

مدل بندی پاسخ های جفت شده ی ترتیبی

انوشیروان کاظم نژاد و فرید زایری

دانشگاه تربیت مدرس

چکیده. حدود ربع قرن از انتشار مقاله ی مکولا [۱۱] در مورد مدل بندی پاسخ های ترتیبی یک متغیره می گذرد. پس از انتشار این مقاله، مدل پیشنهادی وی به تدریج تعمیم یافت، به طوری که ما هم اکنون قادریم پاسخ های چندمتغیره، همبسته و ترتیبی را به کمک مدل های نسبتاً پیچیده اما کارا تحلیل نموده، ارتباط بین این گونه پاسخ ها و متغیرهای کمکی مختلف را مورد بررسی قرار دهیم. در این مقاله قصد داریم پیشرفت های به وجود آمده در زمینه ی مدل بندی پاسخ های ترتیبی همبسته را با تأکید بر پاسخ های دومتغیره ی ترتیبی حاصل از مطالعه ی اندام های جفتی بدن نظیر چشم، گوش، کلیه، دست، پا، و ... بررسی کرده، در پایان، مدلی جدید برای تحلیل داده های ترتیبی همبسته معرفی نماییم. هنگامی که داده های پاسخ مورد بررسی دارای توزیعی دومتغیره و نامتقارن فرض می شوند، این مدل را می توان جایگزینی مناسب برای مدل پروبیت تجمعی دومتغیره محسوب نمود. در پایان، به عنوان مثالی کاربردی، داده های مربوط به وضعیت پرودنتال دانش آموزان دبیرستانی شهر تهران به کمک این مدل، تحلیل و نتایج حاصل با مدل های مشابه مقایسه می شود.

© ۱۳۸۳ پژوهشکده ی آمار. همه ی حقوق محفوظ است.

واژگان کلیدی. پاسخ های ترتیبی همبسته؛ توزیع نهانی دومتغیره؛ معادلات برآوردگر تعمیم یافته؛ مدل های خطی تعمیم یافته.

۱ مقدمه

مدل بندی داده های پاسخ ترتیبی، معمولاً پیچیده تر از پاسخ های دوحالتی و پیوسته به نظر می رسد. به همین دلیل، معمول است که آمارشناسان در ابتدا مدل هایی برای تحلیل داده های دوحالتی یا پیوسته معرفی نموده، سپس این مدل ها را برای تحلیل داده های ترتیبی، تغییر یا تعمیم می دهند. از معروف ترین موارد در این زمینه

می‌توان به مدل لوژستیک معمولی اشاره کرد که در ابتدا فقط برای مدل‌بندی پاسخ‌های دوحالتی یک متغیره به کار برده می‌شد، اما هم‌اکنون تعمیم این مدل را می‌توان برای تحلیل داده‌های چندمتغیره، همبسته و ترتیبی مورد استفاده قرار داد.

با مروری بر مقاله‌های کاربردی منتشر شده در سالیان دور، به‌سادگی می‌توان مشاهده کرد که محققان برای تحلیل داده‌های پاسخ ترتیبی و بررسی ارتباط آن‌ها با متغیر مستقل، بیش‌تر از آزمون‌های آماری استفاده می‌کرده‌اند تا مدل‌سازی آماری. در برخی موارد نیز این پاسخ‌های ترتیبی به پاسخ‌های دیگری (مثلاً دوحالتی) تغییر می‌یافتند و سپس با کمک مدل‌های موجود، مورد بررسی قرار می‌گرفتند. در سال ۱۹۸۰ میلادی، مکولا [۱۱] با انتشار مقاله‌ای در مورد مدل‌بندی پاسخ‌های ترتیبی، تحولی در این زمینه ایجاد کرد. با گذشت حدود ربع قرن از انتشار این مقاله، نام آن را همچنان می‌توان در فهرست مرجع‌های اکثر کتاب‌ها و مقاله‌های منتشر شده در مورد تحلیل داده‌های ترتیبی یک متغیره یا چندمتغیره دید.

با وجود کارایی مناسب مدل مکولا در تحلیل داده‌های ترتیبی یک متغیره، مسئله‌ی مدل‌بندی داده‌های چندمتغیره‌ی ترتیبی و حتی دوحالتی همبسته، همچنان مسئله‌ی پیچیده به نظر می‌رسید. چهار سال بعد، یعنی در سال ۱۹۸۴، رازنر [۱۲] دو مدل مناسب برای تحلیل پاسخ‌های دومتغیره با توزیع دوجمله‌ای و نرمال ارائه کرد. این دو مدل رگرسیون چندگانه، ابزاری مناسب برای تحلیل پاسخ‌های به دست آمده از اندام‌های جفتی بدن در اختیار تحلیل‌گران قرار می‌دادند، که به‌کمک آن‌ها همبستگی موجود بین مشاهدات پاسخ در مدل‌ها لحاظ می‌شد.

روش معادلات برآورده‌گر لیانگ و زیگر [۸] را می‌توان مهم‌ترین تحول در زمینه‌ی تحلیل داده‌های چندمتغیره‌ی همبسته در سالیان اخیر نامید. در این روش، که بعدها به نام روش GEE مشهور شد، نیازی به دانستن توزیع مشترک داده‌های پاسخ چندمتغیره نیست و کافی است که توزیع حاشیه‌ای هر یک از پاسخ‌ها عضوی از خانواده‌ی نمایی فرض شود. این روش به‌طور مستقیم برای تحلیل پاسخ‌های ترتیبی همبسته قابل استفاده نبود، اما لیپ‌شیتس و دیگران [۹] در سال ۱۹۹۴ آن را برای مدل‌بندی داده‌های رسته‌ای همبسته تعمیم دادند. برای برآورد پارامترهای رگرسیونی در هر دو روش یاد شده، از روش شبه‌درست‌نمایی به‌جای روش حد اکثر درست‌نمایی استفاده می‌شود.

با وجود پیچیده بودن تابع درست‌نمایی برای پاسخ‌های چندمتغیره‌ی کیفی همبسته، عده‌ای از آمارشناسان، همچنان بر استفاده از این روش تأکید ورزیده، مدل‌هایی مبتنی بر برآوردهای حداکثر درست‌نمایی را به کار می‌گرفتند. لیپ‌شیتس و دیگران [۱۰] در سال ۱۹۹۰ روشی مبتنی بر حد اکثر درست‌نمایی را برای تحلیل داده‌های جفت‌شده‌ی دوحالتی ارائه کردند. این روش برای تحلیل داده‌های طولی با دو موقعیت اندازه‌گیری یا پاسخ‌های به دست آمده از مطالعات مقطعی کاربرد دارد. با گذشت زمان، مدل‌های مبتنی بر روش حد اکثر درست‌نمایی برای برآورد پارامترها، در تحلیل پاسخ‌های دومتغیره‌ی ترتیبی نیز وارد شدند.

نخستین مقاله در این زمینه در سال ۱۹۹۵ از سوی کیم [۷] ارائه گردید. وی با ابتکاری جالب توجه، یک توزیع دومتغیره‌ی نرمال را توزیع پایه و نهانی برای داده‌های پاسخ دومتغیره فرض کرد و از این توزیع در تشکیل معادلات درست‌نمایی برای برآورد پارامترهای مدل استفاده نمود. مدل کیم به نام مدل پروبیت تجمعی دومتغیره مشهور است.

تحلیل داده‌های چندمتغیره‌ی ترتیبی به این نقطه ختم نشد. یک سال بعد، ویلیامسون و دیگران [۱۴] مدلی دومتغیره برای تحلیل پاسخ‌های ترتیبی همبسته ارائه کردند که در آن از نسبت بخت‌های تعمیم‌یافته به‌عنوان معیار همبستگی بین مشاهدات پاسخ استفاده می‌شد. چند سال بعد، مدل‌های مبتنی بر اثرهای تصادفی [۴] نیز وارد عرصه‌ی تحلیل پاسخ‌های ترتیبی همبسته شدند. طرفداران تحلیل‌های بیزی نیز سرانجام به میدان آمده، مدلی مبتنی بر نمونه‌گیری MCMC و الگوریتم Gibbs برای تحلیل داده‌های ترتیبی دومتغیره ارائه کردند. این مقاله در سال ۲۰۰۲ از سوی بیسواس و داس [۳] منتشر شد. با گسترش روش‌های مختلف برای تحلیل داده‌های ترتیبی چندمتغیره، نیازی مبرم به ایجاد امکانات نرم‌افزاری برای برازش این مدل‌ها احساس شد. ویلیامسون و دیگران [۱۵] در سال ۱۹۹۹ طی مقاله‌ای دو نرم‌افزار رایانه‌ای برای تحلیل داده‌های پاسخ رسته‌ای همبسته معرفی کردند. پس از سال ۲۰۰۰، نرم‌افزار قدرتمند SAS، امکانات لازم را برای برازش گروه عمده‌ای از مدل‌های مناسب برای تحلیل داده‌های چندمتغیره‌ی همبسته شامل برآوردهای GEE و حد اکثر درست‌نمایی فراهم آورد. با این حال، برازش برخی از انواع این مدل‌ها، همچنان نیازمند طراحی برنامه‌های پیچیده‌ی رایانه‌ای در محیط برنامه‌نویسی نرم‌افزاری نظیر SAS، S-Plus، و ... است.

در این مقاله امکان بحث در مورد همه‌ی روش‌های یاد شده وجود ندارد. از آن جایی که مدل ارائه شده از سوی ما، با فرض یک توزیع پیوسته‌ی نهانی برای داده‌های جفت‌شده‌ی ترتیبی بنا می‌شود، در بخش بعد فقط به توصیف مدل‌های مبتنی بر توزیع‌های دومتغیره‌ی نهانی اکتفا می‌کنیم. سپس مهم‌ترین مدل دومتغیره‌ی ارائه شده در سال‌های اخیر (با فرض یک توزیع دومتغیره‌ی نهانی) را مورد بحث قرار می‌دهیم. همان‌طور که پیش از این گفته شد، این مدل را کیم ارائه کرده است و در آن از توزیع نرمال دومتغیره‌ی استاندارد به‌عنوان متغیر نهانی استفاده می‌شود. با وجود مناسب بودن این مدل برای تحلیل طیف وسیعی از داده‌های جفت‌شده‌ی ترتیبی، هنگامی که پاسخ دومتغیره حالت نامتقارنی نداشته باشد، به‌کارگیری این مدل (به‌علت فرم متقارن توزیع نرمال دومتغیره) چندان منطقی به نظر نمی‌آید. برای رفع این نقیصه، در بخش ۳/۲ مدلی ارائه می‌کنیم که در آن، توزیع دومتغیره‌ی نهانی از فرم نامتقارنی برخوردار بوده، نتیجه انتخابی مناسب برای تحلیل برخی از پاسخ‌های دومتغیره‌ی نامتقارن است.

پس از توصیف این مدل، روش‌های مورد نیاز برای برآورد پارامترهای آن بیان شده، سپس داده‌های حاصل از یک مطالعه‌ی اپیدمیولوژیک مربوط به وضعیت پریدنتال دانش‌آموزان دبیرستانی شهر تهران را

به وسیله‌ی این مدل، تحلیل می‌کنیم. در پایان، نتایج حاصل از برازش مدل پیشنهادی، با نتایج حاصل از مدل پروبیت تجمعی دومتغیره‌ی کیم و همچنین مدل مبتنی بر برآورد GEE مقایسه خواهد شد.

۲ مدل‌های رگرسیونی برای پاسخ‌های ترتیبی یک متغیره

در این بخش به‌طور مختصر، مدل‌های رگرسیونی با پاسخ‌های ترتیبی یک متغیره را مورد بحث قرار می‌دهیم. در همه‌ی این مدل‌ها فرض بر این است که یک متغیر نهانی پیوسته وجود دارد که به‌طور مستقیم قابل مشاهده نیست، اما توزیع آن به‌عنوان توزیع پایه‌ی متغیر پاسخ در نظر گرفته می‌شود. فرض کنیم y^* یک متغیر تصادفی پیوسته است که به‌طور مستقیم قادر به مشاهده‌ی آن نیستیم. در عوض، فقط مشاهده می‌کنیم که آیا y^* به فواصل مجاور نظیر (θ_{h-1}, θ_h) ، $h = 1, \dots, H$ ، با فرض $\theta_0 = -\infty$ و $\theta_H = +\infty$ ، تعلق دارد یا خیر. نقاط برش، یعنی $\theta = (\theta_1, \dots, \theta_{H-1})$ ، معمولاً نامعلوم هستند و در عوض، متغیر تصادفی ترتیبی y به‌صورت زیر قابل مشاهده است:

$$\theta_{h-1} < y^* \leq \theta_h \quad \Rightarrow \quad y = h.$$

در حالت کلی فرض کنیم N فرد در مطالعه حضور دارند و y_i^* ، $(i = 1, \dots, N)$ ، نشان‌دهنده‌ی متغیر پیوسته‌ی نهانی پاسخ برای i امین فرد تحت مطالعه با یک رگرسیون خطی روی متغیرهای کمکی \mathbf{x}_i است. به‌طور دقیق‌تر، فرض بر این است که $\frac{y_i^* - \mathbf{x}_i' \beta}{\sigma}$ ، متغیرهای تصادفی با تابع توزیع مشترک F هستند. در این جا σ پارامتر مقیاس و β برداری از ضرایب رگرسیونی فرض می‌شود. بدون از دست رفتن کلیت مطلب، فرض کنیم $\sigma = 1$. در این صورت، متغیر پاسخ رسته‌ای y_i برای i امین فرد، دارای تابع چگالی احتمال زیر خواهد بود:

$$\Pr(y_i = h) = \pi_h(\theta, \beta; \mathbf{x}_i),$$

که در آن:

$$\begin{aligned} \pi_h(\theta, \beta; \mathbf{x}_i) &= \Pr(\theta_{h-1} < y_i^* \leq \theta_h) \\ &= F(\theta_h - \mathbf{x}_i' \beta) - F(\theta_{h-1} - \mathbf{x}_i' \beta). \end{aligned}$$

حال، بر حسب احتمالات تجمعی $\sum_{s=1}^h \pi_s$ ، γ_h می‌توان مدل رگرسیونی را به‌صورت زیر نوشت:

$$F^{-1}\{\gamma_h(\theta, \beta; \mathbf{x}_i)\} = \theta_h - \mathbf{x}_i' \beta.$$

در این جا می‌توان یک توزیع نهانی پیوسته‌ی F به‌عنوان تابع ربط برای احتمال‌های تجمعی فوق فرض کرد. با این کار، در واقع، مدل به‌صورت یک مدل خطی تعمیم‌یافته نوشته خواهد شد. معمول‌ترین انتخاب در

این زمینه، تابع توزیع لوژستیک است که مدل لوژستیک را نتیجه می‌دهد:

$$\log\left(\frac{\gamma_h}{1-\gamma_h}\right) = \theta_h - \mathbf{x}'_i\beta.$$

تابع توزیع گامبل، انتخاب دیگری است که مدل مکمل log-log را به دست می‌دهد:

$$\log\{-\log(1-\gamma_h)\} = \theta_h - \mathbf{x}'_i\beta.$$

همچنین انتخاب تابع توزیع نرمال، مدل پروبیت را نتیجه خواهد داد:

$$\varphi^{-1}(\gamma_h) = \theta_h - \mathbf{x}'_i\beta,$$

که در آن $\varphi^{-1}(\cdot)$ وارون تابع توزیع نرمال استاندارد است.

به‌علت وجود تقارن در توزیع‌های لوژستیک و نرمال، معمولاً مدل‌های لوژستیک و پروبیت، نتایجی مشابه یکدیگر خواهند داشت. اما از مدل‌هایی با تابع پیوند log-log در عمل، هنگامی استفاده می‌شود که داده‌های پاسخ، حالت متقارن نداشته باشند [۲].

پیش از ارائه‌ی روش‌های مناسب برای تحلیل چندمتغیره‌ی پاسخ‌های جفت‌شده‌ی ترتیبی، تحلیل‌گران در ابتدا پاسخ‌های چندمتغیره را با روش‌هایی چون گرفتن میانگین یا میانه، در نظر گرفتن پاسخ ماکسیمم برای هر فرد، یا روش‌های دیگری مانند روش تأکید بر اندام بیمارتر [۵] به پاسخی یک‌متغیره تبدیل کرده، سپس با استفاده از مدل‌های توصیف‌شده در این بخش، داده‌ها را تحلیل می‌کردند. اما به تدریج با مطرح شدن روش‌هایی برای تحلیل پاسخ‌های دومتغیره یا چندمتغیره، مدل‌های یک‌متغیره اهمیت خود را در تحلیل این‌گونه پاسخ‌ها از دست دادند. آنانث و کلاین‌باوم [۱۱] مقاله‌ی مروری مبسوطی در مورد سیر تحولات روش‌های مربوط به تحلیل داده‌های ترتیبی یک‌متغیره ارائه کرده‌اند. اشکال مهم به‌کارگیری مدل‌های یک‌متغیره برای تحلیل پاسخ‌های چندمتغیره‌ی همبسته این است که با تبدیل یک پاسخ چندمتغیره به مشاهده‌ای یک‌متغیره، عملاً امکان احتساب همبستگی بین مشاهدات پاسخ در مدل را از دست می‌دهیم. همچنین واضح است که تبدیل پاسخ‌های چندمتغیره به یک‌متغیره ممکن است موجب بروز تغییرات کلی در ساختار روابط پاسخ مورد بررسی و متغیرهای تبیینی گردد. در بخش‌های بعدی، مدل‌هایی را معرفی می‌کنیم که در آن‌ها امکان تحلیل پاسخ‌های دومتغیره‌ی همبسته و ترتیبی با احتساب همبستگی بین مشاهدات پاسخ وجود دارد.

۳ مدل‌های رگرسیونی برای پاسخ‌های جفت‌شده‌ی ترتیبی

۳/۱ مدل پروبیت تجمعی برای پاسخ‌های دومتغیره‌ی متقارن

همان‌طور که در فصل پیش اشاره شد، مدل پروبیت تجمعی دومتغیره در سال ۱۹۹۵ از سوی کیم ارائه گردید. قبل از توصیف این مدل، بهتر است اشاره‌ای به ساختار داده‌های جفت‌شده در حالت کلی داشته باشیم.

در مطالعه‌ای پزشکی مربوط به اندام‌های جفتی بدن، فرض کنیم y_{i1} و y_{i2} نشان‌دهنده‌ی پاسخ‌های دومتغیره‌ی ترتیبی، به ترتیب برای اندام راست و چپ (یا در مطالعات دندان‌پزشکی، بالا و پایین) فرد i ام باشد. بنا بر این، پاسخ دومتغیره‌ی فرد i ام را می‌توان به صورت $\mathbf{y}_i = (y_{i1}, y_{i2})'$ نمایش داد. همچنین فرض کنیم y_{i1} و y_{i2} مقادیر ۱ و ۲ و ... و H را اختیار کنند. حال همانند حالت یک‌متغیره، متغیره‌های پیوسته و غیر قابل مشاهده‌ی y_{i1}^* و y_{i2}^* را در نظر می‌گیریم. بار دیگر، نقاط برش $\theta_0, \theta_1, \dots, \theta_H$ ، $(\theta_H = +\infty$ و $\theta_0 = -\infty)$ وجود دارند، به طوری که

$$\theta_{h-1} < y_{ij}^* \leq \theta_h \Rightarrow y_{ij} = h.$$

توجه کنید که در این جا نقاط برش یکسانی برای هر دو حاشیه در نظر گرفته شده است. فرض یکسان بودن حاشیه‌ها برای اندام‌های جفتی بدن، فرضی غیرمنطقی به نظر نمی‌رسد. با این حال، چنانچه ماهیت داده‌ها ایجاب کند، می‌توان نقاط برش θ_h ، $(h = 1, \dots, H-1)$ را با θ_{jh} ، $j = 1, 2$ و $(h = 1, \dots, H-1)$ جایگزین نمود.

در مطالعات مربوط به اندام‌های جفتی بدن، دو نوع متغیر کمکی را می‌توان در نظر گرفت: متغیره‌های کمکی مختص فرد و متغیره‌های کمکی مختص اندام. اگر \mathbf{x}_i نشان‌دهنده‌ی برداری p -بعدی از متغیره‌های کمکی مختص فرد i ام، و \mathbf{x}_{i1} و \mathbf{x}_{i2} نشان‌دهنده‌ی بردارهای q -بعدی مختص اندام‌های راست و چپ مربوط به همین فرد باشند، آنگاه ماتریس متغیره‌های کمکی (ماتریس طرح) مربوط به i امین فرد را می‌توان به صورت زیر نمایش داد:

$$\mathbf{X}'_i = \begin{pmatrix} \mathbf{x}'_i & \mathbf{x}'_{i1} \\ \mathbf{x}'_i & \mathbf{x}'_{i2} \end{pmatrix}.$$

با این تعریف، این بار فرض کنیم $\frac{\mathbf{y}_i - \mathbf{X}'_i \beta}{\sigma}$ متغیره‌های تصادفی مستقل با توزیع یکسان و تابع توزیع دومتغیره‌ی مشابه F_ρ باشند، به طوری که پارامتر ρ همبستگی بین y_{i1}^* و y_{i2}^* را مشخص می‌کند. همچنین، $\mathbf{y}_i^* = (y_{i1}^*, y_{i2}^*)'$ نشان‌دهنده‌ی بردار متغیره‌های نهانی مختص فرد i ام است. در این جا، $\beta = (\delta', \epsilon')'$

بردار پارامترهای رگرسیونی مختص فرد و مختص اندام را نمایش می‌دهد و σ همان پارامتر مقیاس معمولی است. بار دیگر فرض کنیم $\sigma = 1$.

از آنجایی که متغیرهای پاسخ مربوط به اندام‌های راست و چپ، هر یک ممکن است در یکی از طبقات ۱ تا H قرار گیرند، پاسخ دومتغیره‌ی هر فرد ممکن است در یکی از خانه‌های یک جدول $H \times H$ قرار گیرد و در نتیجه پاسخ کل این N فرد را می‌توان در این جدول $H \times H$ خلاصه کرد. حال فرض می‌کنیم $\pi_{hg}(\omega; \mathbf{X}_i)$ نشان‌دهنده‌ی احتمال وقوع (h, g) امین خانه‌ی جدول توافقی $H \times H$ یاد شده باشد. به عبارت دیگر،

$$\pi_{hg}(\omega; \mathbf{X}_i) = \text{Pr}(y_{i1} = h, y_{i2} = g)$$

برای هر $h, g = 1, \dots, H$. در این جا $\omega = (\theta', \rho, \delta', \varepsilon')$ بردار $(H + p + q)$ -بعده‌ی مربوط به پارامترهای مدل است. در واقع، در فرایند مدل‌سازی باید $H - 1$ نقطه‌ی برش، یک پارامتر همبستگی، p پارامتر مربوط به متغیرهای کمکی مختص فرد، و q پارامتر مربوط به متغیرهای کمکی مختص اندام برآورد شود.

بار دیگر به‌دلیل ماهیت ترتیبی متغیر پاسخ دومتغیره، بهتر است احتمال‌های تجمعی γ_{hg} را در مدل وارد کنیم. هنگامی که فرض تقارن توزیع پاسخ‌های حاشیه‌ای به دست آمده از هر اندام، فرضی معقول باشد، معمول‌ترین انتخاب ممکن برای متغیر پایه‌ی نهانی، توزیع نرمال استاندارد دومتغیره است [۷]. به عبارت بهتر، می‌توان نوشت:

$$\begin{aligned} \gamma_{hg}(\omega; \mathbf{X}_i) &= F_{\rho}(\theta_h - \mathbf{x}'_i \delta - \mathbf{x}'_{i1} \varepsilon, \theta_g - \mathbf{x}'_i \delta - \mathbf{x}'_{i2} \varepsilon) \\ &= \varphi_{\rho}(\theta_h - \mathbf{x}'_i \delta - \mathbf{x}'_{i1} \varepsilon, \theta_g - \mathbf{x}'_i \delta - \mathbf{x}'_{i2} \varepsilon), \end{aligned}$$

که در آن φ_{ρ} تابع توزیع نرمال دومتغیره‌ی استاندارد است. بنا بر این، احتمال (h, g) امین خانه‌ی جدول $H \times H$ یاد شده را می‌توان به صورت زیر نشان داد:

$$\begin{aligned} \pi_{hg}(\omega; \mathbf{X}_i) &= \varphi_{\rho}(\theta_h - \mathbf{x}'_i \delta - \mathbf{x}'_{i1} \varepsilon, \theta_g - \mathbf{x}'_i \delta - \mathbf{x}'_{i2} \varepsilon) \\ &\quad - \varphi_{\rho}(\theta_h - \mathbf{x}'_i \delta - \mathbf{x}'_{i1} \varepsilon, \theta_{g-1} - \mathbf{x}'_i \delta - \mathbf{x}'_{i2} \varepsilon) \\ &\quad - \varphi_{\rho}(\theta_{h-1} - \mathbf{x}'_i \delta - \mathbf{x}'_{i1} \varepsilon, \theta_g - \mathbf{x}'_i \delta - \mathbf{x}'_{i2} \varepsilon) \\ &\quad + \varphi_{\rho}(\theta_{h-1} - \mathbf{x}'_i \delta - \mathbf{x}'_{i1} \varepsilon, \theta_{g-1} - \mathbf{x}'_i \delta - \mathbf{x}'_{i2} \varepsilon). \end{aligned}$$

به منظور برآورد پارامترهای مدل بالا، تابع لگاریتم درست‌نمایی برای داده‌های دومتغیره‌ی ترتیبی به صورت

زیر خواهد بود:

$$\ell(\omega) = \sum_{i=1}^N \sum_{h=1}^H \sum_{g=1}^H I_{i \setminus h} I_{i \setminus g} \log \pi_{hg}(\omega; \mathbf{X}_i),$$

که در آن، $I_{ijh} = 1$ اگر پاسخ مربوط به زامین اندام فرد i ام در رسته h قرار داشته باشد، و در غیر این صورت $I_{ijh} = 0$. برآورد کردن پارامترها در این مدل، نیازمند روشی تکراری نظیر نیوتون-رافسون برای حل معادله $\frac{\partial \ell}{\partial \omega} = 0$ و ارزیابی ماتریس مشتق‌های جزئی مرتبه دوم (Hessian matrix) به صورت $-\frac{\partial^2 \ell}{\partial \omega \partial \omega}$ برای به دست آوردن خطاهای معیار پارامترهای برآورد شده است.

۳/۲ مدل تجمعی دومتغیره برای پاسخ‌های نامتقارن

فرض متقارن بودن توزیع‌های حاشیه‌ای متغیرهای پاسخ ترتیبی، فرضی است که در عمل، همواره برقرار نخواهد بود. بنا بر این، به‌کارگیری توزیع نرمال دومتغیره به‌عنوان متغیر پایه‌ی نهانی برای انواع داده‌های پاسخ به دست آمده از اندام‌های جفتی بدن، ممکن است در برخی موارد منجر به استنباط‌هایی گمراه‌کننده در باره‌ی متغیر پاسخ مورد بررسی شود.

برای رفع نقیصه‌ی یاد شده، ما مدل دیگری مبتنی بر یک توزیع دومتغیره‌ی نامتقارن نهانی پیشنهاد می‌کنیم. این توزیع در سال ۱۹۷۸ از سوی ساترت ویت و هاچینسون [۱۳] پیشنهاد شده است. در این بخش، نخست این توزیع دومتغیره را به‌طور مختصر مورد بحث قرار داده، سپس مدل پیشنهادی خود را ارائه می‌کنیم.

در سال ۱۹۶۱، گامبل [۶] دو حالت از توزیع لوژستیک دومتغیره ارائه کرد، که یکی از آن‌ها دارای تابع توزیع زیر بود:

$$\psi(x, y) = (1 + e^{-x} + e^{-y})^{-1}.$$

اشکال بزرگ این توزیع دومتغیره، فقدان پارامترهای مرتبط با همبستگی بین دو متغیر تصادفی X و Y است. برای رفع این اشکال، ساترت ویت و هاچینسون توزیع فوق را به‌صورت زیر اصلاح کردند:

$$\psi(x, y) = (1 + e^{-x} + e^{-y})^{-v},$$

که در آن، $v > 0$ پارامتری است که به‌کمک آن می‌توان ضریب همبستگی بین X و Y را محاسبه کرد. این دو آمارشناس در مقاله‌ی خود، خواص مهم این توزیع (نظیر توزیع‌های حاشیه‌ای، امید ریاضی و واریانس، کوواریانس بین X و Y ، و ...) را ارائه نمودند. بر این اساس، آنان ضریب همبستگی بین X و Y را

به صورت زیر به دست آوردند:

$$\rho = \frac{\zeta(2, v)}{\zeta(2, v) + \frac{1}{6}\pi^2},$$

که در آن، $\zeta(s, a) = \sum_{m=0}^{\infty} (m+a)^{-s}$ تابع زتای ریمان (zeta Riemann function) است. بدیهی است در این توزیع، همواره داریم $0 \leq \rho \leq 1$. توزیع فوق بر خلاف توزیع نرمال دومتغیره، دارای شکلی نامتقارن است.

با تعاریفی کاملاً مشابه با مدل پروبیت تجمعی دومتغیره برای احتمال‌های π_{hg} ، می‌توان مدل زیر را با فرض متغیر نهانی توصیف‌شده در این بخش، بر اساس احتمال‌های تجمعی γ_{hg} بنا کرد:

$$\begin{aligned} \gamma_{hg}(\omega; \mathbf{X}_i) &= F_{\rho}(\theta_h - \mathbf{x}'_i \delta - \mathbf{x}'_{i1} \varepsilon, \theta_g - \mathbf{x}'_i \delta - \mathbf{x}'_{i2} \varepsilon) \\ &= \psi_{\rho}(\theta_h - \mathbf{x}'_i \delta - \mathbf{x}'_{i1} \varepsilon, \theta_g - \mathbf{x}'_i \delta - \mathbf{x}'_{i2} \varepsilon), \end{aligned}$$

که در آن، ψ_{ρ} تابع توزیع پیشنهاد شده از سوی ساترتویت و هاجینسون است. در این‌جا نیز $\omega = (\theta', \rho, \delta', \varepsilon')$ بردار پارامترهای مدل را نمایش می‌دهد. پس از برآورد کردن پارامتر همبستگی v در فرآیند مدل‌سازی، به سادگی می‌توان مقدار ضریب همبستگی ρ را به کمک آن برآورد نمود. برای برآورد کردن پارامترهای مدل، با فرض

$$z_{sj} = \theta_s - \mathbf{x}'_i \delta - \mathbf{x}'_{ij} \varepsilon$$

برای $j = 1, 2$ ، می‌توان مدل را در حالت کلی به صورت زیر نوشت:

$$\begin{aligned} \pi_{hg}(\omega; \mathbf{x}_i) &= \psi_v(z_{h1}, z_{g2}) - \psi_v(z_{h1}, z_{(g-1), 2}) \\ &\quad - \psi_v(z_{(h-1), 1}, z_{g2}) + \psi_v(z_{(h-1), 1}, z_{(g-1), 2}). \end{aligned}$$

تشکیل تابع لگاریتم درست‌نمایی، بر حسب مقدار کاملاً مشابه مدل پروبیت تجمعی دومتغیره خواهد بود. برای برآورد کردن پارامترهای این مدل، از محیط برنامه‌نویسی نرم‌افزار S-Plus استفاده کردیم. روال `nlminb` در این نرم‌افزار، ابزاری مناسب برای به دست آوردن برآورد پارامترهای مدل‌های مختلف غیر خطی به کمک روش حد اکثر درست‌نمایی است. برای به دست آوردن برآورد پارامترها، روال `nlminb` از روش تکراری گاوس-نیوتون استفاده می‌کند. رایانه‌های جدید با قدرت پردازش بالا قادرند عملیات فرآیند مدل‌سازی را در طی چند دقیقه انجام دهند.

۴ نتایج حاصل از یک تحقیق دندان پزشکی

در این بخش، داده‌های حاصل از یک مطالعه‌ی اپیدمیولوژیک در مورد وضعیت پریدنتال دانش‌آموزان دبیرستانی شهر تهران به‌کمک روش‌های توصیف‌شده، مورد تحلیل قرار می‌گیرد. نمونه‌ای تصادفی به اندازه‌ی ۸۶۷ دانش‌آموز ۱۵-۱۹ ساله‌ی تهرانی انتخاب شده و وضعیت لثه‌ی آنان به‌وسیله‌ی شاخصی بهداشتی موسوم به شاخص CPI (Community Periodontal Index) مورد بررسی قرار گرفته است. CPI شاخصی ترتیبی با مقادیر زیر است: ۰ = لثه‌ی سالم، ۱ = خون‌ریزی در لثه، ۲ = وجود جرم در لثه، ۳ = وجود پاکت با عمق ۴-۶ میلی‌متر، ۴ = وجود پاکت با عمق بیش از ۶ میلی‌متر.

مقدار CPI معمولاً در شش دندان شاخص، که نشان‌دهنده‌ی شش بخش فک بالا و پایین دهان هر فرد است، اندازه‌گیری می‌شود. این شش بخش (sextant) عبارت‌اند از: سمت راست فک بالا، قسمت میانی فک بالا، سمت چپ فک بالا، سمت راست فک پایین، قسمت میانی فک پایین، و سمت چپ فک پایین. بنا بر این، برای هر فرد، شش مقدار ترتیبی همبسته به‌عنوان پاسخ در اختیار داریم. از آنجایی که هدف ما مدل‌بندی داده‌های ترتیبی دومتغیره است، در ابتدا ضریب همبستگی اسپیرمن برای سمت راست بالا و پایین، میانه‌ی بالا و پایین، و سمت چپ بالا و پایین محاسبه شد. چون بالاترین مقدار ضریب همبستگی اسپیرمن در داده‌های مربوط به بخش میانی فک بالا و پایین به دست آمد ($r = 0.741$ با $P < 0.0001$)، ما تحلیل خود را بر این دو نقطه متمرکز خواهیم کرد. جدول ۱ پاسخ دومتغیره را به‌طور خلاصه نمایش می‌دهد.

جدول ۱. وضعیت پریدنتال (CPI) بخش میانی فک بالا و پایین در نمونه‌ی دانش‌آموزان تهران

فک پایین						فک بالا
سالم	خون‌ریزی	جرم	پاکت کم عمق	پاکت عمیق	جمع	
۲۹۷	۸۸	۲۱	۴	۰	۴۱۰	سالم
۲۷	۱۳۱	۶۵	۹	۰	۲۳۲	خون‌ریزی
۵	۱۴	۱۳۳	۱۷	۰	۱۶۹	جرم
۳	۶	۱۹	۲۷	۱	۵۶	پاکت کم عمق
۰	۰	۰	۰	۰	۰	پاکت عمیق
۳۳۲	۲۳۹	۲۳۸	۵۷	۱	۸۶۷	جمع

از بین متغیرهای کمکی مختلف، ما تأثیر شش متغیر کمکی مختص فرد، و فقط یک متغیر کمکی مختص اندام را بر وضعیت پریدنتال دانش‌آموزان مورد بررسی قرار می‌دهیم. متغیرهای کمکی مختص فرد عبارت‌اند از: جنس (۱ = زن، ۲ = مرد)، تحصیلات پدر (۱ = کم‌تر از دیپلم، ۲ = دیپلم و بالاتر)، تحصیلات

مادر (۱=کم‌تر از دیپلم، ۲=دیپلم و بالاتر)، استفاده از مسواک (۱=هرگز یا به‌طور نامرتب، ۲=به‌طور مرتب، حد اقل روزی یک بار)، استفاده از نخ دندان (۱=هرگز یا به‌طور نامرتب، ۲=به‌طور مرتب، حد اقل روزی یک بار)، ویزیت توسط دندان‌پزشک (۱=فقط در مواقع اضطراری، ۲=به‌طور مرتب به‌منظور پیش‌گیری). همچنین، فقط وجود پلاک‌دندانی قابل مشاهده در قسمت میانی فک بالا یا فک پایین به‌عنوان متغیر کمکی مختص اندام در نظر گرفته شد. در هر یک از فک‌های بالا و پایین، این متغیر به این صورت ثبت شد: پلاک قابل مشاهده (۱=وجود، ۲=نبود).

پیش از برازش مدل‌های دومتغیره به این داده‌ها، مدل‌های لجیت، پروبیت و $\log\text{-}\log$ برازانده شده به پاسخ‌های حاشیه‌ای به دست آمده از فک‌های بالا و پایین را مقایسه کردیم. آماره‌های نیکویی برازش برای مدل با پیوند نامتقارن $\log\text{-}\log$ به‌خصوص برای فک پایین (به‌طور معناداری بهتر از نتایج برازش مدل‌های لجیت و پروبیت به دست آمد. در مرحله‌ی بعد، برازش سه مدل چندمتغیره‌ی زیر را مورد بررسی قرار دادیم:

۱. مدل لجیت تجمعی که پارامترهای آن به روش GEE برآورد می‌شوند؛
۲. مدل پروبیت تجمعی که پارامترهای آن با روش حد اکثر درست‌نمایی برآورد می‌شوند (مدل با تابع پیوند متقارن یا مدل کیم)؛
۳. مدل ارائه شده در بخش ۳/۲ (مدل با تابع پیوند نامتقارن).

به‌علت کوچک بودن اندازه‌ی نمونه برای رسته‌ی دارای $\text{CPI}=4$ ، رسته‌های ۳ و ۴ داده‌های پاسخ را ادغام کردیم. بنا براین، پاسخ‌های حاشیه‌ای به این صورت در مدل وارد شدند: $0 = \text{سالم}$ ، $1 = \text{خون‌ریزی}$ ، $2 = \text{جرم}$ ، $3 = \text{پاکت با عمق بیش از ۴ میلی‌متر}$. جدول ۲ نتایج برازش این سه مدل را نمایش می‌دهد. همان‌طور که از مقادیر منتهای دو برابر لگاریتم درست‌نمایی مشخص است، بهترین برازش، همان‌گونه که انتظار می‌رفت، مربوط به مدل نامتقارن و پس از آن متعلق به مدل کیم و در آخر، مربوط به مدل مبتنی بر برآوردهای GEE است. بدیهی است که برآورد پارامترها، مقادیر خطاهای معیار و پی‌مقدار برای سه مدل، یکسان نخواهد بود. بر طبق نتایج به دست آمده از مدل نامتقارن، کلیه‌ی متغیرهای کمکی دارای تأثیری معنادار بر وضعیت پرودنتال دانش‌آموزان دبیرستانی شهر تهران هستند. مطابق مدل کیم، متغیر تحصیلات پدر، یگانه متغیر کمکی غیر معنادار در مدل است. همچنین بر طبق نتایج حاصل از برازش GEE، دو متغیر تحصیلات پدر و تحصیلات مادر، تأثیر معناداری بر وضعیت پرودنتال دانش‌آموزان تهرانی نداشته‌اند. تفسیر ضرایب رگرسیونی برآورد شده را می‌توان بر حسب مقیاس متغیر نهانی انجام داد. به‌عنوان مثال، در مدل نامتقارن، از آنجایی که اثر استفاده از مسواک به‌صورت $\delta_4 = 17755$ و اثر استفاده از نخ دندان به‌صورت $\delta_5 = 17546$ برآورد شده است، اگر این دو اثر را به‌طور توأم در نظر بگیریم، یعنی $17755 + 17546 = 37301$ ، می‌توان

جدول ۲. نتایج برازش مدل‌های رگرسیونی چندمتغیره برای داده‌های پر یودنتال

پارامتر	مدل نامتقارن			مدل کیم			مدل GEE		
	P	SE	برآورد	P	SE	برآورد	P	SE	برآورد
θ_1	†	۱,۲۳۵	۱۱,۰۵۴	†	۱,۰۹۱	۵,۱۱۲	†	۰,۱۱۵	-۲,۲۲۴
θ_2	†	۱,۵۴۱	۱۳,۱۸۴	†	۰,۱۳۲	۷,۴۹۹	†	۰,۱۰۳	-۱,۲۸۷
θ_3	†	۱,۵۳۸	۱۵,۲۶۹	†	۰,۳۲۵	۹,۴۲۵	†	۰,۱۰۲	۰,۰۹۵
δ_1 : جنس	۰,۰۳۱	۰,۶۳۲	۰,۹۹۹	۰,۰۴۱	۰,۲۷۸	۰,۷۳۳	۰,۰۰۱	۰,۰۷۲	۰,۴۲۰
δ_2 : تحصیلات پدر	۰,۰۳۸	۰,۱۲۷	۰,۳۲۹	۰,۰۶۱	۰,۰۸۳	۰,۱۹۸	۰,۱۹۸	۰,۱۶۰	۰,۲۰۶
δ_3 : تحصیلات مادر	< ۰,۰۰۱	۰,۲۸۹	۰,۹۱	< ۰,۰۰۱	۰,۱۶۰	۰,۵۵۲	< ۰,۰۰۱	۰,۲۳۸	۰,۴۴۶
δ_4 : مسواک	< ۰,۰۰۱	۰,۲۱۲	۱,۷۵۵	< ۰,۰۰۱	۰,۱۱۲	۰,۸۳۰	< ۰,۰۰۱	۰,۰۷۹	۰,۶۵۳
δ_5 : نخ دندان	< ۰,۰۰۱	۰,۱۹۷	۱,۵۴۶	< ۰,۰۰۱	۰,۱۹۹	۰,۹۵۶	< ۰,۰۰۱	۰,۱۴۱	۰,۷۲۹
δ_6 : ویزیت	< ۰,۰۰۱	۰,۲۷۴	۱,۱۳۹	< ۰,۰۰۱	۰,۱۹۴	۰,۷۷۳	< ۰,۰۰۱	۰,۰۷۹	۰,۵۱۳
ϵ_1 : پلاک	< ۰,۰۰۱	۰,۱۱۳	۱,۵۸۱	< ۰,۰۰۱	۰,۰۹۱	۱,۰۰۵	< ۰,۰۰۱	۰,۰۶۴	۱,۱۷۵
پارامتر همبستگی	†	۰,۰۹۸	۰,۴۴۷	†	۰,۱۳۱	۰,۸۰۳	†	†	†
-2log(likelihood)	۳۲۵۶,۲۸۳			۳۳۰۱,۱۶۸			۳۵۲۲,۲۲۷		

† محاسبه نمی‌شود.

چنین استنباط کرد که تغییر از استفاده به عدم استفاده از مسواک و نخ دندان به طور توأم، باعث تغییری به اندازه‌ی ۱/۳۳ در مقیاس متغیر نهانی می‌شود. با توجه به مقادیر عرض از مبدأ برآورد شده، این میزان تغییر می‌تواند وضعیت لثه‌ی یک دانش‌آموز را از حالت خون‌ریزی ($\hat{\theta}_1 \leq y_{ij}^* = ۱۳,۱۸۴$) به حالت وجود پلاک بیش از ۴ میلی‌متر ($\hat{\theta}_3 \leq y_{ij}^* = ۱۵,۲۶۹$) تغییر دهد.

همان‌طور که پیش از این بیان شد، در مدل پیشنهادی ما، مقدار ضریب همبستگی بین پاسخ‌های فک بالا و پایین را می‌توان به وسیله‌ی برآورد پارامتر همبستگی مدل (ρ) محاسبه نمود. با توجه به فرمول ارائه شده در بخش ۳/۲ و مقدار برآورد شده‌ی $\hat{\rho} = ۰/۴۴۷$ ، ضریب همبستگی ρ را می‌توان به صورت زیر برآورد نمود:

$$\hat{\rho} = \frac{\zeta(2, ۰/۴۴۷)}{\zeta(2, ۰/۴۴۷) + \frac{1}{2}\pi^2};$$

یعنی برآورد ضریب همبستگی بین پاسخ‌های مشاهده شده در مدل نامتقارن، در حدود ۰/۷۸۴ بوده است. همچنین، با توجه به مدل کیم، ضریب همبستگی بین پاسخ‌ها را می‌توان مستقیماً با برآورد کردن ضریب همبستگی توزیع نرمال دومتغیره (توزیع نهانی) به دست آورد. این برآورد برای داده‌های مورد مطالعه، در حدود ۰/۸۰ محاسبه شده است، که تفاوت چندانی با برآورد حاصل از مدل نامتقارن ندارد. توجه داریم که مقدار ضریب همبستگی به دست آمده از مدل نامتقارن، به ضریب همبستگی اسپیرمن ($r = ۰/۷۴۱$)

نزدیک‌تر است تا به برآورد حاصل از مدل پروبیت تجمعی دومتغیره.

۵ بحث

هدف ما از ارائه‌ی این مقاله، معرفی روش‌هایی برای مدل‌بندی پاسخ‌های جفت‌شده‌ی ترتیبی بود. بدین منظور، بر مدل‌هایی که یک توزیع دومتغیره‌ی پیوسته‌ی نهانی را به‌عنوان تابع پیوند مورد استفاده قرار می‌دهند، تأکید بیش‌تری شد. اولین پرسش در این جهت ممکن است بدین صورت مطرح شود که آیا این مدل‌ها با توجه به پیچیده و پردردسر بودن برآزاندن‌شان، واقعاً اطلاعاتی بیش از نتایج حاصل از برازش مدل‌های حاشیه‌ای یک‌متغیره به دست می‌دهند؟ پاسخ را به این صورت می‌توان بیان کرد:

۱. هنگامی که مدل‌های حاشیه‌ای را مورد بررسی قرار می‌دهیم، ممکن است رابطه‌ی متغیرهای کمکی مختص فرد یا اندام با پاسخ سمت راست (یا بالا) معنادار و با پاسخ سمت چپ (یا پایین) غیر معنادار باشد. چنین نتایجی ممکن است ما را در استنباط‌های نهایی دچار اشکال کند.
۲. مدل‌های دومتغیره این امکان را برای ما فراهم می‌سازند که برآوردی صریح از ضریب همبستگی بین متغیرهای پاسخ در حضور متغیرهای کمکی مختلف به دست آوریم. این امکان در مدل‌های یک‌متغیره وجود ندارد.
۳. در مدل‌های دومتغیره قادریم اثر عوامل خطر مختص یک اندام را در حالی که همبستگی بین این اندام و اندام دوم به‌طور همزمان در مدل احتساب شده است، مورد بررسی قرار دهیم. مدل‌های یک‌متغیره حاشیه‌ای چنین امکانی را در اختیار ما قرار نمی‌دهند.

با در نظر گرفتن موارد فوق، هنگامی که نمونه‌ای با اندازه‌ی بزرگ در اختیار داریم، بهتر است مزایای مدل‌های دومتغیره را در مقابل عیب بزرگ آن‌ها (پیچیده بودن برازش) قرار داده، در مورد انتخاب مدل مناسب تصمیم‌گیری کنیم.

مطلب دیگری که در این بخش باید به آن اشاره شود، این است که مدل‌های دومتغیره‌ی توصیف‌شده در بخش سوم را به‌سادگی می‌توان برای کاربردهای دیگر تعمیم داد. به‌عنوان مثال، می‌توان پارامترهای رگرسیونی مختلفی برای حاشیه‌ها در نظر گرفت. در این حالت، کافی است ماتریس طرح را به‌صورت زیر تغییر داد:

$$X'_i = \begin{pmatrix} x'_i & x'_{i1} & 0 \\ x'_i & 0 & x'_{i2} \end{pmatrix},$$

پارامترهای رگرسیونی متناظر را به صورت $(\delta, \epsilon_1, \epsilon_2)$ در نظر بگیریم، که در آن، ϵ_1 پارامترهای رگرسیونی مربوط به متغیرهای کمکی به دست آمده از سمت راست، و ϵ_2 پارامترهای رگرسیونی مرتبط با متغیرهای کمکی به دست آمده از سمت چپ هستند. همچنین در صورتی که داده‌ها ایجاب کند، قادریم نقاط برش متفاوتی برای حاشیه‌ها در نظر بگیریم. اگر افراد مورد مطالعه بر حسب دو مقیاس ترتیبی متفاوت دسته‌بندی شده باشند، در نظر گرفتن نقاط برش مختلف برای حاشیه‌ها ضروری به نظر می‌رسد. در مورد داده‌های CPI، از آنجایی که دسته‌بندی متغیر پاسخ برای هر دو حاشیه کاملاً یکسان بوده است، چنین تعمیمی منطقی نخواهد بود.

نکته‌ی قابل ذکر دیگر، انتخاب حدس‌های اولیه‌ی مناسب برای آغاز فرایند برآورد پارامترهای رگرسیونی است. اگر حدس‌های اولیه به‌طور نادرست و دور از مقادیر واقعی انتخاب شوند، ممکن است روش‌های تکراری در ماکسیمم‌های محلی (local) گرفتار شوند و برآوردهای گمراه‌کننده‌ای نتیجه دهند. بر اساس تجربه، ما از حدس‌های اولیه‌ی $\hat{\nu} = 1$ ، $\hat{\delta} = 0$ ، و $\hat{\epsilon} = 0$ برای آغاز این فرایند استفاده کردیم. همچنین برای انتخاب حدس‌های اولیه در مورد نقاط برش مدل، می‌توان از برآورد GEE یا نتایج به دست آمده از برازش مدل‌های حاشیه‌ای استفاده کرد. در این‌جا ذکر این نکته ضروری است که در صورت طراحی یک برنامه‌ی رایانه‌ای مناسب، برآوردهای نهایی با تغییرات جزئی در حدس‌های اولیه، تغییر نخواهند کرد. به‌عنوان مثال، در برنامه‌ی رایانه‌ای طراحی شده توسط ما و کیم، با انتخاب حدس‌های اولیه‌ی مناسب و ایجاد تغییرات جزئی در آن‌ها، مقادیر برآوردهای نهایی دستخوش تغییر نمی‌شوند. انتخاب حدس‌های اولیه‌ی اندکی نامناسب، فقط سبب افزایش تعداد تکرارهای لازم برای رسیدن به همگرایی نهایی، و در نتیجه موجب بالا رفتن مدت زمان پردازش رایانه‌ای خواهد شد.

مطلب دیگری که در این قسمت باید به آن اشاره شود این است که مدل‌های مبتنی بر توزیع‌های نهایی توأم، چه برتری‌ای نسبت به مدل‌های مبتنی بر برآوردهای GEE دارند. همان‌طور که از نتایج مندرج در جدول ۲ مشخص است، در مدل‌های مبتنی بر توزیع‌های نهایی دومتغیره، همبستگی بین مشاهدات پاسخ به‌عنوان یک پارامتر مهم در مدل وارد می‌شود و در پایان، برآوردی صریح از این پارامتر در حضور متغیرهای کمکی مختلف به دست می‌آید. اما در روش GEE، همبستگی بین مشاهدات پاسخ به‌عنوان یک پارامتر مزاحم در نظر گرفته می‌شود و به‌علت به کار بردن توزیع‌هایی حاشیه‌ای به‌جای توزیع توأم متغیرهای پاسخ، در حین فرایند مدل‌سازی، برآورد صریحی از این همبستگی به دست نمی‌آید. این مطلب می‌تواند موجب برازش ضعیف‌تر مدل‌های مبتنی بر روش GEE نسبت به مدل‌هایی با برآورد حد اکثر درست‌نمایی شود.

به‌عنوان توصیه‌ی نهایی، ذکر این نکته الزامی است که مدل‌های مبتنی بر توزیع‌های نهایی که در این مقاله به آن‌ها اشاره شد، بهتر است هنگامی مورد استفاده قرار گیرند که پاسخی دومتغیره و ترتیبی با همبستگی نسبتاً بالا در اختیار داشته باشیم. برای تصمیم‌گیری در مورد انتخاب مدل ما یا مدل پروبیت تجمعی

کیم، می‌توان از برازش‌های حاشیه‌ای مدل‌های یک‌متغیره استفاده کرد. اگر مدل‌های پروبیت یا لوزستیک تجمعی برای پاسخ‌های حاشیه‌ای، برازش‌های بهتری نسبت به مدل‌های غیرمتقارن مانند $\log\text{-}\log$ به دست دهند، انتظار داریم که مدل کیم بهتر از مدل ما عمل کند. در غیر این صورت، مدل ارائه شده از سوی ما احتمالاً (و نه لزوماً) جایگزینی مناسب برای مدل کیم خواهد بود.

مرجع‌ها

- [1] Ananth, C.V.; Kleinbaum, D.G. (1997). Regression models for ordinal responses: a review of methods and applications. *Int. J. Epidemiol.* **26**, 1323-1333.
- [2] Agresti, A. (2002). *Categorical Data Analysis*, 2nd ed. Wiley, New York.
- [3] Biswas, A.; Das, K. (2002). A Bayesian analysis of bivariate ordinal data: Wisconsin epidemiologic study of diabetic retinopathy revisited. *Statist. Med.* **21**, 549-559.
- [4] Dos Santos, D.M.; Berridge, D.M. (2000). A continuation ratio random effects model for repeated ordinal responses. *Statist. Med.* **19**, 3377-3388.
- [5] Gange, S.J.; Linton, K.L.; Scott, A.J.; DeMets, D.L.; Klein, R. (1995). A comparison of methods for correlated ordinal measures with ophthalmic applications. *Statist. Med.* **14**, 1961-1974.
- [6] Gumbel, E.J. (1961). Bivariate logistic distributions. *J. Amer. Statist. Assoc.* **56**, 335-349.
- [7] Kim, K. (1995). A bivariate cumulative probit regression model for ordered categorical data. *Statist. Med.* **14**, 1341-1352.
- [8] Liang, K.Y.; Zeger, S.L. (1986). Longitudinal data analysis using generalized linear models. *Biometrika* **73**, 13-22.
- [9] Lipsitz, S.R.; Kim, K.; Zhao, L. (1994). Analysis of repeated categorical data using generalized estimating equations. *Statist. Med.* **13**, 1149-1163.
- [10] Lipsitz, S.R.; Laird, N.M.; Harrington, D.P. (1990). Maximum likelihood regression methods for paired binary data. *Statist. Med.* **9**, 1517-1525.
- [11] McCullagh, P. (1980). Regression models for ordinal data (with discussion). *J. Roy. Statist. Soc. B* **42**, 109-142.
- [12] Rosner, B. (1984). Multivariate methods in ophthalmology with application to other paired-data situations. *Biometrics* **40**, 1025-1035.
- [13] Satterthwaite, S.P.; Hutchinson, T.P. (1978). A generalisation of Gumbel's bivariate logistic distribution. *Metrika* **25**, 163-170.
- [14] Williamson, J.M.; Kim, K.; Lipsitz, S.R. (1995). Analyzing bivariate ordinal data using a global odds ratio. *J. Amer. Statist. Assoc.* **90**, 1432-1437.

- [15] Williamson, J.M.; Lipsitz, S.R.; Kim, K.M. (1999). GEECAT and GEEGOR: computer programs for the analysis of correlated categorical response data. *Comput. Methods Programs Biomed.* **58**, 25-34.

دریافت: ۱ تیر ۱۳۸۳

آخرین اصلاح: ۲۷ آبان ۱۳۸۳

انتشار: ۲۰ دی ۱۳۸۳

فرید زابری

گروه آمار زیستی، دانشکده‌ی علوم پزشکی،

دانشگاه تربیت مدرس،

پل نصر، بزرگراه جلال آل احمد،

تهران، ایران.

پیام‌نگار: fzayeri@yahoo.com

انوشیروان کاظم نژاد

گروه آمار زیستی، دانشکده‌ی علوم پزشکی،

دانشگاه تربیت مدرس،

پل نصر، بزرگراه جلال آل احمد،

تهران، ایران.

پیام‌نگار: kazem_an@modares.ac.ir

Modeling Paired Ordinal Response Data

Anoshiravan Kazemnejad and Farid Zayeri

Tarbiat Modarres University

Abstract. About 25 years ago, McCullagh proposed a method for modeling univariate ordinal responses. After publishing this paper, other statisticians gradually extended his method, such that we are now able to use more complicated but efficient methods to analyze correlated multivariate ordinal data, and model the relationship between those responses and host of covariates. In this paper, we aim to present the recent progressions in modeling ordinal response data, especially in bivariate ordinal responses that arise from medical studies relating to paired organs such as ophthalmology, otology, nephrology, etc. Additionally, we present a new model for analyzing correlated ordinal response data. This model is an appropriate alternative for bivariate cumulative probit regression model, when joint distribution of response data is not symmetric. Finally, as an applied example, we analyze the obtained data from an epidemiologic study relating to periodontal status among high school students in Tehran using this method and compare the results with the similar models. © 2004 Statistical Research Center. All rights reserved.

Keywords. correlated ordinal responses; bivariate latent distribution; generalized estimating equations; generalized linear models.

