

نمایی بر مدل‌های رده پنهان برای مدل‌بندی توأم اندازه‌گیری‌های طولی و داده‌های بقا

پروانه مهدی‌زاده^۱، تابان باغفلکی^۲، مهدی اسماعیلیان^۱

تاریخ دریافت: ۱۳/۰۶/۱۴۰۰

تاریخ پذیرش: ۰۱/۰۸/۱۴۰۰

چکیده:

مدل‌های توأم در مطالعات پیگیری شونده برای بررسی ارتباط بین اندازه‌گیری‌های نشانگر طولی و یک پیشامد بقا استفاده می‌شود و به وضعیت‌هایی با چند نشانگر طولی و یا مخاطره‌های رقابتی تعمیم‌یافته است. بسیاری از دستاوردهای آماری در زمینه مدل‌بندی توأم در مدل‌های پارامتر مشترک متمرکز شده است که شامل مشخصه‌هایی از نشانگر طولی به‌عنوان متغیرهای تبیینی در مدل بقا در نظر گرفته می‌شود. یک رهیافت کمتر شناخته‌شده، مدل رده پنهان توأم است، این مدل با فرض اینکه ارتباط بین نشانگرهای طولی و خطر رخداد با یک ساختار رده پنهان کاملاً مشخص می‌شود، بنا شده است. مدل رده پنهان به دلیل انعطاف‌پذیری در مدل‌بندی ارتباط بین نشانگرهای طولی و زمان تا رخداد پیشامد و همچنین توانایی در برگرفتن متغیرهای تبیینی به‌ویژه برای پیش‌بینی مناسب است. در این مقاله یک نمای کلی از مدل رده پنهان توأم و تعمیم‌های آن ارائه می‌دهیم، در این راستا، ابتدا مروری بر مدل‌های بحث شده انجام می‌شود و سپس برآورد پارامترهای مدل مورد بحث قرار می‌گیرد. در بخش کاربرد، دو مجموعه‌ی داده‌ی واقعی مورد تحلیل و بررسی قرار می‌گیرند.

واژه‌های کلیدی: الگوریتم EM، برآوردگر درستمایی ماکسیمم، داده‌های طولی، مدل بقا، مدل رده پنهان، مخاطره‌های رقابتی.

۱ مقدمه

کرد [۳۲]. در برخی کاربردهای دیگر، هدف، بررسی روند نشانگر طولی در طی یک بیماری است و مدل‌بندی توأم اریبی ناشی از وقوع رخداد را اصلاح می‌کند. در نهایت، مدل‌های توأم زمانی مورد نیاز است که درک ارتباط بین داده‌های نشانگر طولی و خطر وقوع رخداد از اهمیت بالایی برخوردار باشد. روش مدل‌بندی توأم شامل یک مدل طولی معمولاً یک مدل آمیخته و یک مدل زمان تا رخداد پیشامد معمولاً یک مدل مخاطره متناسب و ارتباط بین این دو مدل است که با استفاده از یک ساختار اشتراکی پنهان، در قالب اثرهای تصادفی، توصیف می‌شود. منابع شایان‌ذکری در این زمینه وجود دارد که هر دو استنباط کلاسیک و بیزی را برای برآورد پارامترهای مدل توأم به‌کار برده‌اند [۱۰، ۴۰، ۱۶، ۱۵، ۱۹، ۳۸، ۸، ۴۱، ۱۷، ۳۶، ۳۹، ۳، ۴، ۵، ۱۲، ۱۳].

چنین مدل توأمی امکان مدل‌بندی داده‌های همگن و تنها با یک الگو را فراهم می‌سازد و استفاده از آن در تحلیل داده‌هایی با زیرگروه‌هایی با الگوهای مختلف از نمیرخ‌های پاسخ^۴ معقول نیست زیرا در چنین شرایطی فرض همگن بودن جامعه مورد مطالعه نمی‌تواند واقع‌گرایانه باشد. برای مثال در کارآزمایی‌های بالینی چندمرکزی انتظار می‌رود مراکز مختلف به دلیل جمعیت‌های

در برخی مطالعات طولی، علاقه‌مند به جمع‌آوری هر دو اندازه‌های مکرر طولی و پیشامد زمان بقا هستیم. مدل‌های توأم از آنجایی که یک تحلیل واحد با مدنظر گرفتن ارتباط بین خطر یک رخداد و تغییرات نشانگر طولی در طی زمان، ارائه می‌دهند، بسیار حائز اهمیت هستند و استفاده از آن‌ها در چنین مطالعاتی متداول است [۳۵، ۳۸]. به‌عنوان مثال، در مطالعات HIV، تعداد $CD4$ ها و زمان ابتلا به ایدز و در مطالعه سرطان پروستات، اندازه‌گیری‌های آنتی‌ژن مخصوص پروستات^۳ (PSA) و خطر عود سرطان بررسی می‌شوند. از آنجاکه در چنین مثال‌هایی، دو برآمد همبسته هستند، مدل‌بندی توأم آن‌ها برای در نظر گرفتن پیوند دو برآمد لازم و ضروری است. همچنین، ممکن است علاقه‌مند به پیش‌بینی نشانگر طولی در طی زمان باشیم، زیرا در برخی کاربردها مانند مطالعه سرطان پروستات، روند PSA بر عود بیماری حائز اهمیت است و پیش‌بینی آن می‌تواند کمک شایان‌ذکری برای درمان بیماری محسوب شود.

با استفاده از مدل‌بندی توأم می‌توان اریبی ناشی از خطاهای اندازه‌گیری تصادفی و اندازه‌گیری متناوب نشانگر طولی را اصلاح

^۱گروه آمار و علوم کامپیوتر، دانشگاه محقق اردبیلی

^۲گروه آمار، دانشگاه تربیت مدرس، (نویسنده مسئول: t.baghfalaki@modares.ac.ir)

^۳Prostate Specific Antigen

^۴Response Profiles

این معنی که ناهمگنی توسط هیچ‌یک از متغیرهای مشاهده‌شده قابل‌دستیابی نیست.

لین و همکاران [۲۴] یک مدل رده پنهان برای تحلیل توأم نشانگرهای طولی و داده‌های بقا پیشنهاد کردند. آن‌ها احتمال عضویت رده را توسط یک مدل لوژیستیک چندجمله‌ای مشخص کردند و همچنین فرض کردند که فرایندهای طولی و پیشامد بقا در هر رده پنهان مستقل هستند. بر این اساس، پارامترهای مدل را با استفاده از روش درست‌نمایی ماکسیمم و از طریق الگوریتم EM برآورد کردند. همچنین، مدل خود را برای مطالعه داده‌های مربوط به سرطان پروستات به کار بردند. مدل آن‌ها چهار زیرگروه را نشان می‌دهد که بیان‌کننده سطوح مختلف نشانگر آنتی‌ژن مخصوص پروستات و خطر ابتلا به سرطان پروستات است.

لارسن [۲۱]، یک مدل رده پنهان را برای تحلیل هم‌زمان پیشامد بقا و چندین نشانگر دودویی در مطالعه سلامتی زنان و پیری^۵ به کار برد. بر اساس این مدل، سه رده پنهان کشف شد که نشان‌دهنده توجه به مسئله تحرک است. در پژوهش [۲۳] یک مدل آمیخته الگوی پنهان^۶ برای داده‌های گمشده متناوب آگاهی‌بخش در مطالعات طولی ارائه شد.

هان و همکاران [۱۴] مدل رده پنهان با یک رخداد که در تحقیق [۲۴] معرفی شده بود را به چند رخداد تعمیم دادند به طوری که مدل‌های آمیخته الگوی پنهان معرفی شده برای داده‌های گمشده متناوب آگاهی‌بخش در منبع [۲۳] را به عنوان حالتی خاص تحت یک الگوی پارامتری در برگرفته است. شین و همکاران [۳۴] مدل رده پنهان تعدیل‌یافته برای بررسی داده‌های رسته‌ای با انصراف آگاهی‌بخش را برای تحلیل داده‌های ترک سیگار ارائه کردند. همچنین، یک مدل توأم رده پنهان برای بررسی ارتباط بین نشانگر طولی آنتی‌ژن مخصوص پروستات و عود سرطان پروستات در پژوهش [۳۳] پیشنهاد شده است.

بیشتر روش‌های استفاده‌شده در مقالات ذکر شده، فرض استقلال شرطی پیشامدهای طولی و بقا برای هر رده پنهان را در نظر گرفته‌اند. چنین فرض استقلال شرطی، ممکن است برای به دست آوردن همبستگی بین دو پیشامد طولی و بقا ناکافی باشد، به این صورت که ممکن است منجر به کم‌برآوردی قدرت ارتباط و عدم قطعیت ساختار پیوند بین این دو فرایند شود [۲۵]. البته یک آزمون امتیاز توسط جاکمین-گادا و همکاران [۲۰] برای ارزیابی اینکه آیا فرض استقلال شرطی در مدل‌های رده پنهان توأم برقرار است یا خیر، پیشنهاد شده است.

گری و همکاران [۱۱] برای مدل‌بندی توأم داده‌های طولی و بقا،

بیمار متفاوت و الگوهای ارجاعی مختلف، عملکردهای متفاوتی نسبت به بقا اولیه داشته باشند. بنابراین، باید زیرگروه‌های بالقوه اندازه‌گیری‌های طولی و نتایج بقا را که ممکن است توزیع‌های مختلف را دنبال کنند، در نظر بگیریم.

یک تعمیم استاندارد از مدل‌های مخاطره نسبی که چنین ساختاری را در برمی‌گیرند، در نظر گرفتن چندین رده است. به طور خاص، فرض بر این است که بیماران به رده‌های مختلفی تقسیم می‌شوند که هر رده تابع مخاطره پایه مربوط به خودش را دارد، اما مقادیر ضرایب رگرسیونی بین رده‌ها مشترک است.

به عبارتی تحت یک مدل رده‌بندی شده، تابع مخاطره بیمار i متعلق به رده k به صورت زیر بیان می‌شود:

$$h_{ik}(t) = h_{\circ k}(t) \exp\{\gamma' w_i + \alpha m_i(t)\}, \quad (1)$$

که در آن $h_{\circ k}(t)$ نشانگر تابع مخاطره پایه رده k ، w_i بردار متغیرهای تبیینی، γ بردار پارامترهای متناظر با w_i و α ضریب رگرسیونی است.

مدل‌های توأم با یک زیر مدل مخاطره نسبی رده‌بندی شده در بسته JM، تحت عملکرد تابع مخاطره پایه تقریبی B-اسپلاین موجود است.

بر اساس فرمول استاندارد مدل‌های مخاطره نسبی رده‌بندی شده، تأثیر هر متغیر در رده‌های مختلف ثابت در نظر گرفته می‌شود. در صورتی که این فرضیه همیشه منطقی نیست، زیرا در بسیاری از موارد برخی از متغیرها ممکن است اثر متفاوتی برای هر رده داشته باشند. به عنوان مثال در یک کارآزمایی بالینی چندمرکزی با جمعیت بیمار مشابه در هر مرکز، ممکن است منطقی باشد که فرض کنیم اثر سن در مراکز مختلف یکسان است، اما نمی‌توان ادعا کرد که اثر درمان در تمام مراکز یکنواخت باشد. برای رفع مشکل در چنین مواردی مدل (۱) را با در نظر گرفتن اثرات متغیرها برای هر رده به صورت زیر تعمیم می‌دهیم:

$$h_{ik}(t) = h_{\circ k}(t) \exp\{\gamma_k' w_i + \alpha_k m_i(t)\}, \quad (2)$$

که در آن نه تنها تابع مخاطره پایه $h_{\circ k}(t)$ بلکه ضرایب رگرسیونی α_k و γ_k هم به رده k بستگی دارد.

مدل‌های توأم رده پنهان که به مدل توأم رده‌بندی شده مرتبط می‌باشد، برای رسیدگی به چنین فرایندهای ناهمگن استفاده می‌شود [۲۳، ۲۴، ۲۹]. انگیزه و هدف این نوع مدل‌های توأم نیز در نظر گرفتن ناهمگنی احتمالی در جمعیت است. باین حال، و برخلاف مدل توأم رده‌بندی شده، مدل‌های توأم رده پنهان فرض می‌کنند که زیرجامعه‌هایی که جامعه را تشکیل می‌دهند، پنهان هستند، به

⁵Women's Health and Aging Study

⁶Latent Pattern Mixture Model

در پژوهش [۲۹]، یک روش رده پنهان غیرخطی برای مدل‌بندی توأم اندازه‌های طولی چند متغیره و زمان بقا پیشنهاد شده است. همچنین پروست-لیما و همکاران [۳۲] یک مقاله مروری از مدل‌های رده پنهان توأم داده‌های طولی و زمان تا رخداد پیشامد ارائه کردند. انتینک و همکاران [۹]، یک مدل آمیخته معرفی کردند که ناهمگنی موجود در مسیرهای توسعه‌ای^{۱۱} را به دست آورد، درحالی‌که توسط مدل پاسخ چند سطحی^{۱۲} و سایر متغیرهای تبیینی به‌طور کامل قابل توصیف نبود. به‌منظور مدل‌بندی یک الگوی غیرخطی برای پیشامد طولی، چن و هوانگ [۷]، با دیدگاه بیزی یک ترکیب متناهی از مدل‌های توأم اثرات آمیخته نیمه پارامتری با توزیع چوله- t برای اندازه‌های طولی به دست آوردند. همچنین، هوانگ و همکاران [۱۸] برای مدل‌بندی داده‌های طولی چندمتغیره ناهمگن غیرخطی، ترکیبی متناهی از مدل‌های اثرات آمیخته غیرخطی برای مدل‌بندی رده پنهان نشانگر طولی در نظر گرفتند. به‌علاوه آن‌ها، مدل‌بندی چندین پیشامد رخداد را با استفاده از مخاطرات متناسب پیشنهاد کردند که تابع مخاطره برای هر رده پنهان به‌صورت تابع پله‌ای تعریف می‌شود. رونت و همکاران [۳۷]، یک مدل رده پنهان توأم بیماری-مرگ برای رخدادهای سانسور شده فاصله‌ای نیمه رقابتی و یک نشانگر طولی پیشنهاد کردند. آن‌ها دو تفسیر از مدل را توسعه دادند، یکی، مارکوفی و دیگری، نیمه مارکوفی و همچنین یک تعمیم برای تحلیل توأم از نشانگرهای طولی چندگانه ارائه کردند.

پروست-لیما و همکاران [۲۷]، نشانگرهای طولی و مخاطرات مرگ طبیعی و زوال عقل را با استفاده از رده‌های پنهان مشترک به هم پیوند زدند. این رده‌های پنهان، زیرجامعه‌های پنهانی که با نمایش معینی از تغییرات برای نشانگر و مخاطره رخدادها مشخص شده‌اند را فرمول‌بندی می‌کند. به‌این ترتیب، علاوه بر جداسازی دو نوع همبستگی و مدل‌بندی انعطاف‌پذیر وابستگی، یک ناهمگنی مورد انتظار در منحنی‌های طولی و خطرات رخداد را نیز مدل می‌کنند.

همچنین، [۳۰، ۳۱] یک بسته نرم‌افزاری R با نام **lemm** معرفی کردند که مجموعه‌ای از توابع را برای برآورد مدل‌های آماری بر اساس تئوری مدل آمیخته خطی فراهم می‌کند. از جمله، تابع **Jointlemm** برای برآورد مدل‌های آمیخته رده پنهان توأم برای یک پیشامد طولی (گوسی یا خم‌خطی^{۱۳}) و یک پیشامد بقا که می‌تواند چپ-بریده شده یا راست-سانسور شده باشد و در یک

یک مدل نقطه تغییر رده پنهان توأم با دو رده پنهان پیشنهاد کردند که در آن پیشامد طولی رده پنهان یک را با استفاده از یک مدل اثرات تصادفی تنها-عرض از مبدأ^{۱۴} و پیشامد طولی رده پنهان دو را با استفاده از یک مدل نقطه تغییر تصادفی تقسیم‌شده مدل‌بندی کردند و یک روش بیزی برای برازش مدل به کار بردند. بینکنز و همکاران [۶]، مدل رده پنهان را برای تحلیل داده‌های طولی ناقص به کار بردند. به‌طوری‌که بدون در نظر گرفتن پیشامد بقا، مدل لوژستیک را برای وضعیت گمشدگی استفاده کردند.

لیو و همکاران [۲۵]، مدل اثرات تصادفی توأم را در چارچوب مدل رده پنهان ادغام کردند تا بهتر بتوانند به ناهمگنی موجود در هر دو فرایند طولی و بقا با استفاده از یک روش متداول دست پیدا کنند. آن‌ها فرض کردند که در هر رده پنهان، یک مدل توأم متمایز از نقاط انتهایی بقا و طولی با اثرات تصادفی مشترک وجود داشته باشد و نشان دادند که مدل آن‌ها جامع‌تر است، از این نظر که می‌تواند، عامل‌هایی که بر عضویت رده پنهان تأثیر می‌گذارند را شناسایی و ارتباط بین اندازه‌های طولی و فرایند بقا را مدل‌بندی کند. علاوه بر این، با فراهم کردن یک ساختار کلی می‌تواند اثرات رده-مشخص متغیرها را روی مدل‌های بقا و طولی در نظر بگیرد. لازم به ذکر است که، در سال‌های اخیر استفاده از مدل رده پنهان برای مدل‌بندی توأم داده‌های طولی و زمان بقا بیشتر مورد توجه قرار گرفته و توسعه یافته است. همچنین، کاربردهایی از این نوع مدل‌بندی در آزمایش‌های بالینی مطرح شده است.

در برخی از آسیب‌شناسی‌ها، علاوه بر تعداد وقایع بالینی، چندین نشانگر یا رخداد در طول زمان جمع‌آوری می‌شوند و معمولاً در مواردی رخ می‌دهند که انجام فرایند طولی موردنظر از طریق نتایج مشاهده‌شده دشوار باشد. برای مثال، در بیماری آلزایمر (AD)^۸ عملکرد شناختی^۹ با آزمایش‌های روان‌سنجی متعدد اندازه‌گیری می‌شود تا شناخت بنیادی بهتر تخمین زده شود [۲۸]. همچنین به‌عنوان مثال دیگر، کیفیت زندگی توسط مجموعه‌ای از نشانگرها اندازه‌گیری می‌شود یا وابستگی عملکردی^{۱۰} توسط مجموعه‌ای از فعالیت‌های روزمره زندگی تعریف می‌شود [۲۷]. از این رو، مدل رده پنهان توأم، ابتدا با تمرکز بر یک نشانگر طولی و یک رخداد بالینی مطالعه شده است و سپس به حالت‌های چندین نشانگر طولی (پیوسته یا گسسته-ترتیبی) و یا چند رخداد به مفهوم علت-مشخص برای مخاطره‌های رقابتی تعمیم یافته است.

⁷Intercept-only random-effects

⁸Alzheimer's Disease

⁹Cognitive Functioning

¹⁰Functional Dependency

¹¹Developmental Trajectories

¹²Multilevel Item Response Model

¹³Curvilinear

کلاسیک که تنها مجموعه اثرات تصادفی b_i را برای هر دو نوع ارتباط در نظر می‌گیرد، فراهم می‌کند.

تحت فرض‌های استقلال شرطی مذکور، یک تعریف کلی از مدل توأم رده پنهان شامل زیرمدل‌های زیر خواهد بود:

$$(۴) \left\{ \begin{array}{l} (Y_i(s)|c_i = g, b_i; \beta_g) = X_i'(s)\beta_g + Z_i'(s)b_{ig} + \epsilon_i(s), \\ \epsilon_i(s) \sim N(0, \sigma^2) \\ h_i(t|c_i = g; \gamma_g) = h_{\cdot g}(t) \exp(w_i' \gamma_g) \\ \pi_{ig} = Pr(c_i = g; \lambda_g) = \frac{\exp(u_i' \lambda_g)}{\sum_{L=1}^G \exp(u_i' \lambda_g)}. \end{array} \right.$$

به طوری که فرض می‌کنیم بیماران در رده‌های پنهان متفاوت هم تکامل طولی متفاوت و هم تابع مخاطره مختلفی برای وقوع یک رخداد دارند. معادله اول رابطه (۴) تحت رده پنهان g ، یک مدل آمیخته خطی استاندارد برای بردار اندازه‌های تکرارشونده نشانگر طولی، $Y_i(s) = (Y_i(s_{i1}), \dots, Y_i(s_{in_i}))'$ ، در زمان‌های متفاوت اندازه‌گیری s_{ij} ، $j = 1, \dots, n_i$ ، را در نظر می‌گیرد که در آن X_i و Z_i به ترتیب ماتریس طرح $p \times n_i$ و $q \times n_i$ بعدی متناظر با بردار p بعدی اثرات ثابت β_g و بردار q بعدی اثرات تصادفی b_{ig} برای فرد i ام در رده g است. به طوری که q -بردار اثرات تصادفی رده مشخص به صورت، $b_{ig} = b_i|_{c_i=g} \sim N_q(\mu_g, D_g)$ ، در نظر گرفته می‌شود که در آن، $D_g = \sigma_g^2 D$. به عبارت دیگر، بردار اثرات تصادفی به صورت $b_i \sim \sum_{g=1}^G \pi_{ig} N_q(\mu_g, D_g)$ مشخص می‌شود که در آن احتمال عضویت در رده پنهان g ام است. ملاحظه می‌شود که وابستگی ماتریس کوواریانس اثرات تصادفی به رده c_i ، تنها از طریق ضرب پارامتر واریانس یعنی σ_g^2 ، بیان می‌شود و به منظور شناساپذیری $\sigma_G = 1$ در نظر گرفته می‌شود.

معادله دوم رابطه (۴)، یک مدل خطر متناسب کاکس با عملکرد خطی و پارامترهای رده-مشخص را توصیف می‌کند که در آن $h_{\cdot g}(t)$ تابع خطر پایه رده g و w_i بردار متغیر کمکی برای فرد i ام در زمان t متناظر با بردار پارامتر γ_g است.

آخرین معادله رابطه (۴)، یک زیر مدل چندجمله‌ای برای احتمال عضویت در رده پنهان g ، $g = 1, \dots, G$ را مشخص می‌کند، که در آن بیانگر برداری از متغیرهای کمکی مرتبط با این احتمالات برای فرد i ام و λ بردار پارامتر رده-مشخص متناظر با w_i هستند. $\lambda' = (\lambda'_1, \dots, \lambda'_G)$ و برای شناساپذیری فرض $\lambda_G = 0$ را در نظر می‌گیریم.

یک نتیجه سودمند حاصل از مدل توأم معرفی شده این است که

ساختار رقابتی تعریف شود. لی و همکاران [۲۲]، یک روش رده پنهان برای مدل‌بندی توأم پیشامد زمان تا رخداد و چند نشانگر طولی پیشنهاد کردند درحالی که اندازه‌گیری‌های نشانگر سانسور شده را به دلیل محدودیت‌های تشخیص^{۱۴} به حساب آورده‌اند. آندریونپولو و همکاران [۲] برای غلبه به محدودیت‌های انتخاب تعداد رده بهینه، با استفاده از یک رویکرد کاملاً بیزی، رده‌های پنهان را در مدل توأم پارامتر مشترک ادغام کردند و برای انتخاب تعداد رده‌های بهینه یک ساختار مدل آمیخته معرفی کردند.

این مقاله به این ترتیب سازمان‌دهی شده است. در بخش دوم از این مقاله، مدل رده پنهان توأم را معرفی می‌کنیم که شامل مدل‌بندی توأم داده‌های طولی و زمان بقا و مدل‌بندی توأم اندازه‌های طولی و داده‌های مخاطره‌های رقابتی است. در بخش بعد به کاربرد مدل توصیف شده خواهیم پرداخت. این مقاله با بحث و نتیجه‌گیری خاتمه می‌یابد.

۲ مدل‌های توأم رده پنهان

۱.۲ مدل‌بندی توأم داده‌های طولی و زمان بقا

برای تعریف مدل‌های توأم رده پنهان، فرض می‌کنیم که G تعداد زیرجامعه‌هایی باشد که جامعه اصلی ما را تشکیل می‌دهند و نشانگر رده مشاهده نشده (پنهان) را به صورت $c_i = g$ ، $g = 1, \dots, G$ تعریف می‌کنیم که عضویت i امین مؤلفه را در رده پنهان g نشان می‌دهد.

با در نظر گرفتن δ_i به عنوان نشانگر رخداد، دوتایی (T_i, δ_i) زمان رخداد مشاهده شده را توصیف می‌کند، همچنین بردار y_i نشان‌دهنده اندازه‌های طولی و θ بردار کل پارامترهای مدل توأم است، در این صورت این مدل‌ها تحت مفروضات استقلال شرطی زیر برقرار هستند:

$$(۳) \quad P(T_i, \delta_i, y_i | c_i = g, b_i; \theta) = P(T_i, \delta_i | c_i = g; \theta) \times P(y_i | c_i = g, b_i; \theta),$$

$$P(y_i | c_i = g, b_i; \theta) = \prod_j P\{y_i(s_{ij}) | c_i = g, b_i; \theta\}.$$

در این مدل‌ها فرض بر این است که همبستگی بین اندازه‌گیری‌های مکرر در فرایند طولی توسط اثرات تصادفی b_i بیان می‌شود، درحالی که، ارتباط بین زمان رخداد و فرایند طولی توسط شاخص رده پنهان مشترک c_i توضیح داده می‌شود. استفاده از دو مؤلفه پنهان از این نظر سودمند است که امکان ایجاد ساختارهای انعطاف‌پذیرتر برای اشتراک دو زیر مدل طولی و زمان بقا، در مقایسه با مدل توأم

¹⁴Detection Limits

علاقه‌مند به ارائه یک تفسیر ساده از مکانیسم ارتباط دو فرایند باشیم [۳۵].

با استفاده از مدل برازش داده‌شده، همچنین می‌توان طبقه‌بندی پسین بیماران را در نمونه به دست آورد که با استفاده از ماکسیم احتمالات پسینی حاصل می‌شود:

$$\hat{\pi}_{ig}^{(Y,T)} = \Pr(c_i = g | T_i, \delta_i, \mathbf{Y}_i; \hat{\theta}) = \frac{\pi_{ig}(\hat{\theta}) h_i(T_i | c_i = g; \hat{\theta})^{\delta_i} S_i(T_i | c_i = g; \hat{\theta}) P(\mathbf{Y}_i | c_i = g; \hat{\theta})}{\sum_{l=1}^G \pi_{il}(\hat{\theta}) h_i(T_i | c_i = l; \hat{\theta})^{\delta_i} S_i(T_i | c_i = l; \hat{\theta}) P(\mathbf{Y}_i | c_i = l; \hat{\theta})}$$

بنابراین، رابطه $\hat{c}_i = \arg \max_g \{ \Pr(c_i = g | T_i, \delta_i, \mathbf{Y}_i; \hat{\theta}) \}$ می‌دهد که فرد i ام در رده g طبقه‌بندی شده است، که از نظر مفهومی مشابه برآورد بیز تجربی اثرات تصادفی است.

۱۰.۱۰۲ استفاده از الگوریتم EM برای برآورد پارامترهای مدل رده پنهان در داده‌های طولی و برآمد بقا

ماکسیم کردن مستقیم رابطه (۵) با جمع رده‌های پنهان در لگاریتم و با صدها پارامتر بالقوه، ممکن است دارای پیچیدگی‌های محاسباتی بسیار باشد، از این رو، بسیاری از محققین حوزه مدل‌بندی توأم داده‌های طولی و زمان بقا بر اساس مدل رده پنهان، از الگوریتم EM برای برآورد پارامترهای مدل توأم استفاده کرده‌اند. بر اساس رابطه (۵)، تابع امتیاز لگاریتم درستنمایی داده‌های مشاهده‌شده را می‌توان به صورت زیر به دست آورد:

$$\begin{aligned} S(\theta) &= \sum_i \frac{\partial}{\partial \theta'} \log P(T_i, \delta_i, \mathbf{Y}_i; \theta) \quad (۶) \\ &= \sum_i \frac{\partial}{\partial \theta'} \log \int \sum_{g=1}^G \pi_{ig}(\theta) p(T_i, \delta_i | c_i = g; \theta) \\ &\quad \times p(\mathbf{Y}_i | c_i = g, \mathbf{b}_i; \theta) p(b_i | c_i = g; \theta) db_i \\ &= \sum_i \frac{1}{P(T_i, \delta_i, \mathbf{Y}_i; \theta)} \frac{\partial}{\partial \theta'} \int \sum_{g=1}^G \pi_{ig}(\theta) p(T_i, \delta_i | c_i = g; \theta) \\ &\quad \times p(\mathbf{Y}_i | c_i = g, \mathbf{b}_i; \theta) p(b_i | c_i = g; \theta) db_i \\ &= \sum_i \frac{1}{p(T_i, \delta_i, \mathbf{Y}_i; \theta)} \int \frac{\partial}{\partial \theta'} \{ \sum_{g=1}^G \pi_{ig}(\theta) p(T_i, \delta_i | c_i = g; \theta) \\ &\quad \times p(\mathbf{Y}_i | c_i = g, \mathbf{b}_i; \theta) p(b_i | c_i = g; \theta) \} db_i \\ &= \sum_i \int \frac{\partial}{\partial \theta'} \log \{ \sum_{g=1}^G \pi_{ig}(\theta) p(T_i, \delta_i | c_i = g; \theta) \\ &\quad \times p(\mathbf{Y}_i | c_i = g, \mathbf{b}_i; \theta) p(b_i | c_i = g; \theta) \} \\ &\quad \times \frac{\sum_{g=1}^G \pi_{ig}(\theta) p(T_i, \delta_i | c_i = g; \theta)}{p(T_i, \delta_i, \mathbf{Y}_i; \theta)} \\ &\quad \times \frac{p(\mathbf{Y}_i | c_i = g, \mathbf{b}_i; \theta) p(b_i | c_i = g; \theta)}{1} db_i \\ &= \sum_i \int \frac{\partial}{\partial \theta'} \log \{ \sum_{g=1}^G \pi_{ig}(\theta) p(T_i, \delta_i | c_i = g; \theta) \\ &\quad \times p(\mathbf{Y}_i | c_i = g, \mathbf{b}_i; \theta) p(b_i | c_i = g; \theta) \} \\ &\quad \times p(\mathbf{b}_i | T_i, \delta_i, \mathbf{Y}_i, c_i = g; \theta) db_i. \end{aligned}$$

می‌توان لگاریتم درستنمایی تحت این مدل را به فرم زیر بیان کرد: (۵)

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^n \log P(T_i, \delta_i, \mathbf{Y}_i) \\ &= \sum_{i=1}^n \log \left\{ \sum_{g=1}^G \Pr(c_i = g; \theta) h_i(t_i | c_i = g; \theta)^{\delta_i} \right. \\ &\quad \times S_i(t_i | c_i = g; \theta) \\ &\quad \times \int [\prod_j P\{Y_i(s_{ij}) | c_i = g, b_i; \theta\} P(b_i | c_i = g; \theta) db_i] \\ &\quad \left. \times \sum_{l=1}^n \log \left\{ \sum_{g=1}^G \pi_{il}(\theta) h_i(t_i | c_i = g; \theta)^{\delta_i} \right. \right. \\ &\quad \left. \left. \times S_i(t_i | c_i = g; \theta) \times P(\mathbf{Y}_i | c_i = g; \theta) \right\} \right\} \end{aligned}$$

که در مقایسه با مدل توأم کلاسیک به راحتی قابل دستیابی است، چون برای محاسبه تابع بقا و درستنمایی نیاز به روش‌های عددی برای محاسبه انتگرال‌ها ندارد. به خصوص، انتگرال موجود در تعریف تابع بقا یک فرم بسته دارد (البته به تابع مخاطره پایه هم‌بستگی دارد)، چون هیچ مؤلفه وابسته به زمان از فرایند طولی را شامل نمی‌شود و همچنین انتگرال نهایی به شرط اثرات تصادفی دارای فرم بسته‌ای است، زیرا تنها شامل مدل طولی شرطی $P(y_i | c_i = g, b_i)$ و چگالی اثرات تصادفی $P(b_i | c_i = g)$ است که تحت فرض نرمال هر دو عبارت منجر به یک توزیع گوسی چندمتغیره می‌شود. با این وجود، مسئله چالش برانگیز در مدل‌های توأم رده پنهان این است که ممکن است لگاریتم تابع درستنمایی منجر به چندین ماکسیم موضعی شود، که در این صورت توصیه می‌شود با چندین بار استفاده از مجموعه‌های متفاوت از مقادیر اولیه مدل را اصلاح کنیم و همگرایی را بررسی کنیم. در حقیقت واضح است که این مشکل، تقریباً مزیت محاسباتی نداشتن تقریب عددی انتگرال‌ها را باطل می‌کند، زیرا اصلاح مجدد مدل، حتی اگر ساده‌تر از برازش مدل توأم کلاسیک باشد، همچنان دارای پیچیدگی محاسباتی است. این مسئله همچنین با توجه به اینکه از ابتدا تعداد مناسب رده‌های پنهان مشخص نیست، افزایش می‌یابد. بنابراین، لازم است چندین مدل با تعداد رده‌های فزاینده برازش داده شوند و از نظر آماری، تعداد رده‌ها را به گونه‌ای انتخاب کنیم که بهترین برازش را روی داده‌ها انجام داده باشد. این انتخاب معمولاً بر اساس معیارهای انتخاب مدل مانند AIC یا BIC انجام می‌شود. مسئله دیگر در مورد مدل‌های توأم رده پنهان این است که تفسیر در مورد ساختار اشتراک زیر مدل طولی و زمان بقا ساده نیست. تحت فرمول‌بندی رده پنهان، هیچ مجموعه‌ای از پارامترها وجود ندارد که به طور مستقیم قدرت پیوند دو فرایند طولی و خطر یک رخداد را اندازه‌گیری کند، که در بسیاری از موارد موضوع اصلی است. بنابراین، در اصل این نوع مدل‌بندی توأم، زمانی که علاقه‌مند به بهبود ناهمگنی جامعه هدف هستیم، مفید است و نه زمانی که

تابع $Q(\theta|\theta^{(r)})$ را به‌عنوان مقدار امید ریاضی تابع لگاریتم درستنمایی داده‌های کامل نسبت به توزیع شرطی داده‌های مشاهده نشده به‌شرط مشاهدات و مقدار پارامتر در مرحله r ام به‌صورت زیر تعریف می‌کنیم:

$$Q(\theta|\theta^{(r)}) = E_{b,c|Y,T,\delta,\check{X};\theta^{(r)}}[\ell_c(\theta; Y, T, \delta, \check{X}, b, c)], \quad (۱۲)$$

بنابراین با جایگذاری معادلات (۸) در تابع $Q(\theta|\theta^{(r)})$ خواهیم داشت:

$$\begin{aligned} Q(\theta|\theta^{(r)}) &= \sum_{i=1}^n E_{b,c|Y,T,\delta,\check{X};\theta^{(r)}} [\ell_{com}(\theta; Y_i, T_i, \delta_i, \check{X}_i, b_i, c_i)] \\ &= \sum_{i=1}^n E_{b,c|Y,T,\delta,\check{X};\theta^{(r)}} \left[\log Pr(c_i | \check{X}_i) \right. \\ &\quad \left. + \log P(b_i | c_i, \check{X}_i) + \log P(Y_i | b_i, c_i, \check{X}_i) \right. \\ &\quad \left. + \log P(T_i, \delta_i | c_i, \check{X}_i) \right] \\ &= \sum_{i=1}^n \left[\sum_{g=1}^G \left\{ E(c_{ig} | Y_i, T_i, \delta_i, \check{X}_i; \theta^{(r)}) \right. \right. \\ &\quad \left. \left. \times \left(\mathbf{u}'_i \boldsymbol{\lambda}_g - \log \left(\sum_{L=1}^G \exp(\mathbf{u}'_i \boldsymbol{\lambda}_g) \right) \right) \right\} \right. \\ &\quad \left. + \sum_{g=1}^G E(c_{ig} | Y_i, T_i, \delta_i, \check{X}_i; \theta^{(r)}) \right. \\ &\quad \left. \times \left(-\frac{q+n_i}{\nu} \log(\nu\pi) - \frac{1}{\nu} \log |D_g| - \frac{1}{\nu} \log(\sigma^2) \right) \right. \\ &\quad \left. + \sum_{g=1}^G \left\{ -\frac{1}{\nu} E(c_{ig} \mathbf{b}'_{ig} D_g^{-1} \mathbf{b}_{ig} - \nu c_{ig} \boldsymbol{\mu}'_g D_g^{-1} \mathbf{b}_{ig} \right. \right. \\ &\quad \left. \left. + c_{ig} \boldsymbol{\mu}'_g D_g^{-1} \boldsymbol{\mu}_g | Y_i, T_i, \delta_i, \check{X}_i; \theta^{(r)}) \right\} \right. \\ &\quad \left. + \sum_{g=1}^G \left\{ -\frac{1}{\nu} E(\sigma^{-\nu} c_{ig} (Y_i(s) - X'_i \boldsymbol{\beta}_g)' (Y_i(s) - X'_i \boldsymbol{\beta}_g) \log P(\mathbf{b}_i | c_i) \right. \right. \\ &\quad \left. \left. - \nu \sigma^{-\nu} c_{ig} (Y_i(s) - X'_i \boldsymbol{\beta}_g)' Z'_i \mathbf{b}_{ig} \right. \right. \\ &\quad \left. \left. + \sigma^{-\nu} c_{ig} \mathbf{b}'_{ig} Z_i Z'_i \mathbf{b}_{ig} | Y_i, T_i, \delta_i, \check{X}_i; \theta^{(r)}) \right\} \right. \\ &\quad \left. + \sum_{g=1}^G E(c_{ig} | Y_i, T_i, \delta_i, \check{X}_i; \theta^{(r)}) (\log h_{\cdot g}(t_i) + \mathbf{w}'_i \boldsymbol{\gamma}_g) \right. \\ &\quad \left. + \log \left\{ \int_{t_i}^{\infty} h_{\cdot g}(t) \exp(\mathbf{w}'_i \boldsymbol{\gamma}_g) dt \right\} \right) + \sum_{g=1}^G E(c_{ig} \\ &\quad \log \left\{ \int_{t_i}^{\infty} h_{\cdot g}(t) \exp(\mathbf{w}'_i \boldsymbol{\gamma}_g) dt \right\} | Y_i, T_i^{obs}, \check{X}_i; \theta^{(r)}) \Big] \log P(Y_i | b_i, c_i, \check{X}_i) \end{aligned}$$

گام E الگوریتم EM: محاسبه امید ریاضی‌های شرطی

با فرض این‌که، $\tilde{x} = E(x | y, T, \delta, \check{X}) \equiv \tilde{E}(x)$ نشان‌دهنده امید ریاضی شرطی (پسینی) متغیر تصادفی x به‌شرط داده‌های مشاهده‌شده باشد، می‌توان امید ریاضی‌های شرطی را به‌صورت زیر به دست آورد.

ابتدا امید شرطی عضویت رده را به‌صورت زیر محاسبه می‌کنیم:

$$\begin{aligned} \tilde{c}_{ig} &= E(c_{ig} | Y_i, T_i, \delta_i, \check{X}_i; \theta^{(r)}) \\ &= \frac{\pi_{ig} P(Y_i | c_{ig} = 1, \check{X}_i; \theta^{(r)}) P(T_i, \delta_i | c_{ig} = 1, \check{X}_i; \theta^{(r)})}{\sum_{g=1}^G \pi_{ig} P(Y_i | c_{ig} = 1, \check{X}_i; \theta^{(r)}) P(T_i, \delta_i | c_{ig} = 1, \check{X}_i; \theta^{(r)})}, \end{aligned} \quad (۱۳)$$

بنابراین، ملاحظه می‌کنیم که بردار امتیاز داده‌های مشاهده‌شده به‌صورت مقدار امید ریاضی بردار امتیاز داده‌های کامل نسبت به توزیع پسین اثرات تصادفی به دست می‌آید. به عبارتی اگر عضویت رده c_i و اثرات تصادفی b_i ها را به‌عنوان داده‌های گمشده در نظر بگیریم، رابطه حاصل برای $S(\theta)$ نقش دوگانه ایفا می‌کند. به‌طور خاص، اگر معادلات امتیاز را نسبت به θ با فرض اینکه $p(b_i | T_i, \delta_i, Y_i, c_i = g; \theta)$ در مقدار تکرار قبلی از θ ثابت باشد، حل کنیم، یک الگوریتم EM خواهیم داشت، درحالی‌که اگر معادلات امتیاز را نسبت به θ و با فرض اینکه $p(b_i | T_i, \delta_i, Y_i, c_i = g; \theta)$ است، حل کنیم، متناظر با ماکسیم سازی مستقیم لگاریتم درستنمایی داده‌های مشاهده‌شده خواهد بود. اگر \check{X}_i نشان‌دهنده بردار ترکیب متغیرهای کمکی Z_i, X_i, u_i و w_i تعریف‌شده در رابطه (۴) باشد، در این صورت لگاریتم درستنمایی داده‌های کامل $(Y_i, T_i, \delta_i, b_i, c_i)$ به‌صورت زیر قابل محاسبه است:

$$\begin{aligned} \ell_{com}(\theta; Y, T, \delta, \check{X}, b, c) &= \sum_{i=1}^n \left\{ \log Pr(c_i | \check{X}_i) + \log P(b_i | c_i) \right. \\ &\quad \left. + \log P(Y_i | b_i, c_i, \check{X}_i) + \log P(T_i, \delta_i | c_i, \check{X}_i) \right\}, \end{aligned}$$

بسط هر یک از عبارت‌های رابطه (۷) را می‌توان به‌صورت زیر محاسبه کرد:

$$\begin{aligned} \log Pr(c_i | \check{X}_i) &= \sum_{g=1}^G c_{ig} \left\{ \mathbf{u}'_i \boldsymbol{\lambda}_g - \log \left[\sum_{L=1}^G \exp(\mathbf{u}'_i \boldsymbol{\lambda}_g) \right] \right\}, \\ \log P(Y_i | b_i, c_i, \check{X}_i) &= \sum_{g=1}^G c_{ig} (\log(\pi_{ig}) + \log \phi_q(\mathbf{b}_{ig} | \boldsymbol{\mu}_g, D_g)) \\ &= \sum_{g=1}^G c_{ig} \left\{ \mathbf{u}'_i \boldsymbol{\lambda}_g - \log \left[\sum_{L=1}^G \exp(\mathbf{u}'_i \boldsymbol{\lambda}_g) \right] \right. \\ &\quad \left. - \frac{q}{\nu} \log(\nu\pi) - \frac{1}{\nu} \log |D_g| \right. \\ &\quad \left. - \frac{(\mathbf{b}_{ig} - \boldsymbol{\mu}_g)' D_g^{-1} (\mathbf{b}_{ig} - \boldsymbol{\mu}_g)}{\nu} \right\}, \end{aligned} \quad (۹)$$

$$\begin{aligned} \log P(T_i, \delta_i | c_i, \check{X}_i) &= \sum_{g=1}^G c_{ig} \log(Y_i(s) | b_i, c_i, \check{X}_i) \\ &= \sum_{g=1}^G c_{ig} \left\{ -\frac{n_i}{\nu} \log(\nu\pi) - \frac{1}{\nu} \log(\sigma^2) \right. \\ &\quad \left. - \frac{1}{\nu} (Y_i(s) - X'_i \boldsymbol{\beta}_g - Z'_i \mathbf{b}_{ig})' \right. \\ &\quad \left. \times \sigma^{-\nu} I_{n_i} (Y_i(s) - X'_i \boldsymbol{\beta}_g - Z'_i \mathbf{b}_{ig}) \right\}, \end{aligned}$$

$$\begin{aligned} \log P(T_i, \delta_i | c_i, \check{X}_i) &= \sum_{g=1}^G c_{ig} \left\{ \delta_i \log h_i(t_i | c_i, \check{X}_i) \right. \\ &\quad \left. + \log S_i(t_i | c_i, \check{X}_i) \right\} \\ &= \sum_{g=1}^G c_{ig} \left\{ \delta_i (\log h_{\cdot g}(t_i) + \mathbf{w}'_i \boldsymbol{\gamma}_g) \right. \\ &\quad \left. + \log \left\{ \int_{t_i}^{\infty} h_{\cdot g}(t) \exp(\mathbf{w}'_i \boldsymbol{\gamma}_g) dt \right\} \right\}. \end{aligned} \quad (۱۱)$$

عبارت‌های زیر با فرم بسته برای پارامترها به دست می‌آید:

که در آن

$$\begin{aligned}
 \hat{D}_g^{(r+1)} &= \arg \max_{D_g} \{Q(\theta|\theta^{(r)})\} \\
 &= \frac{1}{n} \sum_{i=1}^n \left\{ (\tilde{b}_{ig} - \mu_g) (\tilde{b}_{ig} - \mu_g)' \right\} \\
 \hat{\beta}_g^{(r+1)} &= \arg \max_{\beta_g} \{Q(\theta|\theta^{(r)})\} \\
 &= \left(\sum_{i=1}^n X_i X_i' \right)^{-1} \sum_{i=1}^n X_i (Y_i - Z_i \tilde{b}_{ig})' \\
 \hat{\sigma}_g^{(r+1)} &= \arg \max_{\sigma_g^2} \{Q(\theta|\theta^{(r)})\} \\
 &= \frac{1}{\sum_{i=1}^n n_i} \sum_{i=1}^n \left\{ (Y_i - X_i' \hat{\beta}_g)' (Y_i - X_i' \hat{\beta}_g - \gamma_g' Z_i' \tilde{b}_{ig}) \right. \\
 &\quad \left. + \tilde{b}_{ig}' Z_i Z_i' \tilde{b}_{ig} \right\} \\
 \hat{\mu}_g^{(r+1)} &= \arg \max_{\mu_g} \{Q(\theta|\theta^{(r)})\} \\
 &= \sum_{i=1}^n c_{ig} \tilde{b}_{ig} (c_{ig})^{-1},
 \end{aligned}
 \tag{۱۶}$$

$$P(Y_i | \check{X}_i, c_{ig} = 1) = MVN_{n_i} \left(X_i \beta, Z_i \sigma_g^2 D Z_i' + \sigma^2 I_{n_i} \right)$$

و $MVN_{n_i}(mu, \Sigma)$ نشان‌دهنده چگالی نرمال چندمتغیره n_i بعدی با میانگین μ و واریانس Σ است.

همچنین، $P(T_i, \delta_i | c_{ig} = 1, \check{X}_i; \theta^{(r)})$ از رابطه زیر به دست می‌آید:

$$\begin{aligned}
 P(T_i, \delta_i | c_{ig} = 1, \check{X}_i; \theta^{(r)}) &= h_{\cdot g}(t) \exp(w_i' \gamma_g^{(r)}) \\
 &\times \left\{ \int_{t_i}^{\infty} h_{\cdot g}(t) \exp(w_i' \gamma_g^{(r)}) dt \right\}.
 \end{aligned}$$

توزیع پسینی b_{ig} ها که در پیوست محاسبه شده است به صورت نرمال چندمتغیره به ترتیب با میانگین μ_g^* و واریانس V_i به دست می‌آید:

$$b_{ig} | c_{ig}, Y_i, \check{X}_i; \theta \sim MVN_q(\mu_g^*, V_i),$$

که در آن $V_i = \mu_g^* = V_i \left(Z_i (Y_i - X_i' \beta_g) \sigma^{-2} + D_g^{-1} \mu_g \right) \left(Z_i Z_i' \sigma^{-2} + D_g^{-1} \right)^{-1}$.

بر این اساس، سایر امیدهای شرطی را به صورت زیر محاسبه می‌کنیم. لازم به ذکر است که جزئیات بیشتر و نحوه به دست آوردن این امید ریاضی‌های شرطی در پیوست آورده شده است:

$$\begin{aligned}
 \tilde{b}_{ig} &= V_i \left(Z_i (Y_i - X_i' \beta_g^{(r)}) \sigma^{-2(r)} + D_g^{-1} \mu_g^{(r)} \right), \\
 \tilde{b}_{ig}' \tilde{b}_{ig} &= V_i + \tilde{b}_{ig}' \left(\left((Y_i - X_i' \beta_g^{(r)})' Z_i' \sigma^{-2(r)} + \mu_g^{(r)'} D_g^{-1} \right) V_i \right), \\
 c_{ig} \tilde{b}_{ig} &= c_{ig} (\tilde{b}_{ig}), \\
 c_{ig} \tilde{b}_{ig}' D_g^{-1} \tilde{b}_{ig} &= tr \left(\tilde{c}_{ig} D_g^{-1} \left(\tilde{b}_{ig}' \tilde{b}_{ig} \right) \right), \\
 c_{ig} \tilde{b}_{ig}' Z_i Z_i' \tilde{b}_{ig} &= tr \left(\tilde{c}_{ig} Z_i Z_i' \left(\tilde{b}_{ig}' \tilde{b}_{ig} \right) \right),
 \end{aligned}
 \tag{۱۵}$$

در نهایت برای محاسبه

$$E(c_{ig} \log \left\{ \int_{t_i}^{\infty} h_{\cdot g}(t) \exp(w_i' \gamma_g) dt \right\} | Y_i, T_i^{obs}, \delta_i, \check{X}_i; \theta^{(r)})$$

روش‌های عددی اتخاذ می‌گردد که در آن T_i^{obs} نشان‌دهنده زمان‌های مشاهده شده با نشانگر $\delta_i = 1$ است.

گام M الگوریتم EM: بیشینه‌سازی

در این مرحله، بر اساس امیدهای شرطی حاصل در گام قبلی، تابع $Q(\theta|\theta^{(r)})$ را در هر تکرار از الگوریتم بیشینه می‌کنیم و مقدار جدید برآوردگر درستنمایی ما کسیم پارامترها را به دست می‌آوریم.

همچنین، برای پارامترهای زیر فرم بسته وجود ندارد و از روش‌های عددی برای برآورد آن‌ها استفاده می‌شود:

$$\hat{\lambda}_g^{(r+1)} = \arg \max_{\lambda_g} \{Q(\theta|\theta^{(r)})\}, \quad \hat{\gamma}_g^{(r+1)} = \arg \max_{\gamma_g} \{Q(\theta|\theta^{(r)})\}.$$

گام‌های E و M الگوریتم را تا جایی ادامه می‌دهیم که معیار همگرایی زیر برقرار گردد:

$$\ell(\theta^{(r+1)}) - \ell(\theta^{(r)}) < \epsilon,$$

که در آن ϵ مقدار مثبت بسیار کوچکی در نظر گرفته می‌شود و تابع $\ell(\theta)$ لگاریتم درستنمایی داده‌های مشاهده شده می‌باشد.

۲.۲ مدل‌بندی توأم اندازه‌های طولی و داده‌های مخاطره‌های رقابتی

برای تعریف مدل رده پنهان به منظور مدل‌بندی توأم اندازه‌های طولی و داده‌های مخاطره‌های رقابتی، مشابه بخش قبلی فرض می‌کنیم که جامعه اصلی از G زیرجامعه تشکیل شده باشد و نشانگر رده پنهان را به صورت $g = 1, \dots, G, c_i = g$ تعریف می‌کنیم که عضویت i امین مؤلفه را در رده پنهان g نشان می‌دهد. فرض می‌کنیم T_{ik}^* زمان رخداد عامل k ام، $k = 1, \dots, K$ و \tilde{T}_i زمان سانسور برای مؤلفه i ام باشد، در این صورت $T_i = \min(\tilde{T}_i, T_{i1}^*, \dots, T_{iK}^*)$ زمان رخداد مشاهده شده خواهد بود. δ_i را به عنوان نشانگر رخداد در نظر می‌گیریم که $\delta_i = k$ نشان‌دهنده این است که عامل k ام زودتر از بقیه رخ داده است و $\delta_i = 0$ اگر فرد i ام سانسور شده باشد. در ادامه سه زیرمدلی که مشخص‌کننده مدل رده پنهان توأم است را معرفی می‌کنیم. یک مدل آمیخته خطی برای نشانگر طولی پیوسته با پارامترهای رده مشخص

مفروضات استقلال شرطی (۳) که در بخش قبلی معرفی شدند، لگاریتم درست‌نمایی را می‌توان به صورت زیر نوشت:

$$\ell(\theta) = \sum_{i=1}^n \log \sum_{g=1}^G \{P(Y_i | c_i = g; \theta) \times P(T_i, \delta_i | c_i = g; \theta) Pr(c_i = g | \theta)\}, \quad (20)$$

در این رابطه، $Pr(c_i = g | \theta) = \pi_{ig}$ از فرمول (۱۹) به دست می‌آید و تابع چگالی مدل مخاطره‌های رقابتی به صورت زیر بیان می‌شود:

$$P(T_i, \delta_i | c_i = g; \theta) = \exp \left\{ - \sum_{k=1}^K H_k(T_i | c_i = g; \theta) \right\} \times \prod_{k=1}^K h_k(T_i | c_i = g; \theta)^{I(\delta_i=k)},$$

که در آن $h_k(T_i | c_i = g; \theta)$ مخاطره لحظه‌ای علت-مشخص تعریف شده در رابطه (۱۸) و $H_k(T_i | c_i = g; \theta)$ تابع مخاطره تجمعی متناظر با آن است. همچنین،

$$P(Y_i | c_i = g; \theta) = \int \left[\prod_j P(Y_i(s_{ij}) | c_i = g, b_i; \theta) \right] \times P(b_i | c_i = g; \theta) db_i, \quad j = 1, \dots, n_i.$$

همچنین، احتمالات پسینی را می‌توان از مدل توأم به دست آورد. احتمالات پسینی عضویت رده به شرط تمام مشاهدات داده شده برای نشانگر طولی و رخدادها را که به صورت $\pi_{ig}^{(Y,T)}(\theta) = Pr(c_i = g | Y_i, T_i, \delta_i; \theta)$ تعریف می‌کنیم را برای طبقه‌بندی پسینی افراد مورد استفاده قرار می‌دهیم، که در آن

$$\hat{\pi}_{ig}^{(Y,T)} = Pr(c_i = g | Y_i, \bar{T}_i, \delta_i; \hat{\theta}) = \frac{\hat{\pi}_{ig} P(Y_i | c_i = g; \hat{\theta}) P(T_i, \delta_i | c_i = g; \hat{\theta})}{\sum_{l=1}^G \hat{\pi}_{il} P(Y_i | c_i = l; \hat{\theta}) P(T_i, \delta_i | c_i = l; \hat{\theta})},$$

$$\hat{c}_i = \arg \max_g \left\{ Pr(c_i = g | Y_i, T_i, \delta_i; \hat{\theta}) \right\}.$$

۱۰.۲.۲ استفاده از الگوریتم EM برای برآورد پارامترهای مدل رده پنهان در داده‌های طولی و برآمدهای مخاطره‌های رقابتی

مشابه آنچه در زیر بخش قبلی برای برآورد پارامترهای مدل توأم رده پنهان داده‌های طولی و بقا با استفاده از الگوریتم EM شرح دادیم، در اینجا نیز می‌توان با استفاده از این الگوریتم پارامترهای مدل رده پنهان را برای مدل‌بندی توأم اندازه‌های طولی و داده‌های مخاطره رقابتی به دست آورد. تنها با این تفاوت که توزیع زمان رخداد مشاهدات با استفاده از مخاطره‌های رخداد K عامل تعریف شده در مدل به دست می‌آید. بنابراین خواهیم داشت:

$$Q(\theta | \theta^{(r)}) = \sum_{i=1}^n E_{b_i, c_i | Y_i, T_i, \delta_i, \bar{X}_i; \theta^{(r)}} [\ell_{com}(\theta; Y_i, T_i, \delta_i, \bar{X}_i, b_i, c_i)] = \sum_{i=1}^n E_{b_i, c_i | Y_i, T_i, \delta_i, \bar{X}_i; \theta^{(r)}} [\log Pr(c_i | \bar{X}_i) + \log P(b_i | c_i, \bar{X}_i) + \log P(Y_i | b_i, c_i, \bar{X}_i) + \log P(T_i, \delta_i | c_i, \bar{X}_i)] \quad (21)$$

که توزیع شرطی Y_i ها را نسبت به $c_i = g$ توصیف می‌کند را به صورت زیر در نظر می‌گیریم:

$$Y_i(s) = X_i'(s)\beta_g + Z_i'(s)b_{ig} + \epsilon_i(s), \quad (17)$$

$$\epsilon_i(s) \sim N(0, \sigma^2), \quad i = 1, \dots, n,$$

اندازه‌های تکرارشونده نشانگر طولی به صورت $Y_i(s) = (Y_i(s_{i1}), \dots, Y_i(s_{in_i}))'$ در زمان‌های متفاوت اندازه‌گیری $i = 1, \dots, n$ بردار X_i در آن بردار ضرایب متناظر با اثرات ثابت β_g و Z_i بردار ضرایب متناظر با اثرات تصادفی b_{ig} برای فرد i ام در رده g است. بردار دارای توزیع نرمال با میانگین و ماتریس کوواریانس رده مشخص به صورت $b_{ig} \sim N(\mu_g, \sigma_g^2 D)$ ، در نظر گرفته می‌شود و در آن $\sigma_g = 1$. همچنین یک مدل مخاطره متناسب رده-مشخص برای مخاطره‌های رقابتی تعریف می‌کنیم. در این مدل فرض می‌شود که مخاطره رخداد برای هر رده همگن است و بر اساس متغیرهای کمکی شرح داده می‌شود. مخاطره رخداد عامل k در رده g برای فرد i ام به صورت زیر در نظر گرفته می‌شود:

$$h_{ik}(t | c_i = g) = h_{\nu_{kg}}(t, \nu_{kg}) \exp(w_{ik}' \gamma_{kg}), \quad (18)$$

$$k = 1, \dots, K, \quad g = 1, \dots, G,$$

که در آن w_i بردار متغیرهای کمکی متناظر با مخاطره رخداد عامل k در رده g با بردار پارامترهای γ_{kg} است. تابع مخاطره پایه $h_{\nu_{kg}}$ برای هر رده و هر عامل مشخص است و به صورت پارامتری (با ν_{kg}) با استفاده از یک توزیع بقا استاندارد مثل وایبول یا گامپترتز و یا با استفاده از تعداد کمی از توابع پله‌ای یا M -اسپلاین‌ها برای انعطاف بیشتر، مدل‌بندی می‌شود. سرانجام یک مدل رگرسیون چندجمله‌ای بر اساس متغیرهای کمکی برای مدل احتمال عضویت رده پنهان به صورت زیر در نظر می‌گیریم:

$$\pi_{ig} = Pr(c_i = g) = \frac{\exp(u_i' \lambda_g)}{\sum_{L=1}^G \exp(u_i' \lambda_g)}, \quad (19)$$

که در آن u_i برداری از متغیرهای کمکی مرتبط با بردار پارامترهای رده-مشخص λ_g است. u_i شامل عرض از مبدأ است و در عمل، زمانی که هیچ پیش‌بینی‌کننده‌ای برای عضویت رده فرض نمی‌شود به عرض از مبدأ کاهش می‌یابد. برای شناساپذیری، یک رده مرجع لازم است که رده G را انتخاب می‌کنیم، در نتیجه $\lambda_G = 0$.

مدل رده پنهان توأم اندازه‌های طولی و داده‌های مخاطره‌های رقابتی با ماکسیم کردن لگاریتم درست‌نمایی برای تعداد معینی از رده پنهان G برآورد می‌شود، سپس مدلهایی با مقادیر متفاوتی از G توسط برخی معیارها، مانند BIC ، باهم مقایسه می‌شوند. برداری که شامل تمام پارامترهای سه زیر مدل تعریف شده (۱)، (۱۸) و (۱۹) با G رده پنهان است را با θ نمایش می‌دهیم. با در نظر گرفتن

که در آن

$$\log P(T_i, \delta_i | c_i, \check{X}_i) = -\sum_{k=1}^K H_k(T_i | c_i = g; \theta) + \sum_{k=1}^K I(\delta_i = k) \log \{h_k(T_i | c_i = g; \theta)\}$$

و بقیه عبارت‌ها در رابطه (۲۱)، مشابه موارد متناظر در زیر بخش قبلی محاسبه می‌شوند. در نتیجه، برآوردگرهای درست‌نمایی ماکسیمم با استفاده از الگوریتم EM را می‌توان به دست آورد.

$$\hat{\theta}_g^{(r+1)} = \arg \max_{\theta_g} \{Q(\theta | \theta^{(r)})\}.$$

۳ کاربرد

در نرم‌افزار R مدل‌های توأم رده پنهان با استفاده از تابع $\text{Jointlemm}()$ از بسته نرم‌افزاری **lemm** برازش داده می‌شوند [۳۰]. برخلاف تابع $\text{JointModel}()$ که ابتدا لازم است به‌طور جداگانه مدل‌های اثرات آمیخته خطی و بقا را برازش دهیم، تابع $\text{Jointlemm}()$ تنها نیاز به ارائه فرمول جداگانه دارد که قسمت‌های مختلف مدل را مشخص می‌کند.

به‌عنوان مثال، یک مدل رده پنهان از مجموعه داده ایدز برای مدل‌بندی توأم داده‌های طولی و زمان بقا و همچنین یک مدل رده پنهان برای داده‌های صرع به‌منظور مدل‌بندی توأم داده‌های طولی و مخاطره‌های رقابتی را بررسی می‌کنیم.

۱.۳ داده‌های ایدز

در مجموعه داده‌های ایدز، $n = 467$ بیمار آلوده به ویروس نقص ایمنی انسانی (HIV) به‌طور تصادفی تحت درمان با دو داروی دیدانوسین (ddI) و زالتیساین (ddC) قرار می‌گیرند. در این مطالعه، اندازه‌های مکرر تعداد سلول‌های $CD4$ هر دو ماه از ابتدا مطالعه تا ۲۰ ماه ثبت می‌شوند. همچنین، پیشامد زمان تا مرگ نیز در نظر

گرفته می‌شود. در این مطالعه، ارتباط بین تعداد $CD4$ ها و پیشامد مرگ، و همچنین تأثیر داروهای ddI و ddC بر روی این دو پیشامد مورد توجه قرار می‌گیرد. جزئیات بیشتر در مورد طراحی این مطالعه را می‌توان در [۱] یافت.

بنابراین، بر اساس این مجموعه داده، برای بخش طولی، در قسمت ثابت تأثیر اصلی زمان و درمان و در قسمت تصادفی عرض از مبدأ تصادفی و شیب تصادفی را در نظر می‌گیریم:

$$(22)$$

$$y_i(s) | c_i = g = \beta_{0g} + \beta_{1g}s + \beta_{2g}ddI_i + b_{i0} + b_{i1}s + \varepsilon_i(s), \\ i = 1, \dots, n, \quad g = 1, \dots, G.$$

برای سادگی، فرض می‌کنیم که هر دو متغیر در قسمت اثرات ثابت به رده وابسته هستند، از این رو، $(b_{i0}, b_{i1})' \sim N_2(0, D)$ و $\varepsilon_i(s) \sim N(0, \sigma^2)$ در نظر گرفته می‌شوند. برای زیرمدل بقا، فرض می‌کنیم که رده پنهان تابع مخاطره پایه و اثر درمان را مشخص می‌کند:

$$h_i(t | c_i = g) = h_{0g}(t) \exp(\gamma_g ddI_i) \quad (23)$$

توابع خطر پایه رده-مشخص با پارامتر γ_g به‌صورت تکه‌ای با شش گره و در صدک‌های مربوط به توزیع زمان رخداد فرض می‌شود، البته در نرم‌افزار R توابع پایه دیگری نیز موجود است. سرانجام، در مدل چندجمله‌ای عضویت پنهان، فرض می‌کنیم که احتمال پیشین هر بیمار متعلق به یک رده خاص به درمان بستگی دارد، یعنی:

$$\Pr(c_i = g) = \frac{\exp(\lambda_{0g} + \lambda_{1g} ddI_i)}{\sum_{l=1}^G \exp(\lambda_{0l} + \lambda_{1l} ddI_i)}. \quad (24)$$

چهار مدل توأم را به ترتیب با دو، سه، چهار و پنج رده پنهان برازش می‌دهیم. مقادیر AIC و BIC برای $g = 2, \dots, 5$ در جدول ۱ آورده شده است.

جدول ۱: مقادیر لگاریتم درست‌نمایی و معیارهای اطلاع برای تعداد رده پنهان مدل توأم برازش داده‌شده برای داده‌های ایدز

BIC	AIC	loglik	تعداد رده‌ها
۸۶۵۴٫۲۶	۸۵۵۴٫۷۵	-۴۲۵۳٫۳۷	۲
۸۶۵۰٫۴۴	۸۵۰۵٫۳۲	-۴۲۱۷٫۶۶	۳
۸۶۶۹٫۰۳	۸۴۷۸٫۳۰	-۴۱۹۳٫۱۵	۴
۸۷۲۴٫۳۷	۸۴۸۸٫۵۳	-۴۱۸۷٫۲۷	۵

صحيح زیرگروه‌های پنهان را پیشنهاد می‌کند [۲۳، ۲۹]. به دنبال این توصیه، مدل با سه رده پنهان را انتخاب می‌کنیم. همچنین، نتایج مدل رده پنهان توأم با سه رده پنهان برای این مجموعه داده، در جداول ۲ و ۳ خلاصه شده است.

می‌شود.

به‌طور متوسط ۴۶ اندازه‌گیری طولی از ۱ تا ۱۵ برای بیماران ثبت شده است.

مدل رده پنهان توأم را برای تحلیل این داده‌ها به کار می‌گیریم. مدل طولی را به‌صورت زیر در نظر می‌گیریم:

$$y_i(s)|c_i=g = \beta_{0g} + \beta_{1g}s + \beta_2 LG_i + \beta_{3g} LG_i s \quad (25)$$

$$+ b_{i0} + b_{i1}s + \epsilon_i(s);$$

$$i = 1, \dots, n_i; j = 1, \dots, n_i,$$

که در آن $(b_{i0}, b_{i1})' \sim N(0, D)$ (به عبارتی برای سادگی، $\sigma_g^2 =$ و $(1, g = 1, \dots, G)$ و $\epsilon_i \sim N(0, \sigma^2)$. همچنین، LG_i به‌عنوان اثر درمان در مدل به‌صورت دودویی تعریف می‌شود که مقدار ۱ را می‌گیرد اگر بیمار به‌طور تصادفی داروی LG را دریافت کند و مقدار صفر را می‌گیرد اگر بیمار داروی CBZ را دریافت کند. مدل مخاطره‌علت- مشخص را به‌صورت زیر در نظر می‌گیریم:

$$h_{ikg}(t) = h_{0k}(t) \exp\{\gamma_{0kg} + \gamma_{1k} LG_i\}, \quad (26)$$

$$k = 1, 2, g = 1, \dots, G$$

که در آن γ_{11} همان γ_{1SC} و γ_{12} برابر γ_{1UAE} و $\gamma_{0kG} = 0$ توابع خطر پایه به‌صورت تکه‌ای با پنج گره و با پارامتر ν_k در نظر گرفته می‌شوند. همچنین،

$$Pr(c_i = g) = \frac{\exp\{\lambda_g\}}{1 + \sum_{g=1}^{G-1} \exp\{\lambda_g\}}; g = 1, \dots, G. \quad (27)$$

لازم به ذکر است که برای انتخاب تعداد رده‌ها، $G = 1, \dots, 6$ در نظر گرفته شده است. از این رو، با مقایسه مقادیر BIC بهترین را انتخاب می‌کنیم. مقادیر لگاریتم درست‌نمایی، BIC و همچنین احتمال‌های مربوط به هر رده برای مقادیر مختلف G در جدول ۴ ارائه شده است. با توجه به اینکه معیار انتخاب را BIC قرار داده‌ایم، از آوردن مقادیر AIC در جدول خودداری نموده‌ایم. مشاهده می‌شود که بهترین مقدار BIC برای $G = 5$ حاصل شده است. مقادیر برآورد پارامترها برای $G = 5$ نیز در جدول ۵ گزارش شده است.

اختلاف نظر شدیدی بین دو معیار اطلاعات مشاهده می‌شود. به این صورت که، AIC به‌طور پیوسته‌ای کاهش می‌یابد مادامی‌که تعداد رده‌ها افزایش می‌یابد و مدل چهار رده‌ای را توصیه می‌کند، درحالی‌که مدل بهینه مطابق با BIC ، مدل سه رده است. مطالعات تجربی در این زمینه نشان داده است که BIC ، اغلب مدل با تعداد

بر اساس مقادیر آماره والد و پی-مقدار به‌دست‌آمده از جدول ۲، چنانچه پی-مقدار کمتر از ۰.۰۵ و قدر مطلق آماره والد بزرگ‌تر از ۱.۹۶ (مقدار توزیع نرمال استاندارد در سطح ۰.۰۵) باشد، ضریب متغیر مستقل معنی‌دار خواهد بود. بنابراین، مشاهده می‌کنیم که برای برخی پارامترها، مثل پارامتر زمان در قسمت طولی تفاوت بین طبقات قابل‌توجه است، به این صورت که، اثر زمان به جزء رده ۳ همه‌جا معنادار است. همه رده‌ها دارای عرض از مبدأ هستند و اثر درمان در هیچ رده‌ای معنادار نیست. تمام پارامترهای مربوط به تابع خطر پایه مدل بقا در رده ۱ معنادار هستند، از طرفی داروی ddI در هیچ رده‌ای روی بقا مؤثر نیست. همچنین، برای مدل احتمال، عرض از مبدأ و اثر درمان برای رده ۱ معنادار هستند.

۲.۳ داده‌های صرع

داده‌هایی که در این بخش تحلیل می‌شوند، مربوط به مطالعه داروهای استاندارد و جدید ضد صرع ($SANAD$)^{۱۵} است، که در آن، بیماران مبتلا به صرع را به آزمایش AED ^{۱۶} های جدید متقاعد می‌کنند.

در اینجا، تنها بیمارانی که با CBZ (کاربامازپین) یا LTG (لاموتریزین) تحت درمان هستند، در نظر گرفته می‌شوند. زمان تا رخداد عدم موفقیت درمان به‌عنوان پیشامد بقا ثبت می‌شود و تحلیل مخاطره‌های رقابتی برای عدم موفقیت درمان شامل عدم موفقیت درمان ناشی از کنترل نامناسب تشنج (ISC) یا اثرهای جانبی غیرقابل‌قبول (UAE) در نظر گرفته می‌شود [۲۶].

این مجموعه داده شامل ۶۰۵ بیمار است که به‌طور تصادفی ۲۹۲ نفر با CBZ و ۳۱۳ نفر با LTG تحت درمان هستند. ۹۴ بیمار به دلیل UAE و ۱۲۰ نفر به دلیل ISC طی حداکثر زمان پیگیری ۶۶ سال (میانگین ۲۸ سال)، از مطالعه خارج می‌شوند. خروج به دلایل دیگر در این مطالعه سانسور در نظر گرفته می‌شود. برای مقایسه AED ها پس از تنظیم نرخ تیتراسیون، ابتدا کالیراسیون دوز با استاندارد کردن دوز هر دو دارو نسبت به نقطه میانی دامنه نگهداری آن انجام شده است. این دوزهای کالیبره شده به‌عنوان اندازه‌گیری‌های طولی در مدل توأم مخاطره‌های رقابتی در نظر گرفته

¹⁵Standard and New Anti-epileptic Drugs

¹⁶Anti Epileptic Drugs

جدول ۲: نتایج مدل رده پنهان توأم اندازه‌های طولی و داده‌های زمان بقا برای مجموعه داده ایدز با $g = 3$

پارامتر رده برآورد خطای استاندارد آماره والد پی-مقدار					
اثرهای ثابت مدل طولی					
β_{01}	۱ رده	۴,۵۵۲	۰,۲۱۰	۲,۶۹۹	۰,۰۰۰
β_{02}	۲ رده	۱۱,۹۷۳	۰,۶۱۲	۱۹,۵۷۷	۰,۰۰۰
β_{03}	۳ رده	۱۵,۳۸۵	۰,۵۰۴	۳۰,۵۰۳	۰,۰۰۰
β_{11}	۱ رده	-۰,۱۳۵	۰,۰۱۷	-۷,۷۵۰	۰,۰۰۰
β_{12}	۲ رده	-۰,۳۳۴	۰,۰۴۴	-۷,۵۶۰	۰,۰۰۰
β_{13}	۳ رده	۰,۰۰۸	۰,۰۳۲	۲,۲۴۳	۰,۸۰۸
β_{21}	۱ رده	۰,۱۹۲	۰,۲۷۲	۰,۷۰۵	۰,۴۸۱
β_{22}	۲ رده	-۰,۳۶۴	۰,۰۹۰۳	-۰,۴۰۴	۰,۶۸۶
β_{23}	۳ رده	-۱,۰۳۴	۰,۵۷۰	-۱,۸۱۵	۰,۰۶۹
پارامترهای مدل بقا					
ν_{11}	۱ رده	۰,۱۷۵	۰,۰۱۶	۱۰,۶۱۲	۰,۰۰۰
ν_{21}	۱ رده	۰,۱۷۸	۰,۰۱۷	۱۰,۶۲۸	۰,۰۰۰
ν_{31}	۱ رده	۰,۲۰۶	۰,۰۱۹	۱۰,۵۲۴	۰,۰۰۰
ν_{41}	۱ رده	۰,۲۸۴	۰,۰۲۶	۱۰,۱۰۶	۰,۰۰۰
ν_{51}	۱ رده	۰,۲۲۰	۰,۰۲۲	۱۰,۱۰۸	۰,۰۰۰
ν_{12}	۲ رده	۰,۰۷۵	۰,۰۴۳	۱,۷۶۸	۰,۰۷۷
ν_{22}	۲ رده	۰,۰۰۰	۰,۰۷۷	۰,۰۰۴	۰,۹۹۷
ν_{32}	۲ رده	۰,۰۰۰	۰,۰۶۵	۰,۰۰۶	۰,۹۹۵
ν_{42}	۲ رده	۰,۰۰۰	۰,۰۹۴	۰,۰۰۰	۱,۰۰۰
ν_{52}	۲ رده	۰,۱۷۸	۰,۰۵۳	۳,۳۳۳	۰,۰۰۱
ν_{13}	۳ رده	۰,۰۱۴	۰,۰۲۹	۰,۴۶۶	۰,۶۴۱
ν_{23}	۳ رده	-۰,۰۱۲	۰,۰۲۶	-۰,۴۵۹	۰,۶۴۶
ν_{33}	۳ رده	۰,۰۲۲	۰,۰۴۷	۰,۴۷۷	۰,۶۳۳
ν_{43}	۳ رده	۰,۰۰۰	۰,۰۱۴	۰,۰۰۰	۱,۰۰۰
ν_{53}	۳ رده	۰,۰۰۰	۰,۰۱۳	۰,۰۰۰	۱,۰۰۰
γ_1	۱ رده	۰,۲۸۱	۰,۱۵۹	۱,۷۶۹	۰,۰۷۷
γ_2	۲ رده	-۰,۳۹۳	۰,۹۵۱	-۰,۴۱۳	۰,۶۸۰
γ_3	۳ رده	۳,۶۰۳	۴,۲۰۶	۰,۸۵۷	۰,۳۹۲
اثرات ثابت مدل احتمال عضویت رده					
λ_{01}	۱ رده	۱,۸۹۱	۰,۲۳۶	۸,۰۲۶	۰,۰۰۰
λ_{02}	۲ رده	۰,۳۷۸	۰,۳۱۲	۱,۲۱۰	۰,۲۲۶
λ_{11}	۱ رده	-۰,۶۱۷	۰,۳۰۴	-۲,۰۳۲	۰,۰۴۲
λ_{12}	۲ رده	-۰,۸۳۳	۰,۴۳۴	-۱,۹۲۱	۰,۰۵۵

جدول ۳: ادامه نتایج مدل رده پنهان توأم برای $g = 3$ بر اساس داده‌های ایدز

ماتریس واریانس-کواریانس اثرات تصادفی (مستقل از رده‌ها)	
پارامتر	برآورد
D_{11}	۴۹۹۸
D_{12}	-۰٫۲۴۰
D_{22}	۰٫۰۱۷
خطای استاندارد مانده	
پارامتر	برآورد
σ	۱٫۷۳۴
خطای استاندارد	
	۰٫۰۴۷

جدول ۴: مقادیر لگاریتم درستنمایی، معیار اطلاع و احتمال‌های عضویت هر رده، برای ۶ رده پنهان مدل توأم برازش داده‌شده بر اساس داده‌های صرع

تعداد رده (G)	loglik	BIC	% رده ۱	% رده ۲	% رده ۳	% رده ۴	% رده ۵	% رده ۶
۱	-۳۵۷۶٫۵۸	۷۲۶۸٫۴۵	۱۰۰	۰	۰	۰	۰	۰
۲	-۳۴۷۳٫۵۶	۷۱۰۰٫۸۶	۳۲٫۰۷	۶۷٫۹۳	۰	۰	۰	۰
۳	-۳۳۳۸٫۷۰	۷۰۶۹٫۵۵	۷۰٫۵۸	۱۰٫۰۸	۱۹٫۳۴	۰	۰	۰
۴	-۳۴۰۶٫۶۹	۷۰۴۴٫۳۹	۶٫۴۵	۷۴٫۷۱	۱۳٫۸۸	۴٫۹۶	۰	۰
۵	-۳۳۶۰٫۸۹	۶۹۹۰٫۸۰	۶٫۲۸	۱۵٫۲۱	۶۶٫۴۵	۷٫۱۱	۴٫۹۶	۰
۶	-۳۳۹۸٫۲	۷۱۰۳٫۸۵	۰	۷٫۲۷	۷۲٫۲۳	۱۶٫۳۶	۲٫۹۸	۱٫۱۶

نتایج جدول ۵ نشان می‌دهد که همه رده‌ها جزء رده ۱ برای مدل طولی دارای عرض از مبدأ هستند. اثر زمان فقط برای رده‌های ۲ و ۴ معنی‌دار است. اثر متقابل تیمار و زمان برای همه رده‌ها، جزء رده ۱ و ۳ معنی‌دار است. اثر مخاطره رقابتی اول فقط برای رده ۳ معنی‌دار نیست در صورتی که اثر مخاطره رقابتی دوم تنها در رده اول معنی‌دار است. همچنین، اثر درمان مخاطره دوم (γ_{UAE}) معنی‌دار است.

جدول ۵: نتایج مدل رده پنهان توأم اندازه‌های طولی و داده‌های مخاطره رقابتی برای $g = 5$ بر اساس داده‌های صرع

پارامتر	رده	برآورد	خطای استاندارد	آماره والد پی-مقدار
اثرهای ثابت مدل طولی				
β_{01}	۱ رده	۰/۲۵۳	۰/۱۹۶	۱/۲۸۸
β_{02}	۲ رده	۱/۸۴۹	۰/۰۹۷	۱۹/۰۰۴
β_{03}	۳ رده	۱/۸۸۷	۰/۰۵۲	۳۶/۶۴۲
β_{04}	۴ رده	۲/۲۱۲	۰/۱۴۵	۱۵/۲۰۶
β_{05}	۵ رده	۳/۸۷۱	۰/۱۹۴	۱۹/۹۰۹
β_2	-	-۰/۱۶۴	۰/۰۶۱	-۲/۷۰۲
β_{11}	۱ رده	-۰/۴۹۲	۱/۱۹۴	-۰/۴۱۲
β_{12}	۲ رده	۰/۵۱۶	۰/۰۹۱	۵/۷۰۱
β_{13}	۳ رده	۰/۰۳۸	۰/۰۲۵	۱/۴۸۳
β_{14}	۴ رده	۱/۶۷۳	۰/۲۲۴	۷/۴۷۱
β_{15}	۵ رده	۰/۱۰۶	۰/۱۰۲	۱/۰۳۷
β_{21}	۱ رده	۱/۲۳۷	۱/۱۲۲	۱/۱۰۳
β_{22}	۲ رده	۰/۴۳۲	۰/۱۰۴	۴/۱۳۹
β_{23}	۳ رده	۰/۰۶۵	۰/۰۳۴	۱/۹۰۴
β_{24}	۴ رده	۰/۷۶۱	۰/۲۷۰	۲/۸۱۴
β_{25}	۵ رده	۰/۲۶۵	۰/۱۲۳	۲/۱۵۳
مدل مخاطره‌های رقابتی				
ν_{11}	$\log(\text{piecewise}1)$	-۳/۸۲۵	۰/۶۲۶	-۶/۱۱۲
ν_{21}	$\log(\text{piecewise}2)$	-۲/۹۸۳	۰/۵۹۶	-۵/۰۰۷
ν_{31}	$\log(\text{piecewise}3)$	-۲/۴۶۸	۰/۵۷۴	-۴/۳۰۱
ν_{41}	$\log(\text{piecewise}4)$	-۲/۳۳۳	۰/۵۷۵	-۴/۰۶۱
γ_{011}	۱ رده	۳/۶۱۵	۰/۷۴۴	۴/۸۶۱
γ_{012}	۲ رده	۱/۷۱۵	۰/۵۷۸	۲/۹۶۷
γ_{013}	۳ رده	-۱/۲۴۵	۰/۷۰۸	-۱/۷۵۸
γ_{014}	۴ رده	۳/۳۲۷	۰/۶۳۲	۵/۲۶۷
ν_{12}	$\log(\text{piecewise}1)$	-۲/۳۱۶	۰/۷۸۰	-۳/۶۱۴
ν_{22}	$\log(\text{piecewise}2)$	-۲/۹۲۳	۰/۷۷۹	-۳/۷۵۴
ν_{32}	$\log(\text{piecewise}3)$	-۳/۷۱۷	۰/۸۱۷	-۴/۵۵۱
ν_{42}	$\log(\text{piecewise}4)$	-۴/۷۷۹	۱/۰۱۸	-۴/۶۹۷
γ_{021}	۱ رده	۴/۲۱۱	۰/۸۰۷	۵/۲۱۷
γ_{022}	۲ رده	۰/۵۸۴	۱/۲۶۹	۰/۴۶۰
γ_{023}	۳ رده	۰/۴۴۷	۰/۷۹۹	۰/۵۵۹
γ_{024}	۴ رده	۰/۷۷۲	۱/۲۱۰	۰/۶۳۸
γ_{ISE}	-	-۰/۱۴۲	۰/۲۳۲	-۰/۶۰۹
γ_{UAE}	-	-۰/۶۹۳	۰/۲۴۸	-۲/۷۹۴
پارامترهای مدل احتمال				
λ_1	۱ رده	۰/۱۸۱	۰/۳۱۶	۰/۵۶۶
λ_2	۲ رده	۱/۳۳۱	۰/۳۲۱	۴/۱۴۸
λ_3	۳ رده	۲/۴۱۰	۰/۲۷۷	۸/۶۹۸
λ_4	۴ رده	۰/۳۱۹	۰/۳۴۸	۰/۳۵۹

جدول ۶: ادامه نتایج مدل رده پنهان توأم برای $g = 5$ بر اساس داده‌های صرع

ماتریس واریانس-کواریانس اثرات تصادفی (مستقل از رده‌ها)	
پارامتر	برآورد
D_{11}	۰٫۲۸۶
D_{12}	۰٫۰۲۷
D_{22}	۰٫۰۲۸
خطای استاندارد مانده	
پارامتر	برآورد
σ	۰٫۴۳۲
	خطای استاندارد
	۰٫۰۰۷

۴ بحث و نتیجه‌گیری

رخداد پیشامد برای هر رده را در نظر گرفته‌اند در صورتی که لیو و همکاران [۲۵]، از این فرض استفاده نکردند و مدل جدیدی را با ادغام مدل اثرات تصادفی توأم در چارچوب رده پنهان ارائه کردند. بیشترین روش استفاده‌شده در منابع، برای برآورد پارامترهای مدل رده پنهان توأم، روش درستنمایی ماکسیمم و به‌ویژه استفاده از الگوریتم EM بوده است، البته برخی هم روش استنباط بیزی را به‌کاربرده‌اند [۱۱، ۲]. لازم به ذکر است که بیشتر منابع، شامل کاربردهایی از آزمایش‌های بالینی است.

اگرچه در سال‌های اخیر، استفاده از مدل رده پنهان برای مدل‌بندی توأم داده‌های طولی و زمان تا رخداد پیشامد بیشتر موردتوجه قرار گرفته است، اما در نظر گرفتن متغیرهای ورودی بیشتری در زیر مدل بقا به‌منظور لحاظ کردن اطلاعات اضافی در خصوص بیماران و همچنین کاربردهای بیشتر در زمینه پژوهش‌های بالینی، معرفی پکیج‌های کاربردی جهت پردازش سریع، روش‌های انتخاب رده بهینه و استفاده از فرض‌های توزیعی دیگر برای اثرهای تصادفی به‌جای فرض توزیعی نرمال، می‌تواند موضوعات قابل‌بحث در پژوهش‌های آینده باشد.

در سال‌های اخیر، مدل‌های رده پنهان برای تحلیل توأم اندازه‌های طولی و بقا مورد استفاده قرار گرفته است. این مدل‌ها، فرض می‌کنند که جمعیت در G رده پنهان تقسیم می‌شوند که با توابع مخاطره مختلف برای رخداد پیشامد و نمایش‌های متفاوت از نشانگر طولی توسط یک مدل آمیخته برای هر رده قابل توصیف است. در این مقاله، مطالعاتی که در دو دهه اخیر در این زمینه صورت گرفته است را مرور کردیم و یک نمای کلی از مدل‌های رده پنهان توأم ارائه دادیم. همچنین، کاربردهایی از مدل توأم رده پنهان را بر اساس دو مجموعه داده ایدز و داروهای ضد صرع بیان کردیم. ملاحظه کردیم که مدل‌های رده پنهان ابتدا بر اساس یک نشانگر طولی و یک زمان رخداد مطالعه شده‌اند، سپس به حالت‌های پیچیده‌تری تعمیم یافته‌اند، حالت‌هایی که ممکن است چندین نشانگر طولی یا متغیر طولی چندمتغیره شامل متغیر پیوسته، گسسته ترتیبی یا رسته‌ای و یا چند رخداد به مفهوم علت-مشخص برای مخاطره‌های رقابتی را در نظر بگیرند. بیش‌تر روش‌های استفاده‌شده در این زمینه مدل‌بندی، فرض استقلال شرطی بین فرایند طولی و زمان تا

مراجع

- [1] Abrams, D. I., Goldman, A. I., Launer, C., Korvick, J. A., Neaton, J. D., Crane, L. R., and Terry Bein Community Programs for Clinical Research on AIDS. (1994). A comparative trial of didanosine or zalcitabine after treatment with zidovudine in patients with human immunodeficiency virus infection. *New England Journal of Medicine*, **330**, 657-662.
- [2] Andrinopoulou, E. R., Nasserinejad, K., Szczesniak, R., and Rizopoulos, D. (2020). Integrating latent classes in the Bayesian shared parameter joint model of longitudinal and survival outcomes. *Statistical methods in medical research*, **29**, 3294-3307.
- [3] Baghfalaki, T., Ganjali, M. and Berridge, D. (2013). Robust joint modeling of longitudinal measurements and time to event data using normal/independent distributions: a Bayesian approach. *Biometrical Journal*, **55**, 844-865.
- [4] Baghfalaki, T., Ganjali, M., and Hashemi, R. (2014). Bayesian joint modeling of longitudinal measurements and time-to-event data using robust distributions. *Journal of biopharmaceutical Statistics*, **24**, 834-855.
- [5] Baghfalaki, T. and Ganjali, M. (2015). A Bayesian approach for joint modeling of skew-normal longitudinal measurements and time to event data. *REVSTAT-Statistical Journal*, **13**, 169-191.
- [6] Beunckens, C., Molenberghs, G., Verbeke, G., and Mallinckrodt, C. (2008). A latent-class mixture model for incomplete longitudinal Gaussian data. *Biometrics*, **64**, 96-105.
- [7] Chen, J., and Huang, Y. (2015). A Bayesian mixture of semiparametric mixed-effects joint models for skewed-longitudinal and time-to-event data. *Statistics in medicine*, **34**, 2820-2843.
- [8] Diggle, P. J., Sousa, I. and Chetwynd, A. G. (2008). Joint modelling of repeated measurements and time-to-event outcomes: the fourth Armitage lecture. *Statistics in Medicine*, **27**, 2981-2998.
- [9] Entink K. R. H., Fox, J. P., and van den Hout, A. (2011). A mixture model for the joint analysis of latent developmental trajectories and survival. *Statistics in medicine*, **30**, 2310-2325.
- [10] Faucett, C.L and Thomas, D. C. (1996). Simultaneously modeling censored survival data and repeatedly measured covariates: a Gibbs sampling approach. *Statistics in Medicine*, **15**, 1663-1685.
- [11] Garre, F. G., Zwiderman, A. H., Geskus, R. B., and Sijpkens, Y. W. (2008). A joint latent class changepoint model to improve the prediction of time to graft failure. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, **171**, 299-308.
- [12] Gould, A. L., Boye, M. E., Crowther, M. J., Ibrahim, J. G., Quartey, G., Micallef, S. and Bois, F. Y. (2015) Joint modeling of survival and longitudinal non-survival data: current methods and issues: Report of the DIA Bayesian joint modeling working group. *Statist. Med.*, **34**, 2181-2195.
- [13] Guler, I., Faes, C., Cadarso-Suárez, C., Teixeira, L., Rodrigues, A. and Mendonca, D. (2017). Two stage model for multivariate longitudinal and survival data with application to nephrology research. *Biometrical Journal*, **59**, 1204-1220.
- [14] Han, J. , Slate, EH., and Peñ˜sa EA. (2007). Parametric latent class joint model for a longitudinal biomarker and recurrent events. *Statistics in Medicine*, **26**, 5285-5302.
- [15] Henderson, R., Diggle, p. and Dobson, A. (2000). Joint modeling of longitudinal measurements and event time data. *Biostatistics*, **4**, 465-480.

- [16] Hogan, J. and Laird, N. (1997). Mixture models for the joint distribution of repeated measures and event times. *Statistics in Medicine*, **16**, 239 – 258.
- [17] Huang, Y., Dagne, G. and Wu, L. (2011). Bayesian inference on joint models of HIV dynamics for time-to-event and longitudinal data with skewness and covariate measurement errors. *Statistics in Medicine*, **30**, 2930–2946.
- [18] Huang, Y., Lu, X., Chen, J., Liang, J., and Zangmeister, M. (2018). Joint model-based clustering of nonlinear longitudinal trajectories and associated time-to-event data analysis, linked by latent class membership: with application to AIDS clinical studies. *Lifetime data analysis*, **24**, 699-718.
- [19] Ibrahim, J. G., Chen, M. and Sinha, D. (2004). Bayesian methods for joint modeling of longitudinal and survival data with applications to cancer vaccine trials, *Statistica Sinica*, **14**, 863–883.
- [20] Jacqmin-Gadda, H., Proust-Lima, C., Taylor, J. M., and Commenges, D. (2010). Score test for conditional independence between longitudinal outcome and time to event given the classes in the joint latent class model. *Biometrics*, **66**, 11-19.
- [21] Larsen, K. (2004). Joint analysis of time-to-event and multiple binary indicators of latent classes. *Biometrics*, **60**, 85-92.
- [22] Li, M., Lee, C. W., and Kong, L. (2020). A latent class approach for joint modeling of a time-to-event outcome and multiple longitudinal biomarkers subject to limits of detection. *Statistical methods in medical research*, **29**, 1624-1638.
- [23] Lin, H., McCulloch, C., and Rosenheck, R. (2004). Latent pattern mixture models for informative intermittent missing data in longitudinal studies. *Biometrics*, **60**, 295 – 305.
- [24] Lin, H., Turnbull, B. W., McCulloch, C. E., and Slate, E. H. (2002). Latent class models for joint analysis of longitudinal biomarker and event process data: Application to longitudinal prostate-specific antigen readings and prostate cancer. *Journal of the American Statistical Association*, **97**, 53–65.
- [25] Liu, Y., Liu, L., and Zhou, J. (2015). Joint latent class model of survival and longitudinal data: An application to CPCRA study. *Computational Statistics and Data Analysis*, **91**, 40-50.
- [26] Marson, A. G., Al-Kharusi, A. M., Alwaidh, M., Appleton, R., Baker, G. A., Chadwick, D. W., Cramp, C., Cockerell, O. C., Cooper, P. N., Doughty, J., Eaton, B., Gamble, C., Goulding, P. J., Howell, S. J. L., Hughes, A., Jackson, M., Jacoby, A., Kellett, M., Lawson, G. R., Leach, J. P., Nicolaides, P., Roberts, R., Shackley, P., Shen, J., Smith, D. F., Smith, P. E. M., Tudur-Smith, C., Vanoli, A. and Williamson, P. R. (2007). The SANAD study of effectiveness of carbamazepine, gabapentin, lamotrigine, oxcarbazepine, or topiramate for treatment of partial epilepsy: an unblinded randomised controlled trial. *Lancet*, **369**, 1000–1015.
- [27] Proust-Lima, C., Dartigues, J.F., and Jacqmin-Gadda, H. (2016). Joint modeling of repeated multivariate cognitive measures and competing risks of dementia and death: A latent process and latent class approach. *Statistics in Medicine*, **35**, 382–398.
- [28] Proust, C., Jacqmin-Gadda, H., Taylor, J. M., Ganiayre, J., and Commenges, D. (2006). A nonlinear model with latent process for cognitive evolution using multivariate longitudinal data. *Biometrics*, **62**, 1014-1024.
- [29] Proust-Lima, C., Joly, P., Dartigues, J.-F., and Jacqmin-Gadda, H. (2009). Joint modelling of multivariate longitudinal outcomes and a time-to-event: A nonlinear latent class approach. *Computational Statistics and Data Analysis*, **53**, 1142–1154.
- [30] Proust-Lima, C., and Liqueur, B. (2011). LCMM: an R package for estimation of latent class mixed models and joint latent class models. In *The R User Conference, useR! 2011 August 16-18 2011 University of Warwick, Coventry, UK*, (p. 66)

- [31] Proust-Lima, C., Philipps, V., Diakite, A., Liqueet, B., Proust, M. C., and Lima, P. (2021). Package ‘lcm’.
[32] Proust-Lima, C., S’ene, M., Taylor, J. M. G., and Jacqmin- Gadda, H. (2014). Joint latent class models for longitudinal and time-to-event data: A review. *Statistical Methods in Medicine Research*, **23**, 74–90.
[33] Proust-Lima, C., and Taylor, J. M. (2009). Development and validation of a dynamic prognostic tool for prostate cancer recurrence using repeated measures of post treatment PSA: a joint modeling approach. *Biostatistics*, **10**, 535-549.
[34] Qin, L., Weissfeld, L. A., Shen, C., and Levine, M. D. (2009). A two-latent-class model for smoking cessation data with informative dropouts. *Communications in Statistics—Theory and Methods*, **38**, 2604-2619.
[35] Rizopoulos, D. (2012). *Joint Models for Longitudinal and Time-to-Event Data: With Applications in R*, Vol. 6. Boca Raton:Chapman and Hall.
[36] Rizopoulos, D. (2012). Fast fitting of joint models for longitudinal and event time data using a pseudo-adaptive Gaussian quadrature rule. *Computational Statistics and Data Analysis*, **56**, 491 – 501.
[37] Rouanet, A., Joly, P., Dartigues, J. F., Proust-Lima, C., and Jacqmin-Gadda, H. (2016). Joint latent class model for longitudinal data and interval-censored semi-competing events: Application to dementia. *Biometrics*, **72**, 1123-1135.
[38] Tsiatis, A. A. and Davidian, M. (2004). Joint modeling of longitudinal and time-to-event data : An overview. *Statistica Sinica*, **14**, 809–834.
[39] Wu, L., Liu, W., Yi, G. Y. and Huang, Y. (2012). Analysis of longitudinal and survival data: joint modeling, inference methods, and issues. *Journal of Probability and Statistics*.
[40] Wulfsohn, M. S. and Tsiatis, A. A. (1997). A joint model for survival and longitudinal data measured with error. *Biometrics*, **53**, 330-339.
[41] Ye, W., Lin, X. and Taylor, J. (2008). Semiparametric modeling of longitudinal measurements and time-to-event data a two stage regression calibration approach. *Biometrics*, **64**, 1238-1246.

A view on latent class models for joint modeling of longitudinal measurements and survival data

Parvaneh Mehdizadeh¹, Taban Baghfalaki², Mahdy Esmailian¹

Abstract:

Joint models are used in follow-up studies to investigate the relationship between longitudinal marker and a survival outcome and have been generalized to analyze multiple markers or competing risks data. Many statistical achievements in the field of joint modeling focus on shared random effects models which include characteristics of longitudinal markers as explanatory variables in the survival model. A less-known approach is the joint latent class model, assuming that a latent class structure fully captures the relationship between the longitudinal markers and the event risk. The latent class model may be appropriate because of the flexibility in modeling the relationship between the longitudinal marker and the time of event, as well as the ability to include explanatory variables, especially for predictive problems. In this paper, we provide an overview of the joint latent class model and its generalizations. In this regard, first a review of the discussed models is introduced and then the estimation of the model parameters is discussed. In the application section, two real data sets are analyzed.

Keywords: Competing risks, Latent class model, Longitudinal data, Maximum likelihood estimator, Survival model, The EM Algorithm.

¹Department of Statistics and Computer Sciences, Faculty of Sciences, University of Mohagheh Ardabili

²Department of Statistics, Faculty of Mathematical Sciences, Tarbiat Modares University