

مجله علوم آماری، بهار و تابستان ۱۳۸۹

جلد ۴، شماره ۱، ص ۳۵-۵۸

بهبود الگوریتم ساختاری مونت کارلوی زنجیر مارکوف در مدل‌های چندسطحی با متغیر پاسخ نرمال

عاطفه فرخی، موسی گل‌علی‌زاده

گروه آمار، دانشگاه تربیت مدرس

تاریخ دریافت: ۱۳۸۸/۷/۱۶ تاریخ آخرین بازنگری: ۱۳۸۹/۲/۲۳

چکیده: مدل‌های چندسطحی در علوم کاربردی شامل علوم اجتماعی، جامعه‌شناسی، پزشکی و اقتصاد برای تحلیل داده‌های همبسته مورد استفاده قرار می‌گیرند. روش‌های متفاوتی برای برآورد این مدل‌ها با متغیر پاسخ نرمال وجود دارند. در این مقاله برای به کارگیری روش بیزی از تعمیم الگوریتم مونت کارلوی زنجیر مارکوف استفاده می‌شود که قالبی ساده داشته و باعث حذف همبستگی بین نمونه‌های شبیه‌سازی شده برای پارامترهای ثابت و خطاهای منتسب به گروه‌ها می‌شود. چون بعد ماتریس کواریانس بردار خطای جدید افزایش می‌یابد، برای تسریع همگرایی این روش دو راهکار بر مبنای تجزیه چولسکی ماتریس‌های کواریانس پیشنهاد می‌شود. سپس عملکرد این روش‌ها در مطالعه شبیه‌سازی و مثالی کاربردی مورد ارزیابی قرار می‌گیرد.

واژه‌های کلیدی: داده‌های چندسطحی، مدل‌های عرض از مبدا تصادفی، الگوریتم مونت کارلوی زنجیر مارکوف، تجزیه چولسکی.

آدرس الکترونیک مسئول مقاله: موسی گل‌علی‌زاده، gotalizadeh@modares.ac.ir
کد موضوع‌بندی ریاضی (۲۰۰۰): ۶۲F۱۵ و ۶۲J۹۹

معمولاً بررسی‌های آماری روابط بین متغیرهای تبیینی و وابسته که به صورت ترکیب‌های مختلف از مشاهدات کمی و کیفی باشند در قالب رگرسیون، تحلیل واریانس و تحلیل‌های کواریانس صورت می‌گیرد. یکی از فرض‌های اساسی در کاربرد این گونه مدل‌ها استقلال آماری بین مشاهدات است. گاهی اوقات این فرض برای موضوع مورد مطالعه صادق نبوده، در نتیجه به کارگیری مدل‌های رگرسیونی متداول دارای اشکال می‌باشد. به عنوان نمونه پینهریو و بیتس (۲۰۰۰) نشان دادند، نادیده گرفتن فرض همبستگی بین مشاهدات منجر به کم‌برآوردی خطای برآورد ضرایب رگرسیونی می‌شود. مثال‌های زیادی در حوزه علوم کاربردی شامل علوم اجتماعی، علوم پزشکی، جامعه‌شناسی و غیره وجود دارند که در آن‌ها همبستگی بین مشاهدات کاملاً مشهود است. به عنوان مثال، در بررسی میزان کلسترول خون تعدادی بیمار، می‌توان انتظار داشت شباهتی نسبی بین وضعیت بیماران که تحت درمان یک پزشک هستند وجود داشته باشد (توایسک، ۲۰۰۶). مدل مناسب تحلیل داده‌هایی مانند مثال فوق، مدل چندسطحی است. ویژگی اصلی داده‌های چندسطحی خصوصیت گروه‌بندی آن‌ها است که در مدل‌بندی آماری لحاظ می‌شود. روش‌های استنباط آماری زیادی راجع به پارامترهای مدل‌های چندسطحی معرفی شده است (گلدستاین، ۱۹۹۹). طبق معمول چنین استنباط‌هایی به دو شیوه بسامدی و بیزی صورت می‌گیرد. در این بین روش‌های بیزی شامل موضوعات نظری و روش‌های مبتنی بر شبیه‌سازی کامپیوتری و به‌ویژه روش‌های مونت کارلوی زنجیر مارکوف^۱ (MCMC) هستند. مقایسه این روش‌ها توسط براون و درایپر (۲۰۰۰) صورت گرفت. گلفند و همکاران (۱۹۹۵) نیز کاربرد روش‌های مختلف بیزی را در مدل‌های چندسطحی مورد مطالعه قرار دادند.

برای برآورد کردن پارامترهای مدل چندسطحی با روش‌های MCMC و برطرف کردن مشکل همگرایی نمونه‌های تولید شده سرجنت و همکاران (۲۰۰۰) روشی بنام روش ساختاری مونت کارلوی زنجیر مارکوف^۲ (SMCMC) را پیشنهاد کردند.

^۱ Markov Chain Monte Carlo

^۲ Structured Markov Chain Monte Carlo

عاطفه فرخی، موسی گل‌علی‌زاده ۳۷

این روش کارایی بسیار خوبی در مقایسه با روش‌های بازپارامتری داشته و در عین حال از ساختار ساده‌ای برخوردار است، به طوری که اکثر مدل‌های چندسطحی را می‌توان در آن گنجانند. اما نکته قابل تأمل در روش SMCMC، حجم بالای محاسبات ناشی از بزرگ بودن ماتریس کواریانس مدل تغییر یافته است.

در بخش دوم این مقاله روش SMCMC به اختصار معرفی می‌شود. سپس در بخش سوم ابتدا راهکارهایی برای بهبود روش SMCMC پیشنهاد شده و به منظور تشریح بهتر مطالب دو مدل عرض از مبدا تصادفی در نظر گرفته می‌شود. در بخش چهارم با شبیه‌سازی و یک مثال کاربردی روش‌های موجود و پیشنهادی مورد مقایسه قرار می‌گیرند. بحث و نتیجه‌گیری در بخش پنجم ارائه خواهد شد.

۲ روش ساختاری مونت کارلوی زنجیر مارکوف

در روش‌های MCMC تعیین زمان همگرایی، اخذ نمونه‌های تقریباً مستقل، کاهش برآورد خطای برآوردگرها و خصوصیات کارای دیگر زنجیر از موضوعات بسیار اساسی به‌شمار می‌رود. برای آنکه بتوان نمونه‌هایی تقریباً مستقل از توزیع هدف تولید کرد، لازم است نحوه تولید زنجیر مارکوف، زمان همگرایی آن، کاهش همبستگی نمونه‌های حاصل از زنجیر و معیارهای دیگر مورد توجه قرار گیرد. مقاله‌های متعددی در بررسی این مطالب چاپ شده‌اند. به‌عنوان مثال می‌توان به هیل و اسمیت (۱۹۹۲)، گیلکس و رابرت (۱۹۹۶) اشاره کرد. در بین روش‌های MCMC موضوع بهبود زمان دستیابی به همگرایی از طریق کاهش همبستگی نمونه‌های حاصل از زنجیر و در عین حال رسیدن به برآورد دقیق‌انگیزه‌ای شد تا محققین در این زمینه تحقیقات بیشتری به‌عمل آورند. روش‌های بازپارامتری^۳ یکی از نتایج این فعالیت‌ها بوده است. روش‌های بازپارامتری شامل روش‌های متعامدسازی^۴ (هیل و اسمیت، ۱۹۹۲)، روش‌های مرکزی کردن سلسله مراتبی^۵

^۳ Reparameterization

^۴ Orthogonalization

^۵ Hierarchical centering

۳۸ بهبود الگوریتم SMCMC در مدل‌های چندسطحی

(گلفند و همکاران، ۱۹۹۵) و روش‌های بسط پارامتر^۶ (گلمن و همکاران، ۲۰۰۸) است. پاپاس پیلوپوس و همکاران (۲۰۰۷) نیز یک صورت کلی برای پارامتربندی مدل‌های سلسله مراتبی ارائه کردند.

سرجنت و همکاران (۲۰۰۰) به نیت کاهش همبستگی نمونه‌های حاصل از زنجیر مربوط به پارامترهای مدل‌های چندسطحی روش جدیدی تحت عنوان SMCMC ارائه کردند که شباهت زیادی با مدل‌بندی رگرسیون چندگانه دارد. آن‌ها مطالعه خود را روی مدل‌هایی با پارامترهای زیاد متمرکز کردند. مدل‌هایی با پارامتر زیاد شامل مدل‌های سلسله مراتبی، مدل‌های مولفه واریانس و بعضی از مدل‌های فضایی است (بسیگ و همکاران، ۱۹۹۲). قابل ذکر است که روش SMCMC در تعدادی از مدل‌های اقتصادی و همچنین مدل‌های رشد نیز به کار گرفته شد (میرا و سرجنت، ۲۰۰۳). همچنین استفاده عملی آن در مدلی خاص از آمار فضایی و نمایش کارایی آن نسبت به مدل‌گزینی‌های دیگر توسط بانرجی و همکاران (۲۰۰۴) صورت گرفت.

ساده‌ترین مدل سلسله مراتبی، مدلی با اثرات تصادفی یک طرفه متعادل به صورت

$$y_{ij} = \beta_0 + u_j + \epsilon_{ij}, \quad i = 1, \dots, n_j, \quad j = 1, \dots, J \quad (1)$$

است، که در آن اندیس‌های i و j به ترتیب بیانگر سطح اول و سطح دوم، $\sum_{j=1}^J n_j = N$ تعداد کل مشاهدات، β_0 پارامتر ثابت بیانگر عرض از مبدا، و u_j ها و ϵ_{ij} ها مستقل و دارای توزیع نرمال با میانگین صفر و واریانس‌های به ترتیب σ_u^2 و σ_ϵ^2 هستند. مدل (۱) مدل عرض از مبدا تصادفی با متغیر پاسخ پیوسته y_{ij} نامیده می‌شود. نمادگذاری مدل‌های چندسطحی در این مقاله بر اساس نمادگذاری‌های گلداستاین (۱۹۹۹) است که علاوه بر مطالعه آن‌ها برآورد بسامدی پارامترهای مدل را نیز ارائه کرده است. برآورد بیزی پارامترهای این مدل و مدل‌های چندسطحی دیگر در براون (۲۰۰۹) آمده است.

^۶ Parameter expansion

با پیروی از سرجنت و همکاران (۲۰۰۰) و به منظور به کارگیری روش SMCMC قرار داده می‌شود $\beta_{\circ j} = \beta_{\circ} + u_j$. این رابطه را می‌توان به صورت $\circ = -\beta_{\circ j} + \beta_{\circ} + u_j$ نیز نوشت. با توجه به این تساوی و رابطه (۱) داریم

$$\begin{aligned} y_{ij} &= \beta_{\circ j} + \epsilon_{ij} \\ \circ &= -\beta_{\circ j} + \beta_{\circ} + u_j. \end{aligned} \quad (2)$$

حال به ازای $i = 1, \dots, n_j$ و $j = 1, \dots, J$ رابطه (۲) را می‌توان به صورت برداری

$$\begin{pmatrix} y_{11} \\ y_{21} \\ \vdots \\ y_{n_J J} \\ 0_J \end{pmatrix} = \begin{pmatrix} 1_{n_1} & 0_{n_1} & \dots & 0_{n_1} \\ 0_{n_2} & 1_{n_2} & \dots & 0_{n_2} \\ \vdots & \ddots & \vdots & \vdots \\ 0_{n_J} & \dots & 1_{n_J} & 0_{n_J} \\ -I_J & & & 1_J \end{pmatrix} \begin{pmatrix} \beta_{\circ 1} \\ \vdots \\ \beta_{\circ J} \\ \beta_{\circ} \end{pmatrix} + \begin{pmatrix} \epsilon_{11} \\ \vdots \\ \epsilon_{n_J J} \\ u_1 \\ \vdots \\ u_J \end{pmatrix} \quad (3)$$

نوشت، که در آن 0_q بردار ستونی q بعدی از صفرها، 1_q بردار ستونی q بعدی از یک‌ها و I_J ماتریس همانی با بعد J است. واضح است برابری (۳) را می‌توان به صورت ماتریسی

$$Y = X\Theta + E \quad (4)$$

نوشت، که در آن $Y = (y_{11}, \dots, y_{n_J J}, 0_J)^T$ بردار متغیرهای وابسته،

$$X = \begin{pmatrix} \text{diag}\{1_{n_1}, \dots, 1_{n_J}\} & 0_N \\ -I_J & 1_J \end{pmatrix}$$

ماتریس طرح، $\Theta = (\beta_{\circ 1}, \dots, \beta_{\circ J}, \beta_{\circ})^T$ بردار پارامترهای مدل و $E = (\epsilon_{11}, \dots, \epsilon_{n_J J}, u_1, \dots, u_J)^T = (\epsilon, u)^T$ بردار خطاها در ساختار جدید هستند. در ادامه رابطه رگرسیونی جدید (۴) مدل جدید یا مدل ساختاری نامیده می‌شود. واضح است بردار خطاهای E در رابطه (۴) دارای توزیع نرمال چندمتغیری با بردار میانگین 0_{N+J} و ماتریس کواریانس $\Sigma_{N+J} = \Sigma$

$$\Sigma = \begin{pmatrix} \text{Cov}(\epsilon) & 0_{N \times J} \\ 0_{J \times N} & \text{Cov}(u) \end{pmatrix} = \begin{pmatrix} \sigma_{\epsilon}^2 I_N & 0_{N \times J} \\ 0_{J \times N} & \sigma_u^2 I_N \end{pmatrix}.$$

۴۰ بهبود الگوریتم SMCMC در مدل‌های چندسطحی

است (رنچر، ۱۹۹۵). به کمک اطلاعات موجود و در نظر گرفتن پیشین ناآگاهی بخش تخت برای بردار $J + 1$ بعدی پارامتر θ ، یعنی $\pi(\theta) \propto 1$ ، توزیع شرطی کامل آن عبارت است از

$$\theta | (Y, X, \Sigma) \sim MVN \left((X^T \Sigma^{-1} X)^{-1} (X^T \Sigma^{-1} Y), (X^T \Sigma^{-1} X)^{-1} \right). \quad (5)$$

به کمک الگوریتم نمونه‌گیری گیبس می‌توان نمونه‌های تصادفی از توزیع شرطی کامل θ تولید کرد و از طریق آن به استنباط بیزی راجع به پارامترهای β_0 و β_1, \dots, β_J یا به طور معادل u_1, \dots, u_J پرداخت.

روش SMCMC تنها تغییر اساسی را در به‌روز رسانی θ اعمال می‌کند. در حالی که به‌روز نمودن σ_e^2 و σ_u^2 بدون تغییر و با استفاده از روش‌های معمول نمونه‌گیری گیبز و از توزیع شرطی کامل آن‌ها که گاما یا وارون گاما است انجام می‌شود (براون، ۲۰۰۹). به عنوان مثال برای مدل (۱) با انتخاب پیشین $\Gamma(a_u, b_u)$ و $\Gamma(a_e, b_e)$ به ترتیب برای $1/\sigma_e^2$ و $1/\sigma_u^2$ توزیع‌های شرطی کامل به ترتیب $\Gamma(N/2 + a_e, b_e + \sum_{ij} e_{ij}^2/2)$ و $\Gamma(J/2 + a_u, b_u + \sum_{j=1}^J u_j^2/2)$ خواهد بود، که در آن‌ها $e_{ij} = y_{ij} - \beta_0 - u_j$. این موضوع برای مدل‌های دیگری که در ادامه مطالعه می‌شوند نیز صادق است جز این که تغییرات جزئی در پارامترهای توزیع‌های گاما رخ خواهد داد. به منظور جلوگیری از افزایش حجم مطالب مقاله از بیان جزئیات بیشتر راجع به انتخاب توزیع‌های پیشین و تاثیر آن‌ها بر توزیع‌های پسین و شرطی کامل پارامترهای درگیر مدل خودداری می‌شود. خواننده علاقه‌مند به مطالعه بیشتر در این زمینه به براون و درایپر (۲۰۰۰) و براون (۲۰۰۹) ارجاع داده می‌شود.

در روش‌های معمول MCMC پارامترهای β_0, u_1, \dots, u_J مجزا به‌روز می‌شوند و چون در طی تولید زنجیر بین آن‌ها همبستگی به‌وجود می‌آید روش معمول باعث تاخیر در زمان همگرایی می‌شود (براون، ۲۰۰۹). لذا اخذ نمونه‌های تقریباً مستقل از توزیع‌های پسین آن‌ها مستلزم اجرای الگوریتم MCMC برای مدت طولانی است تا بتوان از تاثیر همبستگی بین نمونه‌ها کاست. اما نکته حائز اهمیت در به‌کارگیری روش SMCMC این است که به دلیل به‌روز شدن توام پارامترهای β_0 و β_1, \dots, β_J در یک بلوک زمان همگرایی الگوریتم MCMC کاهش می‌یابد

عاطفه فرخی، موسی گل‌علی‌زاده ۴۱

(سرچنت و همکاران، ۲۰۰۰). به عبارتی دیگر با به‌روز نمودن توام این پارامترها همبستگی درونی بین آن‌ها در کل زنجیر کاهش می‌یابد و لذا زنجیر بهتر ترکیب می‌شود. واضح است که این امر خود موجب افزایش اندازه نمونه موثر^۷ (ESS) می‌شود. این نکته در ادامه در مثال‌های شبیه‌سازی و کاربردی منعکس خواهد شد.

تعمیم مدل (۱)، مدل عرض از مبدا تصادفی همراه با حضور متغیرهای مستقل است. مدل چند سطحی عرض از مبدا تصادفی با حضور تنها یک متغیر مستقل x_{ij} و با متغیر پاسخ پیوسته y_{ij} عبارت است از

$$y_{ij} = \beta_0 + u_j + x_{ij}\beta_1 + \epsilon_{ij} \quad (۶)$$

فرض‌های مدل (۶) همان فرض‌های مدل (۱) می‌باشند جز این‌که در مدل اخیر پارامتر ثابت اضافی β_1 بیانگر شیب ثابت مدل است (گلداستاین، ۱۹۹۹). برای بکارگیری روش SMCMC ابتدا رابطه (۶) را به صورت دو معادله

$$\begin{aligned} y_{ij} &= \beta_{0j} + x_{ij}\beta_1 + \epsilon_{ij} \\ 0 &= -\beta_{0j} + \beta_0 + u_j. \end{aligned}$$

نوشته و سپس به ازای کلیه مشاهدات آن‌را به صورت برداری

$$\begin{pmatrix} y_{11} \\ \vdots \\ y_{nJJ} \\ 0_J \end{pmatrix} = \begin{pmatrix} 1_{n1} & 0_{n1} & \dots & 0_{n1} & x_1 & 0_{n1} \\ 0_{n2} & 1_{n2} & \dots & 0_{n2} & x_2 & 0_{n2} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0_{nJ} & \dots & 1_{nJ} & 0_{nJ} & x_J & 0_{nJ} \\ & & -I_J & & 0_J & 1_J \end{pmatrix} \begin{pmatrix} \beta_{01} \\ \vdots \\ \beta_{0J} \\ \beta_1 \\ \beta_0 \end{pmatrix} + \begin{pmatrix} \epsilon_{11} \\ \vdots \\ \epsilon_{nJJ} \\ u_1 \\ \vdots \\ u_J \end{pmatrix}$$

می‌نویسیم، که در آن $x_q = (x_{1q}, x_{2q}, \dots, x_{nqq})^T$. واضح است این تساوی نیز می‌تواند به صورت ماتریسی $Y = X\Theta + E$ با ماتریس طرح جدید نوشته شود. لذا توزیع شرطی کامل Θ همان توزیع (۵) خواهد بود، جز اینکه ماتریس طرح X تغییر خواهد کرد. به عنوان مثال، بنا به مدل جدید و رابطه (۵) معکوس ماتریس

^۷ Effective Sample Size

کواریانس توزیع شرطی Θ به شرط سه‌تایی (X, Y, Σ) به صورت زیر خواهد بود

$$X^T \Sigma^{-1} X = \begin{pmatrix} \frac{1}{\sigma_u^2} + \frac{n_1}{\sigma_e^2} & \dots & \circ & \frac{\sum x_{i1}}{\sigma_e^2} & \frac{-1}{\sigma_u^2} \\ \vdots & \ddots & \vdots & \vdots & \vdots \\ \circ & \circ & \frac{1}{\sigma_u^2} + \frac{n_J}{\sigma_e^2} & \frac{\sum x_{iJ}}{\sigma_e^2} & \frac{-1}{\sigma_u^2} \\ \frac{\sum x_{i1}}{\sigma_e^2} & \dots & \frac{\sum x_{iJ}}{\sigma_e^2} & \frac{\sum_{j=1}^J \sum_{i=1}^{n_j} x_{ij}}{\sigma_e^2} & \circ \\ \frac{-1}{\sigma_u^2} & \dots & \frac{-1}{\sigma_u^2} & \circ & \frac{J}{\sigma_u^2} \end{pmatrix}. \quad (V)$$

مجدداً برآورد بردار پارامترهای $\Theta = (\beta_{\circ 1}, \dots, \beta_{\circ J}, \beta_{\circ})^T$ به کمک نمونه‌گیری گیبز از توزیع (۵) با ماتریس طرح جدید X به دست می‌آید. همچنین σ_u^2 و σ_e^2 با روش معمول نمونه‌گیری گیبس و به کمک نمونه‌های تولید شده از توزیع وارون گاما برآورد می‌شوند.

به منظور به کارگیری الگوریتم SMCMC برای مدل‌های عرض از مبدا تصادفی و همچنین مدل‌های چندسطحی دیگر ضروری است معادله مدل به صورت معادله رگرسیون $Y = X\Theta + E$ نوشته شود. سپس به کمک رابطه (۵) برآورد تعدادی از پارامترهای مدل به کمک روش‌های شبیه‌سازی MCMC انجام می‌گیرد. در روش SMCMC به دلیل تعریف مدل ساختاری همبستگی بین پارامترهای ثابت و خطاهای مربوط به سطوح بالا کم خواهد شد. اما نکته مهم روش SMCMC این است که معمولاً بعد ماتریس کواریانس مدل جدید خیلی زیاد است و به دلیل این که به روز نمودن این ماتریس در هر گام از الگوریتم MCMC ضروری است، به کارگیری این روش زمان همگرایی زنجیر را افزایش خواهد داد. نکته دیگر این که با افزایش تعداد گروه‌ها و همچنین ورود متغیرهای مستقل و لحاظ نمودن مدل‌های مختلف چندسطحی بعد ماتریس $(X^T \Sigma^{-1} X)^{-1}$ بزرگ‌تر شده و محاسبه آن حجم وسیعی از محاسبات مربوط به تولید نمونه برای پارامترها را به خود اختصاص خواهد داد.

در مقاله حاضر، برای رفع مشکلاتی شبیه آن‌چه در بالا به آن‌ها اشاره شد دو روش پیشنهاد و کارایی آن‌ها نسبت به روش SMCMC مورد ارزیابی قرار می‌گیرد. روش‌های پیشنهادی را بهبودهای روش SMCMC می‌نامیم. بیان این نکته ضروری است که توجه این مقاله تنها معطوف به مدل عرض از مبدا تصادفی بوده و مدل‌های دیگر چندسطحی مثل مدل شیب تصادفی مورد مطالعه قرار نگرفته است. یکی از

دلایل این امر عدم بهبود مناسب روش‌های پیشنهادی برای مدل‌های دیگر شامل مدل شیب تصادفی بوده است. اما امید می‌رود در تحقیقات آتی این موضوع نیز با راهکارهای دیگر مورد مطالعه قرار گیرد.

۳ بهبود روش ساختاری مونت کارلوی زنجیر مارکوف

به دلیل افزایش بیش از حد بعد ماتریس کواریانس به روش SMCMC انتظار می‌رود دقت برآورد پارامترها در مسائل کاربردی کاهش یابد. این امر از تقریب وارون ماتریس کواریانس و ترکیب آن با ماتریس طرح ناشی می‌شود. یک راه حل مناسب برای بهبود روش SMCMC استفاده از تجزیه ماتریس است.

روش‌های مختلفی برای تجزیه یک ماتریس وجود دارد که هر کدام دارای محاسن و معایب خاص خود هستند. اما تجزیه چولسکی به دلیل سرعت بالای محاسبات نسبت به تجزیه‌های دیگر ارجحیت دارد (بیکر، ۲۰۰۳). اگر A ماتریسی متقارن معین مثبت باشد، تجزیه چولسکی آن به صورت $A = LL^T$ است که در آن L ماتریسی بالا مثلثی است. توجه کنید برای مدل‌های این مقاله شاید تجزیه‌های دیگر ساختار مناسب‌تری برای معکوس ماتریس کواریانس ارائه دهند ولی بنا به بیکر (۲۰۰۳) استفاده از آن‌ها بهبود قابل ملاحظه‌ای در دقت و زمان همگرایی روش SMCMC ایجاد نخواهند کرد.

بیان این نکته ضروری است که برای پیاده‌سازی روش‌های بهبود یافته SMCMC، ابتدا با نرم‌افزار Maple نسخه ۱۱ تجزیه چولسکی معکوس ماتریس‌های مورد نیاز محاسبه و سپس مدل‌بندی و دستیابی به برآوردهای مورد نیاز با نرم‌افزار R نسخه ۲.۹.۰، انجام شده‌اند. همچنین برای برآزش مدل‌های چندسطحی تعدادی از کتابخانه‌های ضروری شامل MASS، lme4 و nlme در نرم‌افزار آماری R به خدمت گرفته شده است.

۱.۳ تجزیه چولسکی و ارون ماتریس کواریانس پارامترهای بلوک بندی شده

در این بخش برای بهبود روش SMCMC، از تجزیه چولسکی ماتریس کواریانس استفاده می‌شود. از آنجا که محاسبه ماتریس کواریانس پارامترهای بلوک‌بندی شده زمان‌بر است، با استفاده از تجزیه چولسکی مشکل افزایش بعد ماتریس کواریانس مدل ساختاری و تاثیر آن در زمان همگرایی روش SMCMC مرتفع می‌شود. روش پیشنهادی برای دو مدل مطرح شده در بخش قبل به کار گرفته خواهد شد.

با توجه به مدل عرض از مبدا تصادفی (۱) و بنا به رابطه (۵) تجزیه چولسکی و ارون ماتریس کواریانس مدل ساختاری یعنی $X^T \Sigma^{-1} X$ عبارت است از

$$\Sigma^* \Sigma^{*T} = \begin{pmatrix} aI_{J+1} & \frac{-1}{a\sigma_u^2} \\ 0_J & \sqrt{\frac{n_J}{n\sigma_u^2 + \sigma_e^2}} \end{pmatrix}^T \begin{pmatrix} aI_{J+1} & \frac{-1}{a\sigma_u^2} \\ 0_J & \sqrt{\frac{n_J}{n\sigma_u^2 + \sigma_e^2}} \end{pmatrix} \quad (۸)$$

که در آن

$$a = \sqrt{\frac{n\sigma_u^2 + \sigma_e^2}{\sigma_u^2 \sigma_e^2}}$$

با توجه به این تساوی و رابطه (۵) به نظر می‌رسد محاسبات مربوط به اجرای الگوریتم MCMC خیلی ساده‌تر از حالت معمول آن (عدم استفاده از تجزیه چولسکی) باشد. به ویژه با جایگذاری تجزیه چولسکی $\Sigma^* \Sigma^{*T}$ به جای $X^T \Sigma^{-1} X$ در رابطه (۵)، ماتریس کواریانس توزیع شرطی کامل Θ به صورت $(\Sigma^{*T})^{-1} (\Sigma^*)^{-1}$ خواهد بود. بنابراین توزیع شرطی کامل Θ به صورت

$$\Theta | (Y, X, \Sigma) \sim MVN \left((\Sigma^* \Sigma^{*T})^{-1} (X^T \Sigma^{-1} Y), (\Sigma^* \Sigma^{*T})^{-1} \right)$$

تغییر خواهد کرد، که با توجه به رابطه (۸) محاسبه واریانس ساده‌تر خواهد شد. از آنجا که عبارت $X^T \Sigma^{-1} Y$ برابر حاصل ضرب تعدادی مقادیر معلوم در $1/\sigma_e^2$ است، با محاسبه مقادیر ثابت آن خارج از تکرار الگوریتم MCMC می‌توان زمان همگرایی را بیش از پیش کاهش داد. به‌عنوان مثال برای مدل عرض از مبدا تصادفی و بدون حضور متغیر مستقل (رابطه (۱) را ببینید) کمیت $X^T \Sigma^{-1} Y$ عبارت است از

$$\frac{1}{\sigma_e^2} \left(\sum_{i=1}^{n_1} y_{ij}, \sum_{i=1}^{n_2} y_{ij}, \dots, \sum_{i=1}^{n_J} y_{ij} \right)^T$$

حال نحوه به کارگیری روش SMCMC در مدل عرض از مبدا تصادفی و با حضور یک متغیر تبیینی تشریح می‌شود. واضح است روش‌های پیشنهادی برای وقتی که تعداد متغیرهای تبیینی بیشتر باشند نیز قابل تعمیم است. توجه کنید تجزیه چولسکی ماتریس $X^T \Sigma^{-1} X$ در مدل چندسطحی عرض از مبدا تصادفی و با حضور متغیر تبیینی در مدل‌هایی که تعداد واحدهای درون سطوح نابرابر باشند دارای فرم استاندارد نیست. به عبارتی دیگر تجزیه چولسکی بسته به تعداد واحد هر گروه دارای ساختار متفاوتی است. لذا نمی‌توان یک دستورالعمل سراسری را برای چنین مدل‌هایی پیشنهاد نمود. از اینرو برای آشنایی بیشتر با نحوه اجرای روش SMCMC و در عین حال سادگی محاسبات آتی واحدهای سطوح متفاوت مدل چندسطحی عرض از مبدا تصادفی و با حضور متغیر تبیینی به صورت برابر در نظر گرفته می‌شود. چنین طرحی به طرح متعادل معروف است.

وارون ماتریس کواریانس مدل چندسطحی عرض از مبدا تصادفی و با حضور متغیر مستقل در (۷) آمده است. با لحاظ نمودن فرض طرح متعادل، تجزیه چولسکی ماتریس (۷) عبارت است از

$$\Sigma^* \Sigma^{*T} = \begin{pmatrix} \frac{f_0 I_J}{\sqrt{f_0 \sigma_e^2}} & \cdots & \frac{\sum x_{iJ}}{\sqrt{f_0 \sigma_e^2}} & 0_{J \times 2} & \circ \\ \frac{-1}{\sqrt{f_0 \sigma_u^2}} & \cdots & \frac{-1}{\sqrt{f_0 \sigma_u^2}} & f_1 & \circ \\ & & & f_2 & f_2 \end{pmatrix} \times \begin{pmatrix} \frac{f_0 I_J}{\sqrt{f_0 \sigma_e^2}} & \cdots & \frac{\sum x_{iJ}}{\sqrt{f_0 \sigma_e^2}} & 0_{J \times 2} & \circ \\ \frac{-1}{\sqrt{f_0 \sigma_u^2}} & \cdots & \frac{-1}{\sqrt{f_0 \sigma_u^2}} & f_1 & \circ \\ & & & f_2 & f_2 \end{pmatrix}^T$$

که در آن f_0, f_1, f_2 و f_2 به صورت زیر تعریف می‌شوند

$$f_0^2 = (n\sigma_u^2 + \sigma_e^2) / (\sigma_u^2 \sigma_e^2),$$

$$f_1^2 = \left(\sigma_e^2 \sum x_{ij}^2 + (n-1) \sum x_{ij}^2 \sigma_u^2 - 2\sigma_u^2 \sum x_{ij} x_{mj} \right) / (f_0 \sigma_e^2),$$

$$f_{\gamma}^{\gamma} = \frac{((nm - 1)\sigma_e^{\gamma} + nm(n - 1)\sigma_u^{\gamma}) \sum x_{ij}^{\gamma} - 2(nm\sigma_u^{\gamma} + \sigma_e^{\gamma}) \sum x_{ij}^{\gamma} x_{mj}^{\gamma}}{f_0 \sigma_e^{\gamma} \sigma_u^{\gamma} \sum x_{ij}^{\gamma} + (n - 1) \sum x_{ij}^{\gamma} \sigma_u^{\gamma} - 2\sigma_u^{\gamma} \sum x_{ij}^{\gamma} x_{mj}^{\gamma}},$$

$$f_{\gamma}^{\gamma} = \frac{\sum x_{ij}^{\gamma} \sqrt{\sigma_e^{\gamma}}}{(f_0 \sigma_e^{\gamma} \sigma_u^{\gamma} \sum x_{ij}^{\gamma} + (n - 1) \sum x_{ij}^{\gamma} \sigma_u^{\gamma} - 2\sigma_u^{\gamma} \sum x_{ij}^{\gamma} x_{mj}^{\gamma})}.$$

فرآیند به کارگیری این تجزیه در به روز کردن پارامتر θ مشابه وضعیتی است که برای مدل عرض از مبدا تصادفی بدون حضور متغیر تبیینی توضیح داده شد. تنها تغییر جزئی ورود متغیر تبیینی x_{ij} در مدل و در نتیجه در تجزیه چولسکی است. با این حال همان گونه که روابط مربوط به ساختار تجزیه چولسکی برای مدل عرض از مبدا تصادفی و با حضور یک متغیر تبیینی نشان می‌دهد با اضافه نمودن متغیرهای مستقل دیگر تجزیه چولسکی وضعیتی به مراتب پیچیده‌تر به خود خواهد گرفت. لذا اعمال آن در فرآیند به روز کردن پارامترها تأثیر بسزایی در کاهش زمان همگرایی الگوریتم MCMC نخواهد داشت.

روش پیشنهادی به کار گرفته شده در این بخش را روش SMCMC1 نامیده و با یک مثال کاربردی و مطالعه شبیه‌سازی دقت برآوردهای حاصل و همچنین زمان همگرایی سیستم با استفاده از این روش با روش‌های دیگر مقایسه می‌شوند. اما قبل از آن روش پیشنهادی دوم این مقاله که در مورد نحوه به کارگیری روش تجزیه چولسکی ماتریس کواریانس مدل اولیه برای دو مدل عرض از مبدا و عرض از مبدا در حضور متغیر مستقل است مورد مطالعه قرار می‌گیرد.

۲.۳ تجزیه چولسکی ماتریس کواریانس مدل ساختاری

پیشنهاد دیگر برای بهبود روش SMCMC استفاده از تجزیه چولسکی برای ماتریس کواریانس مدل ساختاری است. به عبارتی دیگر در این روش ابتدا تجزیه چولسکی معکوس ماتریس کواریانس مدلی که به صورت رابطه (۴) نوشته شده باشد را محاسبه و سپس با دخالت آن در رابطه (۵) مدلی ساده‌تر به دست آورده می‌شود. این روش را SMCMC2 نامگذاری می‌نماییم.

برای مدل عرض از مبدا تصادفی (۱) ماتریس کواریانس مدل ساختاری (۴) ماتریسی قطری متشکل از N تا σ_e^2 و J تا σ_u^2 است. لذا ماتریس بالا مثلثی حاصل از تجزیه چولسکی، که با Σ_1 نمایش داده می‌شود، ماتریسی قطری بترتیب شامل N مقدار σ_e^{-1} و J مقدار σ_u^{-1} است. حال با جای‌گذاری این تجزیه در عبارت $X^T \Sigma^{-1} X$ نمونه‌گیری از توزیع شرطی کامل Θ به صورت

$$\Theta | (Y, X, \Sigma) \sim MVN \left((X^T \Sigma_1 \Sigma_1^T X)^{-1} X^T \Sigma_1^T \Sigma_1 Y, (X^T \Sigma_1 \Sigma_1^T X)^{-1} \right). \quad (9)$$

خواهد بود. اکنون پیش‌بینی می‌شود با اختیار متغیر جدید $Z = X \Sigma_1^{-1}$ روش SMCMC، بهبود بیشتری نسبت به روش SMCMC1 داشته باشد. به کمک ماتریس جدید، میانگین و واریانس توزیع شرطی کامل Θ در (۹) به ترتیب برابر $(Z^T Z)^{-1} Z^T \Sigma_1 Y$ و $(Z^T Z)^{-1}$ خواهند بود. به عبارتی دیگر به منظور به‌روز کردن Θ ضروری است از توزیع نرمال چندمتغیره با میانگین $(Z^T Z)^{-1} Z^T \Sigma_1 Y$ و واریانس $(Z^T Z)^{-1}$ نمونه‌گیری کرد. لذا محاسبه Z و ذخیره‌سازی آن به جای محاسبه مستقیم $(X^T \Sigma_1 \Sigma_1^T X)^{-1}$ بر سرعت همگرایی خواهد افزود.

شکل ماتریسی تجزیه چولسکی ماتریس کواریانس به روش SMCMC2 در مدل عرض از مبدا تصادفی و در حضور تنها یک متغیر تبیینی مانند مدل عرض از مبدا تصادفی و بدون حضور متغیر مستقل است. تنها تفاوت این دو مدل در ماتریس طرح آن‌ها خواهد بود. در این حالت اگر تجزیه چولسکی ماتریس کواریانس مدل ساختاری مربوطه با Σ_* نشان داده شود رابطه (۵) به صورت

$$\Theta | (Y, X, \Sigma) \sim MVN \left((X^T \Sigma_* \Sigma_*^T X)^{-1} X^T \Sigma_*^T \Sigma_* Y, (X^T \Sigma_* \Sigma_*^T X)^{-1} \right).$$

تغییر خواهد کرد. همان‌طور که ملاحظه می‌شود تعدادی از کمیت‌های موجود در ماتریس کواریانس در بردار میانگین نیز ظاهر شده‌اند. لذا با ذخیره اولیه تعدادی از کمیت‌های ثابت که در مراحل شبیه‌سازی SMCMC نیازی به به‌روز سازی آن‌ها نیست بر سرعت الگوریتم SMCMC افزوده خواهد شد.

۴ تحلیل یک مثال کاربردی و مطالعه شبیه‌سازی

در این بخش ابتدا یک مثال واقعی مورد مطالعه قرار خواهد گرفت. سپس با شبیه‌سازی داده‌هایی که تا حدودی مشابه مثال کاربردی بخش ابتدایی هستند تحلیل آماری مربوطه صورت می‌گیرد. به ویژه اجرای الگوریتم‌های پیشنهادی با روش‌های موجود قبلی مورد مقایسه قرار می‌گیرد.

۱.۴ مثال کاربردی

پایگاه <http://www.emgo.nl/quality-of-our-research/research-tools/multilevel> شامل داده‌هایی مربوط به ۴۴۱ بیمار در محدوده سنی ۴۴ تا ۸۶ سال است که تحت درمان ۱۲ پزشک بوده و میزان کلسترول خون آن‌ها بر حسب میلی مول در لیتر ثبت شده است. بنا به نمادگذاری مدل چندسطحی گروه‌بندی بیماران بر اساس پزشکان است. از آن جا که از نقطه نظر کلسترول خون بیماران تحت پزشک خاص تا حدودی با هم شباهت دارند، برای تحلیل مشاهدات نمی‌توان از روش‌های متداول رگرسیون که همبستگی بین بیماران را نادیده می‌گیرد استفاده کرد. از این رو از مدل دو سطحی، که سطح اول بیماران تحت درمان و سطح دوم پزشکان هستند، استفاده می‌شود. با توجه به داده‌ها تعداد بیماران هر گروه، یعنی n_j به ازای $j = 1, \dots, 12$ به ترتیب عبارتند از ۳۶، ۳۶، ۳۹، ۳۶، ۳۶، ۳۹، ۳۶، ۳۶، ۳۹، ۳۶، ۳۶، ۳۶ و ۳۶. این داده‌ها با ترکیب‌های مختلفی از مدل‌های چندسطحی با روش کمترین توان‌های دوم تعمیم یافته بازگشتی مورد تحلیل قرار گرفت (توایسک، ۲۰۰۶). اما برآورد بیزی پارامترهای مدل و به‌ویژه با روش MCMC تعیین نشده است. در ادامه ابتدا پارامترهای مدل چندسطحی با استفاده از روش‌های MCMC، SMCMC و بهبودهای روش SMCMC برآورد خواهند شد. سپس زمان همگرایی الگوریتم در روش‌های متفاوت مورد بررسی قرار خواهد گرفت. هنگام اجرای روش‌های MCMC، SMCMC و بهبودهای آن زنجیر ۱۰۰۰۰۰ بار تکرار کرده و برای محاسبه میانگین‌های توزیع‌های شرطی کامل، به‌عنوان برآورد پارامترها، گام داغیدن برابر ۱۰۰۰ در نظر گرفته شده است. همچنین توزیع‌های پیشین پارامترهای

مدل به صورت

$$\pi(\beta_1) \propto 1, \quad \pi(\beta_0) \propto 1, \quad \pi(1/\sigma_u^2) \sim \Gamma(a_u, b_u), \quad \pi(1/\sigma_e^2) \sim \Gamma(a_e, b_e)$$

اختیار شده است به طوری که مقادیر ابر پارامترها همان مقادیر پیشنهاد شده توسط براون و درایپر (۲۰۰۰) و براون (۲۰۰۹) یعنی $a_u = a_e = b_u = b_e = 0.001$ هستند. توجه شود که انتخاب پیشین‌های یکنواخت برای پارامترهای اثر ثابت منجر به تساوی برآوردگرهای بیزی و ماکسیمم درست‌نمایی خواهد شد. برای برآورد پارامترها در مدل‌های چندسطحی ابتدا پارامترهای مدل ساده عرض از مبدا تصادفی و بدون حضور متغیر مستقل برآورد می‌شوند. سپس مدل عرض از مبدا تصادفی و با حضور متغیر مستقل سن مورد بررسی قرار می‌گیرد.

برآورد پارامترها همراه با خطای برآورد عرض از مبدا در یک مدل دوسطحی عرض از مبدا تصادفی در جدول ۱ آمده است. برآورد پارامترهای مدل عرض از مبدا تصادفی با روش‌های *MCMC*، *SMCMC*، *SMCMC1* و *SMCMC2* محاسبه شده است. تغییرات کم برآورد β_0 حاکی از دقت روش‌های برآورد است. نتایج نشان می‌دهند که با استفاده از روش‌های متفاوت، تفاوت زیادی بین برآورد پارامترها وجود ندارد. اما در ادامه نشان داده خواهد شد که روش پیشنهادی از دو مزیت عمده برخوردار است.

جدول ۱: برآورد و خطای استاندارد برآورد پارامترهای مدل عرض از مبدا تصادفی با روش‌های مختلف *MCMC*.

$\hat{\sigma}_e^2$	$\hat{\sigma}_u^2$	$\hat{\beta}_0 (SE(\hat{\beta}_0))$	روش
۰/۵۴۱	۰/۵۲۲	۵/۹۰۰ (۰/۰۷۲)	<i>MCMC</i>
۰/۵۳۶	۰/۵۲۳	۵/۹۸۱ (۰/۰۸۱)	<i>SMCMC</i>
۰/۵۳۶	۰/۵۲۸	۶/۰۲۹ (۰/۰۸۲)	<i>SMCMC1</i>
۰/۵۳۶	۰/۵۲۳	۵/۹۸۱ (۰/۰۸۱)	<i>SMCMC2</i>

تجزیه چولسکی ماتریس کواریانس ساختاری برای مدل عرض از مبدا تصادفی و در حضور متغیر تبیینی در حالتی که تعداد واحدها در گروه‌های متفاوت یکسان نباشد روند منظمی را دنبال نمی‌کند. اما برای انجام یک مقایسه مناسب

۵۰ بهبود الگوریتم SMCMC در مدل‌های چندسطحی

تعداد مشاهدات هر گروه ۳۶ در نظر گرفته شد و برآورد پارامترهای مدل به روش SMCMC1 انجام شده است. قابل ذکر است متغیر مستقل مدل سن بیماران تحت درمان است. نتایج آن همراه با روش‌های MCMC، SMCMC و SMCMC2 در جدول ۲ ارائه شده است. نتایج نشان می‌دهد بین برآورد پارامترها با روش‌های متفاوت تنها تغییرات جزئی وجود دارد. اما برآورد پارامترها در این مدل با برآورد پارامترها در مدل عرض از مبدا تصادفی و بدون حضور متغیر مستقل متفاوت است. به علاوه مشاهده می‌شود مقدار عرض از مبدا (β_0) نسبت به قبل کاهش چشمگیری داشته و در واریانس‌های مدل نیز تغییراتی ایجاد شده است. همان‌گونه که در

جدول ۲: برآورد و خطای استاندارد برآورد پارامترهای مدل عرض از مبدا تصادفی با حضور متغیر تبیینی سن با روش‌های مختلف MCMC.

روش	$\hat{\beta}_0 (SE(\hat{\beta}_0) \times 10^5)$	$\hat{\beta}_1 (SE(\hat{\beta}_1))$	$\hat{\sigma}_{u_0}^2$	$\hat{\sigma}_e^2$
MCMC	۲/۹۲۹(۹/۸۱۲)	۰/۰۴۹(۰/۰۷۷)	۰/۴۹۱	۰/۳۲۵
SMCMC	۲/۹۴۸(۳/۲۴۱)	۰/۰۴۹(۰/۰۹۲)	۰/۴۹۳	۰/۳۳۴
SMCMC1	۲/۹۳۰(۹/۴۳۹)	۰/۰۴۹(۰/۰۷۸)	۰/۴۹۴	۰/۳۳۰
SMCMC2	۲/۹۰۱(۱/۱۵۱)	۰/۰۵۰(۰/۰۸۸)	۰/۴۹۸	۰/۳۳۲

دو مدل مورد بررسی ملاحظه شد، روش SMCMC و بهبودهای آن برآوردهایی حدوداً مشابه روش‌های موجود قبلی ارائه می‌نمایند. اما نکته قابل تامل در مورد روش‌های SMCMC و بهبودهای آن نتایج حاصل از مقدار ESS و میزان زمان مصرفی سیستم (به ثانیه) در ارائه برآورد پارامترهای مدل‌های چندسطحی متفاوت شامل مدل عرض از مبدا تصادفی و عرض از مبدا تصادفی در حضور متغیر تبیینی سن است. جدول ۳ مقدار ESS و میزان زمان مصرفی سیستم در دستیابی به برآورد پارامترهای مدل عرض از مبدا تصادفی و بدون حضور متغیر مستقل به روش‌های MCMC، SMCMC، SMCMC1 و SMCMC2 را نشان می‌دهد. همان‌گونه که ملاحظه می‌شود هم زمان محاسبه برآورد پارامترهای مدل مورد نظر با روش تجزیه چولسکی ماتریس کواریانس به روش SMCMC2 خیلی کمتر است و هم ESS آن از همه بیشتر است. گرچه زمان کاربر در روش SMCMC2 تنها

عاطفه فرخی، موسی گل‌علی‌زاده ۵۱

بیشتر از روش SMCMC است، اما این خسارت ناشی از تعدادی دستورات اضافی ضروری در این روش بوده و زیاد چشمگیر نمی‌باشد.

جدول ۳: مقدار ESS و زمان دستیابی برآورد پارامترها در مدل عرض از مبدأ تصادفی با روش‌های مختلف MCMC.

روش	ESS	زمان کاربر (به ثانیه)	زمان سیستم (به ثانیه)
MCMC	۳۷۴۲۱	۱۴۴/۰۷	۰/۱۱
SMCMC	۹۵۷۴۱	۱۰۸/۷۸	۰/۱۷
SMCMC1	۹۵۸۸۲	۱۹۹/۹۲	۰/۱۳
SMCMC2	۹۶۳۹۶	۱۱۰/۱۵	۰/۰۶

زا ضریب عدم دقت اهرت‌ماراپ در برابر روش‌های MCMC، SMCMC، SMCMC1 و SMCMC2 نسبت به روش‌های دیگر کمتر است. همچنین روش‌های SMCMC1 و SMCMC2 نسبت به روش‌های دیگر دارای کمترین خطای استاندارد است. از آنجا که روش‌های SMCMC1 و SMCMC2 نسبت به روش‌های دیگر دارای کمترین خطای استاندارد است، بنابراین روش‌های SMCMC1 و SMCMC2 را می‌توان به عنوان بهترین روش‌ها برای برآورد پارامترها در مدل عرض از مبدأ تصادفی در نظر گرفت. همچنین روش‌های SMCMC1 و SMCMC2 نسبت به روش‌های دیگر دارای کمترین خطای استاندارد است، بنابراین روش‌های SMCMC1 و SMCMC2 را می‌توان به عنوان بهترین روش‌ها برای برآورد پارامترها در مدل عرض از مبدأ تصادفی در نظر گرفت. همچنین روش‌های SMCMC1 و SMCMC2 نسبت به روش‌های دیگر دارای کمترین خطای استاندارد است، بنابراین روش‌های SMCMC1 و SMCMC2 را می‌توان به عنوان بهترین روش‌ها برای برآورد پارامترها در مدل عرض از مبدأ تصادفی در نظر گرفت.

جدول ۴: مقدار ESS و زمان (برحسب ثانیه) دستیابی به برآورد پارامترها در مدل عرض از مبدا تصادفی با حضور متغیر تبیینی سن با روش‌های مختلف MCMC.

روش	ESS	زمان کاربر	زمان سیستم
MCMC	۲۵۸۷۶	۱۸۰/۰۷	۲۳/۴
SMCMC	۸۸۷۳۵	۱۲۴/۱۶	۰/۳۶
SMCMC۱	۸۸۴۳۲	۱۶۷/۶۱	۰/۱۶
SMCMC۲	۸۸۹۳۲	۱۵۲/۴۴	۰/۱۱

۲.۴ مطالعه شبیه‌سازی

در این بخش بر اساس داده‌های شبیه‌سازی از مدل‌های تحت بررسی، روش‌های پیشنهادی مقایسه می‌شوند. بنا به خلاصه‌های آماری مربوط به متغیرهای کلاسترول خون و سن بیماران در بخش قبل، سعی شد داده‌های شبیه‌سازی تا حدودی مشابه آن‌ها باشند. با استفاده از ۴۴۱ بار شبیه‌سازی از توزیع‌های $U(۳, ۹)$ و $N(۶, ۱)$ برای متغیر وابسته و از توزیع‌های $U(۴۴, ۸۶)$ و $N(۶۵, ۸۱)$ برای متغیر تبیینی به تحلیل آماری روش‌های پیشنهادی پرداخته شد. برای اجرای روش‌های SMCMC و بهبودهای متناظر با آن زنجیر ۱۰۰۰۰۰ بار تکرار و در محاسبه میانگین‌های توزیع‌های شرطی کامل، به عنوان برآورد پارامترها، ۱۰۰۰ مرحله داغیدن در نظر گرفته شده‌اند. فرض‌های زیر بنایی مدل همانند بخش‌های قبل در نظر گرفته شده است. برآورد پارامترهای مدل دوسطحی عرض از مبدا تصادفی به روش‌های مختلف SMCMC و تحت شبیه‌سازی از دو توزیع مشروحه فوق در جدول ۵ آمده است.

همان‌طور که ملاحظه می‌شود، برآورد پارامترها تحت توزیع نرمال برای متغیر وابسته تفاوت چندانی با برآوردهای متناظر آن هنگام کار با داده‌های واقعی ندارد. با این حال در هنگام استفاده از توزیع یکنواخت برآورد واریانس خطای سطح اول خیلی کم شده است. این موضوع ناشی از تغییرپذیری زیاد مشاهدات هنگام استفاده از توزیع یکنواخت است. در نتیجه استواری برآوردها نسبت به تغییرپذیری توزیع‌ها کم است. نکته مهم شباهت نسبی برآوردها تحت روش‌های مختلف اما با فرض

جدول ۵: برآورد و خطای استاندارد برآورد پارامترهای مدل عرض از مبداء تصادفی با روش‌های مختلف MCMC.

توزیع	روش	$\hat{\beta}_0 (SE(\hat{\beta}_0))$	$\hat{\sigma}_u^2$	$\hat{\sigma}_e^2$
یکنواخت	SMCMC	۶/۰۲۹ (۰/۰۴۷)	۰/۰۷۷	۲/۹۵۵
	SMCMC۱	۶/۱۲۹ (۰/۰۴۷)	۰/۰۸۴	۲/۹۵۹
	SMCMC۲	۶/۰۲۹ (۰/۰۴۶)	۰/۰۷۶	۲/۹۵۵
نرمال	SMCMC	۶/۰۱۰ (۰/۰۵۳)	۰/۰۵۸	۵/۰۶۴
	SMCMC۱	۶/۰۷۸ (۰/۰۵۳)	۰/۰۱۹	۶/۰۷۸
	SMCMC۲	۵/۹۸۲ (۰/۰۵۲)	۰/۰۱۶	۶/۰۶۴

یک توزیع خاص برای متغیرهای مورد مطالعه است. این موضوع مقایسه روش‌های برآورد براساس مقدار ESS و زمان مصرفی سیستم را میسر می‌سازد.

با فرض ثابت بودن شیب خط رگرسیون برای گروه‌های متفاوت متغیر تبیینی را وارد مدل نموده و با روش SMCMC و بهبود آن پارامترهای مدل برآورد شد. نتایج در جدول ۶ ارائه شده‌اند. همان‌طور که ملاحظه می‌شود بین برآورد پارامترها هنگام استفاده از توزیع‌های یکنواخت و نرمال و برآوردهای متناظر آن‌ها هنگام استفاده از داده‌های واقعی تفاوت عمده‌ای وجود دارد. بویژه برآورد واریانس خطای سطح اول در هر دو حالت کاهش یافته ولی مقدار آن در سطح دوم افزایش یافته است. به علاوه، مثل حالت مدل عرض از مبداء تصادفی، مقایسه برآوردهای حاصل از یک توزیع خاص، مثلاً نرمال، نشانگر شباهت نسبی بین برآوردها و خطای برآورد است. در نتیجه می‌توان آن‌ها را بر اساس مقدار ESS و زمان مصرفی سیستم مورد مقایسه قرار داد.

علاوه بر بررسی دقت و نیرومندی برآوردها با استفاده از شبیه‌سازی مدل، لازم است میزان بهبود زمان اجرای الگوریتم‌های برآورد بر اساس روش‌های موجود مورد ارزیابی قرار گیرد. برای این منظور مقدار ESS، زمان‌های کاربر و سیستم برای هر کدام از وضعیت‌ها و مدل‌ها در جدول‌های ۷ و ۸ ارائه شده‌اند. نتایج حاکی از آن است که روش تجزیه چولسکی و ارون ماتریس کواریانس مدل اولیه کمترین زمان سیستم (به ثانیه) در دستیابی به برآورد پارامترها را به خود اختصاص داده

است. به علاوه مقدار ESS حاصل از این روش در هر دو وضعیت شبیه‌سازی داده‌ها بیشترین مقدار ممکن را دارد. در نتیجه روش تجزیه چولسکی معکوس ماتریس کواریانس مدل اولیه روش برتر در بین روش‌های مورد بررسی است.

برای مدل عرض از مبدا تصادفی همراه با حضور یک متغیر تبیینی نیز مقدار ESS و مقایسه زمان سیستم در تولید نمونه‌های کارا مورد محاسبه قرار گرفت. همان‌گونه که در جدول ۸ ملاحظه می‌شود تجزیه چولسکی وارون ماتریس کواریانس مدل اولیه و سپس اجرای الگوریتم SMCMC از زمان مصرفی کمتری برخوردار است. به علاوه مقدار ESS حاصل از این روش در هر دو وضعیت شبیه‌سازی داده‌ها بیشترین مقدار ممکن را دارد. بنابراین اجرای الگوریتم SMCMC همراه با تجزیه چولسکی ماتریس کواریانس مدل ساختاری بهبود قابل ملاحظه‌ای در کاهش زمان مصرفی سیستم برای تولید نمونه‌های کارا اعمال می‌نماید. جالب اینکه کاهش زمان مصرفی متأثر از توزیع‌های آماری مربوط به متغیرهای مورد مطالعه نیست. اجرای شبیه‌سازی یک نتیجه فرعی دیگری را نیز ارائه نموده است. پایداری برآوردها شدیداً تحت تاثیر توزیع‌ها است. در حالی که با فرض نرمال بودن متغیرهای تبیینی و وابسته برآوردها مورد اعتمادند، اما این وضعیت با فرض توزیع یکنواخت از اعتبار لازم برخوردار نیست.

بحث و نتیجه‌گیری

با گسترش توانایی کامپیوتر برای برآورد بیزی پارامترها در مدل‌های چندسطحی محققین تعمیم‌های متفاوتی از الگوریتم MCMC را معرفی نمودند. یکی از آن روش‌ها روش SMCMC است. در این مقاله نشان داده شد که اگرچه روش SMCMC به کمک تعریف ساختار جدید باعث حذف همبستگی بین پارامترهای ثابت و خطای تصادفی مربوط به سطوح بالا می‌شود اما به دلیل افزایش بعد ماتریس کواریانس مدل ساختاری و به روز نمودن آن در هر گام MCMC زمان همگرایی روش SMCMC کاهش می‌یابد. برای رفع این مشکل دو روش بر اساس تجزیه چولسکی ماتریس کواریانس پیشنهاد شد. همچنین کارکرد روش‌های پیشنهادی بر

جدول ۶: برآورد و خطای استاندارد برآورد پارامترهای مدل عرض از مبداء تصادفی در حضور متغیر تبیینی با روش‌های مختلف MCMC.

$\hat{\sigma}_e^2$	$\hat{\sigma}_u^2$	$\hat{\beta}_1 (SE(\hat{\beta}_1))$	$\hat{\beta}_0 (SE(\hat{\beta}_0) \times 10^5)$	روش	توزیع
۲/۹۵۰	۰/۰۱۳	۰/۰۱۰ (۰/۱۵۵)	۶/۰۰۴ (۳/۲۴)	SMCMC	
۲/۹۴۹	۰/۰۱۲	۰/۰۱۰ (۰/۲۰۰)	۶/۰۰۶ (۵/۵۹)	SMCMC۱	یکنواخت
۲/۹۴۹	۰/۰۱۳	۰/۰۱۰ (۰/۲۴۲)	۶/۰۰۶ (۵/۵۹)	SMCMC۲	
۶/۰۷۶	۰/۰۱۸	۰/۰۰۷ (۰/۱۷۴)	۵/۵۰۸ (۲/۹۰)	SMCMC	
۶/۰۵۰	۰/۰۱۸	۰/۰۰۷ (۰/۱۸۴)	۵/۵۰۹ (۳/۵۰)	SMCMC۱	نرمال
۶/۰۴۹	۰/۰۱۸	۰/۰۰۷ (۰/۱۸۳)	۵/۵۰۷ (۳/۳۰)	SMCMC۲	

جدول ۷: مقدار ESS و زمان (بر حسب ثانیه) محاسبه برآورد پارامترها در مدل عرض از مبداء تصادفی با روش‌های مختلف MCMC.

زمان سیستم	زمان کاربر	ESS	روش	توزیع
۰/۰۴	۱۱۱/۱۷	۶۴۳۲۱	SMCMC	
۰/۰۷	۱۶۶/۹۷	۸۴۳۲۹	SMCMC۱	یکنواخت
۰/۰۳	۱۱۱/۰۰	۸۶۵۶۷	SMCMC۲	
۰/۱	۱۰۸/۵۸	۵۷۳۰۰	SMCMC	
۰/۱۲	۱۶۲/۴۸	۷۴۲۸۷	SMCMC۱	نرمال
۰/۰	۱۰۹/۷۰	۷۶۹۳۱	SMCMC۲	

جدول ۸: مقدار ESS و زمان (بر حسب ثانیه) محاسبه برآورد پارامترها در مدل عرض از مبداء تصادفی در حضور متغیر تبیینی با روش‌های مختلف MCMC.

زمان سیستم	زمان کاربر	ESS	روش	توزیع
۰/۰۲	۱۲۵/۴۳	۵۴۶۵۴	SMCMC	
۰/۰۴	۱۳۰/۸۸	۷۴۵۰۰	SMCMC۱	یکنواخت
۰/۰۱	۱۲۰/۷۸	۷۶۵۰۰	SMCMC۲	
۰/۸۴	۱۵۵/۵۸	۴۹۱۰۰	SMCMC	
۰/۸۶	۱۶۳/۷۸	۷۰۵۰۱	SMCMC۱	نرمال
۰/۰۷	۱۵۲/۰۲	۷۳۸۹۳	SMCMC۲	

اساس یک مثال کاربردی و یک مطالعه شبیه‌سازی مورد ارزیابی قرار گرفت و تجزیه ماتریس کواریانس مدل ساختاری و سپس اجرای روش شبیه‌سازی MCMC به عنوان یک راه‌کار مفید و موثر به منظور تولید نمونه‌های موثر مستقل‌تر و همچنین کاهش زمان اجرای الگوریتم SMCMC پیشنهاد گردید. به علاوه نشان داده شد که کارایی روش SMCMC با فرض نرمال بودن متغیرهای درگیر بیشتر از فرض توزیع یکنواخت است و پیش‌بینی می‌شود این مطلب با فرض توزیع‌های دیگر نیز صادق باشد.

تقدیر و تشکر

نویسندگان از داوران محترم که پیشنهادهای سازنده‌ای برای بهبود این مقاله ارائه کردند تقدیر و تشکر می‌نمایند.

مراجع

- Baker, A. J. (2003), *Matrix Groups: An Introduction to Lie Group Theory*, New York, Springer-Verlag.
- Banerjee, S, Carlin, B. P. and Gelfand, A. E. (2004), *Hierarchical Modelling and Analysis for Spatial Data*, London, Chapman and Hall.
- Besag, J., York, J. and Mollie, A. (1992), Bayesian Image Restoration with Two Applications in Spatial Statistics, *Annals of Institute of Statistical Mathematics*, **43**, 1-59.
- Browne, W. J. (2009), *MCMC Estimation in MLwiN (Version 2.10)*, Center for Multilevel Modelling, University of Bristol.
- Browne, W. J. and Draper, D. (2000), Implementation Issues in Bayesian Fitting of Multilevel Models, *Computational Statistics*, **15**, 391-420.

- Gelfand, A. E., Sahu, S. K. and Carlin, B. P. (1995), Efficient Parameterisations for Linear Mixed Models, *Biometrika*, **83**, 479-488.
- Gelman, A., van Dyk, D. A., Huang, Z.Y. and Boscardin, W. J. (2008), Using Redundant Parameterization to Fit Hierarchical Models, *Journal of Computational and Graphical Statistics*, **17**, 95-122.
- Gilks, W.R. and Robert, G. O. (1996), Strategies for Improving MCMC, In: Gilks, W. R., Richardson, S. and Spiegelhalter, D. J. (eds), *Markov Chain Monte Carlo in Practice*, 89-140, London, Chapman and Hall.
- Goldstein, H. (1999), *Multilevel Statistical Models*, London, Institute of Education, Multilevel Project.
- Hills, S. E. and Smith, A. F. M. (1992), Parameterization Issues in Bayesian Inference, In Bernardo, J., Berger, J., Dawid, A. and Smith, A. *Bayesian Statistics*, **4**, 227-246.
- Mira. A. and Sargent. D. J. (2003), A New Strategy for Speeding Markov Chain Monte Carlo Algorithms. *Statistical Methods and Applications*, **12**, 49-60.
- Papaspiliopoulos, O., Roberts, G. O. and Skold, M. (2007), A General Framework for Parametrisation of Hierarchical Models, *Statistical Sciences*, **22**, 59-73.
- Pinheiro, J. C. and Bates, D. M. (2000), *Mixed-Effects Models in R and S-plus*, New York, Springer-Verlag.
- Rencher, A. C. (1995), *Methods of Multivariate Analysis*, New York, Wiley.

۵۸ بهبود الگوریتم SMCMC در مدل‌های چندسطحی

Sargent, D. J., Hodges, J. S. and Carlin, B. P. (2000), Structured Markov Chain Monte Carlo, *Journal of Computational and Graphical Statistics*, **9**, 217-234.

Twisk, J. W. R. (2006), *Applied Multilevel Analysis: Practical Guides to Biostatistics and Epidemiology*, Cambridge, Cambridge University Press.

Improving of Structured Markov Chain Monte Carlo Algorithm in Multilevel Models

Farokhi, A. and Golalizadeh, M.

Department of Statistics, Tarbiat Modares University, Tehran, Iran.

Abstract: The multilevel models are used in applied sciences including social sciences, sociology, medicine, economic for analysing correlated data. There are various approaches to estimate the model parameters when the responses are normally distributed. To implement the Bayesian approach, a generalized version of the Markov Chain Monte Carlo algorithm, which has a simple structure and removes the correlations among the simulated samples for the fixed parameters and the errors in higher levels, is used in this article. Because the dimension of the covariance matrix for the new error vector is increased, based upon the Cholesky decomposition of the covariance matrix, two methods are proposed to speed the convergence of this approach. Then, the performances of these methods are evaluated in a simulation study and real life data.

Keywords: Multilevel data, Random intercept models, MCMC algorithm, Cholesky decomposition.

Mathematics Subject Classification (2000): 62F15, 62J99