

مجله علوم آماری، پاییز و زمستان ۱۳۹۳

جلد ۸ شماره ۲، ص ۱۸۵-۲۰۵

چگالی‌های پیشین با ساختار معین برای تحلیل بیزی جدول‌های پیشاپردازی کامل و ناقص

سید کامران قربیشی

گروه آمار، دانشگاه قم

تاریخ دریافت: ۱۳۹۲/۱۲/۵ تاریخ آخرین بازنگری: ۱۳۹۳/۸/۱۹

چکیده: در تحلیل جدول‌های پیشاپردازی، اغلب محققان از چگالی‌های پیشین خاص برای پارامترهای مدل لگ-خطی یا احتمال خانه‌های جدول استفاده می‌کنند. اما در عمل گاهی اطلاعات با ارزشی، ترجیح‌آمیز، در خصوص نسبت بخت‌های (تعمیم یافته) وجود دارد. لذا محقق نیازمند به رهیافت قوی‌تری است که بتواند باور پیشین خود را روی نسبت بخت‌های تعمیم یافته قرار دهد. از این توزیع‌های پیشین به عنوان چگالی‌های پیشین با ساختار معین یاد خواهد شد. در این مقاله ابتدا الگوی کلی چگالی‌های پیشین با ساختار معین معرفی خواهند شد. سپس به دلیل کاربرد وسیع این پیشین‌ها در آزمایه‌های بالینی و به‌ویژه در جدول‌های پیشاپردازی کامل و ناقص 2×2 ، پیشین‌های متناظر این حالت تحت سه شرط متفاوت به دست آورده می‌شوند.

واژه‌های کلیدی: تحلیل بیزی، جدول‌های پیشاپردازی، چگالی پیشین، مدل‌های نسبیتی، نسبت بخت‌ها، نسبت مخاطره.

آدرس الکترونیک مسئول مقاله: سید کامران قربیشی، atty_ghoreishi@yahoo.com

کد موضوع بنای ریاضی (۲۰۱۰): ۶۲H۱۷، ۶۲H۱۲

چگالی‌های پیشین با ساختار معین

جدول‌های پیشایندی به طور وسیعی در زمینه‌های مختلف پزشکی، جامعه‌شناسی، علوم رفتاری و غیره کاربرد دارند. تاکنون آمارشناسان فراوانی به تحلیل بسامدی و بیزی این جدول‌ها پرداخته‌اند. با این وجود به نظر می‌رسد تحلیل بیزی، در مواردی که استفاده از چگالی پیشین خاصی مورد نظر است، از اقبال کمتری برخوردار بوده باشد. شاید دلیل اصلی این موضوع را بتوان در فقدان یک الگوی مناسب دانست که محقق را قادر سازد تا هر اعتقاد پیشین خود را در قالب یک مدل کارآمد به احتمال خانه‌های جدول پیشایندی، اعم از کامل یا ناقص، تحمیل نماید.

مرور منابع موجود در خصوص تحلیل بیزی جدول‌های پیشایندی نشانگر آن است که اولین توزیع پیشین برای پارامترهای توزیع چند جمله‌ای توزیع دریکله بوده است. (لیندلی، ۱۹۶۴؛ گود، ۱۹۶۵). ابتدا سوال مهمی که در این زمینه مطرح گردید این بود که آیا این توزیع پیشین از چنان انعطاف‌پذیری برخوردار است که بتواند هر باور پیشینی را در برگیرد؟ آمارشناسانی نظیر لئونارد (۱۹۷۳)، آیتیجینسون (۱۹۸۵)، جوتیس (۱۹۹۳) و فورستر و اسکین (۱۹۹۴) خاطر نشان کردند که به دلیل دارا بودن تعداد پارامترهای کم، توزیع دریکله از انعطاف‌پذیری مناسب برای کاربرد وسیع در دنیای واقعی برخوردار نیست. یکی از محدودیت‌های این چگالی این است که به کمک آن نمی‌توان هر ساختار پیوندی معینی را، که باور پیشین محقق برآن دلالت دارد، در تحلیل بیزی شرکت داد. به همین دلیل آمارشناسان مذکور استفاده از چگالی پیشین نرم‌ال چند متغیره را پیشنهاد کردند. آن‌ها از این چگالی برای پارامترهای مدل لوجیت چند جمله‌ای استفاده نمودند. اگرچه چگالی پیشین نرم‌ال چندمتغیره نسبت به چگالی دریکله از انعطاف‌پذیری قابل قبولی برخوردار بود لیکن تمرکز هر دو چگالی پیشین بر برآورد بیزی احتمال خانه‌های جدول استوار بود. این در حالی است که به منظور مطالعه ساختار پیوند در تحلیل جدول‌های پیشایندی، بیشترین تمرکز بر پارامترهای مدل لگ-خطی معطوف است تا احتمال خانه‌های جدول. به همین دلیل لئونارد (۱۹۷۵) از چگالی پیشین نرم‌ال برای پارامترهای مدل لگ-خطی استفاده نمود. در ادامه کار او لرد (۱۹۷۸) احتمال

خانه‌های جدول را با استفاده از مدل لگ-خطی و رهیافت بیز تجربی برآورد کرد. تاکنون کارهای گوناگون دیگری شامل: استفاده از توزیع پیشین پایا، مدل‌های دیریکله سلسله مراتبی، تحلیل بیزی مدل‌های پیوند و مدل‌های توافقی در این رابطه صورت گرفته که در اینجا از ذکر جزئیات آنها خودداری شده به منبع بسیار غنی اگرستی و هیچکوک (۲۰۰۵) ارجاع داده می‌شود.

اگر چه امروزه استفاده از چگالی پیشین مبتنی بر توزیع نرمال، برای پارامترهای پیوند، در مدل‌های لگ-خطی، مبنای تحلیل بیزی جدول‌های پیشاپندی است، اما باید توجه داشت که مبنا و شالوده این مدل‌ها نسبت بختهای مجاور است که در یک جدول دو بعدی $J \times I$ تعداد آنها برابر $(J - 1)(I - 1)$ بوده و به صورت

$$\theta_{ij} = \frac{p_{ij}p_{i+1,j+1}}{p_{i+1,j}p_{i,j+1}}; \quad i = 1, \dots, I - 1, \quad j = 1, \dots, J - 1 \quad (1)$$

تعریف می‌شوند (بیشап و همکاران، ۱۹۷۵)، که در آن p_{ij} احتمال خانه (j, i) ام جدول است. با لگاریتم‌گیری از (۱) نمایش معادل عبارت است از:

$$\log \theta_{ij} = \log p_{ij} + \log p_{i+1,j+1} - \log p_{i+1,j} - \log p_{i,j+1}. \quad (2)$$

حالات مختلف نسبت بختهای مجاور منجر به مدل‌های پیوند مختلف و در نتیجه مدل‌های لگ-خطی گوناگون می‌شود (بیشап و همکاران، ۱۹۷۵؛ گودمن، ۱۹۷۲). بنابراین طبق رابطه (۲) توزیع پیشین روی پارامترهای مدل لگ-خطی معادل توزیع پیشین روی لگاریتم نسبت بختهای در نتیجه روی خود نسبت بختهای است. اکنون سوال بسیار مهمی که ممکن است مطرح شود این است که آیا هر باور پیشین روی نسبت بختهای یا لگاریتم آنها را می‌توان در قالب مدل‌های لگ-خطی وارد تحلیل بیزی کرد؟ به عنوان مثال، اگر در یک جدول $J \times I$ باور ذهنی محقق این باشد که نسبت بختهای مجاور برای چهار خانه، γ برابر نسبت بختهای مجاور برای چهار خانه دیگر از همان جدول است که در آن γ دارای توزیع پیشین داده شده است، آیا می‌توان مدل لگ-خطی متناظر و در نتیجه چگالی پیشین، برای احتمال خانه‌های جدول، متناظر با آن را تعیین نمود؟

لازم به ذکر است که این قبیل سوالات جزء لاینفک مسایل مطرح در زمینه‌های جامعه‌شناسی، روانشناسی و پژوهشکی است. علاوه بر آن که لزومی ندارد این سوال

۱۸۸ چگالی‌های پیشین با ساختار معین

صرفاً به نسبت بخت‌های مجاور محدود شود بلکه برای هر چهار خانه دلخواه از جدول قابل بررسی است.

سوال مشابه را می‌توان در جدول‌های پیش‌سازنده ناقص مطرح نمود. برای پرداختن دقیق به آن فرض کنید که با جدول ناقص 2×2 زیر مواجه‌ایم.

p_{11}	p_{12}
-	p_{22}

طبق تعریف تنگ و جیانگ (۲۰۱۱)، نسبت مخاطره^۱ برای این جدول ناقص به صورت

$$\phi = \frac{\frac{p_{11}}{p_{11} + p_{12}}}{p_{11} + p_{12}} = \frac{p_{11}}{(p_{11} + p_{12})^2} \quad (3)$$

تعریف می‌شود. برای ملاحظه یک مجموعه داده واقعی و کاربرد نسبت مخاطره به اگرستی (۲۰۰۲) مراجعه شود. حال فرض کنید محقق بخواهد باور پیشین خود در خصوص تحلیل بیزی این جدول ناقص را با فرض یک چگالی پیشین برای نسبت مخاطره ϕ وارد تحلیل بیزی نماید. چگونه این کار عملی خواهد بود؟ آیا چگالی پیشین معرفی شده توسط لرد (۱۹۷۸) که بر مبنای مدل لگ-خطی است از چنان قابلیتی برخوردار است که بتواند جواب مناسب را برای این سؤال فراهم آورد؟

تمایز آشکار رهیافت ارائه شده در این مقاله در مقایسه با متون موجود در این زمینه آن است که به جای در نظر گرفتن چگالی برای احتمال خانه‌ها یا پارامترهای مدل لگ-خطی، چگالی پیشین برای نسبت بخت‌ها یا در حالت کلی برای نسبت بخت‌های تعمیم یافته در نظر گرفته می‌شود. به نظر نویسنده و تجربه شخصی او، اغلب مدل‌بندی باور محقق برای نسبت بخت‌ها یا نسبت بخت‌های تعمیم یافته به واقعیت نزدیکتر است تا مدل‌بندی پارامترهای مدل لگ-خطی یا حتی احتمال خانه‌های جدول. نقطه قوت دیگر این رهیافت آن است که محقق را قادر می‌سازد تا بتواند باور ذهنی خود را، حتی زمانی که این باور به بخشی از خانه‌های جدول دلالت دارد، در قالب یک چگالی پیشین به مساله تحمیل نماید.

^۱ Risk ratio

در بخش‌های بعدی به طور مفصل در خصوص روابط (۲) و (۳) بحث خواهد شد. با این وجود در بخش ۲ مفاهیم و برخی تعاریف ارائه خواهند شد. بخش ۳ به معرفی چگالی‌های پیشین با ساختار معین اختصاص دارد. بخش ۴ به تعیین چگالی‌های پیشین مورد نظر این مقاله برای جدول‌های ناقص و کامل 2×2 اختصاص یافته است. در نهایت در بخش ۵ دو مجموعه داده‌های واقعی تحلیل می‌شوند.

۲ مفاهیم اولیه و تعاریف

فرض کنید جدول پیشاندی تحت مطالعه از k متغیر رسته‌ای X_1, \dots, X_k به ترتیب با I_1, \dots, I_k سطح باشند که در آن $\chi_1, \dots, \chi_k, I_i \geq 2; i = 1, \dots, k$. فضای $\chi_1 \times \dots \times \chi_k$ نمونه‌ای متناظر با این متغیرها تعریف می‌شوند. نقطه $(x_1, \dots, x_k) \in \chi_1 \times \dots \times \chi_k$ یک خانه جدول نامیده می‌شود هرگاه دارای احتمال غیر صفر برای داشتن بسامد مشاهده شده بوده و در غیر این صورت خانه‌ای با صفر ساختاری در نظر گرفته می‌شوند. کلیه خانه‌های جدول پیشاندی $I = I_1 \times I_2 \times \dots \times I_k$ است که در ادامه از آن برای اشاره به تعداد خانه‌های با ساختار غیر صفر استفاده می‌شود.

فرض کنید خانه‌های جدول به صورت فرهنگ لغتنامه‌ای مرتب شده باشند. در این صورت از زیرنویس τ برای اشاره به خانه τ ام و از p_i برای اشاره به احتمال این خانه اشاره خواهد شد. در این حالت فرض می‌شود بسامد مشاهده شده جدول از توزیع چند جمله‌ای پیروی می‌کنند. با تبدیل p_i به c_i ، بسامد مورد انتظار خانه τ ام در توزیع پواسون، نتایج برای توزیع حاصل ضرب پواسون نیز برقرار خواهد بود. با توجه به فرضیات فوق n را به عنوان بسامد مشاهده شده خانه τ ام تعریف کرده از $n = \sum n_i$ به عنوان جمع کل بسامدهای مشاهده شده یاد خواهد شد.

تعریف ۱ (کلیسموا و همسکاران، ۲۰۱۲): ترکیب خطی $\sum_{i=1}^I c_i \log p_i$ ، لگاریتم نسبت بخت‌های تعمیم‌یافته تعریف می‌کنیم که در آن ضرایب c_i مقادیر حقیقی مقدار معلوم‌اند که همگی همزمان صفر نیستند. اگر شرط $\sum_{i=1}^I c_i = 0$ نیز برقرار باشد از آن به عنوان لگاریتم نسبت بخت‌های تعمیم‌یافته همگن یا مقابله

۱۹۰ چگالی‌های پیشین با ساختار معین

تعمیم یافته یاد می‌شود.

پر واضح است که فضای تولید شده از لگاریتم نسبت بخت‌های تعمیم یافته نسبت به عمل جمع و تفریق بسته است. به این معنی که جمع دو لگاریتم نسبت بخت تعمیم یافته و همچنین تفریق آنها یک لگاریتم نسبت بخت‌های تعمیم یافته است. این خاصیت برای کلیت بخشی به رهیافت ما در تعیین چگالی‌های پیشین با ساختار معین مهم و حیاتی است. دو مثال زیر به درک بیشتر موضوع کمک می‌کند.

مثال ۱ : یک جدول 3×3 کامل را در نظر بگیرید. از رابطه (۱) دو نسبت بخت‌های (تعمیم یافته)

$$\theta_{11} = \frac{p_{11}p_{22}}{p_{21}p_{12}}, \quad \theta_{22} = \frac{p_{22}p_{33}}{p_{32}p_{23}}$$

را در نظر بگیرید. با لگاریتم گیری داریم:

$$\log \theta_{11} = \log p_{11} + \log p_{22} - \log p_{21} - \log p_{12},$$

$$\log \theta_{22} = \log p_{22} + \log p_{33} - \log p_{32} - \log p_{23}.$$

مالحظه می‌شود که هر دو در تعریف لگاریتم نسبت بخت‌های تعمیم یافته صدق می‌کنند. حال فرض کنید $\theta_{11} = \gamma\theta_{22}$. در این صورت با لگاریتم گیری و قرار دادن لگاریتم احتمال‌ها در یک طرف، یک لگاریتم نسبت بخت‌های تعمیم یافته همگن جدید حاصل می‌شود که برابر ثابت $\log \gamma$ است.

مثال ۲ : جدول پیشایندی ناقص 2×2 توضیح داده شده در بخش قبل را مجدداً در نظر بگیرید. فرض کنید:

$$\rho = \frac{p_{11}p_{22}}{p_{12}} = \frac{\phi(1-a)^2}{(1-a\phi)^2} = \gamma \quad (4)$$

که در آن $p_{11} + p_{12} = a$. با لگاریتم گیری مجدداً یک لگاریتم نسبت بخت‌های تعمیم یافته جدیدی حاصل می‌شود که برابر ثابت $\log \gamma$ است. یعنی:

$$\log p_{11} - 2 \log p_{12} + 2 \log p_{22} = \log \gamma. \quad (5)$$

دلیل استفاده از رابطه (۴) برای جدول ناقص قضیه زیر است که به خاطر سرراست بودن برهان آن تنها به ارائه صورت آن بسنده می‌شود.

قضیه ۱ : کمیت ρ در (۴) تابعی صعودی از نسبت مخاطره ϕ در (۳) بوده و برای آن شرایط زیر برقرار است:

$$\text{الف) } \rho = 1 \Leftrightarrow \phi = 1$$

$$\text{ب) } 0 < \phi < \frac{1}{p_{11} + p_{12}}$$

بنابراین به نظر می‌رسد فرضیاتی که بر اساس لگاریتم نسبت بخت‌های تعمیم یافته نظیر $\delta = \sum_{i=1}^I c_i \log p_i$ در نظر گرفته می‌شوند، تا حد قابل قبولی بتوانند هر باور محقق در خصوص ساختار چگالی پیشین را مدل‌بندی کنند.

در بعضی موارد طبیعی خواهد بود که باور ذهنی محقق به دو یا چند لگاریتم نسبت بخت‌های تعمیم یافته منجر شود. این حالت وقتی اتفاق می‌افتد که اعتقاد پیشین محقق به طور مجزا به دسته‌های گوناگون (r تا I) از خانه‌های جدول مربوط باشد. این فرضیات، با قبول استقلال خطی بردار ضرایب $c_j = (c_{j1}, c_{j2}, \dots, c_{jI})'$ که شرط اصلی برای داشتن r فرض مجزا است، به صورت زیر داده می‌شود:

$$\begin{cases} \sum_{i=1}^I c_{1i} \log p_i = \delta_1 \\ \vdots \\ \sum_{i=1}^I c_{ri} \log p_i = \delta_r. \end{cases}$$

نمایش ماتریسی این دستگاه معادلات به صورت

$$C \log p = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1I} \\ c_{21} & c_{22} & \dots & c_{2I} \\ \vdots & \vdots & \ddots & \vdots \\ c_{r1} & c_{r2} & \dots & c_{rI} \end{pmatrix} \begin{pmatrix} \log p_1 \\ \log p_2 \\ \vdots \\ \log p_I \end{pmatrix} = \begin{pmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_r \end{pmatrix} = \delta \quad (6)$$

است.

۳ چگالی پیشین با ساختار معین

در بخش قبل باور ذهنی محقق در خصوص لگاریتم نسبت بخت‌ها در قالب مدل کلی (۶) ارائه شد. در این مدل سطرهای ماتریس C بردارهایی در R^I می‌باشند. به دلیل آن که رابطه (۶) از r فرض مستقل خطی و مجزا تشکیل شده است می‌توان نتیجه گرفت که ماتریس C یک ماتریس پر رتبه سط्रی است. مانند آن‌چه در مدل‌های نسبتی^۲ متداول است (کلیمووا و همکاران، ۲۰۱۲) همواره می‌توان ماتریس پر رتبه سطري دیگری نظیر A از رتبه $I \times (I - r)$ یافت به‌طوری که شرایط زیر برقرار باشند:

$$\cdot AC' = \circ \quad \cdot CA' = \circ$$

$$2. [C' A'] \text{ یک ماتریس پر رتبه و از مرتبه } I \times I \text{ است.}$$

در اینجا لازم به توضیح است که تعیین ماتریس‌های A و C از مهمترین چالش‌های این مساله بوده و در کلیمووا و همکاران (۲۰۱۲) تاحدی در خصوص چگونگی تعیین آن‌ها بحث شده است. با توجه به شرط تعامد دو ماتریس A و C ، افزایش

$$\log p = A'\beta + C'\theta \quad (7)$$

برای بردار I بعدی $\log p$ برقرار است که در آن بردار β به عنوان پارامتر کم اهمیت تر دارای نقش معینی نبوده و به‌همین دلیل اغلب از چگالی پیشین ناسره برای آن استفاده می‌شود. این نکته که با بردار β به عنوان پارامتر مزاحم رفتار شود از محاسن مدل (۷) است. زیرا در غیر این صورت همواره می‌توان ماتریس دیگری نظیر A^* یافت که $A^*/\beta^* = A'/\beta$ و در نتیجه تحلیل بیزی مخدوش می‌شود.
با ضرب طرفین رابطه (۷) در C و استفاده از خاصیت تعامد داریم:

$$C \log p = CC'\theta. \quad (8)$$

در این صورت با مقایسه روابط (۶) و (۸) خواهیم داشت:

$$\theta = (CC')^{-1}\delta. \quad (9)$$

رابطه (۹) به این معنی است که هر باور پیشین روی نسبت بخت‌های تعمیم‌یافته و به طور معادل روی δ از طریق بردار پارامتر θ در مدل (۷) قابل تبیین است. همان‌طور که از مباحث بالا بر می‌آید به نظر توانسته باشیم تا هر باور ذهنی در خصوص نسبت بخت‌های تعمیم‌یافته را از طریق مدل (۷) به تحلیل بیزی القاء کنیم. حال بدون کاستن از کلیت مساله، در ادامه از چگالی پیشین نرمال چندمتغیره برای پارامترهای مورد نظر و پارامترهای مزاحم استفاده خواهد شد. برای این منظور داریم:

- برای ماتریس معلوم W از بعد $r \times r$ ، فرض می‌شود δ دارای توزیع نرمال چند متغیره با میانگین \circ و ماتریس کوواریانس W^{-1} باشد که در آن ثابت q به عنوان ابر پارامتر به کار می‌رود. در اینجا فرض می‌شود ماتریس همیشه مشبّت W حاوی تمام باور پیشین محقق در خصوص چگونگی ساختار ارتباطی نسبت بخت‌های تعمیم‌یافته است.
- فرض می‌شود بردار پارامترهای مزاحم β دارای توزیع نرمال با میانگین صفر و ماتریس واریانس کواریانس J_{I-r}^{-1} است، که در آن η عددی بسیار کوچک نظیر 0.01 ، 0.02 در نظر گرفته می‌شود. در اینجا ماتریس J_{I-r} ، ماتریس واحد و از مرتبه $I - r$ است. یعنی:

$$\beta \sim N(\circ, \eta^{-1} J_{I-r}). \quad (10)$$

با توجه به فرضیات فوق و با استفاده از رابطه (۹) توزیع پیشین θ به صورت

$$\theta \sim N((CC')^{-1}\delta_0, q^{-1}(CC')^{-1}W(CC')^{-1}) \quad (11)$$

خواهد بود. روابط (۷)، (۱۰) و (۱۱) چگالی پیشین مورد نظر برای احتمال خانه‌های جدول و در حالت کلی برای بسامدهای مورد انتظار را فراهم می‌آورند. یکی از مشکلات استفاده از چگالی پیشین فوق وجود ابر پارامتر q است. برای مشارکت دادن اثر آن در تحلیل بیزی ممکن است علاقمند به استفاده از یک چگالی گاما، به عنوان چگالی پیشین سلسله مراتبی برای آن باشیم. با این وجود یک برآورده نسبتاً مناسب از آن به شرح زیر داده می‌شود:

۱۹۴ چگالی‌های پیشین با ساختار معین

فرض کنید S ماتریس واریانس برآوردهای $\log p$ باشد که در آن احتمال‌ها با نسبت‌های نمونه‌ای $\frac{n_i}{n}$ جای گزین شده باشند، (اگرستی، ۲۰۰۲). از رابطه (۸) و با استفاده از چگالی پیشین (۱۰) داریم:

$$CSC' \approx q^{-1}W.$$

با ضرب بردار واحد ۱ از دو طرف در رابطه فوق، برآورد زیر برای q حاصل می‌شود:

$$\tilde{q} = \frac{\mathbf{1}'W\mathbf{1}}{\mathbf{1}'CSC'\mathbf{1}}. \quad (12)$$

برای انجام استنباط بیزی لازم است توزیع پسین متناظر را به دست آورد. چون توزیع حاصل دارای نمایش بسته نیست با استفاده از روش متروپولیس-هاستینگ می‌توان از آن نمونه‌گیری کرد. بر اساس این روش با در نظر گرفتن تابع احتمال $f(n|p)$ و چگالی پیشین (p, π) ، احتمال پذیرش نمونه جدید i' برای مؤلفه n ام بردار p برابر است با:

$$\alpha(p_i, p'_i) = \min\left\{1, \frac{f(n|(p_1, p_2, \dots, p'_i, \dots, p_I))\pi(p_1, p_2, \dots, p'_i, \dots, p_I)}{f(n|(p_1, p_2, \dots, p_i, \dots, p_I))\pi(p_1, p_2, \dots, p_i, \dots, p_I)}\right\}.$$

۴ چگالی‌های پیشین برای جداول 2×2

در این بخش چگالی پیشین با ساختار معین برای جداول 2×2 معرفی خواهد شد.

۱.۴ چگالی‌های پیشین برای جداول ناقص

در این زیربخش سه حالت مختلف، که ممکن است محقق را به سوی تعیین پیشین خاص رهنمون سازد، در نظر گرفته می‌شود.

• حالت اول جدول ناقص 2×2

را در نظر بگیرید. هر باور پیشین روی نسبت مخاطره یا روی لگاریتم نسبت بخت‌های تعمیم یافته، ارائه شده در (۵)، منجر به ماتریس‌های A و C و

p_{11}	p_{12}
-	p_{22}

چگالی پیشین به صورت زیر می شود:

$$A = \begin{pmatrix} \circ & 1 & 1 \\ 2 & 1 & \circ \end{pmatrix}, \quad C = \begin{pmatrix} 1, & -2, & 2 \end{pmatrix},$$

$$\theta \sim N\left(\frac{1}{9}\delta_0, \frac{q^{-1}}{81}W\right).$$

در نتیجه چگالی پیشین نرمال برای بردار لگاریتم احتمال خانه‌های، $(\log p_{11}, \log p_{12}, \log p_{22})$ با توجه به روابط (۷)، (۱۰) و (۱۱) حاصل می‌گردد.

• **حالت دوم** در این حالت k جدول ناقص 2×2 به صورت را در نظر

p_{111}	p_{121}
-	p_{221}
p_{112}	p_{122}
-	p_{222}

بگیرید. طبق رابطه (۳) فرض کنید $\phi_k, \phi_1, \dots, \phi_n$ به ترتیب نسبت مخاطره آنها و ρ_1, \dots, ρ_k تبدیل‌های یک به یک از آن‌هاست، (نتیجه قضیه ۱). در این حالت فرض می‌شود ϕ_i ها و در نتیجه ρ_i ها نمونه‌هایی تصادفی از یک توزیع معین، مثلاً نرمال باشند. بنابر این ماتریس‌های A و چگالی پیشین عبارتند از:

$$A = \begin{pmatrix} \circ & 1 & 1 \\ 2 & 1 & \circ \end{pmatrix} \otimes J_k, \quad C = \begin{pmatrix} 1 & -2 & 2 \end{pmatrix} \otimes J_k,$$

$$\theta \sim N_k\left(\frac{1}{9}\delta_0, \frac{q^{-1}}{81}W\right)$$

که در آن \otimes نماد ضرب کرونکر است. توزیع پیشین معین برای بردار لگاریتم احتمال خانه‌های جدول

$$(\log p_{111}, \log p_{121}, \log p_{221}, \dots, \log p_{11k}, \log p_{12k}, \log p_{22k})$$

۱۹۶ چگالی‌های پیشین با ساختار معین

با توجه به روابط (۷)، (۱۰) و (۱۱) حاصل می‌شوند.

- حالت سوم با در نظر گرفتن فرضیات حالت قبل، فرض کنید باور پیشین ما در این حالت فرض برابری نسبت مخاطره‌ها و در نتیجه تبدیلات آن‌هاست. یعنی:

$$\rho_1 = \dots = \rho_k$$

که آن را می‌توان به صورت

$$\rho_1 = \rho_2, \dots, \rho_1 = \rho_k \quad (۱۳)$$

بازنویسی کرد. در این صورت

$$C = \begin{pmatrix} 1 & -2 & 2 \end{pmatrix} \otimes \begin{pmatrix} 1_{k-1}, & -J_{(k-1) \times (k-1)} \end{pmatrix}$$

و

$$A = \left(\begin{array}{c|c} \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 2 & 1 & 0 \end{pmatrix}, & \circ^{3 \times 3(k-2)} \\ \left(\begin{array}{ccc} 2 & 1 & 0 \end{array} \right) \otimes 1_{2(k-1)}, & \left(\begin{array}{ccc} 2 & 1 & 0 \\ 0 & 1 & 1 \end{array} \right) \otimes J_{(k-1) \times (k-1)} \end{array} \right),$$

علاوه بر این که بردارهای δ و در نتیجه θ نیز مساوی صفر هستند که در نتیجه بخش دوم سمت راست مدل (۷) حذف می‌شود. رابطه (۱۳) را می‌توان به راحتی به حالت کلی تر

$$\rho_1 = \gamma_2 \rho_2, \dots, \rho_1 = \gamma_k \rho_k,$$

نیز توسعه داد، که در آن $\log \gamma_i$ مستقل و دارای توزیع نرمال‌اند.

۲.۴ چگالی‌های پیشین برای جداول کامل

در این زیر بخش نیز سه حالت مختلف در نظر می‌گرفته می‌شود:

p_{11}	p_{12}
p_{21}	p_{22}

• حالت اول برای جدول کامل 2×2

هر باور پیشین روی نسبت بخت‌ها یا روی لگاریتم آن‌ها منجر به ماتریس‌های

$$A = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{pmatrix}, \quad C = \begin{pmatrix} 1 & -1 & -1 & 1 \end{pmatrix}$$

می‌شود. با فرض این که (δ_0, q^{-1}) چگالی پیشین نرمال برای θ عبارت است از:

$$\theta \sim N\left(\frac{\delta_0}{4}, \frac{q^{-1}}{16}W\right).$$

در این صورت چگالی پیشین بردار $(\log p_{11}, \log p_{12}, \log p_{21}, \log p_{22})$, با توجه به روابط (۷)، (۱۰) و (۱۱) به دست می‌آید.

• حالت دوم برای k جدول کامل 2×2

p_{111}	p_{121}	p_{112}	p_{122}	\dots	p_{11k}	p_{12k}
p_{211}	p_{221}	p_{212}	p_{222}		p_{21k}	p_{22k}

فرض کنید $\theta_k^*, \theta_1^*, \dots, \theta_i^*$ به ترتیب نسبت بخت‌های در k جدول باشند. با آن که θ_i^* ها نمونه‌هایی تصادفی از یک توزیع معین، مثلاً نرمال باشند، داریم:

$$\delta_i \sim N(\delta_0, q^{-1}), \quad i = 1, \dots, k,$$

$$A = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{pmatrix} \otimes J_k,$$

$$C = \begin{pmatrix} 1 & -1 & -1 & 1 \end{pmatrix} \otimes J_k,$$

$$\theta \sim N_k\left(\frac{\delta_0}{4}, \frac{q^{-1}}{16}W\right).$$

در نهایت چگالی پیشین بردار

$$(\log p_{111}, \log p_{121}, \log p_{211}, \log p_{221}, \dots, \log p_{11k}, \log p_{12k}, \log p_{21k}, \log p_{22k})$$

با توجه به روابط (۷)، (۱۰) و (۱۱) به دست می‌آید.

- **حالت سوم** با در نظر گرفتن فرضیات حالت قبل، فرض کنید باور پیشین برابری نسبت بخت‌ها باشد. یعنی:

$$\theta_1^* = \dots = \theta_k^*$$

که آن را می‌توان به صورت

$$\theta_1^* = \theta_2^*, \dots, \theta_1^* = \theta_k^* \quad (۱۴)$$

بازنویسی کرد. در این صورت

$$C = \begin{pmatrix} 1 & -1 & -1 & 1 \end{pmatrix} \otimes \begin{pmatrix} 1_{k-1}, & -J_{(k-1) \times (k-1)} \end{pmatrix}$$

و

$$A = \left(\begin{array}{c|c} \begin{pmatrix} 1, 1, 0, 0, 1, 1, 0, 0 \\ 1, 1, 0, 0, 1, 0, 1, 0 \\ 1, 0, 1, 0, 1, 1, 0, 0 \\ 1, 0, 1, 0, 1, 0, 1, 0 \end{pmatrix}, & \circ^{4 \times 4(k-2)} \\ \hline \begin{pmatrix} 1, 1, 1, 1 \\ 1, 1, 0, 0 \\ 1, 0, 1, 0 \end{pmatrix} \otimes 1_{2(k-1)}, & \begin{pmatrix} 1, 1, 1, 1 \\ 1, 1, 0, 0 \\ 1, 0, 1, 0 \end{pmatrix} \otimes J_{(k-1) \times (k-1)} \end{array} \right).$$

علاوه بر این که بردارهای δ و در نتیجه θ نیز مساوی صفر است بخشن دوم سمت راست مدل (۷) نیز حذف می‌شود. رابطه (۱۴) را می‌توان به راحتی به حالت کلی تر

$$\theta_1^* = \gamma_2 \theta_2^*, \dots, \theta_1^* = \gamma_k \theta_k^*$$

نیز توسعه داد، که در آن $\log \gamma_i$ مستقل و دارای توزیع نرمال‌اند.

۵ کاربرد

در این بخش به تحلیل بیزی دو مجموعه داده واقعی پرداخته می‌شود. این دو مثال به روشنی قابلیت بالای رهیافت معرفی شده در این مقاله را نشان می‌دهند. برای تحلیل این داده‌ها از نرم افزار WINBUGS14 استفاده شده است.

مثال ۳ : مجموعه داده‌ها از یک مطالعه درمانی دو مرحله‌ای تشکیل شده است. در اینجا با سه جدول ناقص 2×2 روبرو هستیم که تحت سه شرط مختلف شدت بیماری حاصل شده‌اند. به عبارتی بیماران در سه گروه مختلف با شدت بیماری: «خفیف»، «متوسط» و «شدید» طبقه‌بندی شده، سپس در دو مرحله مورد مداوا قرار گرفته‌اند. جدول ۱ داده‌ها را نشان می‌دهد (تنگ و جیانگ، ۲۰۱۱).

تابع احتمال برای این داده‌ها حاصل ضرب چند جمله‌ای با ضابطه

$$f(n|p) = \prod_{i=1}^k n_i! \prod_{j=1}^3 \frac{p_{ij}^{n_{ij}}}{n_{ij}!},$$

است. این داده‌ها توسط آمارشناسان مذکور مورد تحلیل بسامدی قرار گرفت تا درستی فرض برابری نسبت مخاطره‌ها، برای بیماران با شدت مختلف بیماری، را بیازمایند. در اینجا از رهیافت ارائه شده در این مقاله برای تحلیل بیزی داده‌ها استفاده می‌شود. فرض پیشین در اینجا فرض «برابری نسبت مخاطره‌ها برای انواع شدت بیماری» بود. در این تحلیل $3 = k$ و بر اساس رابطه (۱۲)، $q = 0 / 123 = 0.0100$ به دست آمد. از توزیع‌های پیشین $(0.0, 100)$ برای پارامترهای مزاحم، از توزیع $\Gamma(1, 0.01)$ برای ابرپارامتر q و از توزیع $(0, q^{-1})$ برای پارامترهای اصلی θ استفاده شده است. احتمال خانه‌های سه جدول به صورت فرهنگ لغت نامه‌ای مرتب و در یک بردار از بعد ۹ قرار گرفتند. با تعیین ماتریس‌های A و C و نیز بر اساس نتایج حاصل از حالت دوم در بخش ۴-۱، قبول فرض پیشین‌های ناسره برای پارامترهای مزاحم و استفاده از تکنیک نمونه‌گیری MCMC با 10000 بار تکرار، نتایج به شرح جدول‌های ۲ و ۳ به دست آمد. جدول ۲ حاوی برآوردهای بیزی پارامترهای θ_1 ، θ_2 و θ_3 است. برای ارزیابی این که آیا چگالی پسیون مؤید باور پیشین مبنی بر برابری نسبت مخاطره‌ها برای انواع شدت بیماری است، از ملاک

احتمال پسین

$$P(|\theta_i - \theta_j| > \epsilon | D); \quad i \neq j,$$

برای بررسی اختلاف θ ها، مشروط بر مشاهدات موجود (D) استفاده می‌شود. نتایج حاصل از تعیین این احتمال‌ها بر اساس نمونه‌گیری از توزیع پسین برای مقادیر مختلف ϵ در جدول ۳ آمده است. همان طور که ملاحظه می‌شود احتمال پسین اختلاف‌های دو به دوی کمیت‌های θ_1 , θ_2 و θ_3 مبین وجود عدم اختلاف معنی دار بین θ_1 و θ_2 است. به عبارتی بر اساس این احتمال‌ها پسین می‌توان نسبت مخاطره برای شدت بیماری «خفیف» و «متوسط» را تقریباً برابر فرض نمود. این در حالی است که اندازه نسبت مخاطره برای شدت بیماری «شدید» متفاوت از بقیه و کمتر از آن‌هاست. این نتایج با قبول تقریب نرمال برای چگالی پسین(با توجه به حجم نسبتاً بزرگ بسامدی‌های مشاهده شده) و ساختن بازه اطمینان بیزی برای اختلاف θ ها نیز حاصل می‌شود. بازه اطمینان‌های ۹۵ درصد به روش MCMC برای این اختلاف‌ها عبارتند از:

$$\theta_1 - \theta_2 : -0/072 \pm 0/206,$$

$$\theta_1 - \theta_3 : 0/258 \pm 0/253,$$

$$\theta_2 - \theta_3 : 0/332 \pm 0/276.$$

جدول ۱: داده‌های یک مطالعه درمانی دو مرحله‌ای تحت سه شرط شدت بیماری

شدت بیماری						
شدید		متوسط			خفیف	
		مرحله دوم				
مرحله اول	عدم بهبود	عدم بهبود	بهبود	بهبود	عدم بهبود	بهبود
۲۱	۶	۳۷	۱۶	۸۳	۴۶	عدم بهبود
۴۳	-	۹۱	-	۱۷۶	-	بهبود

جدول ۲: خروجی الگوریتم MCMC با ۱۰۰۰ بار تکرار

پارامتر	میانگین انحراف معیار	میانه	
۱/۵۹۳	۰/۰۶۶	۱/۵۹۲	θ_1
۱/۶۶۵	۰/۰۸۰	۱/۶۶۶	θ_2
۱/۳۳۵	۰/۱۱۶	۱/۳۳۴	θ_3

جدول ۳: برآورد احتمال‌های پسین اختلاف θ ‌ها

۰/۵۰	۰/۴۰	۰/۳۰	۰/۲۰	۰/۱۰	ϵ
۰/۰۰۰۲	۰/۰۰۱۱	۰/۰۱۶۵	۰/۱۳۰۰	۰/۴۴۷۸	$\hat{P}(\theta_1 - \theta_2 > \epsilon D)$
۰/۰۳۴۵	۰/۱۴۱۷	۰/۳۷۲۷	۰/۶۶۲۳	۰/۸۸۶۲	$\hat{P}(\theta_1 - \theta_2 > \epsilon D)$
۰/۱۱۵۸	۰/۳۱۱۷	۰/۵۸۶۲	۰/۸۳۶۸	۰/۹۵۲۳	$\hat{P}(\theta_2 - \theta_3 > \epsilon D)$

مثال ۴: جدول ناقص ۴ وضعیت تغییر شغلی پاسخ دهنده‌گان نسبت به پدرشان را نشان می‌دهد (کلیممو و همکاران، ۲۰۱۲). در این تحقیق از ۳۷۶۷۷ نفر در خصوص شغل شخص و نیز شغل پدر او سوال شده است. هدف اصلی مطالعه تغییر شغلی فرزندان نسبت به پدران است. جدول ۴ داده‌های این مثال را نشان می‌دهد. در اینجا به ترتیب سه رده شغلی «کارمندی»، «کارگری» و «کشاورزی» در نظر گرفته شده است. فرض کنید باور پیشین عبارت است از فرض‌های

جدول ۴: جدول داده‌های تغییر شغلی پاسخ دهنده‌گان نسبت به پدران

وضعیت شغل فرزند		
شغل پدر	پسرفت	بدون تغییر
-	۶۳۱۳	۲۷۷۶
۶۳۲۱	۱۰۸۸۳	۲۹۴
۸۶۱۹	۲۴۷۱	-
		کارمندی
		کارگری
		کشاورزی

و $\gamma = ۳\gamma = \frac{p_{۱۱}p_{۲۲}}{p_{۱۲}p_{۲۱}}$. با لگاریتم‌گیری از این روابط و قبول شرط $\delta = \log \gamma$ فرض‌های معادل عبارتند از

$$\begin{cases} \log p_{۱۱} + \log p_{۲۲} - \log p_{۱۲} - \log p_{۲۱} = \delta, \\ \log p_{۲۱} + \log p_{۱۲} - \log p_{۲۲} - \log p_{۱۱} = \delta + \log ۳. \end{cases}$$

این فرضیات به معنی آن است که صرف نظر از مشاهده، معتقد دیم نسبت بخت پسرفت شغلی برای فرزندان کارمندان نسبت به فرزندان کارگران یک سوم همین نسبت در رده شغلی کارگران به کشاورزان است. با فرض توزیع چند جمله‌ای برای بسامد های مشاهده شده و بردار احتمال های ماتریس های متعامد متناظر A و C به صورت زیر به دست می‌آیند:

$$A = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 1 \end{pmatrix},$$

و

$$C = \begin{pmatrix} 1 & -1 & 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & -1 & 1 \end{pmatrix}.$$

برای تحلیل بیزی به روش ارائه شده در این مقاله، از توزیع های پیشین $n(0, 100)$ برای پارامتر های مراحم، از توزیع $(1, 0/0)\Gamma$ برای ابر پارامتر q^{-1} و از توزیع $(1, q^{-1})n$ برای پارامتر های اصلی θ استفاده شده است. چون در تحلیل این داده ها فرض های پیشین با الهام از مشاهدات به دست آمدند، این انتظار وجود دارد که برآوردهای بیزی مؤید باور پیشین باشند. برای ارزیابی این که چگالی پیشین تا چه حد باور پیشین را تأیید می کند، از برآورد پسین فاصله کمیت

$$\kappa = \frac{\frac{p_{12}p_{23}}{p_{12}p_{22}}}{\frac{p_{21}p_{23}}{p_{21}p_{22}}} = \frac{p_{12}p_{23}p_{21}}{p_{21}p_{22}p_{13}},$$

و مقدار آن تحت فرض پیشین، $\frac{1}{3}$ ، استفاده شده است. این معیار عبارت است از:

$$\hat{P}(|\kappa - \frac{1}{3}| > \epsilon | D).$$

با استفاده از ۱۰۰۰۰ بار تکرار الگوریتم MCMC نتایج در جدول ۵ آمده است. چون در این جا احتمالات پسین پیشامد «وجود اختلاف» برای مقادیر نسبتاً کوچک ϵ در مقایسه با متمم آن کوچکتر است پس با مبنای قرار دادن ملاک احتمال پسین می توان نتیجه گرفت که احتمال پسین تمایل به تأیید درستی باور پیشین دارد.

جدول ۵: برآورد احتمال‌های پسین اختلاف κ و $\frac{1}{\epsilon}$			
ϵ	$0/01$	$0/05$	$0/50$
$\hat{P}(\kappa - \frac{1}{\epsilon} > \epsilon D)$	$0/9437$	$0/3977$	$0/0202$

بحث و نتیجه‌گیری

در این مقاله عدم کارآمدی رهیافت‌های موجود برای القاء هر باور پیشین به احتمال خانه‌های جدول پیشایندی تشریح شد. سپس تأکید شد که همواره هر باور پیشین در خصوص احتمال خانه‌های جدول را می‌توان در قالب یک یا چند نسبت بخت تعمیم یافته وارد مدل نمود. با این وجود نکته اساسی این است که چگونه محقق می‌تواند این باور پیشین را از نسبت بخت تعمیم یافته به احتمال خانه‌های جدول انتقال دهد؟ با معرفی دو ماتریس متعامد A و C و با الهام از مدل‌های نسبتی، برای تجزیه بردار احتمال خانه‌ها به دو مؤلفه متعامد، توانستیم برای این سؤال مهم پاسخ مناسب فراهم آوریم. ویژگی برجسته رهیافت ارائه شده در این مقاله این است که حتی شامل مواردی است که اعتقاد محقق نه به همه خانه‌های جدول بلکه به بخشی از آن‌ها، که در نسبت بختهای تعمیم یافته شرکت دارند، نیز تعلق دارد. علاوه بر این چگالی‌های پیشین معرفی شده برای نسبت بختهای تعمیم یافته همگن، از خاصیت پایایی جالبی برخوردار است. به این معنی که برای هر ثابت حقیقی مقدار λ از رابطه (۸) و با استفاده از خواص لگاریتم داریم:

$$C \log \lambda p = CC' \theta \Leftrightarrow C \log p = CC' \theta.$$

این رابطه نشان می‌دهد که با ضرب احتمال خانه‌های جدول در هر عدد ثابت و دلخواه، ساختار چگالی پیشین تغییر نمی‌کند. این نتیجه همواره برقرار بوده و متأثر از پارامترهای مزاحم β ، در رابطه (۷)، نمی‌باشد. این مهم از شرط تعامد ماتریس‌های A و C ناشی می‌شود. زیرا علی‌رغم وجود و نقش این پارامترهای در (۷)، در رابطه (۸) نقشی نداشته و در نتیجه خاصیت پایایی مطلوب همواره برقرار خواهد بود.

تقدیر و تشکر

نویسنده از داوران محترمی که با ارائه نقطه نظرات سازنده خود باعث بهبود کیفیت مقاله شده‌اند، قدردانی می‌نماید.

مراجع

- Agresti, A. (2002), *Categorical Data Analysis*, John Wiley, New York.
- Agresti, A. and Hitchcock, D. B. (2005), Bayesian Inference for Contingency Data Analysis: A Survey, *Statistical Methods and Applications*, **14**, 297-330.
- Aitchinson, J. (1985), Practical Bayesian Problems in Simplex Sample Spaces, *Bayesian Statistics*, **2**, 15-31.
- Bishop, Y. M. M., Fienberg, S. E. and Holland, P. W. (1975), *Discrete Multivariate Analysis: Theory and Practice*, MIT Press.
- Forster, J. J. and Skene, A. M. (1994), Calculation of Marginal Densities for Parameters of Multinomial Distribution, *Statistics and Computing*, **4**, 279-286.
- Good, I. (1965), *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*, MIT Press, Cambridge.
- Goodman, L. O. (1972), Some Multiplicative Models for the Analysis of Cross-Classified Data in: *Proceedings of the Sixth Berkely Symposium of Mathematical Statistics and Probability*.
- Goutis, C. (1993), Bayesian Estimation Methods for Contingency Tables, *Journal of the Italian Statistical Society*, **2**, 35-54.

سید کامران قریشی ۲۰۵

Klimova, A., Rudas, T. and Dobra, A. (2012), Relational Models for Contingency Tables, *Journal of Multivariate Analysis*, **104**, 159-173.

Larid, N. M. (1978), Empirical Bayes Methods for Two-Way Contingency Tables, *Biometrika*, **65**, 581-590.

Leonard, T. (1973), A Bayesian Method for Histograms, *Biometrika*, **60**, 297-308.

Leonard, T. (1975), Bayesian Estimation Methods for Two-way Contingency Tables, *Journal of the Royal Statistical Society, B*, **37**, 23-37.

Lindly, D. V. (1964), The Bayesian Analysis of Contingency Tables, *The Annals of Mathematical Statistics*, **35**, 1622-1643.

Tang, N. S. and Jiang, S. P. (2011), Testing Equality of Risk Ratios in Multiple 2×2 Tables with Structural Zero, *Computational Statistics and Data Analysis*, **55**, 1273-1284.