

تحلیل پاسخ‌های آمیخته‌ی گسسته و پیوسته‌ی رده‌بندی شده

احسان بهرامی‌سامانی^۱ و حمیدرضا عالی‌پور

بخش آمار، دانشگاه شهید بهشتی

تاریخ دریافت: ۹۲/۴/۳۱ تاریخ پذیرش: ۹۳/۳/۲۰

چکیده: در این مقاله، هدف ما رده‌بندی پاسخ‌هایی است که به‌صورت آمیخته‌ای از پاسخ‌های گسسته و پیوسته هستند. برای این‌کار، ابتدا باید تابع توزیع توأم چنین پاسخ‌هایی را داشته باشیم. بنابراین، مدل مکانی عام جدید را برای دستیابی به تابع توزیع توأم پاسخ‌های آمیخته‌ی گسسته و پیوسته معرفی کرده‌ایم و به‌صورت اجمالی با مدل مکانی عام دِ لئون و کریر [۱] مورد مقایسه قرار داده‌ایم. رویکرد مورد استفاده برای رده‌بندی چنین پاسخ‌هایی، استفاده از قاعده‌ی رده‌بندی ولج [۲] به‌همراه ارائه‌ی احتمال‌های بد-رده‌بندی و نرخ واقعی خطا برای این مدل است. همچنین نتایج رده‌بندی مدل مکانی عام جدید با نتایج رده‌بندی مدل مکانی عام مورد مقایسه قرار گرفته است [۱]. لازم به ذکر است قواعد رده‌بندی نمونه‌ای برای مدل مکانی عام جدید با استفاده از برآوردهای به‌دست آمده از طریق ماکسیمم‌سازی تابع درست‌نمایی مطرح می‌شود. در پایان، برای نشان دادن قابلیت قواعد به‌دست آمده، مطالعه‌ی شبیه‌سازی و تحلیل داده‌های واقعی با استفاده از مدل مکانی عام جدید انجام شده است.

واژه‌های کلیدی: قواعد رده‌بندی، مدل مکانی عام، پاسخ‌های آمیخته‌ی پیوسته و گسسته، تابع درست‌نمایی، احتمال بد-رده‌بندی.

رده بندی ریاضی: ۶۲J۰۵، ۶۲J۱۲.

۱- مقدمه

تفکیک و رده‌بندی داده‌ها، این‌که کدام داده به کدام جامعه تعلق دارد، همواره از مسائل مورد توجه در زمینه‌ی علوم پزشکی و سایر علوم بوده است. در این بین مواردی نیز پیش می‌آید که در آن پاسخ‌های به‌دست آمده در زمینه‌ای خاص، صرفاً از یک نوع، یعنی فقط پیوسته یا فقط گسسته، نیستند بلکه آمیخته‌ای از هر دو نوع، یعنی متشکل از پاسخ‌های گسسته و پاسخ‌های پیوسته هستند. مسئله وقتی جذاب‌تر می‌شود که بدانیم بین پاسخ‌ها همبستگی نیز وجود دارد. بنابراین، دیگر استفاده از روش‌های استاندارد قبلی کاری نادرست بوده و نتایج به‌دست آمده از آن‌ها غیر از آن نتایجی است که باید به‌دست آیند. این نوع پاسخ‌ها خود به انواع مختلفی تقسیم‌بندی می‌شوند که از آن جمله می‌توان به پاسخ‌های آمیخته‌ی پیوسته و اسمی، پاسخ‌های پیوسته، اسمی و ترتیبی و غیره اشاره کرد. آنچه در این بین پژوهش‌گر را به چالش می‌کشانند، مسئله‌ی تعیین توزیع توأم این نوع پاسخ‌ها برای دستیابی هر چه آسان‌تر به محاسبه‌ی مقدار همبستگی و نیز سایر خصیصه‌های مورد نیاز وی در این‌گونه پاسخ‌ها است. رویکردهای گوناگونی برای این منظور و توسط افراد مختلف ارائه شده است ولی اکثر آن‌ها دارای انعطاف‌پذیری کافی نیستند و شکل توزیعی استاندارد ندارند. تیت از جمله‌ی پیشگامان در این زمینه است [۳].

یکی از پیشنهادهای قابل‌توجه در این زمینه، رویکرد مستقیم توزیع توأم پاسخ‌های آمیخته است. به این صورت که توزیع توأم این نوع پاسخ‌ها، از طریق تجزیه به توزیع حاشیه‌ای مجموعه‌ای از پاسخ‌ها و توزیع شرطی سایر پاسخ‌ها به شرط این مجموعه به‌دست می‌آید. دو نوع مدل می‌توان برای این روش در نظر گرفت. در مدل اول که به مدل مکانی عام (GLOM)^۱ معروف است، توزیع توأم پاسخ آمیخته از طریق تجزیه به توزیع حاشیه‌ای پاسخ گسسته و توزیع شرطی پاسخ پیوسته به شرط پاسخ گسسته مهیا می‌شود و فرض بر این است که پاسخ‌های پیوسته به شرط پاسخ‌های گسسته دارای توزیع نرمال چند متغیره با میانگین‌های متفاوت و ماتریس کوواریانس یکسان و پاسخ‌های گسسته دارای توزیع چندجمله‌ای حاشیه‌ای باشند. در مدل دیگر که به مدل پیوسته‌ی گروه‌بندی شده‌ی شرطی (CGCM)^۲ معروف است، توزیع توأم با استفاده از تجزیه به توزیع حاشیه‌ای پاسخ پیوسته و توزیع شرطی پاسخ گسسته به شرط پاسخ پیوسته به‌دست می‌آید. در زمینه‌ی

1- General Location Model

2- Conditional Grouped Continuous Model

بررسی این نوع مدل، می‌توان به کاکس و ویرموس اشاره کرد [۴]. همچنین در این زمینه، کاتالانو و رایان در [۵] مدل داده‌های آمیخته بر پایه‌ی متغیر پنهان را مورد بررسی قرار داده‌اند. از سوی دیگر، دِ لئون و کریر در [۱]، یک مدل مکانی عام آمیخته با شامل متغیرهای اسمی و پاسخ‌های پیوسته و ترتیبی معرفی نمودند به طوری که در این مدل، جدول پیشابندی از تقاطع سطوح متغیرهای اسمی حاصل می‌شود و پاسخ‌های پیوسته و ترتیبی درون خانه‌های این جدول قرار می‌گیرند و توزیع توأم پاسخ‌های آمیخته‌ی پیوسته، ترتیبی و اسمی از طریق تجزیه به توزیع پاسخ‌های پیوسته و ترتیبی به شرط متغیرهای اسمی و توزیع حاشیه‌ای متغیرهای اسمی به دست می‌آید.

در ادامه دِ لئون و همکاران در [۶] برای اولین بار به بررسی رده‌بندی داده‌های آمیخته در مدل دِ لئون و کریر در [۱] پرداختند. اما آنچه در این مقاله مطرح می‌شود، پیشنهاد مدل جدیدی است که در آن جدول پیشابندی از تقاطع هم‌زمان سطوح پاسخ‌های اسمی و ترتیبی تشکیل می‌شود و این پاسخ‌های پیوسته هستند که درون خانه‌های جدول پیشابندی قرار می‌گیرند. تابع توزیع توأم نیز از طریق تجزیه به تابع توزیع شرطی پاسخ‌های پیوسته به شرط پاسخ‌های ترتیبی، تابع توزیع شرطی پاسخ‌های ترتیبی به شرط پاسخ‌های اسمی و تابع توزیع حاشیه‌ای پاسخ‌های اسمی به دست می‌آید. رویکرد رده‌بندی در این مقاله نیز قاعده‌ی رده‌بندی ولج [۲] خواهد بود و همان‌طور که اشاره خواهد شد، رده‌بندی این نوع پاسخ‌ها با استفاده از مدل جدید، نتایج بهتری در مقایسه با مدل دِ لئون و کریر [۱] ارائه می‌کند. برای این منظور، ابتدا احتمال‌های بد-رده‌بندی حاصل از مدل مکانی عام جدید با مورد مشابه برای مدل دِ لئون و کریر [۱] با استفاده از مطالعه‌ی شبیه‌سازی مورد مقایسه قرار گرفته و سپس برای تشریح سودمندی رده‌بندی مدل مکانی عام جدید، مطالعه‌ی روی بیماری کبد چرب و رده‌بندی بیماران صورت گرفته است.

این مقاله به این صورت سازماندهی شده است که در بخش دوم، مروری کوتاه بر مدل مکانی عام دِ لئون و کریر [۱] داشته، سپس به معرفی مدل مکانی عام جدید و نیز ارائه تابع درستیابی و قواعد رده‌بندی این مدل می‌پردازیم. در بخش سوم برای مقایسه‌ی مدل مکانی عام جدید و مدل مکانی عام دِ لئون و کریر [۱] از یک مطالعه‌ی شبیه‌سازی استفاده‌شده و در نهایت در بخش چهارم، قواعد رده‌بندی مربوط به مدل مکانی عام دِ لئون و کریر [۱] و مدل مکانی عام جدید بر روی بیماری کبد چرب بیماران بیمارستان طالقانی تهران پیاده‌سازی شده است.

۲- مدل‌های مربوط به پاسخ‌های آمیخته‌ی همبسته‌ی پیوسته و گسسته

در پژوهش‌های آماری مواردی اتفاق می‌افتد که پژوهش‌گر با دو نوع پاسخ پیوسته و گسسته به صورت هم‌زمان و در یک مدل روبرو می‌شود و پیچیدگی چنین مدل‌هایی وقتی ملموس‌تر می‌شود که بین چنین پاسخ‌هایی همبستگی نیز وجود داشته باشد. در این‌گونه موارد، به تنهایی استفاده از مدل‌های پاسخ‌های پیوسته یا حتی مدل‌های پاسخ‌های گسسته که در ادبیات آماری به وفور یافت می‌شوند، عملی ناشیانه است و نتایج اشتباهی را در بر دارد. آنچه که ماهیت این نوع پاسخ‌ها ایجاد می‌کند، یافتن روش مناسب برای تحلیل چنین پاسخ‌های آمیخته است. در سال‌های اخیر و با مطرح‌شدن تحلیل داده‌های آمیخته، فعالیت‌های پژوهشی گسترده‌ای در این زمینه صورت گرفته است و ماحصل آن بیان سه نوع رویکرد مختلف در حالت‌ها و مدل‌های مختلف است. لازم به ذکر است که در مدل‌هایی با پاسخ‌های آمیخته‌ی گسسته و پیوسته، بیش‌ترین چالش فراروی پژوهش‌گران آماری، نحوه‌ی به‌دست آوردن توزیع توأم این نوع پاسخ‌هاست. اولین رویکردی که برای تحلیل پاسخ‌های توأم پیوسته و گسسته استفاده می‌شود و از شیوه‌ای مستقیم برای یافتن تابع توزیع توأم این نوع پاسخ‌ها استفاده می‌کند، مدل تجزیه به عامل‌ها^۱ است. این رویکرد نیز خود به دو روش مدل مکانی عام (GLOM) و مدل پیوسته‌ی گروه‌بندی شده‌ی شرطی (CGCM) تقسیم می‌شود. در مدل مکانی عام، تابع توزیع توأم پاسخ‌های آمیخته‌ی گسسته و پیوسته از طریق تجزیه به توزیع حاشیه‌ای برای پاسخ‌های گسسته و توزیع شرطی پاسخ‌های پیوسته به شرط پاسخ‌های گسسته به‌دست می‌آید. در مدل پیوسته‌ی گروه‌بندی شده‌ی شرطی، توزیع توأم پاسخ‌های آمیخته‌ی گسسته و پیوسته از طریق تجزیه‌ی آن به توزیع حاشیه‌ای برای پاسخ‌های پیوسته و توزیع شرطی برای پاسخ‌های گسسته به شرط پاسخ‌های پیوسته به‌دست می‌آید. رویکرد دوم که به رویکرد غیرمستقیم شهرت دارد و به دو دسته‌ی مفصل‌ها [۷] و اثرهای تصادفی [۸] تقسیم می‌شود به همراه رویکرد سوم که شامل مدل‌های پلاکت-دیل [۹] و توزیع‌های نمایی درجه‌ی دو [۱۰] است، در این مقاله مورد مطالعه قرار نمی‌گیرند.

۲-۱- مدل مکانی عام با پاسخ‌های آمیخته‌ی پیوسته، ترتیبی و اسمی رده‌بندی شده

در این بخش، ابتدا مدل مکانی عام با پاسخ‌های آمیخته‌ی پیوسته و اسمی، آمیخته‌ی پیوسته و ترتیبی و مدل مکانی عام با پاسخ‌های آمیخته‌ی پیوسته، ترتیبی و اسمی معرفی می‌گردد [۱]. سپس مدل مکانی عام جدید با پاسخ‌های آمیخته‌ی پیوسته، ترتیبی و اسمی معرفی و مسئله‌ی رده‌بندی روی این مدل بیان می‌گردد. لازم به ذکر است مسئله‌ی رده‌بندی روی مدل مکانی عام جدید، توسط محققان دیگر مورد بررسی قرار نگرفته است.

۲-۱-۱- مقدمه‌ای بر مدل مکانی عام

فرض کنیم \mathbf{u} و \mathbf{n} به ترتیب به‌عنوان بردار پاسخ‌های پیوسته و اسمی باشند. از تقاطع سطوح متغیرهای اسمی، جدول پیش‌بینی D خانه به‌وجود می‌آید و بردار متغیرهای پیوسته درون خانه‌های این جدول قرار می‌گیرند. دِ لئون و کریر [۱] با تعمیم این مدل، مدل جدیدی را برای سه متغیر پیوسته، ترتیبی و اسمی آمیخته ارائه کرده‌اند. در این مدل همانند مدل قبلی، متغیرهای اسمی تشکیل جدول پیش‌بینی می‌دهند و متغیرهای پیوسته به‌همراه متغیرهای ترتیبی بر پایه‌ی متغیر پنهان درون خانه‌های این جدول قرار می‌گیرند.

مثال ۱. متغیر اسمی N_1 (۱ یا ۰) و متغیر پیوسته‌ی U را در نظر بگیریم. با توجه به دو سطحی بودن متغیر اسمی، جدول پیش‌بینی با دو خانه تشکیل می‌شود که متغیر پیوسته درون خانه‌های این جدول قرار می‌گیرد و با فرض ثابت بودن واریانس درون این دو خانه، میانگین آن از خانه‌ای به خانه‌ی دیگر تغییر می‌کند. این میانگین‌ها در شکل ۱ با دو نماد μ_0 و μ_1 نشان داده شده است.

N_1	
۰	۱
μ_0	μ_1

شکل ۱: مدل مکانی عام دِ لئون و کریر [۱] برای پاسخ‌های پیوسته و اسمی

مثال ۲. فرض کنیم متغیر پیوسته‌ی U و متغیرهای ترتیبی O_1 و O_2 به‌ترتیب با دو و سه سطح موجود باشند.

$$O_r = \begin{cases} 1 & U_r^* \leq \theta_1 \\ 2 & \theta_1 < U_r^* \leq \theta_2 \\ 3 & U_r^* > \theta_2 \end{cases} \quad \text{و} \quad O_1 = \begin{cases} 1 & U_1^* < \lambda \\ 2 & U_1^* \geq \lambda \end{cases}$$

که در آن U_1^* و U_r^* به ترتیب متغیرهای پنهان مربوط به متغیرهای ترتیبی O_1 و O_r و λ ، θ_1 و θ_2 نقاط آستانه‌ای مربوط به متغیرهای ترتیبی هستند. جدول پیشابندی حاصل از تقاطع سطوح متغیرهای ترتیبی، دارای شش خانه خواهد بود که متغیر پیوسته درون خانه‌های این جدول قرار می‌گیرد، منتهی میانگین این متغیر با فرض ثابت بودن واریانس آن از خانه‌ای به خانه‌ی دیگر فرق می‌کند که در شکل ۲ با نماد μ_j ($j=1, \dots, 6$) نشان داده شده است.

$U_1^* \geq \lambda$	$U_1^* < \lambda$	
μ_4	μ_1	$U_r^* \leq \theta_1$
μ_6	μ_2	$\theta_1 < U_r^* \leq \theta_2$
μ_5	μ_3	$U_r^* > \theta_2$

شکل ۲: مدل مکانی عام د لئون و کریر [۱] برای پاسخ‌های پیوسته و ترتیبی

مثال ۳. فرض کنیم متغیر پیوسته‌ی U_1 و U_r ، متغیرهای ترتیبی O_1 و O_r (متغیرهای پنهان متناظر با آن‌ها به ترتیب U_1^* و U_r^* هستند) به ترتیب با دو و سه سطح (مانند مثال ۲) و متغیرهای اسمی N_1 و N_r هر کدام با دو سطح موجود باشند.

N_1					
۱		۰			
U_1	U_1^*	U_1	U_1^*		
U_r	U_r^*	U_r	U_r^*	۰	
U_1	U_1^*	U_1	U_1^*	N_r	
U_r	U_r^*	U_r	U_r^*	۱	

شکل ۳: مدل مکانی عام د لئون و کریر [۱] برای پاسخ‌های پیوسته، اسمی و ترتیبی

جدول پیشابندی حاصل از تقاطع سطوح متغیرهای اسمی با چهارخانه تشکیل خواهد شد که متغیرهای پیوسته و ترتیبی درون خانه‌های آن قرار خواهند گرفت.

۲-۱-۲- مدل و تابع درست‌نمایی مدل مکانی عام جدید

در مدل جدید، در حالت وجود متغیرهای اسمی، ترتیبی و پیوسته جدول پیشابندی برخلاف مدل دِ لئون و کریر [۱] که از تقاطع متغیرهای اسمی تشکیل می‌شود، هم‌زمان از تقاطع سطوح متغیرهای اسمی و ترتیبی تشکیل می‌شود و متغیرهای پیوسته درون خانه‌های جدول پیشابندی قرار می‌گیرند. این مطلب با فرض وجود متغیرهای مثال قبلی در شکل ۴ نشان داده شده است.

N_r					
۱			۰		
$U_r^* > \theta_r$	$\theta_r < U_r^* \leq \theta_r$	$U_r^* \leq \theta_r$	$U_r^* > \theta_r$	$\theta_r < U_r^* \leq \theta_r$	$U_r^* \leq \theta_r$
U_1	U_1	U_1	U_1	U_1	U_1
U_r	U_r	U_r	U_r	U_r	U_r
U_1	U_1	U_1	U_1	U_1	U_1
U_r	U_r	U_r	U_r	U_r	U_r
U_1	U_1	U_1	U_1	U_1	U_1
U_r	U_r	U_r	U_r	U_r	U_r
U_1	U_1	U_1	U_1	U_1	U_1
U_r	U_r	U_r	U_r	U_r	U_r

شکل ۴: مدل مکانی عام جدید برای پاسخ‌های پیوسته، اسمی و ترتیبی

نکته‌ی حائز اهمیت این است که با فرض وجود بردار پاسخ‌های آمیخته‌ی اسمی (\mathbf{n}) ، ترتیبی (\mathbf{o}) و پیوسته (\mathbf{u}) ، تابع چگالی توأم آن‌ها در مدل مکانی عام دِ لئون و کریر [۱] به صورت $f(\mathbf{u}, \mathbf{o} | \mathbf{n})f(\mathbf{n}, \mathbf{o})$ و در مدل مکانی عام جدید به صورت $f(\mathbf{u} | \mathbf{n}, \mathbf{o})f(\mathbf{n}, \mathbf{o})$ تجزیه می‌شود و تعداد خانه‌های جدول مدل مکانی عام جدید بیشتر از تعداد خانه‌های جدول پیشابندی مدل مکانی عام دِ لئون و کریر [۱] است.

در حالت کلی، فرض کنیم \mathbf{n} و \mathbf{o} به ترتیب بردارهایی $S \times 1$ و $Q \times 1$ از متغیرهای اسمی و ترتیبی هستند، که زامین مؤلفه‌های آن‌ها دارای d_j و d_i^* حالت ممکن می‌باشند. زوج (\mathbf{n}, \mathbf{o}) تشکیل یک جدول پیشایندی با $D = \prod_{i=1}^Q \prod_{j=1}^S d_i^* d_j$ حالت ممکن یا به عبارتی D خانه را می‌دهند. بردار اسمی مانند $\mathbf{x} = (X_1, \dots, X_D)$ را می‌توان به گونه‌ای تعریف کرد که اگر (\mathbf{n}, \mathbf{o}) در حالت d ، $d = 1, 2, \dots, D$ ، قرار گیرد، X_d برابر یک و در غیر این صورت برابر صفر باشد.

بنابراین برای فرد خاصی، \mathbf{x} فقط شامل یک عدد ۱ خواهد بود. تحت مدل مکانی عام، \mathbf{x} دارای توزیع چندجمله‌ای است؛

$$\mathbf{x} | \boldsymbol{\pi} \sim \prod_{d=1}^D \pi_d^{X_d} \quad (1)$$

که در آن $\boldsymbol{\pi} = (\pi_1, \dots, \pi_D)^T$ برداری از احتمال‌های مربوط به خانه‌های جدول پیشایندی است. متغیرهای نشانگر را برای فرد k -ام $(k = 1, \dots, n)$ به صورت زیر تعریف می‌کنیم:

$$I_{n_k} = \begin{cases} 1 & N_k = n_k \\ 0 & o.w. \end{cases}, I_{u_k^*} = \begin{cases} 1 & \lambda_k < U_k^* \leq \lambda_{k+1} \\ 0 & o.w. \end{cases}, I_{u_k} = \begin{cases} 1 & U_k \leq u_k \\ 0 & o.w. \end{cases}$$

که در آن λ_k نقطه‌ی برشی برای متغیر پنهان U_k^* است. ارتباط بین متغیر پنهان و متغیر ترتیبی به وسیله‌ی مدل آستانه‌ای تعریف می‌شود، یعنی،

$$O_q = \lambda_q^{l_q} \Leftrightarrow \lambda_q^{l_q-1} < U_q^* \leq \lambda_q^{l_q}, \quad q = 1, \dots, Q$$

که در آن $\lambda_q^{l_q}$ نقاط آستانه‌ای نام دارند و مجموعه‌ی $\{\lambda_q^1, \dots, \lambda_q^{l_q}\}$ نقاط نامعلوم هستند. همچنین $\lambda_q^0 = -\infty$ و $\lambda_q^{l_q+1} = +\infty$. فرض کنیم متغیرهای ترتیبی O_1 و O_p و متغیرهای پنهان U_1^* و U_p^* به ترتیب متناظر با O_1 و O_p و متغیر پیوسته U_k موجود باشند. تابع درست‌نمایی مدل عبارت است از:

$$\begin{aligned}
& L(\boldsymbol{\eta} | \mathbf{u}, O_1, O_\tau, \mathbf{n}) \\
&= \prod_{k=1}^n \frac{\partial}{\partial u_k} \left\{ \Pr(U_k \leq u_k \mid \lambda_{k_1}^{l_{k_1}-1} < U_{k_1}^* \leq \lambda_{k_1}^{l_{k_1}}, \lambda_{k_\tau}^{l_{k_\tau}-1} < U_{k_\tau}^* \leq \lambda_{k_\tau}^{l_{k_\tau}}, n_k, \boldsymbol{\Sigma}) \right\} \\
&\quad \times \Pr(\lambda_{k_\tau}^{l_{k_\tau}-1} < U_{k_\tau}^* \leq \lambda_{k_\tau}^{l_{k_\tau}} \mid \lambda_{k_1}^{l_{k_1}-1} < U_{k_1}^* \leq \lambda_{k_1}^{l_{k_1}}, n_k, \lambda_q^1, \lambda_q^2, \dots, \lambda_q^{l_q}) \\
&\quad \times \Pr(\lambda_{k_1}^{l_{k_1}-1} < U_{k_1}^* \leq \lambda_{k_1}^{l_{k_1}} \mid n_k, \lambda_q^1, \lambda_q^2, \dots, \lambda_q^{l_q}) \times \Pr(N_k = n_k)
\end{aligned}$$

که در آن $\boldsymbol{\Sigma}$ ماتریس کوواریانس است. همچنین احتمال مربوط برای خانه‌ی d ام جدول پیشابندی برای $d = 1, 2, \dots, D$ را که در آن D تعداد کل خانه‌ها است، به‌صورت زیر تعریف می‌کنیم:

$$\pi_d = \Pr(\lambda_{k_1}^{l_{k_1}-1} < U_{k_1}^* \leq \lambda_{k_1}^{l_{k_1}}, \lambda_{k_\tau}^{l_{k_\tau}-1} < U_{k_\tau}^* \leq \lambda_{k_\tau}^{l_{k_\tau}} \mid N_k = n_k) \times \Pr(N_k = n_k).$$

با استفاده از تقریب توزیع نرمال چند متغیره که توسط جوی [۱۱] ارائه شده است، تابع درستنمایی به‌صورت زیر بیان می‌شود:

$$\begin{aligned}
L(\boldsymbol{\eta} | \mathbf{u}, O_1, O_\tau, \mathbf{n}) &\cong \prod_{k=1}^n \left[(\nu\pi)^{-\frac{1}{\nu}} |\boldsymbol{\Sigma}|^{-\frac{1}{\nu}} \exp \left\{ -\frac{1}{\nu} (\mathbf{u}_k - \boldsymbol{\mu}_d)^T \boldsymbol{\Sigma}^{-1} (\mathbf{u}_k - \boldsymbol{\mu}_d) \right\} \right. \\
&\quad \left. + \frac{\partial}{\partial u_k} \boldsymbol{\omega}_{\nu 1}^T \boldsymbol{\omega}_{\nu 1}^{-1} \left[1 - E(I_{u_{k_1}}^*), 1 - E(I_{u_{k_\tau}}^*), 1 - E(I_{n_k}) \right]^T \right] \pi_d
\end{aligned} \quad (2)$$

که در آن

$$E(I_{u_i}^*) = P(\lambda_i < U_i^* \leq \lambda_{i+1}), \lambda_q^{l_q+1} = +\infty, \lambda_q^0 = -\infty$$

$$\boldsymbol{\omega}_{\nu 1} = \left(\boldsymbol{\Sigma}_{I_{u_k}, I_{u_{k_1}}}, \boldsymbol{\Sigma}_{I_{u_k}, I_{u_{k_\tau}}}, \boldsymbol{\Sigma}_{I_{u_k}, I_{n_k}} \right)$$

$$\boldsymbol{\omega}_{\nu 1}^{-1} = \begin{pmatrix} \boldsymbol{\Sigma}_{I_{u_{k_1}}^*, I_{u_{k_1}}^*} & \boldsymbol{\Sigma}_{I_{u_{k_1}}^*, I_{u_{k_\tau}}^*} & \boldsymbol{\Sigma}_{I_{u_{k_1}}^*, I_{n_k}} \\ & \boldsymbol{\Sigma}_{I_{u_{k_\tau}}^*, I_{u_{k_\tau}}^*} & \boldsymbol{\Sigma}_{I_{u_{k_\tau}}^*, I_{n_k}} \\ & & \boldsymbol{\Sigma}_{I_{n_k}, I_{n_k}} \end{pmatrix}$$

که در آن‌ها،

$$\begin{aligned} \Sigma_{I_{u_i^*}, I_{u_j^*}} &= \text{cov}(I_{u_i^*}, I_{u_j^*}) = P(\lambda_i < U_i^* \leq \lambda_{i+1}, \lambda_j < U_j^* \leq \lambda_{j+1}) \\ &- P(\lambda_i < U_i^* \leq \lambda_{i+1})P(\lambda_j < U_j^* \leq \lambda_{j+1}). \end{aligned}$$

۲-۱-۳- رده‌بندی مدل مکانی عام جدید

در این بخش، ابتدا تابع چگالی توأم پاسخ‌های پیوسته، اسمی و ترتیبی را که از طریق مدل مکانی عام جدید به دست آمده است، مطرح نموده، سپس با استفاده از قاعده‌ی رده‌بندی ولج [۲]، قاعده‌ی رده‌بندی برای آن ارائه می‌شود. برای این منظور، فرض کنیم \mathbf{u} ، \mathbf{n} ، ترتیبی و پنهان باشند، با استفاده از مدل مکانی عام جدید، تابع چگالی توأم آن‌ها به صورت زیر بیان می‌شود ($d = 1, 2, \dots, D$):

$$\begin{aligned} f_{\mathbf{O}, \mathbf{U}, \mathbf{N}}(\mathbf{o}, \mathbf{u}, \mathbf{n}) &= f(\mathbf{u} | \mathbf{o}, \mathbf{n}) f(\mathbf{o}, \mathbf{n}) \\ &\equiv \left\{ (\nu \boldsymbol{\pi})^{-1/\Delta} |\boldsymbol{\Sigma}|^{-1/\Delta} \exp \left\{ -\frac{1}{\nu} (\mathbf{u} - \boldsymbol{\mu}_d)^T \boldsymbol{\Sigma}^{-1} (\mathbf{u} - \boldsymbol{\mu}_d) \right\} \right. \\ &\quad \left. + \frac{\partial}{\partial \mathbf{u}} \boldsymbol{\omega}_{\nu, \boldsymbol{\omega}_{\nu}^{-1}} \left(1 - E(I_{u_i^*}), 1 - E(I_{u_j^*}), 1 - E(I_n) \right) \right\} \boldsymbol{\pi}_d \end{aligned} \quad (3)$$

که در آن بردار میانگین \mathbf{u} برای خانه‌ی d ام جدول پیشابندی و $\boldsymbol{\Sigma}$ ماتریس کوواریانس \mathbf{u} است.

فرض کنیم $\boldsymbol{\Pi}^{(1)}$ و $\boldsymbol{\Pi}^{(\nu)}$ دو جامعه با احتمال‌های پیشین، هزینه‌های بد-رده‌بندی و ماتریس کواریانس برابر باشند. مشاهده‌ی آمیخته‌ی $\mathbf{w}' = \{\mathbf{o}, \mathbf{u}, \mathbf{n}\}$ طبق قاعده‌ی رده‌بندی ولج [۲] و با استفاده از تابع چگالی توأم پاسخ‌های آمیخته‌ی به دست آمده در رابطه‌ی (۳) به جامعه‌ی $\boldsymbol{\Pi}^{(1)}$ اختصاص داده می‌شود هرگاه

$$\mathcal{R}_\nu: \frac{\left[a \exp \left\{ -\frac{1}{\nu} (\mathbf{u} - \boldsymbol{\mu}_d^{(1)})^T \boldsymbol{\Sigma}^{-1} (\mathbf{u} - \boldsymbol{\mu}_d^{(1)}) \right\} + t \right]}{\left[a \exp \left\{ -\frac{1}{\nu} (\mathbf{u} - \boldsymbol{\mu}_d^{(\nu)})^T \boldsymbol{\Sigma}^{-1} (\mathbf{u} - \boldsymbol{\mu}_d^{(\nu)}) \right\} + t \right]} \geq \frac{\pi_d^{(\nu)}}{\pi_d^{(1)}} \quad (4)$$

و در غیر این صورت به جامعه‌ی $\Pi^{(r)}$ اختصاص داده می‌شود که در آن مقادیر $t = \frac{\partial}{\partial \mathbf{u}} \boldsymbol{\omega}_r \boldsymbol{\omega}_r^{-1} \left(1 - E(I_{\mathbf{u}_r^*}), 1 - E(I_{\mathbf{u}_r^*}), 1 - E(I_{\mathbf{n}}) \right)^T$ و $a = (\sqrt{2\pi})^{-d} |\boldsymbol{\Sigma}|^{-1/2}$ ثابتی هستند. $\boldsymbol{\mu}_d^{(g)}$ نیز میانگین متغیر پیوسته‌ی \mathbf{u} در خانه‌ی d ام جدول پیشابندی مربوط به جامعه‌ی g ام ($g = 1, 2$) است. به عبارت دیگر، \mathbf{w}' به جامعه‌ی $\Pi^{(1)}$ اختصاص داده می‌شود هرگاه $\mathbf{w}' \in \mathcal{R}_1$ و به جامعه‌ی $\Pi^{(2)}$ اختصاص داده می‌شود هرگاه $\mathbf{w}' \in \mathcal{R}_2$. لازم به ذکر است که با استفاده از نمونه‌ی مستقل به اندازه‌ی $N^{(1)}$ و $N^{(2)}$ به ترتیب از جامعه‌های $\Pi^{(1)}$ و $\Pi^{(2)}$ ، با جایگذاری برآورد پارامترها (به دست آمده از تابع درستنمایی مربوطه) \mathcal{R}_1 به صورت زیر برآورد می‌شود.

$$\hat{\mathcal{R}}_1 : \frac{\left[(\sqrt{2\pi})^{-d} |\mathbf{S}_p|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{u} - \bar{\mathbf{u}}_d^{(1)})^T \mathbf{S}_p^{-1} (\mathbf{u} - \bar{\mathbf{u}}_d^{(1)}) \right\} + \hat{t} \right]}{\left[(\sqrt{2\pi})^{-d} |\mathbf{S}_p|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{u} - \bar{\mathbf{u}}_d^{(2)})^T \mathbf{S}_p^{-1} (\mathbf{u} - \bar{\mathbf{u}}_d^{(2)}) \right\} + \hat{t} \right]} \geq \frac{\hat{\pi}_d^{(2)}}{\hat{\pi}_d^{(1)}} \quad (5)$$

که در آن ماتریس کواریانس نمونه‌ای ادغام‌شده، $\bar{\mathbf{u}}_d^{(g)}$ بردار میانگین نمونه‌ای متغیر پیوسته برای خانه‌ی d ام در جامعه‌ی g ام ($g = 1, 2$) هستند. بنابراین، می‌توان چنین بیان کرد که مشاهده‌ی آمیخته‌ی $\mathbf{w}' = \{\mathbf{o}, \mathbf{u}, \mathbf{n}\}$ به جامعه‌ی $\Pi^{(1)}$ اختصاص داده می‌شود هرگاه رابطه‌ی (5) برقرار باشد در غیر این صورت این مشاهده به جامعه‌ی $\Pi^{(2)}$ اختصاص داده خواهد شد. قاعده‌ی رده‌بندی بهینه‌ی مدل مکانی جدید را می‌توان برای موارد بیش از دو جامعه تعمیم داد. با فرض این‌که $\Pi^{(1)}, \dots, \Pi^{(G)}$ ، $G \geq 2$ ، جوامعی باشند که فقط در بردار پارامترهای مکانی باهم اختلاف دارند، آنگاه طبق قاعده‌ی رده‌بندی بهینه‌ی ولج [۲] و بر پایه‌ی مدل مکانی جدید، با فرض برابری هزینه‌های بد-رده‌بندی برای تمام جوامع، مشاهده‌ی آمیخته‌ی $\mathbf{w}' = \{\mathbf{o}, \mathbf{u}, \mathbf{n}\}$ به جامعه‌ای که برای آن

$$\delta^{(g)} = \log \alpha^{(g)} + \log \pi_d^{(g)} + \log \left[a \exp \left\{ -\frac{1}{2} (\mathbf{u} - \boldsymbol{\mu}_d^{(g)})^T \boldsymbol{\Sigma}^{-1} (\mathbf{u} - \boldsymbol{\mu}_d^{(g)}) \right\} + t \right], \quad g = 1, 2, \dots, G \quad (6)$$

بیشترین مقدار را در بین سایر جوامع داشته باشد، اختصاص داده خواهد شد.

محاسبه‌ی احتمال بد-رده‌بندی. با فرض $G=2$ و برابری احتمال‌های پیشین و هزینه‌ی بد-رده‌بندی برای هر دو جامعه‌ی $\Pi^{(1)}$ و $\Pi^{(2)}$ ، احتمال بد-رده‌بندی بر پایه‌ی مدل مکانی عام جدید به صورت زیر بیان می‌شود:

$$\Pr(g | g') = \sum_d \sum_{\ell} \left(\int_{\mathcal{R}_1} I(g'=1) + \int_{\mathcal{R}_1} I(g'=2) \right) \left(\int \{\phi_c^{(g')}(\mathbf{u}) + t\} \pi_d^{(g')} d\mathbf{u} \right) \quad (7)$$

که در آن $\phi_c(\mathbf{u})$ تابع چگالی نرمال C -متغیره با بردار میانگین $\boldsymbol{\mu}_d$ و ماتریس کوواریانس Σ و $I(\cdot)$ تابع نشان‌گر است.

محاسبه‌ی نرخ واقعی خطا. نرخ واقعی خطا (AER)^۱ به صورت زیر بیان می‌شود:

$$AER = \frac{P(1|2) + P(2|1)}{2}. \quad (8)$$

برای برآورد میانگین رابطه‌ی (۸)، ابتدا برآورد ماکسیمم درست‌نمایی پارامترهای موجود در مدل مکانی عام جدید با استفاده از دستور "nlminb" در نرم‌افزار R که برای می‌نیمم کردن توابع غیرخطی و غیره استفاده می‌شود، به دست می‌آید و سپس انتگرال‌های موجود در (۷) با استفاده از روش‌های مونت کارلو تعیین می‌شوند. همچنین شیوه‌ی دیگری نیز به نام «برآورد جدا نگه داشته شده»^۲ مطرح است، برای اطلاعات بیشتر به [۱۲] رجوع شود.

۳- مطالعه‌ی شبیه‌سازی

در این بخش، با استفاده از یک مطالعه‌ی شبیه‌سازی، قواعد رده‌بندی مدل مکانی عام د لئون و کریر [۱] و مدل مکانی عام جدید مورد بررسی قرار گرفته، نرخ خطای بد-رده‌بندی مدل مکانی عام د لئون و کریر [۱] را با نرخ خطای بد-رده‌بندی مدل مکانی عام جدید مورد مقایسه قرار داده‌ایم.

به منظور شبیه‌سازی، دو جامعه‌ی $\Pi^{(1)}$ و $\Pi^{(2)}$ با یک متغیر پیوسته، یک بردار اسمی با دو حالت و یک متغیر ترتیبی با دو سطح در نظر گرفته شده است. پارامترهای مدل نیز به-

1- Actual Error Rate

2- Holdout Method

صورت $\Theta_1 = (\pi, \mu_1, \mu_2, \lambda)^T$ و $\Theta_2 = (\sigma^2, \rho, \beta)^T$ نشان داده شده‌اند که در آن μ_s بیانگر میانگین متغیر پیوسته‌ی U در خانه‌ی s ، λ نقطه‌ی آستانه‌ای برای متغیر پنهان U^* تحت متغیر ترتیبی O است و همچنین ρ ضریب همبستگی بین U و U^* است. متغیر ترتیبی O به صورت زیر تعریف می‌شود:

$$O = \begin{cases} 1 & U^* \leq \lambda \\ 2 & U^* > \lambda \end{cases}$$

که در آن λ نقطه‌ی برشی برای متغیر پنهان U^* است. مشابه دِ لئون و کریبر [۱] چهار حالت زیر در نظر گرفته شده و مورد بررسی قرار می‌گیرد:

(I) هر دو جامعه، در بردار اسمی با هم اختلاف دارند؛

(II) هر دو جامعه، در متغیر پیوسته با هم اختلاف دارند؛

(III) هر دو جامعه، در متغیر ترتیبی با هم اختلاف دارند؛

(IV) جوامع نسبت به همه متغیرها با هم اختلاف دارند.

برای اجرای قاعده‌ی رده‌بندی، نمونه‌های مستقل به اندازه‌های

$$(N^{(1)}, N^{(2)}) = (50, 100), (100, 100), (200, 250)$$

برای دو جامعه‌ی $\Pi^{(1)}$ و $\Pi^{(2)}$ تولید می‌شود. برای تولید داده‌های مورد نظر، ابتدا متغیر اسمی X را از توزیع برنولی و سپس متغیرهای پیوسته و متغیر پنهان مربوط به متغیر ترتیبی را از توزیع نرمال دو متغیره با استفاده از مقادیر پارامتری زیر برای چهار حالت تولید می‌نماییم.

$$\Theta_2^{(g)} = (\sigma^2, \rho, \beta)^T = (1, 0, 5, 0, 58)^T, \quad g = 1, 2 \quad \text{(I-IV)}$$

$$\Theta_1^{(2)} = (0, 1, 0, 2, -1/15)^T \quad \text{و} \quad \Theta_1^{(1)} = (0, 3, 0, 2, -1/15)^T \quad \text{(I)}$$

$$\Theta_1^{(2)} = (0, 5, 4, 0, 2/31)^T \quad \text{و} \quad \Theta_1^{(1)} = (0, 5, 3, -1, 2/31)^T \quad \text{(II)}$$

$$\Theta_1^{(2)} = (0, 5, 0, 0, 5, -2/0.2)^T \quad \text{و} \quad \Theta_1^{(1)} = (0, 5, 0, 0, 5, -0/29)^T \quad \text{(III)}$$

$$\Theta_1^{(2)} = (0, 1, 0, 0, -1/12)^T \quad \text{و} \quad \Theta_1^{(1)} = (0, 3, 2, 5, 1/1, 0/81)^T \quad \text{(IV)}$$

ابتدا دو مدل مکانی عام دِ لئون و کرییر [۱] و مدل مکانی عام جدید را روی داده‌های تولیدشده در چهار حالت برآزش داده، برآورد پارامترهای هر یک از این مدل‌ها را به دست می‌آوریم. سپس تولید داده‌ها و برآورد پارامترهای هر دو مدل را ۱۰۰۰ بار تکرار می‌کنیم. در نهایت، متوسط تعداد نمونه‌های بد-رده‌بندی شده در هر دو جامعه را برای هر دو مدل با استفاده از احتمال بد-رده‌بندی شده به دست می‌آوریم. نتایج، در جدول (۱) آمده است.

جدول (۱) بیان‌گر متوسط تعداد افرادی است که در هر دو نمونه‌ی $N^{(1)}$ و $N^{(2)}$ بد-رده‌بندی شده‌اند. مثلاً برای نمونه‌ای به اندازه‌ی ۴۵۰ (جمع اندازه‌ی نمونه‌ی هر دو جامعه) در حالت (IV)، عدد ۲۳/۱۵ مربوط به مدل مکانی جدید در مقایسه با مدل دِ لئون و کرییر [۱]، (عدد ۲۴/۳۳) نشان دهنده‌ی رده‌بندی بهتر (متوسط تعداد بد-رده‌بندی‌های پایین) است. این نتیجه برای چهار حالت بیان شده در نمونه‌های مختلف نیز صدق می‌کند و نشان دهنده‌ی این است که مدل مکانی عام جدید بهتر از مدل مکانی عام دِ لئون و کرییر [۱] عمل می‌کند.

جدول ۱: متوسط تعداد نمونه‌های بد-رده‌بندی شده در هر دو جامعه

	مدل مکانی عام جدید	مدل مکانی عام دِ لئون و کرییر [۱]	$N^{(1)}$	$N^{(2)}$	
	۳۸/۴۵	۴۰/۴۳	۵۰	۱۰۰	
	۴۷/۵۶	۴۸/۵۹	۱۰۰	۱۰۰	I
	۱۱۳/۶۶	۱۱۵/۹۰	۲۰۰	۲۵۰	
	۴۱/۱۲	۴۳/۰	۵۰	۱۰۰	
	۶۰/۸۷	۶۱/۴۴	۱۰۰	۱۰۰	II
	۱۳۳/۹۹	۱۳۴/۷۷	۲۰۰	۲۵۰	
	۳۵/۹۶	۳۷/۰۱	۵۰	۱۰۰	
	۱۱/۱۲	۱۲/۵۲	۱۰۰	۱۰۰	III
	۲۵/۴۵	۲۷/۴۹	۲۰۰	۲۵۰	
	۱۰/۱۴	۱۱/۳۱	۵۰	۱۰۰	
	۱۱/۳۴	۱۲/۳۲	۱۰۰	۱۰۰	IV
	۲۳/۱۵	۲۴/۳۳	۲۰۰	۲۵۰	

دلیل اصلی آن، این است که در مدل عام جدید، ترتیبی بودن متغیر در مدل وارد شده و سبب ایجاد خانه‌های بیشتر در مدل مکانی عام شده است در صورتی که در مدل مکانی عام دِ لئون و کرییر [۱]، ترتیبی بودن نقشی همانند متغیر پیوسته در جدول پیشابندی بازی می‌کند و نسبت به مدل مکانی عام جدید اطلاعات کم‌تری را در خود دارد و می‌توان

این‌گونه استنباط کرد که شیوه‌ی رده‌بندی با افزایش اطلاعات (در اینجا با اضافه کردن خانه‌های جدول پیش‌بینی) بهبود می‌یابد. در نتیجه در مسئله‌ی رده‌بندی مشاهدات، مدل مکانی عام جدید در مقایسه با مدل دِ لئون و کریبر [۱]، توصیه می‌گردد.

۴- کاربرد روی داده‌های آمیخته‌ی بیماری کبد چرب

در این مقاله، از داده‌های بیماری کبد چرب که از مطالعه‌ای بر روی بیماران بیمارستان طالقانی تهران جمع‌آوری شده بود، استفاده شده است. بیماری کبد چرب، از تجمع زیاد چربی در اطراف کبد حادث می‌شود و سبب مختل شدن کار کبد در بدن می‌گردد. از سویی این بیماری در ایران بسیار شایع بوده و علائم خاصی نیز ندارد. به همین منظور با استفاده از متغیرهای موثر روی کبد چرب و با استفاده از مدل مکانی عام جدید به دنبال رده‌بندی درست بیماران مبتلا به کبد چرب هستیم. متغیرهای موثر در نظر گرفته شده در بیماری کبد چرب عبارتند از: میزان چاقی که بیان‌کننده نسبت وزن به مجذور قد افراد است (متغیر پیوسته)، وجود/عدم وجود علائم بیماری کبد چرب (متغیر اسمی، ۱: وجود علائم، ۲: عدم وجود علائم) و شدت بیماری کبد چرب (متغیر ترتیبی، با سه سطح ۱: فرد کبد چرب ندارد، ۲: فرد کبد چرب متوسطی دارد و ۳: فرد کبد چرب شدیدی دارد). بیش‌ترین تمرکز این مطالعه، بررسی و شناسایی در تفکیک بیماران نیازمند بستری و درمان سرپایی است. متغیرهای آمیخته براساس $N^{(1)} = 146$ بیمار سرپایی (بیمارانی که ترخیص شدند، جامعه‌ی اول) و $N^{(2)} = 38$ بیمار بستری (بیمارانی که در بیمارستان بستری شدند، جامعه‌ی دوم) مورد بررسی قرار گرفته است.

با استفاده از قاعده‌ی رده‌بندی حاصل از مدل مکانی عام جدید (رابطه‌ی ۴) و جایگزین کردن برآورد ماکسیمم درست‌نمایی پارامترهای مدل مکانی عام جدید، نتایج مربوط به قاعده‌ی رده‌بندی برای بیمار مراجعه‌کننده، به‌دست می‌آید. برای مثال، قاعده‌ی رده‌بندی برای بیمار مراجعه‌کننده به این صورت خواهد بود که اگر یک بیمار بدون شدت کبد چرب و با وجود علائم بیماری کبد چرب باشد، آنگاه برای بستری پذیرفته می‌شود هرگاه رابطه‌ی ۴ برقرار باشد در غیر این صورت بیمار ترخیص می‌شود. قواعد رده‌بندی برای همه‌ی وضعیت‌ها و سطوح متغیر کبد چرب و متغیر وجود/عدم وجود علائم بیماری کبد چرب برآورد می‌شوند. لازم به ذکر است که این نتایج برای استفاده در موقعیت‌های بالینی که نیاز به تصمیم‌گیری برای بستری یا ترخیص احساس می‌شود، می‌توانند به‌کاربرده شوند. نتایج رده‌بندی‌های صورت گرفته براساس مدل مکانی عام جدید و مدل مکانی دِ لئون و

کریر [۱]، با فرض برابری احتمال‌های پیشین و همچنین برابری هزینه‌های بد-رده‌بندی برای هر دو گروه بیمارها در جدول (۲) نشان داده شده است.

جدول ۲: تعداد (درصد) نمونه‌های بد-رده‌بندی شده مدل مکانی عام جدید و مدل دلئون و کریر [۱] برای داده‌های کبد چرب

گروه	اندازه‌ی نمونه	مدل مکانی جدید	مدل مکانی عام دلئون و کریر [۱]
بیماران بستری شده	۳۸	۱۰ (۲۶/۳۱)	۱۳ (۳۴/۲۱)
بیماران سرپایی	۱۴۶	۴۵ (۳۰/۸۲)	۴۸ (۳۲/۸۷)
جمع	۱۸۴	۵۵ (۲۹/۸۹)	۶۱ (۳۳/۱۵)

نتایج جدول (۲) حاکی از آن است که در مدل مکانی دلئون و کریر [۱]، ۱۳ بیمار از گروه بیماران بستری شده در گروه بیماران سرپایی بد-رده‌بندی شده و به اشتباه ترخیص شده‌اند. با توجه به مقدار ۲۶/۳۱ درصد بیماران ترخیص شده که در نوع خود عدد بزرگی است برای بد-رده‌بندی بیمارانی که باید بستری شوند و نشده‌اند، باید به دنبال راهکارهای کاهش خطای بد-رده‌بندی با افزایش میزان بازدهی کارکنان و مسئولین بهداشتی و پزشکان ناظر باشیم. همچنین ملاحظه می‌شود در اینجا نیز مدل مکانی عام جدید، بهتر از مدل مکانی عام دلئون و کریر [۱]، مشاهدات آمیخته را رده‌بندی می‌نماید.

۵- بحث و نتیجه‌گیری

در این مقاله، با استفاده از مدل مکانی عام جدید توانستیم تفکیک و رده‌بندی داده‌های آمیخته را نسبت به مدل دلئون و کریر [۱]، بهبود بخشیم. رده‌بندی بهتر این مدل به خاطر افزایش خانه‌های جدول پیش‌بینی حاصل از ورود متغیر ترتیبی در ساختن جدول پیش‌بینی، ایجاد شده است. از سویی دیگر این موضوع تاکنون مورد بررسی قرار نگرفته است و می‌توان آن را به عنوان موضوع جدیدی در حوزه‌های زیر مورد توجه قرار داد: الف) برای پاسخ‌های دارای گم‌شدگی غیرقابل چشم‌پوشی؛ ب) برای پاسخ‌های طولی.

مراجع

[1] De Leon, A.R. and Carrière, K.C. (2007). General mixed-data model: Extension of general location and grouped continuous models. *The Canadian Journal of Statistics*, **35**, 533 - 548.

- [2] Welch, B.L. (1939). Note on discriminant functions. *Biometrika*, **31**, 218–220.
- [3] Tate, F.R. (1954). Correlation between discrete and continuous variables: Point-Biserial correlation. *Ann. Math. Statist.*, **25**, 603–607.
- [4] Cox, D.R. and Wermuth, N. (1992). Response models for mixed binary and quantitative variables. *Biometrika*, **79**(3), 441–461.
- [5] Catalano, P. and Ryan, L.M. (1992). Bivariate latent variable models for clustered discrete and continuous outcomes. *Journal of the American Statistical Association*, **87**, 651–658.
- [6] De Leon, A.R., Soo, A. and Williamson, T. (2011). Classification with discrete and continuous variables via general mixed-data models. *Journal of Applied Statistics*, **5** (38), 1021–1032.
- [7] Trivedi, P. and Zimmer, D. (2006). *Copula modelling in econometrics: Introduction to practitioners*. Wiley, New York.
- [8] Faes, C., Aerts, M., Molenberghs, G., Geys, H., Teuns, G. and Bijnes, L. (2008). A high-dimensional joint model for longitudinal outcomes of different nature. *Stat. Med.*, **27**, 4408–4427.
- [9] Molenberghs, G. and Verbeke, G. (2005). *Models for Discrete Longitudinal Data*. Springer-Verlag, New York.
- [10] Sammel, M.D. Ryan, L.M. and Legler J.M. (1997). Latent variable models for mixed discrete and continuous outcomes. *Journal of the Royal Statistical Society*, **B**, 59, 667–678.
- [11] Joe, H. (1995). Approximations multivariate normal rectangle probabilities based on conditional expectation. *Journal of the American Statistical Association*, **90**: 957–967.
- [12] Johnson, R.A. and Wichern, D.W. (2007). *Applied multivariate statistical analysis*. 6th ed., Prentice Hall, New York.

Analysis of Mixed Discrete and Continuous Classified Responses

Ehsan Bahrami Samani and Hamid Reza Aalipour

Department of Statistics, Shahid Beheshti University, Tehran, Iran

Abstract

In this paper, classification of mixed discrete and continuous responses is our goal. To do this, first we have the joint distribution function of these data. Therefore, we propose new general location model to obtain the joint distribution function of mixed discrete and continuous responses and briefly compare it with the De Leon and Carrière's general location model [1]. The classification approach is based on Welch's classification rules [2] by providing misclassification probabilities and actual error rate. In addition, the results of classification of the new general location model are compared with the same one of the De Leon and Carrière's general location model [1]. It is noted that, the sample classification rules of new general location model are provided through obtaining parameters estimation of the maximum likelihood function. Finally, to show the performance of the rules, simulation study and analysis of a real data set using new general location model, is performed.

Keywords: Classification rule; General location model; Continuous and Discrete mixed data; Likelihood Function; Misclassification Probability.

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