

نتایج جدیدی درباره‌ی شناساپذیری مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی

الهام تبریزی، احسان بهرامی سامانی^۱ و ناصح جعفری

بخش آمار، دانشگاه شهید بهشتی

تاریخ دریافت: ۱۳۹۳/۷/۱۷ تاریخ پذیرش: ۱۳۹۴/۴/۱۰

چکیده: شناساپذیری یکی از ویژگی‌های لازم برای کفایت یک مدل آماری است. وقتی مدلی شناساپذیر نباشد، با هیچ اندازه‌ای از نمونه، نمی‌توان پارامتر حقیقی مدل را تعیین کرد. در این مقاله، مروری بر مفهوم مشهور شناساپذیری و ویژگی‌های آن شده است. به‌علاوه از آنجایی که مشکل شناسا-ناپذیری در مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی بسیار رایج است، تمرکز اصلی ما بر روی این‌گونه از مدل‌ها بوده است. از سوی دیگر، معمولاً نرم‌افزارهای آماری، بعد از برازش مدل شناساناپذیر در آن‌ها، اشاره‌ای به این مسئله نکرده و خروجی‌های غیر معتبر ارائه می‌دهند. بنابراین یافتن راهی برای بررسی شناساپذیری مدل، قبل از برازش آن، خالی از فایده نخواهد بود. در این راستا، قضایای جدیدی در رابطه با شناساپذیری مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی بیان شده است. همچنین برای تشریح سودمندی قضایای مطرح‌شده، چند مطالعه‌ی شبیه‌سازی روی مدل‌های شناساناپذیر خطی تعمیم‌یافته و خطی تعمیم‌یافته با اثرهای تصادفی انجام شده و مشکل‌های حاصل از برازش آن‌ها مورد بررسی قرار گرفته است.

واژه‌های کلیدی: شناساپذیری، مدل‌های آمیخته‌ی خطی، مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی.

رده‌بندی ریاضی (۲۰۱۰): ۶۲J۰۵، ۶۲J۱۲.

۱- مقدمه

در استنباط‌های آماری، یک مجموعه از داده‌ها، مشاهده‌هایی از یک آزمایش تصادفی تعریف شده روی فضای احتمال (Ω, F, P) است که اندازه‌ی احتمال P ، جامعه نامیده می‌شود [۱]. هر مدل آماری، مجموعه‌ای از فرض‌ها را روی جامعه‌ی P اعمال می‌کند و بیان می‌دارد

۱- آدرس الکترونیکی نویسنده مسئول مقاله: ehsan_bahrami_samani@yahoo.com

که جامعه‌ی P به خانواده‌ای پارامتری مثل $\{P_{\theta} : \theta \in \Theta\}$ تعلق دارد. شناساپذیری، یک ویژگی در مور توزیع احتمال مشاهده‌ها و یا مدل است که در صورت برقرار نبودن این ویژگی، برآورد پارامترهای مدل و تحلیل آن بی‌معنی خواهد بود. این ویژگی زمانی برقرار است که نگاشت $\theta \rightarrow P_{\theta}$ ، نگاشتی یک‌به‌یک است. در واقع این خاصیت، یکتایی توزیع به‌عنوان تابعی از پارامترهای مدل را بیان می‌کند. از سوی دیگر، برآورد پارامترهای هر مدل شناساپذیر، امری غیرمعقول بوده و در صورت برآورد کردن آن‌ها حتی با در دسترس داشتن تعداد نامتناهی از مشاهده‌ها (کل جامعه)، به مقدار حقیقی پارامتر مدل دست نخواهیم یافت.

بررسی شناساپذیری مدل با استفاده از تعریف اصلی آن و قبل از برازش مدل، توسط [۲] انجام شده است. یکی از شایع‌ترین حالت‌هایی که با مشکل شناساپذیری مدل مواجه می‌شویم، هنگام برازش مدل‌های آمیخته‌ی خطی است. محقق‌های زیادی مانند [۳ و ۴] مدل‌های آمیخته‌ی خطی را مورد مطالعه قرار داده‌اند.

شناساپذیری مدل‌های آمیخته‌ی خطی وقتی که یکی از اثرهای تصادفی و خطاها باهم آمیخته شده‌اند نیز به‌صورت مختصر توسط [۵] بررسی شده است. همچنین در مدل‌های آمیخته‌ی خطی، می‌توان ساختارهای متفاوتی برای ماتریس کوواریانس پاسخ در نظر گرفت. این ماتریس کوواریانس، حاصل مجموع دو ماتریس کوواریانس مربوط به بردار اثرهای تصادفی و خطاها خواهد بود. از این‌رو، انتخاب ساختار کوواریانس مناسب برای بردار اثرهای تصادفی و خطاها امری اجتناب‌ناپذیر است. لیست جامعی از این ساختارها توسط [۶ و ۷] ارائه شده است. در این مدل‌ها، حتی اگر ماتریس کوواریانس پاسخ، بیش پارامتری شده نباشد (بیش پارامتری شدن ماتریس کوواریانس به این معنا است که تعداد پارامترهای نامعلوم موجود در این ماتریس بیش‌تر از تعداد عناصر قطری و بالا مثلی آن شود)، آنگاه به دلیل تلفیقی بودن این ماتریس از دو ماتریس دیگر، ممکن است مشکل شناساپذیری حاصل شود. در صورت برازش چنین مدل‌هایی، نرم‌افزارها، خروجی‌های قابل‌اطمینانی برای برآورد پارامترها نداشته و معمولاً اشاره‌ای به این مشکل نخواهند کرد.

لازم به ذکر است، مطالعات انجام‌شده توسط محققان علوم آمار کاربردی به بررسی شناساپذیری در مدل‌های آمیخته‌ی خطی خلاصه‌شده است در صورتی‌که در این مقاله به بررسی شناساپذیری در مدل‌های آمیخته‌ی خطی تعمیم‌یافته با اثرهای تصادفی پرداخته می‌شود که تاکنون مورد توجه محققان دیگر قرار گرفته نشده است. قضایای جدید ارائه‌شده در این مقاله، اولاً شناساپذیری مدل توأم خطی تعمیم‌یافته با اثرهای تصادفی را موقوف به شناساپذیری مدل تکی کرده‌اند و ثانیاً شروط لازم برای شناساپذیری پارامترها در این مدل‌ها را ارائه می‌دهند. سابقاً شروط لازم شناساپذیری مدل‌های آمیخته‌ی خطی، توسط [۸] بررسی شده است اما موارد استفاده از آن‌ها، محدود به همین مدل‌هاست. بنابراین در این مقاله سعی بر یافتن ارتباطی بین

فضای ارائه‌شده در [۸] با فضای معادل آن‌ها در مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی شده است.

همچنین این مقاله به این صورت سازمان‌دهی گردیده است که در بخش دوم مروری کوتاه بر مفهوم شناساپذیری شده است و تعریف خاصی از شناساپذیری، برگرفته از تعریف اصلی آن، در مدل‌های خطی تعمیم‌یافته بیان شده است. در ادامه به چگونگی بررسی شناساپذیری مدل‌های خطی تعمیم‌یافته با استفاده از این تعریف، پرداخته شده است. در بخش سوم مروری بر مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی خواهد شد و با بیان قضایا و برهان آن‌ها در مورد شناساپذیری در مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی، راهکارهای جدیدی برای بررسی شناساپذیری در این مدل‌ها ارائه می‌شود. در نهایت در بخش چهارم، مشکل‌های حاصل از برآزش مدل شناساناپذیر با استفاده از مطالعه‌های شبیه‌سازی، نشان داده می‌شود. مشکل‌های نرم‌افزاری ایجادشده در این مدل‌ها، علاوه بر مشکلاتی که در [۸] برای مدل‌های آمیخته‌ی خطی بیان شده است، موارد قابل تأمل دیگری مانند تغییر برآورد پارامترها با تغییر مقدار اولیه و یا حتی تعویض نرم‌افزار مورد استفاده را دارند. در این بخش به بررسی برخی از این مشکلات خواهیم پرداخت.

۲- شناساپذیری در مدل‌های خطی تعمیم‌یافته

۲-۱- مفهوم شناساپذیری

تعریف ۱. تعریف شناساپذیری بر اساس خانواده‌ی پارامتری

خانواده‌ی پارامتری $\{P_{\theta} : \theta \in \Theta\}$ شناساپذیر است اگر و تنها اگر $P_{\theta_1} = P_{\theta_2}$ دلالت بر $\theta_1 = \theta_2$ بنماید [۱]. به عبارت بهتر، این خانواده‌ی پارامتری و به اصطلاح، Θ شناساپذیر است اگر نگاشت $\theta \rightarrow P_{\theta}$ نگاشتی یک‌به‌یک در نظر گرفته شود

تعریف ۲. تعریف شناساپذیری بر اساس خانواده‌ی از توزیع‌های آماری

مدل آماری (۱)، تعریف‌شده به‌واسطه‌ی خانواده‌ی از توزیع‌ها بر پایه‌ی بردار تصادفی Y که تحت بردار پارامتری θ پارامتری شده است را در نظر بگیرید:

$$P = \{P_{\theta} : \theta \in \Theta\} \quad (1)$$

که در آن P_{θ} توزیع بردار تصادفی Y است که به‌واسطه‌ی بردار پارامتری θ پارامتری شده است. همچنین Θ اشاره به فضای پارامتر دارد. این مدل روی فضای پارامتری Θ ، شناساپذیر است اگر $P_{\theta_1} = P_{\theta_2}$ ، آن‌گاه $\theta_1 = \theta_2$ [۳]. تعریف شناساپذیری (۲)، تعریفی بسیار عمومی از

شناساپذیری به شمار می‌آید و در واقع، مفهوم شناساپذیری را بیان می‌کند. حال برای بررسی شناساپذیری پارامترهای یک مدل، اگر بتوان توابعی از پارامترهای مدل را چنان یافت که تابع چگالی احتمال مدل، به‌واسطه‌ی این توابع کاملاً تعیین شود، آن‌گاه می‌توان تعریف جدیدی برای شناساپذیری بر پایه‌ی تعریف (۳) بیان کرد. این روش برای ساده‌سازی بررسی شناساپذیری پارامترهای مدل، استفاده می‌شود.

۲-۲- چگونگی بررسی شناساپذیری مدل‌های خطی تعمیم‌یافته

مفهوم شناساپذیری با اندکی تغییرات در مدل‌های خطی تعمیم‌یافته نیز کاربرد دارد [۱۰]. در مدل‌های خطی تعمیم‌یافته نیز توزیع بردار پاسخ \mathbf{Y} به‌واسطه‌ی $E(\mathbf{Y})$ و $E(\mathbf{Y})$ و پارامترهای دیگری مثل ϕ که به $E(\mathbf{Y})$ مرتبط نیستند، کاملاً تعیین می‌شود. هر مدل خطی تعمیم‌یافته به‌صورت زیر در نظر گرفته می‌شود:

$$E(\mathbf{Y}) = \mathbf{h}(\mathbf{X}\boldsymbol{\beta}) \quad (۲)$$

که در آن $\mathbf{h}(\cdot)$ هر تابع یکنوا و مشتق‌پذیر است. معکوس این تابع، تابع ربط نامیده می‌شود. در مدل‌سازی‌ها، معمولاً ترجیح می‌دهند که از تابع‌های توزیع تجمعی به‌عنوان تابع $\mathbf{h}(\cdot)$ استفاده کنند اما هر تابع دیگری که در این شرایط صدق کند، نیز می‌تواند به‌عنوان تابع $\mathbf{h}(\cdot)$ مورد استفاده قرار گیرد [۱۱]. همچنین $\mathbf{Y} = (Y_1, \dots, Y_n)'$ بردار پاسخ، $\mathbf{X}_{n \times m}$ ماتریس طرح اثرهای ثابت و $\boldsymbol{\beta}$ بردار پارامتری اثرهای ثابت است. بنابراین می‌توان تعریف (۲) را به‌صورت زیر، برای مدل‌های خطی تعمیم‌یافته، تعمیم داد:

تعریف ۳. تعریف شناساپذیری در مدل‌های خطی تعمیم‌یافته

بردار پارامتری $\boldsymbol{\beta}$ شناساپذیر است اگر برای هر بردار $\boldsymbol{\beta}_1$ و $\boldsymbol{\beta}_2$ ، $\mathbf{h}(\mathbf{X}\boldsymbol{\beta}_1) = \mathbf{h}(\mathbf{X}\boldsymbol{\beta}_2)$ دلالت بر $\boldsymbol{\beta}_1 = \boldsymbol{\beta}_2$ کند. لازم به ذکر است اگر بردار $\boldsymbol{\beta}$ شناساپذیر باشد، آنگاه پارامتریدن $\mathbf{h}(\mathbf{X}\boldsymbol{\beta})$ (مدل) شناساپذیر است. به‌علاوه $\mathbf{g}(\boldsymbol{\beta})$ تابعی از R^p به R است که در آن p بعد بردار پارامتری $\boldsymbol{\beta}$ خواهد بود. شناساپذیر است اگر $\mathbf{h}(\mathbf{X}\boldsymbol{\beta}_1) = \mathbf{h}(\mathbf{X}\boldsymbol{\beta}_2)$ دلالت بر $\mathbf{g}(\boldsymbol{\beta}_1) = \mathbf{g}(\boldsymbol{\beta}_2)$ کند. اگر نوعی پارامتریدن شناساپذیر نباشد اما تابعی از پارامترها چنان وجود داشته باشد که شناساپذیر است، آنگاه گفته می‌شود این نوع پارامتریدن، شناساپذیر جزئی است.

نتیجه ۱. با توجه به این‌که در ویژگی‌های تابع برداری $\mathbf{h}(\cdot)$ ، معکوس‌پذیر بودن آن نیز لحاظ شده است، از تعریف (۳) می‌توان نتیجه گرفت که بردار پارامتری $\boldsymbol{\beta}$ در مدل‌های خطی تعمیم‌یافته‌ی (۲) شناساپذیر است اگر ماتریس طرح \mathbf{X} ، پر رتبه باشد.

قضیه ۱. تابع $g(\beta)$ در مدل خطی تعمیم‌یافته‌ی (۲) شناساپذیر است اگر و تنها اگر $g(\beta)$ تابعی از $h(X\beta)$ است [۱۰].

نتیجه ۲. تابع $g(\beta)$ در مدل خطی تعمیم‌یافته‌ی (۲) شناساپذیر است اگر و تنها اگر $g(\beta)$ تابعی از $X\beta$ است.

۲-۳- چند مثال در شناساپذیری مدل‌های خطی تعمیم‌یافته

در این قسمت به بررسی شناساپذیری در مدل‌هایی با پاسخ‌های گسسته‌ی دودویی همچون رگرسیون لوژستیک و پروبیت در دو حالت با وجود مفهوم متغیر پنهان و بدون وجود مفهوم متغیر پنهان پرداخته می‌شود و در صورت شناساپذیری مدل راه‌کارهایی برای رفع مشکل بیان خواهد شد.

مدل (۲) را مدنظر قرار دهید. هرگاه متغیر پاسخ Y_i در این مدل، متغیر گسسته‌ی دودویی باشد، آنگاه مدل حاصل، رگرسیون دودویی خواهد بود. بنابراین در این مدل‌ها، $E[Y_i | x_i] = h(x_i'\beta)$ که در آن بردار متغیرهای تبیینی مربوط به مشاهده i ام و در واقع سطر i ام از ماتریس طرح X است. به‌طور معمول، مؤلفه‌ی اول بردار x_i ، ۱ و مؤلفه‌ی اول بردار پارامتری β ، پارامتر عرض از مبدأ در نظر گرفته می‌شود. به‌علاوه $h(\cdot)$ هر تابع پیوسته اکیداً صعودی و دو بار مشتق‌پذیر است. برای مطالعه بیشتر در مورد شرایط این تابع به [۳] مراجعه کنید.

۲-۳-۱- شناساپذیری در مدل رگرسیون لوژستیک

اگر در مدل $E[Y_i | x_i] = h(x_i'\beta)$ ، تابع $h(\cdot)$ ، تابع توزیع تجمعی لوژستیک با پارامتر مکان صفر و پارامتر مقیاس یک برگزیده شود، آنگاه مدل حاصل، مدل رگرسیون لوژستیک نامیده می‌شود. از سوی دیگر، برای استفاده از تعریف (۳) و قضیه‌ی (۱) مناسب است اشاره‌ای به تابع برداری $h(\cdot)$ بنماییم. توجه کنید که در این مدل:

$$\pi(x_i) = E[Y_i | x_i] = \frac{\exp\{x_i'\beta\}}{1 + \exp\{x_i'\beta\}}, \quad i = 1, \dots, n.$$

بنابراین $h(s) = \frac{e^s}{1+e^s}$ و تابع برداری $h(s)$ تابعی است که مؤلفه‌ی i ام آن به‌صورت $h(s_i)$ است که در آن مؤلفه‌ی i ام بردار s خواهد بود. در نتیجه:

$$E[Y | X] = h(X\beta)$$

که در آن $\mathbf{Y} = (Y_1, \dots, Y_n)'$ و $\mathbf{h}(\cdot)$ تابع برداری از R^n به R^n تعریف شده است. حال با استفاده از تعریف (۲) می‌توان شناساپذیری در مدل رگرسیون لوژستیک را به صورت زیر مورد بررسی قرار داد:

ابتدا لگاریتم تابع درست‌نمایی بردار پاسخ \mathbf{Y} به صورت زیر به دست می‌آید:

$$\ln f_Y(\mathbf{y}; \boldsymbol{\beta}) = \mathbf{y}'\mathbf{X}\boldsymbol{\beta} - \sum_{i=1}^n \ln(1 + e^{x_i'\boldsymbol{\beta}}).$$

حال به راحتی می‌توان ثابت کرد که پر رتبه بودن ماتریس طرح \mathbf{X} شرط لازم و کافی برای شناساپذیری بردار پارامتری $\boldsymbol{\beta}$ است. فرض کنید ماتریس طرح \mathbf{X} پر رتبه نیست. در این صورت بردار $\mathbf{a} \neq \mathbf{0}$ چنان وجود دارد که

$$\mathbf{X}\mathbf{a} = \begin{bmatrix} \mathbf{x}'_1\mathbf{a} \\ \vdots \\ \mathbf{x}'_n\mathbf{a} \end{bmatrix} = \mathbf{0}$$

بنابراین به ازای هر بردار پارامتری $\boldsymbol{\beta}$ ، بردار پارامتری $\boldsymbol{\beta}^* = \boldsymbol{\beta} + \mathbf{a}$ چنان وجود دارد که به ازای هر \mathbf{y} موجود در تکیه‌گاه، $f_Y(\mathbf{y}; \boldsymbol{\beta}) = f_Y(\mathbf{y}; \boldsymbol{\beta}^*)$ در صورتی که $\boldsymbol{\beta} \neq \boldsymbol{\beta}^*$.

از طرف دیگر اگر بردار پارامتری $\boldsymbol{\beta}$ شناساپذیر در نظر گرفته نشود، آن‌گاه به ازای هر بردار پارامتری $\boldsymbol{\beta}_1$ ، بردار پارامتری $\boldsymbol{\beta}_2$ چنان وجود دارد که اولاً مخالف با $\boldsymbol{\beta}_1$ است و دوماً $f_Y(\mathbf{y}; \boldsymbol{\beta}_1) = f_Y(\mathbf{y}; \boldsymbol{\beta}_2)$ و این یعنی $\mathbf{X}(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2) = \mathbf{0}$ ؛ بنابراین ماتریس طرح \mathbf{X} پر رتبه نیست. همچنین با استفاده از تعریف (۳)، برای هر دو بردار $\boldsymbol{\beta}_1$ و $\boldsymbol{\beta}_2$ که $\mathbf{h}(\mathbf{X}\boldsymbol{\beta}_1) = \mathbf{h}(\mathbf{X}\boldsymbol{\beta}_2)$ ، نتیجه می‌شود که $\mathbf{X}\boldsymbol{\beta}_1 = \mathbf{X}\boldsymbol{\beta}_2$ زیرا $\ln \frac{\pi(\mathbf{x}_i)}{1 - \pi(\mathbf{x}_i)}; \mathbf{x}'_1\boldsymbol{\beta}_1 = \mathbf{x}'_1\boldsymbol{\beta}_2 = \ln \frac{\pi(\mathbf{x}_i)}{1 - \pi(\mathbf{x}_i)}$ ؛ بنابراین بردار پارامتری $\boldsymbol{\beta}$ در مدل لوژستیک شناساپذیر است اگر و تنها اگر ماتریس طرح \mathbf{X} پر رتبه در نظر گرفته شود.

۲-۳-۲- شناساپذیری در مدل رگرسیون پروبیت بر پایه متغیر پنهان

اگر در مدل $E[Y_i | \mathbf{x}_i] = \mathbf{h}(\mathbf{x}'_i\boldsymbol{\beta})$ به جای تابع $\mathbf{h}(\cdot)$ از Φ که تابع توزیع تجمعی نرمال استاندارد است، استفاده شود، مدل پروبیت حاصل می‌شود. این تابع زمانی در تحلیل پاسخ‌های دودویی به کار می‌رود که قضاوت تحلیل‌گر بر اساس یک متغیر پنهان باشد.

بررسی شناساپذیری مدل رگرسیون پروبیت بر اساس تعریف (۲) به صورت زیر انجام می‌شود:

فرض کنید متغیر تصادفی گسسته دودویی Y_i بر پایه‌ی متغیر پنهان Y_i^* ، به این صورت تعریف شده باشد:

$$Y_i = \begin{cases} 1 & Y_i^* > 0 \\ 0 & Y_i^* \leq 0 \end{cases} \Rightarrow Y_i \sim \text{Ber}(P(Y_i^* > 0)).$$

اگر برای Y_i^* ، مدلی رگرسیونی به صورت زیر در نظر گرفته شود:

$$Y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2),$$

که در آن بردار متغیرهای تبیینی و بردار پارامتری اثرهای ثابت است، آنگاه در این مدل بردار پارامتر یعنی $\boldsymbol{\theta} = (\boldsymbol{\beta}', \sigma^2)'$ شناساپذیر نیست. برای بررسی این موضوع به [۱۲] مراجعه کنید. بررسی شناساپذیری مدل رگرسیون پروبیت بر اساس تعریف (۳) به صورت زیر انجام می‌شود: مدل رگرسیون پروبیت به صورت زیر بازنویسی می‌شود:

$$E[Y | \mathbf{X}] = \mathbf{h}_{\sigma^2}(\mathbf{X}\boldsymbol{\beta}),$$

که در آن تابع برداری $\mathbf{h}_{\sigma^2}(\cdot)$ تابعی است که هر مؤلفه‌ی آن، تابع توزیع تجمعی نرمال با میانگین صفر و واریانس σ^2 است. از سوی دیگر تابع $\mathbf{h}(\cdot)$ معرفی شده در مدل (۲) باید تابعی معلوم باشد؛ بنابراین مجدداً مدل به صورت زیر بازنویسی می‌شود:

$$E[Y | \mathbf{X}] = \mathbf{h}(\mathbf{X}\boldsymbol{\beta}) = \Phi\left(\mathbf{X} \frac{\boldsymbol{\beta}}{\sigma}\right)$$

که در آن تابع برداری $\mathbf{h}(\cdot)$ ، تابعی است که هر مؤلفه‌ی آن، تابع توزیع تجمعی نرمال استاندارد خواهد بود. حال طبق تعریف (۳) پر رتبه بودن ماتریس طرح \mathbf{X} شرط لازم و کافی برای شناساپذیری بردار پارامتری $\boldsymbol{\beta}^* = \frac{\boldsymbol{\beta}}{\sigma}$ است. توجه کنید که توزیع بردار تصادفی \mathbf{Y} متعلق به خانواده‌ی نمایی بوده و به واسطه‌ی دو گشتاور اول تعیین می‌شود و تحت همین شرط است که مجاز به استفاده از تعریف (۳) خواهیم بود. برای رفع این مشکل که $\boldsymbol{\beta}^*$ شناساپذیر است و نه $\boldsymbol{\beta}$ و σ (بنابراین نرم‌افزارها خروجی معقول را برای برآورد $\boldsymbol{\beta}^*$ نمایش خواهند داد) می‌توان قیدهای روی مدل اعمال کرد. برای مثال، قید $\sigma^2 = 1$ از متداول‌ترین این قیود است. لازم به ذکر است که اگر این رگرسیون بدون استفاده از متغیرهای پنهان برازش داده شود و در واقع از ابتدای مدل‌سازی، تابع $\mathbf{h}(\cdot)$ تابع توزیع نرمال استاندارد در نظر گرفته شود، آنگاه پر رتبه بودن ماتریس طرح \mathbf{X} شرط لازم و کافی برای شناساپذیری مدل و به تبع بردار پارامتر $\boldsymbol{\beta}$ خواهد

بود. حال، اگر متغیر تصادفی گسسته دودویی Y_i بر پایه‌ی متغیر Y_i^* به صورت زیر تعریف شده باشد:

$$Y_i = \begin{cases} 1 & Y_i^* > \theta \\ 0 & Y_i^* \leq \theta \end{cases}$$

که در آن θ پارامتری نامعلوم (پارامتر آستانه‌ای) است؛ بنابراین برای نمونه‌ای تصادفی به حجم n ، مدل به صورت زیر است:

$$E[\mathbf{Y} | \mathbf{X}] = \mathbf{h}_{\theta, \sigma^2}(\mathbf{X}\boldsymbol{\beta}), \quad (۴)$$

که در آن $\mathbf{h}_{\theta, \sigma^2}(\cdot)$ تابع برداری با مؤلفه‌های تابع توزیع نرمال با میانگین θ و واریانس σ^2 است. مجدداً مدل فوق را برای رسیدن به تابع $\mathbf{h}(\cdot)$ مطلوب و مستقل از پارامتر، بازنویسی می‌کنیم:

$$E[\mathbf{Y} | \mathbf{X}] = \mathbf{h}\left(\mathbf{X}\left(\frac{\boldsymbol{\beta}}{\sigma}\right) - \frac{\theta}{\sigma}\right). \quad (۵)$$

برای استفاده از تعریف (۳)، مدل (۵) مطلوب نیست؛ بنابراین تعریف می‌کنیم

$\mathbf{X}^* = [\mathbf{X} \quad \mathbf{1}_n]$ که در آن $\mathbf{1}_n$ بردار یک‌ها با طول n و $\boldsymbol{\beta}^* = \left(\frac{\boldsymbol{\beta}'}{\sigma}, \frac{\theta}{\sigma}\right)'$ مدل (۵)

به صورت $E[\mathbf{Y} | \mathbf{X}] = \mathbf{h}(\mathbf{X}^* \boldsymbol{\beta}^*)$ بازنویسی می‌شود. بنابراین مانند قبل اولاً باید، قیدی برای مشکل وجود σ در مدل اعمال کرد و ثانیاً برای شناساپذیری $\boldsymbol{\beta}^*$ حتی بعد از اعمال قید روی σ ، شرط پر رتبه بودن برای ماتریس طرح \mathbf{X}^* لازم خواهد بود و اگر بردار پارامتری $\boldsymbol{\beta}$ شامل عرض از مبدأ β_0 در نظر گرفته شود، آنگاه ستون اول ماتریس \mathbf{X} ، $\mathbf{1}_n$ و بنابراین ماتریس طرح \mathbf{X}^* پر رتبه نخواهد شد. راهی دیگر برای برطرف کردن این مشکل، برآزش مدل بدون عرض از مبدأ β_0 است؛ بنابراین در بردار $\boldsymbol{\beta}$ ، عرض از مبدأ β_0 حذف خواهد شد.

۳- شناساپذیری در مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی

در این بخش ابتدا مروری بر مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی خواهیم کرد و در ادامه با ارائه‌ی چند قضیه‌ی جدید در مورد شناساپذیری مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی، راهکارهای آسان‌تری را برای بررسی شناساپذیری در این مدل‌ها بیان می‌کنیم.

۳-۱- مروری بر مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی

مدل آمیخته‌ی خطی زیر را در نظر بگیرید:

$$\mathbf{Y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i, \quad i = 1, \dots, N \quad (۶)$$

که در آن، \mathbf{Y}_i بردار پاسخ $(n_i \times 1)$ مربوط به خوشه (فرد) i ام، \mathbf{X}_i ماتریس طرح اثرهای ثابت در خوشه i ام، $\boldsymbol{\beta}$ بردار ضرایب اثرهای ثابت، \mathbf{u}_i بردار اثرهای تصادفی و \mathbf{Z}_i ماتریس طرح مربوط به اثرهای تصادفی در خوشه‌ی i ام هستند. $\boldsymbol{\varepsilon}_i$ هم بردار خطای تصادفی خوشه i ام است. در این مدل $\mathbf{u}_i \sim (\mathbf{0}, \Sigma_{\mathbf{u}})$ و $\boldsymbol{\varepsilon}_i \sim (\mathbf{0}, \Sigma_{\boldsymbol{\varepsilon}_i})$ هستند و \mathbf{u}_1 تا \mathbf{u}_N و $\boldsymbol{\varepsilon}_1$ تا $\boldsymbol{\varepsilon}_N$ توأماً مستقل در نظر گرفته می‌شوند. به مدل (۶) به ازای همه i ها مدل توأم و به ازای یک i ثابت، مدل تکی می‌گویند. بررسی شناساپذیری چنین مدلی تحت ساختارهای رایج ماتریس کواریانس توسط [۸] صورت گرفته است. در واقع همان‌گونه که بردار پاسخ پیوسته‌ی \mathbf{Y}_i در مدل (۶)، می‌تواند توسط متغیرهای تبیینی و اثرهای تصادفی مدل شود و شناساپذیری پارامترهای آن مورد بررسی قرار گیرد، $\mathbf{Y}_i = [Y_{i1}, \dots, Y_{in_i}]'$ بردار پاسخ گسسته مربوط به فرد یا خوشه‌ی i ام نیز می‌تواند بر اساس متغیرهای تبیینی و اثرهای تصادفی مدل شده و شناساپذیری پارامترهای آن بررسی شود. در این راستا فرض کنید \mathbf{X}_i ، $\boldsymbol{\beta}$ ، \mathbf{Z}_i و \mathbf{u}_i بردارها و ماتریس‌های تعریف‌شده در رابطه‌ی (۶) هستند؛ بنابراین، هر مدل خطی تعمیم‌یافته‌ی اثرهای تصادفی با تابع ربط برداری $\mathbf{g}(\cdot)$ به صورت زیر خواهد بود:

$$\mathbf{g}(E(\mathbf{Y}_i | \mathbf{u}_i, \mathbf{X}_i)) = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{u}_i, \quad i = 1, \dots, N. \quad (۷)$$

مشخصه‌های عمومی مدل‌های خطی تعمیم‌یافته‌ی اثرهای تصادفی به صورت زیر است:

۱. متغیرهای پاسخ Y_{i1}, \dots, Y_{in_i} به شرط \mathbf{u}_i دوه‌دو مستقل هستند و توزیع آن‌ها عضوی از خانواده‌ی پراکندگی نمایی به صورت زیر است:

$$f_{Y_{ij}}(y_{ij} | \mathbf{u}_i) = \exp \left\{ \frac{(y_{ij} \theta_{ij}) - b(\theta_{ij})}{a(\phi)} + c(y_{ij}, \phi) \right\}$$

که در آن $a(\cdot)$ ، $b(\cdot)$ و $c(\cdot)$ تابع‌های حقیقی مقدار، ϕ و θ_{ij} به ترتیب پارامتر پراکنش و پارامتر مکان هستند [۱۱].

۲. گشتاورهای شرطی $\mu_{ij} = E[Y_{ij} | \mathbf{u}_i] = b'(\theta_{ij})$ و $v_{ij} = \text{Var}[Y_{ij} | \mathbf{u}_i] = \frac{b''(\theta_{ij})}{\phi}$ که در آن $b'(\cdot)$ و $b''(\cdot)$ مشتق‌های مرتبه‌ی اول و دوم تابع $b(\cdot)$ نسبت به θ_{ij} هستند در رابطه‌های زیر صدق می‌کنند:

$$\theta_{ij} = g(\mu_{ij}) = \mathbf{x}'_{ij}\beta + \mathbf{z}'_{ij}\mathbf{u}_i, \quad v_{ij} = v(\mu_{ij})\phi,$$

که در آن $g(\cdot)$ و $v(\cdot)$ به ترتیب تابع ربط (مؤلفه‌های تابع ربط برداری $g(\cdot)$) و تابع واریانس هستند (برای مطالعه بیشتر در این زمینه به [۱۱ و ۱۳] مراجعه کنید).

۳. دارای توزیع چندمتغیره مشخص است که معمولاً نرمال با بردار میانگین $\mathbf{0}$ و ماتریس کوواریانس Σ_u ، پارامتری شده به واسطه‌ی بردار پارامتری θ_u برای آن رایج است.

از سوی دیگر برای بررسی شناساپذیری، به تابع چگالی توأم برای آزمودنی i ام نیاز داریم؛ بنابراین طبق مشخصه‌ی ۱، این تابع به صورت زیر محاسبه می‌شود:

$$f_{Y_i}(y_{i1}, \dots, y_{in_i}; \beta, \phi, \theta_u) = \int \prod_{j=1}^{n_i} f_{Y_i}(y_{ij} | \mathbf{u}_i; \beta, \phi) f_{U_i}(\mathbf{u}_i; \theta_u) d\mathbf{u}_i \quad (۸)$$

که در آن وابستگی به \mathbf{x}_{ij} برای سادگی در نوشتار حذف شده است. در ضمن انتگرال فوق، یک انتگرال k گانه است.

۳-۲- چند نتیجه‌ی جدید درباره‌ی شناساپذیری مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی

در این زیر بخش، چند قضیه‌ی جدید در شناساپذیری مدل‌های خطی تعمیم‌یافته ارائه خواهد شد. با استفاده از این قضایا بررسی شناساپذیری در این مدل‌ها، ساده‌تر خواهد شد.

قضیه ۲. مدل توأم (۷) شناساپذیر است اگر و تنها اگر حداقل یکی از مدل‌های تکی آن شناساپذیر در نظر گرفته شود.

برهان. به پیوست مراجعه شود.

قضیه‌ی فوق بررسی شناساپذیری مدل توأم را به بررسی شناساپذیری مدل‌های تکی موکول می‌کند؛ بنابراین برای سهولت در نگارش، چون ما دائماً با یک i ثابت از مدل توأم، کار خواهیم کرد. زیروند i را از مدل (۷) برمی‌داریم و با مدل تکی کار می‌کنیم.

در ادامه قضیه‌ای مطرح خواهد شد که بیان آن مستلزم ارائه‌ی تعریفی برای شناساپذیری ماتریس کواریانس Σ_u بر پایه‌ی تعریف (۲) است. این تعریف توسط [۸] ارائه شده است.

قضیه ۳. پر رتبه بودن ماتریس طرح \mathbf{X} و شناساپذیری ماتریس کواریانس Σ_u در مدل تکی (۷)، شرطی لازم برای شناساپذیری بردار پارامتری θ است.

برهان. به پیوست مراجعه شود.

البته توجه به این نکته ضروری است که پارامتر پراکنش ϕ ، یک در نظر گرفته شده است و توزیع به شکل توزیع طبیعی خانواده‌ی نمایی تبدیل شده است.

ملاحظه ۱. یک شرط لازم برای شناساپذیری پارامترهای مدل، یکتایی دو گشتاور اول به‌عنوان تابعی از پارامترهای مدل است [۳].

$$E_{\theta_1}(\mathbf{Y}) = E_{\theta_2}(\mathbf{Y}), \text{ CoV}_{\theta_1}(\mathbf{Y}) = \text{CoV}_{\theta_2}(\mathbf{Y}) \Rightarrow \theta_1 = \theta_2.$$

توجه به این نکته ضروری است که این شرط لازم است و نه کافی. زمانی که توزیع متغیر پاسخ، به‌واسطه‌ی دو گشتاور اول کاملاً تعیین شود، این شرط لازم و کافی خواهد بود. به همین دلیل در مدل‌های رگرسیونی با فرض نرمالیتی، شرط ذکرشده، برای شناساپذیری پارامترهای مدل، لازم و کافی خواهد بود.

قضیه ۴. مدل (۷) را مدنظر قرار دهید. اگر به ازای هر بردار پارامتری θ_1 عضو مجموعه‌ی Θ ، بردار پارامتری θ_2 در مجموعه‌ی Θ و مخالف با θ_1 چنان وجود داشته باشد که تساوی

$$\left| \Sigma_{\theta_{u_1}}^{-1} \right| \mathbf{h}(\mathbf{X}\beta_{\theta_1} + \mathbf{Z}\Sigma_{\theta_{u_1}}^{-1}(\theta_{u_1})\mathbf{w}) = \left| \Sigma_{\theta_{u_2}}^{-1} \right| \mathbf{h}(\mathbf{X}\beta_{\theta_2} + \mathbf{Z}\Sigma_{\theta_{u_2}}^{-1}(\theta_{u_2})\mathbf{w}) \quad (9)$$

به ازای هر \mathbf{w} در R^n برقرار باشد، آن‌گاه مدل (۷) شناساناپذیر است.

برهان. به پیوست مراجعه شود.

نکته ۱. گاهی اوقات برخلاف مدل (۷)، بیش از یک مؤلفه‌ی تصادفی (\mathbf{u}) در مدل وجود دارد. نتایج بیان‌شده در این مقاله برای برخی از این گونه مدل‌ها نیز کاربرد دارند. برای مثال، مدل زیر با دو مؤلفه‌ی تصادفی را مدنظر قرار دهید:

$$\mathbf{g}(E(\mathbf{Y}_i | \mathbf{u}_{i1}, \mathbf{u}_{i2}, \mathbf{X}_i)) = \mathbf{X}_i\beta + \mathbf{Z}_i\mathbf{u}_{i1} + \mathbf{Z}_i\mathbf{u}_{i2}, \quad i = 1, \dots, N. \quad (10)$$

که در آن $\mathbf{u}_1 \sim (\mathbf{0}, \Sigma_{u_1})$ ، $\mathbf{u}_2 \sim (\mathbf{0}, \Sigma_{u_2})$. اگر ستون‌های ماتریس \mathbf{Z}_1 و ماتریس \mathbf{Z}_2 در مدل (۱۰) مستقل خطی باشند، آنگاه به‌راحتی مدل (۱۰) را به‌صورت زیر به مدل (۷) تبدیل کرده و از مطالب این مقاله استفاده می‌کنیم:

$$\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2), \quad \mathbf{u} = \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix}, \quad \Sigma_u = \begin{bmatrix} \Sigma_{u_1} & \Sigma_{u_1, u_2} \\ \Sigma_{u_2, u_1} & \Sigma_{u_2} \end{bmatrix}.$$

اما در حالتی که ستون‌های ماتریس \mathbf{Z}_1 و ماتریس \mathbf{Z}_2 مستقل خطی نباشند، تنها در حالتی که \mathbf{u}_1 و \mathbf{u}_2 مستقل هستند، می‌توان از این نتایج استفاده کرد. در این حالت اگر فرض کنیم

$\varepsilon = Z_p u_p$ ، آنگاه می‌توان از تمام قضایای [۸] که تحت ساختارهای رایج ماتریس کواریانس بیان شده است، برای بررسی شناساپذیری ماتریس کواریانس اثرهای تصادفی مدل استفاده کرد.

۴- چند مطالعه‌ی شبیه‌سازی

در این بخش، طی چند مطالعه شبیه‌سازی به بررسی اهمیت شناساپذیری پارامترهای مدل‌های خطی تعمیم‌یافته و مدل خطی تعمیم‌یافته با اثرهای تصادفی می‌پردازیم.

۴-۱- مطالعه‌ی شبیه‌سازی روی مدل‌های رگرسیون پروبیت بر پایه‌ی متغیر پنهان

در این بخش ابتدا، طی یک مطالعه‌ی شبیه‌سازی از مدل رگرسیون پروبیت بر پایه‌ی متغیر پنهان و بدون اثرهای تصادفی، به بررسی مشکل‌های حاصل‌شده از شناساپذیری این مدل پرداخته و سپس قیده‌های بیان‌شده برای شناساپذیری مدل را اعمال کرده و مجدداً به بررسی وضعیت مدل شناساشده و مقایسه‌ی آن با مدل قبل می‌پردازیم.

الف) مدل رگرسیون پروبیت بر پایه‌ی متغیر پنهان

مدل (۴)، حالتی که تنها یک متغیر تبیینی در مدل است، در نظر گرفته شود:

$$Y_i = \begin{cases} 1 & Y_i^* > \theta \\ 0 & Y_i^* \leq \theta \end{cases}, \quad Y_i^* = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n,$$

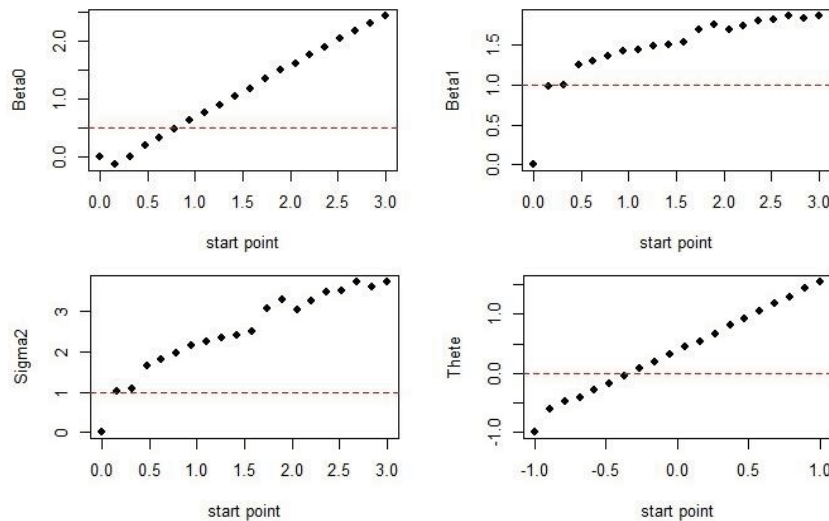
که در آن $\varepsilon_i \sim N(0, \sigma^2)$ ؛ بنابراین $Y_i \sim \text{Ber}(p(Y_i^* > \theta))$. بردار پارامتری $\beta' = (\beta_0, \beta_1, \sigma^2, \theta)$ شناساپذیر نبوده و تمام شش مشکل اشاره شده در بخش بحث و نتیجه‌گیری را خواهد داشت. به این منظور طی یک مطالعه‌ی شبیه‌سازی برای $n = 100$ ، ۱۰۰۰ بار از مدل فوق شبیه‌سازی کرده و تعداد دفعاتی که الگوریتم عددی برای محاسبه‌ی پارامترها همگرا نشده را مد نظر قرار داده‌ایم. برای شبیه‌سازی مدل فوق $\beta' = (0/5, 1, 1, 0)$ لحاظ شده است. در این مطالعه در هر ۱۰۰۰ دفعه شبیه‌سازی، الگوریتم همگرا شد و بنابراین نشانی از شناساپذیری به‌واسطه‌ی همگرا نشدن، در این مدل یافت نمی‌شود. همچنین برای بررسی خطای استانداردهای محاسبه شده برای پارامترهای مدل، جدول (۱) تهیه شده است. همان‌گونه که در جدول (۱) ملاحظه می‌شود، تقریباً برای همه‌ی پارامترها به‌جز β_1 ، چندک‌های ۵ درصد به بالای خطای استاندارد، مقدارهای بسیار بزرگ و غیرمعمولی نسبت به مقدار واقعی خود پارامتر دارند. برای پارامتر β_1 چندک ۵۰٪ به بالا نامعقول و غیرواقعی است.

در هر مدل خطی با پاسخ پیوسته یا گسسته، برآورد پارامترهای مدل نباید به مقدار اولیه‌ای که به‌عنوان نقطه‌ی ابتدایی به تابع برآورد کننده‌ی پارامترها داده می‌شود، وابسته باشد. در واقع

الگوریتم‌های عددی محاسبه‌کننده‌ی برآورد پارامترها در مدل‌های بدون مشکل، وابسته به مقادیر اولیه نیستند. به این منظور، یک‌بار از مدل (۴) شبیه‌سازی کرده و سپس طی ۲۰ دفعه تغییر مقادیرهای اولیه، پارامترهای مدل را برآورد کردیم. در شکل (۱) سعی شده است تغییر برآورد پارامترها با تغییر مقدار اولیه در تابع پروبیت ۱، نوشته‌شده با استفاده از بسته‌ی کمینه‌سازی لگاریتم تابع درست‌نمایی، نمایش داده شود. این دستور، برآورد پارامترهای مدل به روش ماکسیمم درست‌نمایی را در خروجی ارائه می‌دهد. دامنه‌ی تغییرات پارامترها β_0 ، β_1 ، σ^2 و θ به ترتیب $[0, 3]$ ، $[0, 3]$ ، $[0, 3]$ و $[-1, 1]$ لحاظ شده است. شکل (۱) تغییر نامعقول برآورد پارامترها متناسب با تغییر مقادیرهای اولیه در این مدل شناساناپذیر را نشان می‌دهد.

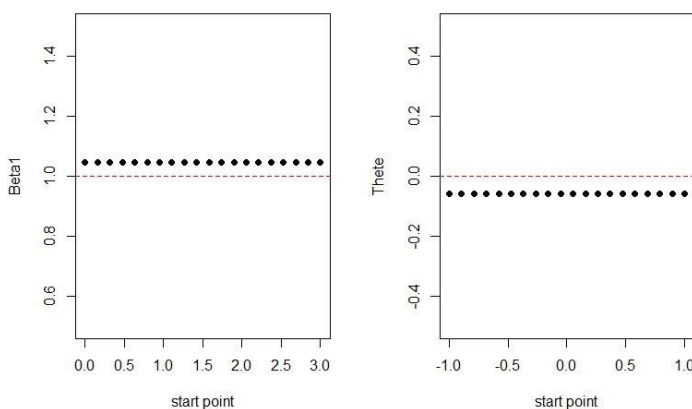
جدول ۱: چندک‌های خطای استاندارد برای پارامترهای مدل شناساناپذیر، به روش ماکسیمم درست‌نمایی طی هزار بار شبیه‌سازی

پارامتر	۰٪	۱٪	۵٪	۵۰٪	۶۰٪	۸۰٪	۹۰٪	۹۵٪	۹۹٪	۱۰۰٪
$\beta_0 = 0.5$	۲/۱۸۶	۶/۰۷	۱۵/۲۳	۴۲/۲۵	۵۱/۵۵	۸۱/۰۶	۱۱۷/۷۹	۱۸۳/۷۴	۳۹۵/۴۲	۵۷۶/۵۷
$\beta_1 = 1$	۰/۱۲	۲/۴۳	۴/۹۱	۲۱/۸	۲۷/۷	۴۵/۴۳	۶۶/۰۷	۹۸/۱۵	۲۱۱/۶۶	۳۳۳/۴۴
$\sigma^2 = 1$	۰/۰۳	۵/۲۷	۱۰/۶۷	۴۶/۱۹	۶۰/۹۶	۱۰۱/۳۹	۱۴۸/۱۶	۲۱۹/۷	۴۰۵/۷۷	۵۹۲/۸۲
$\theta = 0$	۱/۷۱	۲/۱۸۸	۷/۳۸	۳۳/۹۸	۴۳/۸۰	۷۸/۶۱	۱۰۸/۸۲	۱۸۱/۱۲۳	۳۹۹/۴۶	۶۲۱/۲۹



شکل ۱: نمودار پراکنش مقادیرهای اولیه تابع پروبیت ۱ در مقابل برآورد پارامترهای مدل شناساناپذیر

حال با اعمال قیدهای $\sigma^2 = 1$ و حذف β_0 از بردار پارامتری β ، مدل، شناساپذیر خواهد شد. همین مطالعه‌های شبیه‌سازی را برای مدل شناساشده انجام داده و نتایج در شکل (۲) و جدول (۲) نمایش داده شده‌اند. خطای استاندارد برای پارامترهای مدل شناساپذیر فوق، مقدارهای قابل قبول و معقولی اختیار می‌کنند. برآورد خطاهای استاندارد برای پارامتر β_1 در این مطالعه شبیه‌سازی متعلق به بازه‌ی $[0, 0/39]$ و برای پارامتر θ در ۹۹٪ اوقات $[0, 0/17]$ است.



شکل ۲: نمودار پراکنش مقدارهای اولیه تابع پروبیت ۱ در مقابل برآورد پارامترهای مدل شناساپذیر

جدول ۲: چندک‌های خطای استاندارد برای پارامترهای مدل شناساپذیر، به روش ماکسیمم درستمایی طی هزار بار شبیه‌سازی

پارامتر	۰٪	۱٪	۵٪	۵۰٪	۶۰٪	۸۰٪	۹۰٪	۹۵٪	۹۹٪	۱۰۰٪
$\beta_1 = 1$	۰	۰/۱۵	۰/۱۶	۰/۱۹	۰/۲۰	۰/۲۲	۰/۲۴	۰/۲۵	۰/۳	۰/۳۹
$\theta = 0$	۰	۰/۱۲	۰/۱۳	۰/۱۴۶	۰/۱۴۸	۰/۱۵۳	۰/۱۵۸	۰/۱۶	۰/۱۷	۳/۲

با وجود شناساپذیر شدن مدل پروبیت بعد از اعمال قیدهای فوق، یکسان بودن برآورد پارامترها با وجود تغییر در مقدار اولیه ورودی تابع پروبیت ۱، امری طبیعی است که شکل (۲) نیز به تأیید آن می‌پردازد.

(ب) مدل رگرسیون پروبیت بر پایه‌ی متغیر پنهان و با حضور اثر تصادفی

به مدل (۷) توجه کنید. مجدداً حالتی را که تنها یک متغیر تبیینی در مدل است، در نظر بگیرید:

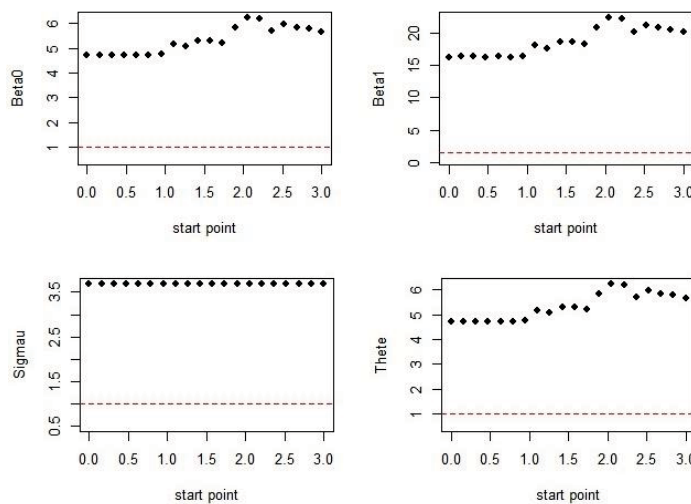
$$Y_{it} = \begin{cases} 1 & Y_{it}^* > \theta \\ 0 & Y_{it}^* \leq \theta \end{cases}, \quad Y_{it}^* = \beta_0 + \beta_1 x_{it} + u_i + \varepsilon_{it}, \quad i = 1, \dots, n, \quad t = 1, \dots, T,$$

که در آن $\varepsilon_{it} \sim N(0, 1)$ ، $u_i \sim N(0, \sigma_u^2)$. همچنین u_i ها و ε_{it} ها توأمأً مستقل هستند. طبق قضیه‌ی (۳) بردار پارامتری $\beta' = (\beta_0, \beta_1, \sigma_u^2, \theta)$ شناساپذیر نخواهد بود. برای بررسی مشکل‌های حاصل از برازش، مشابه بخش قبل، طی یک مطالعه‌ی شبیه‌سازی برای $n = 100$ ، ۱۰۰۰ بار از مدل فوق شبیه‌سازی کرده و تعداد دفعاتی که الگوریتم عددی برای محاسبه‌ی پارامترهای مدل، همگرا نشده است را مدنظر قرار داده‌ایم. در این مطالعه $\beta' = (1, 1/5, 1, 1)$ لحاظ شده است. این بار در ۱۰۰ مورد از ۱۰۰۰ بار، الگوریتم همگرا نشد و در بقیه موارد موجب نمایش خروجی شده که به دلیل شناسا نبودن مدل معتبر نیست، ولیکن نرم‌افزار هیچ اشاره‌ای به معتبر نبودن خروجی نمی‌کند. جدول (۳) چندک‌های خطای استانداردهای محاسبه‌شده برای پارامترهای مدل را نشان می‌دهد. چندک‌های انحراف استاندارد پارامترها، بیش از حد معمول بزرگ برآورد شده‌اند. پارامترهای β_0 و θ که دلیل ایجاد مشکل شناساپذیری مدل هستند از چندک یک درصد به بعد مقادیر نامعقول اتخاذ کرده‌اند.

شکل (۳) تغییر نامعقول برآورد پارامترها با تغییر مقادیر اولیه در این مدل شناساپذیر را نمایش می‌دهد. مدل‌های شناساپذیر تقریباً علائم مشابه دارند اما این که همواره این علامت‌ها در خروجی مشاهده شوند، قطعی نیست و همین مطلب ما را مشتاق به بررسی شناساپذیری مدل قبل از برازش آن می‌کند.

جدول ۳: چندک‌های خطای استاندارد برای پارامترهای مدل شناساپذیر، به روش ماکسیمم درست‌نمایی طی ۱۰۰۰ بار شبیه‌سازی

پارامتر	۰٪	۱٪	۵٪	۵۰٪	۶۰٪	۸۰٪	۹۰٪	۹۵٪	۹۹٪	۱۰۰٪
$\beta_0 = 1$	۰/۱۷	۲۶/۵۰	۱۳۱/۸۴	۲۹۱۰/۸۶	۵۳۴۸/۸	۵۴۳۱۷	۱۱۶۷۱۴	۲۳۲۸۹۱	۳۲۵۸۳۳	۳۴۹۰۶۹
$\beta_1 = 1/5$	۰/۶۸	۳/۴۵	۱۴/۵۰	۱۴۸۸۹	۲۸۰۰۲	۱۱۲۴۶۲	۳۳۳۹۵۶	۵۴۱۰۹۲	۷۰۶۸۰۱	۷۴۸۲۲۸
$\sigma_u^2 = 1$	۰/۶۸	۰/۹۱	۱/۸۱	۲۱۴۱۴	۴۵۳۸۳	۱۲۸۵۱۳	۴۶۱۴۳۶	۵۳۳۰۳۷	۵۹۰۳۱۸	۶۰۴۶۳۸
$\theta = 1$	۰/۱۷	۲۵/۸۱	۱۲۸/۳۶	۲۹۱۰	۵۳۴۸	۵۴۳۱۷	۱۱۶۷۱۴	۲۳۲۸۹۱	۳۲۵۸۳۳	۳۴۹۰۶۹



شکل ۳: نمودار پراکنش مقادیر اولیه تابع پروبیت ۲ در مقابل برآورد پارامترهای مدل شناساناپذیر

۵- بحث و نتیجه‌گیری

گاهی اوقات نرم‌افزارهای آماری در صورت برازش مدل‌های شناساناپذیر، پیام‌هایی مثل همگرا نبودن الگوریتم‌های عددی را گزارش می‌کنند. در ادامه برخی از مشکل‌های حاصل از برازش هر مدل شناساناپذیر ذکر می‌شود. **الف)** نمایش خطای استاندارد خیلی بزرگ و یا صفر بودن مقدار خطای استاندارد برای هر یک از پارامترهای مدل **ب)** ناتوان بودن نرم‌افزار در ساخت فاصله‌های اطمینان **ج)** فاصله‌های اطمینان ساخته‌شده غیرمعقول برای پارامترها. از سوی دیگر در برخی از اوقات، ممکن است هیچ‌کدام از مشکل‌های فوق در خروجی ظاهر نشوند و کاربر بدون اطلاع از این مسئله، از خروجی‌های غیر معتبر استفاده کرده و تفسیرهای نادرستی ارائه دهد. با این اوصاف، بررسی شناساپذیری هر مدل آماری، قبل از برازش آن، امری واجب خواهد بود. از سوی دیگر، اگر مدل خیلی پیچیده نباشد، آنگاه محاسبه‌ی مستقیم عناصر ماتریس کوواریانس پاسخ و استفاده از تعریف اصلی شناساپذیری، برای بررسی آن، ممکن است راهی پیش روی محقق قرار دهد. این کار توسط [۵] انجام‌شده است. بررسی مشکل‌های نرم‌افزاری حاصل از برازش مدل‌های شناساناپذیری توسط [۲ و ۶] صورت گرفته است. تمرکز اصلی [۶] بر روی مشکل‌های حاصل از شناساناپذیری در یکی از نرم‌افزارهای آماری بوده است، در صورتی‌که [۲] به مقایسه‌ی خروجی‌های حاصل از مدل‌های شناساناپذیر در چهار نرم‌افزار آماری پرداخته است.

در این مقاله، بر پایه‌ی تعریف اصلی شناساپذیری، قضایای دیگری که بررسی شناساپذیری مدل‌های خطی تعمیم‌یافته‌ی اثرهای تصادفی را ساده می‌کنند، بیان شد. موارد زیر برای

تحقیقات بیشتر پیشنهاد می‌شود: الف) ارائه شروط لازم و کافی در شناساپذیری پارامترهای ماتریس کوواریانس اثرهای تصادفی و اثرهای ثابت در مدل‌های خطی تعمیم‌یافته با اثرهای تصادفی: برای بررسی شروط کافی شناساپذیری در مدل‌های فوق، توجه به بحث زیر خالی از فایده نخواهد بود. در راستای بررسی کافی بودن شرط پر رتبه بودن ماتریس طرح و شناساپذیری ماتریس Σ_u در قضیه‌ی (۳)، باید به‌صورت زیر عمل کرد:

فرض می‌شود بردار پارامتری θ شناساپذیر نیست و یا به عبارتی به ازای هر بردار پارامتری θ_1 که عضوی از مجموعه‌ی Θ است، بردار پارامتری θ_2 در مجموعه‌ی Θ چنان وجود دارد که اولاً مخالف با θ_1 است و دوماً $f_Y(y; \theta_1) = f_Y(y; \theta_2)$. توجه کنید که $\{\theta_1 \neq \theta_2\}$ معادل با $\{\beta_1 \neq \beta_2, \theta_{u1} = \theta_{u2}\}$ ، $\{\beta_1 = \beta_2, \theta_{u1} \neq \theta_{u2}\}$ و یا $\{\beta_1 \neq \beta_2, \theta_{u1} \neq \theta_{u2}\}$ است. در نتیجه:

$$\int \prod_{j=1}^n f_{Y_j}(y_j | u; \beta_1) f_u(u; \theta_{u1}) du = \int \prod_{j=1}^n f_{Y_j}(y_j | u; \beta_2) f_u(u; \theta_{u2}) du. \quad (11)$$

از رابطه‌ی فوق، نتیجه می‌شود:

$$\int K(u, \beta_1, \theta_{u1}) - K(u, \beta_2, \theta_{u2}) du = 0$$

اگر تابع تحت انتگرال، به ازای هر بردار u (عضوی از R^k است)، همواره مثبت یا همواره منفی در نظر گرفته شود (و یا به ازای هر مجموعه‌ی B عضو سیگما میدان برل R^k ، مقدار انتگرال فوق روی B برابر صفر باشد)، آنگاه این تابع تقریباً همه‌جا برابر صفر می‌شود و یا به عبارتی $K(u, \beta_1, \theta_{u1}) = K(u, \beta_2, \theta_{u2})$. حال سه حالت پیش می‌آید: آ) فرض کنید $\{\beta_1 \neq \beta_2, \theta_{u1} = \theta_{u2}\}$. در نتیجه $f_Y(y | u; \beta_1) = f_Y(y | u; \beta_2)$. مجدداً یادآوری می‌کنیم که توزیع بردار تصادفی $Y | u$ عضوی از خانواده‌ی نمایی بوده و به‌واسطه‌ی دو گشتاور اول کاملاً تعیین می‌شود؛ بنابراین از (۷) داریم:

$$E_{\beta_1}(Y | u, X) = h(X\beta_1 + Zu) = h(X\beta_2 + Zu) = E_{\beta_2}(Y | u, X)$$

حال از معکوس‌پذیر بودن تابع $h(\cdot)$ نتیجه می‌شود که $X(\beta_2 - \beta_1) = 0$ و عبارت قبل نشان‌دهنده‌ی پر رتبه نبودن ماتریس طرح X است؛ اما این مسئله که تابع تحت انتگرال مذکور، تقریباً همه‌جا صفر می‌شود یا خیر، باز است. ب) فرض کنید که $\{\beta_1 = \beta_2, \theta_{u1} \neq \theta_{u2}\}$. در این صورت $f_U(u; \theta_{u1}) = f_U(u; \theta_{u2})$ و این معادل با شناساپذیری ماتریس کواریانس اثرهای تصادفی است. پ) فرض کنید $\{\beta_1 \neq \beta_2, \theta_{u1} \neq \theta_{u2}\}$. در این حالت از تبدیل $w = \Sigma_u^{-1/2} u$ استفاده می‌کنیم. بنابراین w

دارای توزیع نرمال با بردار میانگین صفر و ماتریس کواریانس \mathbf{I} (همانی) است. حال رابطه‌ی (۱۱) به صورت زیر بازنویسی می‌شود:

$$\begin{aligned} & \int \left| \Sigma^{\dagger}(\theta_{u1}) \right| f_Y(y | \Sigma^{\dagger}(\theta_{u1})\mathbf{w}; \beta_1, \theta_{u1}) f_w(\mathbf{w}) d\mathbf{w} \\ &= \int \left| \Sigma^{\dagger}(\theta_{u2}) \right| f_Y(y | \Sigma^{\dagger}(\theta_{u2})\mathbf{w}; \beta_2, \theta_{u2}) f_w(\mathbf{w}) d\mathbf{w}. \end{aligned} \quad (۱۲)$$

مجدداً اگر توابع تحت انتگرال باهم برابر شوند، یعنی:

$$\left| \Sigma^{\dagger}(\theta_{u1}) \right| f_Y(y | \Sigma^{\dagger}(\theta_{u1})\mathbf{w}; \beta_1, \theta_{u1}) = \left| \Sigma^{\dagger}(\theta_{u2}) \right| f_Y(y | \Sigma^{\dagger}(\theta_{u2})\mathbf{w}; \beta_2, \theta_{u2}),$$

آنگاه شناسانپذیری هر دو بردار پارامتری β و θ_u حاصل نخواهد شد؛ بنابراین کافی بودن شرایط قضیه‌ی (۳) برای بررسی شناسانپذیری، مسئله‌ای باز است.

ب) بررسی شناسانپذیری پارامترهای مدل‌های هم‌زمان آمیخته با پاسخ‌های گسسته و پیوسته تحت ساختارهای گوناگون ماتریس کواریانس

پ) گاهی اوقات برخلاف مدل تکی (۶)، بیش از دو مؤلفه‌ی تصادفی $(\varepsilon, \mathbf{u})$ در مدل وجود دارند. بررسی شناسانپذیری در این مدل‌ها، مسئله‌ای باز است.

ت) بررسی شناسانپذیری در مدل‌های خطی تعمیم‌یافته‌ی چند متغیره و مدل‌های خطی تعمیم‌یافته‌ی چند متغیره با اثرهای تصادفی

۶- پیوست

برهان قضیه (۲): (\Leftarrow) اگر مدل توأم شناسانپذیر در نظر گرفته شود، آن‌گاه بدیهی است که حداقل یکی از مدل‌های تکی شناسانپذیر است. در غیر این صورت اگر فرض کنیم همه مدل‌های تکی شناسانپذیر نیستند، $(\theta = (\beta', \theta_u)')$ آن‌گاه $\theta_1 \neq \theta_2$ چنان وجود دارند که در هر مدل تکی، $f_{Y_i}(y_i; \theta_1) = f_{Y_i}(y_i; \theta_2)$ ؛ بنابراین بدیهی است که تابع چگالی توأم که حاصل ضرب توابع چگالی کناری است نیز به ازای این دو بردار پارامتری نابرابر، برابر خواهد شد.

(\Rightarrow) حال اگر مدل تکی \hat{I} ام شناسانپذیر در نظر گرفته شود، آن‌گاه $f_{Y_i}(y_i; \theta_1) = f_{Y_i}(y_i; \theta_2)$ دلالت بر $\theta_1 = \theta_2$ می‌کند. از آنجایی بردار پارامتر در مدل‌های تکی، وابسته به اندیس نیست، بنابراین مدل تکی \hat{I} ام نیز شناسا خواهد بود. پس مدل توأم شناسانپذیر است (باید توجه شود که اگر مدل تکی \hat{I} ام شناسانپذیر در نظر گرفته شود، به ازای "هر" $\theta_1 \neq \theta_2$ نتیجه می‌شود $f_{Y_i}(y_i; \theta_1) \neq f_{Y_i}(y_i; \theta_2)$).

برهان قضیه (۳): آ) فرض کنید ماتریس \mathbf{X} پر رتبه نیست، در این صورت بردار $\beta^* \neq 0$ چنان وجود دارد که $\mathbf{X}\beta^* = 0$ و یا به ازای هر z از 1 تا n ، $\mathbf{x}'_j \beta^* = 0$ ، که در آن \mathbf{x}'_j سطر z ام ماتریس طرح \mathbf{X} است. حال به تابع چگالی توأم (۸) توجه کنید. از آنجایی که $Y_j | \mathbf{u}$ به ازای z از 1 تا n ، عضوی از خانواده‌ی پراکندگی نمایی است؛ بنابراین توزیع آن به‌واسطه‌ی دو گشتاور اول مشخص خواهد شد. با توجه به (۷) می‌دانیم $E(Y_j | \mathbf{u}, \mathbf{x}_j) = h(\mathbf{x}'_j \beta + \mathbf{z}'_j \mathbf{u})$. پس اگر به ازای هر بردار پارامتری θ_1 قرار دهیم $\theta_1 = ((\beta + \beta^*)', \theta'_{u1})'$ آنگاه $f_Y(\mathbf{y}; \theta_1) = f_Y(\mathbf{y}; \theta_2)$

ب) فرض کنید ماتریس Σ_{u1} شناسا نیست؛ بنابراین به ازای هر بردار پارامتری θ_{u1} عضو مجموعه‌ی Θ ، بردار پارامتری θ_{u2} در مجموعه‌ی Θ_{u1} و مخالف با θ_{u1} چنان وجود دارد که $\Sigma(\theta_{u1}) = \Sigma(\theta_{u2})$. حال چون توزیع احتمال بردار اثرهای تصادفی \mathbf{u} نرمال در نظر گرفته شده است، توزیع آن به‌واسطه‌ی دو گشتاور اول مشخص خواهد شد و بنابراین $f_U(\mathbf{u}; \theta_{u1}) = f_U(\mathbf{u}; \theta_{u2})$ ؛ بنابراین با توجه به (۸)، اگر به ازای هر بردار پارامتری θ_1 قرار دهیم $\theta_1 = (\beta', \theta'_{u1})'$ آنگاه $f_Y(\mathbf{y}; \theta_1) = f_Y(\mathbf{y}; \theta_2)$

پ) اگر ماتریس \mathbf{X} پر رتبه و ماتریس کواریانس Σ_{u1} شناساپذیر در نظر گرفته نشوند، آنگاه به‌راحتی با تعریف بردار پارامتری $\theta_1 = ((\beta + \beta^*)', \theta'_{u1})'$ ، شناساپذیری مدل نتیجه می‌شود.

برهان قضیه (۴): فرض کنید به ازای هر بردار پارامتری θ_1 عضو مجموعه‌ی Θ ، بردار پارامتری θ_2 در مجموعه‌ی Θ و مخالف با θ_1 چنان وجود دارد که رابطه‌ی (۹) برقرار است. حال طرفین رابطه فوق را در $\int \mathbf{w} f_{\mathbf{w}}(\mathbf{w}) d\mathbf{w}$ دارای توزیع نرمال با بردار میانگین صفر و ماتریس کواریانس همانی است) ضرب کرده و انتگرال می‌گیریم:

$$\int \left| \Sigma^{\dagger}(\theta_{u1}) \right| h(\mathbf{X}\beta_1 + \mathbf{Z}\Sigma^{\dagger}(\theta_{u1})\mathbf{w}) f_{\mathbf{w}}(\mathbf{w}) d\mathbf{w} = \int \left| \Sigma^{\dagger}(\theta_{u2}) \right| h(\mathbf{X}\beta_2 + \mathbf{Z}\Sigma^{\dagger}(\theta_{u2})\mathbf{w}) f_{\mathbf{w}}(\mathbf{w}) d\mathbf{w}.$$

در عبارات سمت چپ و راست رابطه‌ی فوق، به ترتیب، تغییر متغیرهای $\mathbf{u} = \Sigma^{\dagger}(\theta_{u1})\mathbf{w}$ و $\mathbf{u} = \Sigma^{\dagger}(\theta_{u2})\mathbf{w}$ را اعمال می‌کنیم تا به برابری زیر برسیم:

$$\int h(\mathbf{X}\beta_1 + \mathbf{Z}\mathbf{u}) f_U(\mathbf{u}; \theta_{u1}) d\mathbf{u} = \int h(\mathbf{X}\beta_2 + \mathbf{Z}\mathbf{u}) f_U(\mathbf{u}; \theta_{u2}) d\mathbf{u}. \quad (۱۳)$$

با توجه به رابطه‌ی (۷)، $E_{\theta_0}(\mathbf{Y} | \mathbf{u}, \mathbf{X}) = h(\mathbf{X}\beta + \mathbf{Z}\mathbf{u})$ ، لذا از برابری (۱۳) نتیجه می‌شود:

$$E_{\theta_0}(\mathbf{Y}) = E(E_{\theta_0}(\mathbf{Y} | \mathbf{u}, \mathbf{X})) = E(E_{\theta_1}(\mathbf{Y} | \mathbf{u}, \mathbf{X})) = E_{\theta_1}(\mathbf{Y})$$

با استفاده از ملاحظه (۱)، شناسانپذیری مدل (۷) ثابت می‌شود.

مراجع

- [1] Shao, J. (2003). *Mathematical Statistics* (2nd ed.), Springer.
- [2] West, B., Welch, K.B. and Galecki, A. T. (2007). *Linear Mixed Models: A Practical Guide Using Statistical Software*. New York: Chapman & Hall, CRC.
- [3] Demidenko, E. (2004). *Mixed Models: Theory and Applications*. John Wiley & Sons.
- [4] Verbeke, G. and Molenberghs, G. (2009). *Linear Mixed Model for Longitudinal Data*. Springer.
- [5] Littell, R. C., Henry, P. R. and Ammerman, C. B. (1998). Statistical analysis of repeated measures data using sas procedures. *Journal of Animal Science*, **76**, 1216-1231.
- [6] Pinheiro, J. and Bates, D. (2009). *Mixed-Effects Models in S and S-PLUS*. Springer.
- [7] Wolfinger, R. (1993). Covariance structure selection in general mixed models. *Communications in Statistics: Simulation and Computation*, **22**(4), 1079-1106.
- [8] Wang, W. (2013). Identifiability of linear mixed effects models. *Electronic Journal of Statistics*, **7**, 244-263.
- [9] San Martin, E. and Quintana, E. (2002). Consistency and identifiability revisited. *Brazilian Journal of Probability and Statistics*, **16**, 99-106.
- [10] Christensen, R. (2011). *1 Plane Answers to Complex Questions: The Theory of Linear Models* (Fourth ed.). Springer.
- [11] Agresti, A. (2002). *Categorical Data Analysis*, p. 133. John Wiley & Sons.
- [۱۲] طهماسبی‌نژاد، ژاله (۱۳۹۰). بررسی مدل‌های همبسته آمیخته با پاسخ‌های پیوسته و گسسته دودویی. پایان‌نامه کارشناسی ارشد، دانشگاه علم و صنعت ایران، تهران.
- [۱۳] قهرودی، زهرا رضایی (۱۳۸۷). کاربردهای مدل‌های انتقالی برای تحلیل داده‌های طولی با پاسخ رسته‌ای با و بدون مقادیر گم‌شده، پایان‌نامه دکترا، دانشگاه شهید بهشتی، تهران.

New Results on Identifiability of Generalized Linear Models with Random Effects

Elham Tabrizi, Ehsan Bahrami Samani and Naseh. Jafari

Department of Statistics, Shahid Beheshti University, Tehran, Iran

Abstract

Identifiability is a necessary property for the adequacy of a statistical model. When a model is not identifiable, no amount of data cannot determine true parameter. In this article, well-known concept of identifiability and its properties is reviewed. Moreover, since non-identifiability problem in linear mixed effects models and generalized linear models with random effects is very common, our main focus is on these models. On the other hand, statistical software, after fitting non-identifiable models, don't usually indicate the problem and show invalid outputs. Consequently, it is useful to have a way to check model identifiability before fitting. In this regard, some new theorems to check identifiability in generalized linear models with random effects are presented. data from non-identifiable models are simulated and problems with model non-identifiability are listed for showing advantages of the mentioned theorems.

Keywords: Identifiability, Linear Mixed Effects Models, Generalized Linear Random Effects Models.

Mathematics Subject Classification (2010): 62J12, 62J05.