۱۰/۲	۱۴
۲۶۰۰۰	۲	۱۵۷۰۹	۳۳	۱۶۰۰۰	۲	۱۵۵۶۰	۱۵
۱۷۸۶۶/۸	۱	۱۵۶۴۶/۲۸	۳۵	۱۲۹۸۵/۲۶	۲	۱۵۵۸۲/۸۶	۱۶
۱۸۸۰۰	۲	۱۵۶۵۲	۳۵	۱۹۲۰۰	۱	۱۵۷۰۶/۷۲	۱۷
۴۳۵۹۰	۳	۱۵۲۴۰/۵۴	۳۶	۳۴۲۱۲	۲	۱۵۹۷۸	۱۸



شکل ۱. بافت‌نگار هزینه خالص خانوارهای شهری.

DIC بیزی است. سه معیار اول به‌ترتیب به‌صورت

$$AIC = -2\ell(\hat{\theta}) + 2k,$$

$$BIC = -2\ell(\hat{\theta}) + k \log n,$$

$$HQC = -2\ell(\hat{\theta}) + 2k \log(\log n),$$

محاسبه می‌شوند، که در آن‌ها k تعداد پارامترهای هر مدل و n اندازه نمونه را نشان می‌دهد. معیار DIC بیزی به‌صورت

$$DIC = \bar{D} - D(\bar{\theta}), \quad (12)$$

تعریف می‌شود، که در آن $D(\theta) = -2\ell(\theta)$ را انحراف گویند، $\bar{D} = E(D(\theta)|\mathbf{y})$ و $\bar{\theta} = E(\theta|\mathbf{y})$. برای ملاحظه جزئیات بیشتر درباره نحوه محاسبه و تعبیر معیار DIC ، اشپیگل‌هاتر و همکاران (۲۰۰۲)

را ببینید. جدول ۶ مقادیر معیارهای اطلاع مختلف را برای مدل چوله نرمال آمیخته به ازای تعداد مولفه‌های مختلف، نشان می‌دهد. همان‌طور که ملاحظه می‌شود، مدل چوله نرمال آمیخته دومولفه‌ای، بر اساس همه معیارهای اطلاع یادشده، در مقایسه با سایر مدل‌ها از حمایت بیشتری از سوی مشاهدات برخوردار است. از این‌رو، برای مدل‌بندی متغیر پاسخ از یک توزیع چوله نرمال آمیخته دومولفه‌ای به صورت

$$Y_i | \mathbf{x}_i \sim \eta f_{SN}(\beta_{11}x_{i1} + \beta_{12}x_{i2}, \omega_1^2, \alpha_1) + (1 - \eta) f_{SN}(\beta_{21}x_{i1} + \beta_{22}x_{i2}, \omega_2^2, \alpha_2), \quad (13)$$

$i = 1, 2, \dots, n$ استفاده شده است. پارامترهای این مدل بر اساس رهیافت بیزی پیشنهادی برآورد و در جدول ۷ ارائه شده‌اند. مدل فراوانی‌گرایانه نرمال آمیخته نیز به‌عنوان مدل سنتی رقیب به داده‌ها برازش داده شده است، تا اثر عدم تقارن مشاهدات بر کارایی مدل‌های رقیب به صورت دقیق‌تر مورد مطالعه قرار بگیرد. به منظور ارزیابی مدل بیزی پیشنهادی و مقایسه آن با مدل رقیب فراوانی‌گرا، مقادیر معیارهای اطلاع AIC ، BIC و HQC که ذاتاً معیارهایی فراوانی‌گرا هستند، نیز برای مدل‌های مورد بحث محاسبه و در جدول ۸ ارائه شده‌اند. ملاحظه می‌شود که مدل رگرسیونی چوله نرمال آمیخته بیزی دومولفه‌ای حتی بر اساس معیارهای فراوانی‌گرایانه نیز، در مقایسه با مدل فراوانی‌گرایانه بیشتری دارد. بعلاوه، در رهیافت بیزی پیشنهادی، مشکلات معمول مدل رقیب فراوانی‌گرا مانند کران‌دار نبودن تابع درستنمایی، عدم همگرایی الگوریتم EM و بهینه نبودن نتایج برای نمونه‌های کوچک، بروز پیدا نمی‌کند. برنامه‌های نوشته‌شده برای برازش مدل رگرسیونی بیزی چوله نرمال آمیخته متناهی، در نشانی اینترنتی <https://sites.google.com/site/datasetandprograms> در دسترس هستند.

بحث و نتیجه‌گیری

تحلیل رگرسیونی معمولاً با فرض متقارن بودن متغیر پاسخ صورت می‌گیرد. این در حالی است که در بسیاری از کاربردها، مشاهدات انحراف زیادی از توزیع نرمال داشته و دارای ساختار چندنمایی و چوله هستند. از این‌رو، روش‌های سنتی تحلیل رگرسیونی به نتایج گمراه‌کننده‌ای منجر می‌شوند. در چنین شرایطی، توزیع‌های آمیخته متناهی می‌توانند به خوبی ناهمگنی و تمایز بین زیرجوامع را مدل‌بندی نمایند. علی‌رغم ویژگی‌های جالب توزیع چوله نرمال، برآورد پارامتر چولگی در رهیافت فراوانی‌گرا با مشکلات خاصی همراه است. با توجه به دشواری‌های ذاتی برآورد پارامترها در توزیع‌های آمیخته و این که خواص بهینه برآوردگرهای ماکسیم

جدول ۶. مقادیر معیارهای انتخاب مدل مختلف برای مدل چوله‌نرمال آمیخته بیزی با تعداد مولفه‌های متفاوت.

DIC	HQC	AIC	BIC	تعداد پارامتر	تعداد مولفه‌ها
۲۱۴/۶۷	۲۸۳/۲۴	۲۳۶/۷۴	۲۵۱/۰۰	۹	۲
۲۱۵/۷۵	۳۲۴/۶۱	۲۵۲/۲۷	۲۷۴/۴۴	۱۴	۳
۲۳۰/۸۵	۳۶۴/۱۸	۲۶۴/۰۱	۲۹۶/۰۹	۱۹	۴

جدول ۷. مقادیر برآورد پارامترهای دو مدل چوله‌نرمال آمیخته بیزی و نرمال آمیخته.

مدل آمیخته	β_{11}	β_{12}	β_{21}	β_{22}	ω_1^2	ω_2^2	α_1	α_2
چوله‌نرمال	۱/۶۲۴۰	۰/۱۹۰۴	۲/۱۲۳۹	۰/۶۰۸۹	۵۸/۴۲	۰/۳۴	۲/۲۸	۰/۰۹
نرمال	۲/۰۰۸۰	۰/۲۵۶۴	۲/۴۷۴۳	۰/۶۰۱۲	۳۷/۴۰	۸/۵۳	-	-

جدول ۸. مقادیر معیارهای انتخاب مدل مختلف برای مدل چوله‌نرمال آمیخته به همراه مقادیر متناظر برای مدل نرمال آمیخته.

مدل آمیخته	تعداد مولفه‌ها	تعداد پارامترها	BIC	AIC	HQC
چوله نرمال	۲	۹	۲۵۱/۰۰	۲۳۶/۷۴	۲۸۳/۲۴
	۲	۷	۲۵۸/۸۲	۲۴۷/۷۴	۲۸۳/۹۱
نرمال	۳	۱۱	۲۶۸/۵۲	۲۵۱/۱۰	۳۰۷/۹۳
	۴	۱۵	۲۶۵/۸۳	۲۴۲/۰۷	۳۱۹/۵۸

درست‌نمایی تنها در حالت مجانبی تحقق می‌یابند، رهیافت فراوانی‌گرایانه تحلیل رگرسیونی تحت توزیع‌های آمیخته از کارایی لازم برخوردار نیست. نتایج حاصل از مطالعه شبیه‌سازی و مثال کاربردی نشان می‌دهد که میزان کارایی رهیافت بیزی پیشنهادی، دست‌کم برای نمونه‌های کوچک و متوسط، به‌صورت قابل توجهی از مدل فراوانی‌گرا بیشتر است. افزون بر این، رهیافت بیزی پیشنهادی از نظر محاسباتی بسیار سراسر است بوده و مشکلات مختلف رهیافت فراوانی‌گرا مانند کران‌دار نبودن تابع درست‌نمایی، عدم همگرایی الگوریتم عدم همگرایی الگوریتم EM و بهینه نبودن نتایج برای نمونه‌های کوچک، در این رهیافت وجود ندارند.

تقدیر و تشکر

نویسندگان از داوران و سردبیر گرامی مجله برای ارائه پیشنهادهای سازنده، در راستای بهبود این پژوهش، سپاسگزاری می‌کنند.

مراجع

- Aitkin, M. and Wilson, G. T. (1980), Mixture Models, Outliers and EM Algorithm. *Technometrics*, **22**, 325-331.
- Arellano-Valle, R. B., Castro, L. M., Genton, M. G., and Gomez, H. M. (2008), Bayesian Inference for Shape Mixtures of Skewed Distributions with Application to Regression Analysis. *Bayesian Analysis*, **3**, 513-540.
- Azzalini, A. (1985), A Class of Distribution which Includes the Normal Ones, *Scandinavian Journal of Statistics*, **12**, 171-178.
- Barr, R., Donald, E. and Sherril. (1999), Mean and Variance of Truncated Normal Distribution, *The American Statistician*. **53**, 357-361.
- Bernardi, M., Maruotti, A. and Lea, P. (2012), Skew Mixture Models for Loss Distributions: A Bayesian Approach, *The American Statistician*. **53**, 357-361.
- Cancho, V. G., Dey, K. D., Lachos, V. H., and Andrade, M. (2010a), Bayesian Nonlinear Regression Models with Scale Mixtures of Skew Normal Distributions: Estimation and Case influence diagnostics *Computational Statistics and Data Analysis*, **55**, 588 -602.
- Cancho, V. G., Lachos, V. H., and Ortega, E. M. M. (2010b), A Nonlinear Regression Model with Skew-Normal Errors, *Statistical Papers*, **51**, 547-558.

- Chiogna, M. (2005), A Note on the Asymptotic Distribution of the Maximum Likelihood Estimator for the Scalar for the Skew-Normal Distribution, *Statistical Methods and Applications*, **14**, 331 -347.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977), Maximum Likelihood with Incomplete data via the E-M algorithm. *Journal of the Royal Statistical Society*, **39**, 1-38.
- DeSarbo, W. S. and Corn, W. L. (1988), A Maximum Likelihood Methodology for Clusterwise linear regression. *Journal of Classification*, **5**, 249-282.
- Dey, D. (2010), Estimation of the Parameters of Skew-Normal Distribution by Approximating the Ratio of the Normal Density and Distribution Functions, PhD Thesis in *Applied Statistics*, University of California.
- Figuerase, G., Puig, P. and Pewsey, A. (2007), Goodness-of-Fit Tests for the Skew Normal Distribution when the Parameters are Estimated from the Data, *Communications in Statistics: Theory and Methods*, **36**, 1735-1755.
- Gelman, A., Robert, G. and Gilks, W. (1996), Efficient Metropolis Jumping Rules, In *Bayesian Statistics 5* (Edited by J. M. Bernardo, J. O. Berger, A. P. Dawid and A. F. M. Smith), *Oxford University Press*, New York.
- Gelman, A., Jakulin, M. Grazia Pittau, A. and Su, A. (2008), A Weakly Informative Default Prior Distribution for Logistic and Other Regression Models. *The Annals of Applied Statistics*, **4**, 1360-1383.
- Gelman A. and Rubin D. B. (1992), Inference from Iterative Simulation Using Multiple Sequences, *Statistical Science*, **7**, 457-511.
- Genton, M. G. (2004), *Skew Elliptical Distribution and Their Applications: A Journey Beyond Normality*, USA, Chapman Hall/CRC.
- Jasra, A., Holmes, C. C. and Stephens, D. A. (2000), Markov Chain Monte Carlo Methods and the Label Switching Problem in Bayesian Mixture Modeling. *Statistical Science*, **20**, 1, 50-67.
- 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.
- Lachos, V., Ghosh, P. and Arellano-Valle, R. (2010). Likelihood Based Inference for Skew Normal/Independent Linear Mixed Model. *Statistica Sinica*, **20**, 303-322.
- Lin, T. I., Lee, J. C. and Yen, S. Y. (2007), Finite Mixture Modeling Using the Skew Normal Distribution. *Statistica Sinica*, **17**, 909-927.

- Liu, M., Hancock, G. R., and Haring, J. R. (2011), Using Finite Mixture Modeling to Deal with Systematic Measurement Error: A Case Study. *Journal of Modern Applied Statistical*, **10**, 249-261.
- Bernardi, M., Mariuotti, A. and Lea, P. (2012), Skew Mixture Models for Loss Distributions: A Bayesian Approach, *Insurance: Mathematics and Economics*, **51**, 3, 617-623.
- Papastamoulis, P. (2015), LabelSwitching: An R Package for Dealing with the Label Switching Problem in MCMC Outputs, *Journal of Statistical Software*, **vv**, Code Snippet II, <http://www.jstatsoft.org>.
- Pewsey, A. (2000), Problems of Inference for Azzalini Skew - Normal Distribution. *Journal of Applied Statistics.*, **27**, 859-870.
- Quandt, R. E. (1972), A New Approach to Estimating Switching Regressions. *Journal of the American Statistical Association*, **67**, 306-310.
- Quandt, R. E., and Ramsey, J. B. (1978), Estimating Mixtures of Normal Distributions and Switching Regressions, *Journal of the American Statistical Association*, **73**, 730-738.
- R Core Team (2013), R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, ISBN 3-900051-07-0, URL <http://www.R-project.org>.
- Plummer, M., Best, N. Cowles, K. and Vines, K. (2006), CODA: Convergence Diagnosis and Output Analysis for MCMC, *R News*, **6**, 1, 7-11.
- Richardson, S. and P. J. Green (1997), On Bayesian Analysis of Mixtures with an Unknown Number of Components (With Discusson). *Journal of the Royal Statistical Society, Series B*, **59**, 4, 731-792.
- Sartori, M. (2006), Bias Prevention of Maximum Likelihood Estimation: Skew-Normal and t Distribution. *Journal of Statistical Planning and Inference*, **136**, 4259-4275.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P. and Van Der Linde, A. (2002), Bayesian Measures of Model Complexity and Fit, *Journal of Royal Statistical Society, Series B*, **59**, 731-792.
- Stephens, M. (2000), Dealing with Label Switching in Mixture Models, *Journal of Royal Statistical Society, Series B*, **82**, 4, 795-809.
- Turner, T. R. (2000), Estimating the Propagation Rate of a Viral Infection of Potato Plants via Mixtures of Regressions, *Journal of Applied Statistics*, **49**, 371-384.