

رگرسیون ضرایب متغیر طولی حاشیه‌ای

حسین بهرامی چشمه‌علی، آرش اردلان

گروه ریاضی، دانشگاه یاسوج

تاریخ دریافت: ۱۳۹۵/۱۱/۱۴ تاریخ آخرین بازنگری: ۱۳۹۷/۱/۲۶

چکیده: مدل‌های رگرسیونی ناپارامتری و نیمه‌پارامتری در زمینه داده‌های مستقل توسعه چشمگیری پیدا کرده‌اند، اما رشد آن‌ها در زمینه داده‌های طولی، محدود به چند سال اخیر است. از آنجا که روش‌های رگرسیونی معمول برای داده‌های همبسته نسبت به داده‌های مستقل توانایی کمتری دارند، باید از مدل‌هایی استفاده شود، که همبستگی بین داده‌ها را نیز در نظر بگیرند. در این میان مدل‌های آمیخته و حاشیه‌ای که عامل همبستگی بین داده‌ها را نیز در نظر می‌گیرند، مدل‌هایی هستند که برای برازش داده‌های طولی مورد استفاده قرار می‌گیرند. همچنین با توجه به انعطاف‌پذیری مدل‌های نیمه‌پارامتری نسبت به مدل‌های پارامتری و ناپارامتری، مدل رگرسیون نیمه‌پارامتری طولی حاشیه‌ای با برآوردهای اسپلاین تاوانیده مدل مناسبی برای تحلیل داده‌های طولی است. در این مقاله رگرسیون نیمه‌پارامتری با ضرایب متغیر که در آن ارتباط بین متغیر پاسخ و یک متغیر پیش‌بین بر مبنای متغیر پیش‌بین دیگر مشخص می‌شود، بررسی شده است. همچنین استنباط بیزی برای مدل ناپارامتری روی داده‌های شبیه‌سازی شده و برای مدل نیمه‌پارامتری طولی حاشیه‌ای روی داده‌های واقعی، با نرم‌افزارهای استاندارد انجام شده است که نشان‌دهنده عملکرد قابل قبول این استنباط است.

واژه‌های کلیدی: نمونه‌گیری گیبز، مدل آمیخته، مدل حاشیه‌ای، مدل گرافیکی، مدل بیزی سلسله مراتبی.

۱ مقدمه

مدل رگرسیونی یکی از پرکاربردترین ابزارهای آماری است که به دلیل تنوع زیاد مطالعات، روش‌های گوناگونی برای انجام آن پدید آمده است. از یک منظر رگرسیون را می‌توان به سه بخش پارامتری، ناپارامتری، نیمه‌پارامتری تقسیم کرد. رگرسیون پارامتری بیشتر وقتی مورد استفاده قرار می‌گیرد که اطلاعات تحقیق به قدر کافی زیاد باشد تا آنجایی که بتوان ارتباط بین متغیر پاسخ و پیش‌بین را به وسیله یک مدل خطی یا قابل‌تبدیل به خطی بیان کرد. رگرسیون ناپارامتری در مواقعی که نتوان نوع ارتباط بین متغیرهای پیش‌بین و پاسخ را با یک شکل تابعی معلوم و مشخص بیان کرد، مورد استفاده قرار می‌گیرد، که در آن می‌توان از برآوردهای موضعی مانند اسپلاین استفاده کرد. مدل‌های نیمه‌پارامتری که در این پژوهش تمرکز بیشتر بر آنها است، ترکیبی از مدل‌های پارامتری و ناپارامتری است، به این معنی که بعضی از متغیرها دارای روند معلوم تغییرات نسبت به متغیر پاسخ بوده و برای بعضی از متغیرها نمی‌توان شکل تابعی خاصی را در نظر گرفت. مدل‌های نیمه‌پارامتری از آنجایی که شکل تابعی پیچیده‌تری نسبت به مدل‌های پارامتری و ناپارامتری دارند، مسئله برآورد یابی برای آن‌ها نیز پیچیدگی بیشتری دارد که برای مثال می‌توان به آمیخته شدن این مسئله با مفاهیم محض ریاضی همچون فضای هیلبرت اشاره کرد. از آنجا که فرم‌های تابعی نامعلوم در مدل‌های ناپارامتری و نیمه‌پارامتری اغلب با روش اسپلاین برآورد می‌شوند و نوع داده‌ها نیز طولی است، یک روش منطقی برای بررسی این مدل‌ها استفاده از ارتباط مدل‌های آمیخته و مدل‌های اسپلاین است. گرچه این مسئله تا حدودی برآوردیابی را در حوزه مدل‌های نیمه‌پارامتری آسان می‌کند اما باید توجه داشت که برآوردیابی برای مدل‌هایی با فرض‌های ناگوسی با این روش نیز کمی دچار مشکل است و امکان دارد برای حجم نمونه پایین اطلاعات حاصل اریب نیز باشد. در علوم نوین آمار و همچنین در آمار بیزی برای حل این مشکل از الگوریتم‌های نمونه‌گیری مانند مونت‌کارلوی زنجیر مارکوفی^۱ (MCMC) استفاده می‌کنند. در این مقاله نیز، ماتریس کوواریانس مدل حاشیه‌ای نیمه‌پارامتری به این روش برآورد شده است. یک مسئله قابل توجه ادغام مدل‌های حاشیه‌ای و نیمه‌پارامتری برای بررسی داده‌های طولی است که از مزایای آن می‌توان به موارد زیر اشاره کرد:

۱- برآورد انعطاف‌پذیر از توابع رگرسیونی که با مدل‌های ساده رگرسیونی امکان‌پذیر نیست.

۲- برآورد حاشیه‌ای از ماتریس کوواریانس، که در بحث داده‌های طولی دارای اهمیت زیادی است.

این مسئله با استفاده از روش برآوردیابی اسپلاین تاوانیده، توسط الخدیری و همکاران (۲۰۱۰) مورد بررسی

¹Markov Chain Monte Carlo

قرارگرفته است. همچنین زگر و دیگل (۱۹۹۴) مرجع اولیه برای رگرسیون ناپارامتری طولی حاشیه‌ای فراهم ساخته‌اند. کارل و همکاران (۲۰۰۹) بر روی تأثیر برآوردهای حاشیه‌ای نیمه‌پارامتری برای مدل خطی جمعی جزئی در داده‌های طولی کار کرده‌اند. فان و هوانگ (۲۰۰۷) تحلیل داده‌های طولی با برآورد نیمه‌پارامتری تابع کوواریانس را مورد مطالعه قرار داده‌اند. وانگ و کارل (۲۰۰۵) کارآیی برآورد حاشیه‌ای نیمه‌پارامتری بر روی داده‌های دسته‌بندی شده طولی را مطالعه کرده‌اند. ولهام و همکاران (۲۰۰۷) مقایسه‌ای از مدل‌های آمیخته با اسپلاین تاوانیده برای منحنی برازش داده‌شده انجام داده‌اند. ژائو و همکاران (۲۰۰۶) در مورد طرح‌های بیزی برای مدل‌های آمیخته‌ی خطی به بحث کرده‌اند.

در بخش ۲ مروری بر دیدگاه بیز سلسله مراتبی خواهد شد. در بخش ۳ به معرفی مدل ناپارامتری طولی حاشیه‌ای و داده‌های شبیه‌سازی شده برای آن خواهیم پرداخت. در بخش ۴ مدل طولی حاشیه‌ای ضرایب متغیر که اساس کار این پژوهش است، مورد بررسی قرار می‌گیرد. در بخش ۵ برآورد ماکسیمم درست‌نمایی برای برآورد ضرایب، مدل را مورد مطالعه قرار خواهیم داد و در بخش ۶ بایان محدودیت‌های این برآورد، مقدمه‌ای برای استفاده از استنباط بیزی را فراهم خواهیم آورد. در بخش‌های ۷ و ۸ ابتدا با استفاده از مدل ناپارامتری و داده‌های شبیه‌سازی شده، دقت برآورد بیزی مؤلفه‌های واریانس ارزیابی می‌شود، سپس مدل نیمه‌پارامتری طولی حاشیه‌ای تحلیل خواهد شد.

۲ مدل بیز سلسله مراتبی و مدل گرافیکی

مدل رگرسیون خطی ساده به صورت

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2) \quad (1)$$

را در نظر بگیرید. اگر بخواهیم اطلاعات اولیه مدل (۱) برای یک استنباط بیز سلسله مراتبی بیان شود، می‌توان آن را به صورت

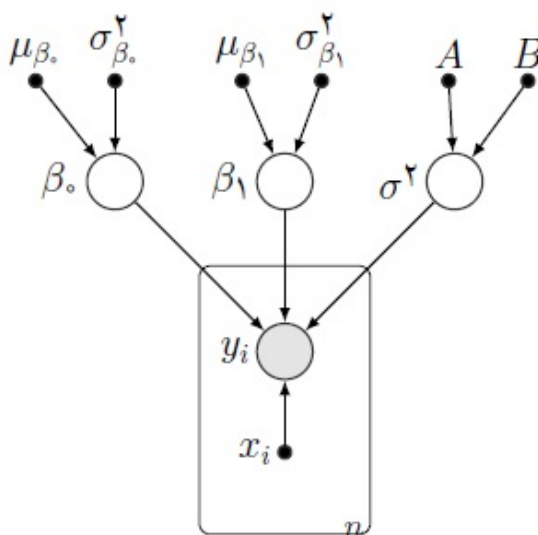
$$\begin{aligned} [y_i | \beta_0, \beta_1, \sigma^2] &\sim N(\beta_0 + \beta_1 x_i, \sigma^2) \\ [b_0] &\sim N(\mu_{\beta_0}, \sigma_{\beta_0}^2) \quad [b_1] \sim N(\mu_{\beta_1}, \sigma_{\beta_1}^2) \\ [\sigma^2] &\sim IG(A, B) \end{aligned}$$

نوشت، که در آن $\{\mu_{\beta_0}, \sigma_{\beta_0}^2, \mu_{\beta_1}, \sigma_{\beta_1}^2, A, B\}$ ابر پارامترهای مدل هستند. با افزایش پارامترهای مدل، مدل بیز سلسله مراتبی شکل پیچیده‌تری پیدا می‌کند. برای سهولت در فهم ارتباطات مدل‌های بیز سلسله مراتبی می‌توان از مدل‌های گرافیکی بهره برد. مدل‌های گرافیکی شاخه‌ای در علوم کامپیوتر و ریاضیات به حساب می‌آیند که ترکیبی از نظریه گراف و احتمال حاصل می‌شود. این مدل‌ها تجسم مدل احتمالی را ساده‌تر کرده و برای تعیین روابط بین متغیرها و بیان کارآمدی مدل، مورد استفاده قرار می‌گیرند. دو روش معمول مدل‌های گرافیکی عبارتند از:

(۱) گراف غیر مدور مستقیم^۲ (DAGs) که با شبکه بیزی نیز شناخته می‌شود.

(۲) گراف غیرمستقیم که با میدان مارکوف تصادفی نشان داده می‌شود.

برای مدل‌های سلسله مراتبی بیزی اغلب از مدل گرافیکی DAGs استفاده می‌شود، که فرم آن در شکل ۱ نشان داده شده است.



شکل ۱. مدل گرافیکی برای مدل بیزی سلسله مراتبی (۲)

²Directed acyclic graph

دایره‌های توپر نشان‌دهنده مقادیر ثابت مدل، در اینجا x_i ها و ابر پارامترها، هستند. دایره سایه زده شده مربوط به متغیر y ، نشان دهنده این موضوع است که این متغیر دارای مقادیر مشاهده شده است. در نهایت نیز دایره‌های توخالی، نشان‌دهنده پارامترهای مدل هستند که بر اساس روش‌های بیزی مورد استنباط قرار می‌گیرند. همچنین n نشان‌دهنده تعداد مشاهدات برای این مدل است. در ادامه از این مدل‌ها برای روشن نمایی مدل‌های بیز سلسله مراتبی در مدل‌های طولی حاشیه‌ای ناپارامتری و نیمه پارامتری استفاده خواهیم کرد.

۳ مدل ناپارامتری طولی حاشیه‌ای

برای $m \leq i \leq n$ آزمودنی، $1 \leq j \leq n$ ($n \ll m$) مشاهده متناظر با متغیر پاسخ y_{ij} و متغیر پیش‌بین x_{ij} را در نظر بگیرید. همچنین فرض کنید y_i و x_i به ترتیب بردار پاسخ و بردار پیش‌بین‌ها برای آزمودنی i ام است. یک رگرسیون ناپارامتری طولی حاشیه‌ای^۳ به صورت

$$E(Y_{ij}) = f(x_{ij}), \quad Cov\{\mathbf{Y}_i | f(\mathbf{x}_i)\} = \Sigma, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n \quad (۲)$$

تعریف می‌شود، که در آن f تابع هموار و Σ ماتریس کوواریانس است. شکل ۲ نشان‌دهنده داده‌های شبیه‌سازی شده برای مدل (۲)، با $m = 100$ ، $n = 2$ و

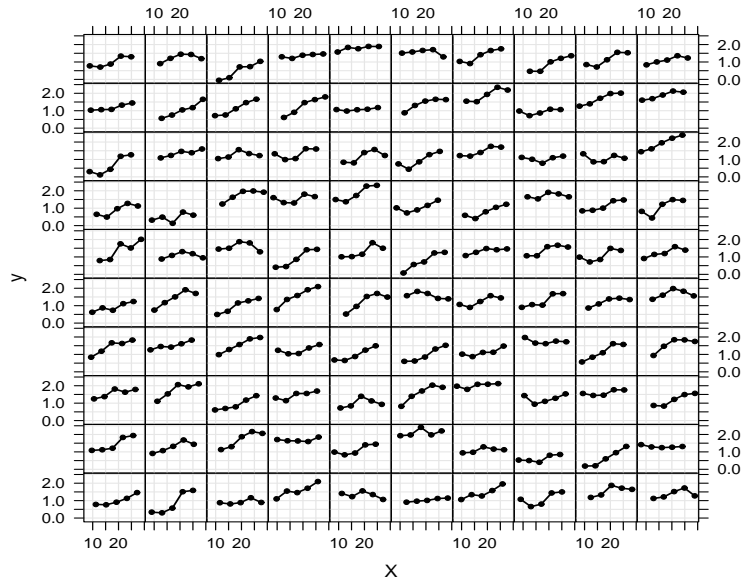
$$f(x) = 1 + \frac{1}{4} \phi((2x - 36)/5) \quad (۳)$$

و ماتریس کوواریانس

$$\Sigma = \begin{pmatrix} 0/122 & 0/098 & 0/073 & 0/063 & 0/05 \\ 0/098 & 0/122 & 0/098 & 0/073 & 0/063 \\ 0/073 & 0/098 & 0/122 & 0/098 & 0/073 \\ 0/063 & 0/078 & 0/098 & 0/122 & 0/098 \\ 0/05 & 0/063 & 0/078 & 0/098 & 0/122 \end{pmatrix} \quad (۴)$$

است، که در آن ϕ تابع چگالی نرمال استاندارد است.

^۳Marginal longitudinal nonparametric regression



شکل ۲. داده‌های شبیه‌سازی شده برای مدل (۲) با در نظر گرفتن فرم تابعی (۳) و ماتریس کوواریانس (۴)

مسئله اصلی پیدا کردن برآوردی برای f و در نهایت برای Σ است. یک روش مناسب برای برآورد f ، استفاده از اسپلاین به صورت

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^K u_k Z_k(x) \quad (5)$$

است، که در آن Z_1, \dots, Z_k یک مجموعه مناسب از توابع پایه‌ای اسپلاین هستند. یک پایه ساده با قرار دادن

$$Z_k(x) = (x - c_k)_+ = \begin{cases} x - c_k & \text{اگر } x - c_k > 0 \text{ یا } x > c_k \\ 0 & \text{اگر } x - c_k < 0 \text{ یا } x < c_k \end{cases}$$

حاصل می‌شود، به طوری که c_1, \dots, c_k یک نمایش از مکان گره‌ها روی دامنه x_i ها هستند که در بیشتر اوقات $k = 25$ برای مدل کفایت می‌کند (واند و اورمروود، ۲۰۰۸). برای پرهیز از بیش برآزشی احتیاج به جریمه ضرایب مدل اسپلاین داریم. یک جریمه مناسب در نظر گرفتن ضرایب اسپلاین، یعنی u_k ها،

به‌عنوان متغیرهای تصادفی از یک توزیع نرمال با میانگین صفر و واریانس σ^2 است. حال با اعمال این جریمه بر روی معادله (۵)، یک نمایش ماتریسی از مدل (۲) را می‌توان به صورت

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon} \quad (۶)$$

نشان داد، که در آن

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_m \end{pmatrix}, \quad \mathbf{Z} = \begin{pmatrix} Z_1(x_1) & \dots & Z_K(x_1) \\ \vdots & \ddots & \vdots \\ Z_1(x_m) & \dots & Z_K(x_m) \end{pmatrix}$$

$$\mathbf{u} = \begin{pmatrix} u_1 \\ \vdots \\ u_K \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_1 \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_m \end{pmatrix} \quad (۷)$$

بردارهای تصادفی در سمت راست معادله (۶) دارای میانگین صفر و ماتریس کوواریانس زیر است:

$$\text{Cov} \begin{pmatrix} \mathbf{u} \\ \epsilon_1 \\ \vdots \\ \epsilon_m \end{pmatrix} = \begin{pmatrix} \sigma^2 I & \circ & \circ & \dots & \circ \\ \circ & \Sigma & \circ & \dots & \circ \\ \circ & \circ & \Sigma & \dots & \circ \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \circ & \circ & \circ & \dots & \Sigma \end{pmatrix} = \begin{pmatrix} \sigma^2 I & \circ \\ \circ & I_m \otimes \Sigma \end{pmatrix}$$

که در آن نماد \otimes نشان دهنده ضرب کرونگر است، همچنین اگر σ^2 و Σ معلوم در نظر گرفته شوند، می‌توان از روش بهترین پیش‌بینی نااریب خطی^۴ (BLUP) برای برآورد $\boldsymbol{\beta}$ و \mathbf{u} استفاده کرد. در عمل σ^2 و Σ نیز هر دو باید برآورد شوند و فرض مناسب برای رسیدن به این هدف این است که

$$\begin{pmatrix} \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} \sim N \left(\begin{pmatrix} \circ \\ \circ \end{pmatrix}, \begin{pmatrix} \sigma^2 I & \circ \\ \circ & I_m \otimes \Sigma \end{pmatrix} \right)$$

^۴Best linear unbiased prediction

۴ مدل ضرایب متغیر طولی حاشیه‌ای

یک نوع از رگرسیون نیمه‌پارامتری با پیش‌بین‌های چندگانه یک مدل با ضرایب گوناگون^۵ است. فرض کنید برای هر y_{ij} پیش‌بین‌های جداگانه‌ای موجود باشد. در این قسمت با حفظ نمادهای بخش قبل، مدل را به دو متغیر کمی برای j امین تکرار از آزمودنی i ام که به صورت x_{ij} و s_{ij} نمایش داده می‌شود، محدود می‌کنیم. مدل ضرایب متغیر طولی حاشیه‌ای برای داده‌های تعریف‌شده در بخش قبل به صورت

$$E(Y_{ij}) = f_0(s_{ij}) + f_1(s_{ij})x_{ij}, \quad Cov\{\mathbf{Y}_i | f_0(\mathbf{s}_i), f_1(\mathbf{s}_i)\} = \Sigma \quad (۸)$$

تعریف می‌شود، که در آن f_0 و f_1 توابع هموار هستند. همان‌طور که از ساختار مدل مشخص است، این مدل به‌گونه‌ای است که برای هر مقدار ثابت s ، یک ارتباط خطی بین y و x بیان می‌کند. حال اگر توابع f_0 و f_1 با روش اسپلاین تخمین زده شود، مدل (۸) به صورت

$$y_{ij} = \alpha_0 + \alpha_1 s_{ij} + \sum_{k=1}^{K_1} u_{1k} (s_{ij} - c_k)_+ + \{\beta_0 + \beta_1 s_{ij} + \sum_{k=1}^{K_2} u_{2k} (s_{ij} - c_k)_+\} x_{ij} + \epsilon_{ij} \quad (۹)$$

حاصل می‌شود. اگر u_1 و u_2 متغیرهای تصادفی با توزیع‌های

$$u_{1k} \stackrel{iid}{\sim} N(0, \sigma_1^2), \quad u_{2k} \stackrel{iid}{\sim} N(0, \sigma_2^2)$$

در نظر گرفته شوند، می‌توان مدل (۹) را به صورت مدل آمیخته

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$$

در نظر گرفت، که در آن

$$\mathbf{u} = [u_{11}, \dots, u_{1K_1}, u_{21}, \dots, u_{2K_2}]^T, \quad \mathbf{X} = [1, \mathbf{s}, \mathbf{x}, \mathbf{s}\mathbf{x}], \quad \boldsymbol{\beta} = [\alpha_0, \alpha_1, \beta_0, \beta_1]^T$$

⁵Varying coefficient models

$$\mathbf{Z} = [(s_{ij} - c_k)_{\substack{1 \leq i \leq m, \\ 1 \leq k \leq K_1}} + (s_{ij} - c_k)_{\substack{1 \leq i \leq m, \\ 1 \leq k \leq K_2}} x_{ij}]_{\substack{1 \leq i \leq m, \\ 1 \leq j \leq n}}$$

ماتریس کوواریانس برای ضرایب اسپلاین و خطای مدل به صورت

$$\begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \boldsymbol{\epsilon} \end{pmatrix} \sim N \left(\begin{pmatrix} \circ \\ \circ \\ \circ \end{pmatrix}, \begin{pmatrix} \sigma_1^2 I & \circ & \circ \\ \circ & \sigma_2^2 I & \circ \\ \circ & \circ & I_m \otimes \boldsymbol{\Sigma} \end{pmatrix} \right)$$

تعریف می‌شود، که در آن \mathbf{u}_1 یک بردار $1 \times K_1$ بعدی شامل u_{1k} ها است و \mathbf{u}_2 نیز به صورت مشابه تعریف می‌شود.

۵ برآورد ماکسیمم درست‌نمایی و بهترین پیش‌بینی

در مبحث مدل‌های حاشیه‌ای مسئله اصلی برآورد مؤلفه‌های واریانس است که اغلب به روش ماکسیمم درست‌نمایی مقید انجام می‌شود. در حالت کلی تعمیم مدل‌های (۲) و (۸) به مدل‌هایی با d تابع هموارساز را می‌توان به مدل آمیخته گاوسی به صورت

$$\mathbf{y} | \mathbf{u} \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, I_m \otimes \boldsymbol{\Sigma}), \quad \mathbf{u} \sim N(\circ, \text{blockdiag}(\sigma_r^2 I_{k_r}))_{\substack{1 \leq r \leq d}} \quad (10)$$

تبدیل کرد، که در آن K_r متناظر با تعداد توابع پایه اسپلاین مورد کاربرد در r امین تابع هموارساز برآورد شده است. فرض کنید $\boldsymbol{\sigma}^2 = (\sigma_1^2, \dots, \sigma_d^2)$ بردار پارامترهای واریانس باشد. لگاریتم درست‌نمایی مدل (۱۰) به صورت

$$\ell(\boldsymbol{\beta}, \boldsymbol{\sigma}^2, \boldsymbol{\Sigma}) = -\frac{1}{2} \{n \log(2\pi) + \log |V| + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T V^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\}$$

حاصل می‌شود، به طوری که

$$V = \text{Cov}(\mathbf{Y}) = \sum_{r=1}^d \sigma_r^2 \mathbf{Z}_{[r]} \mathbf{Z}_{[r]}^T + I_m \otimes \boldsymbol{\Sigma}$$

که در آن $Z_{[1]}, \dots, Z_{[d]}$ افزایشی از Z هستند که متناظر با توابع پایه برای هر تابع هموارساز است. با توجه به برآوردیابی در مدل‌های آمیخته برای مقادیر ثابت σ^2 و Σ ، برآورد ضرایب ثابت به صورت

$$\tilde{\beta}(\sigma^2, \Sigma) = (X^T V^{-1} X)^{-1} X^T V^{-1} y$$

حاصل می‌شود. بر اساس این برآورد لگاریتم درست‌نمایی نیم‌رخ به صورت

$$\ell_p(\sigma^2, \Sigma) = -\frac{1}{2} \left\{ \log |V| + y^T V^{-1} \{ I - X(X^T V^{-1} X)^{-1} X^T V^{-1} \} y \right\} - \frac{n}{2} \log(2\pi)$$

خواهد شد. البته در عمل اغلب لگاریتم درست‌نمایی مقید به صورت

$$\ell_R(\sigma^2, \Sigma) = \ell_p(\sigma^2, \Sigma) - \frac{1}{2} \log |X^T V^{-1} X|$$

مورد استفاده قرار می‌گیرد، زیرا برآورد حاصل از روش ماکسیمم درست‌نمایی مقید (REML) این نکته که β ها برآورد شده‌اند را نیز در نظر می‌گیرد (روپرت و همکاران، ۲۰۰۳). پس مسئله برآوردیابی برای پارامترهای مدل، یعنی β و σ^2 و Σ به شرح زیر است:

۱- به دست آوردن برآورد REML برای σ^2 و Σ با استفاده از بیشینه کردن ℓ_R .

۲- به دست آوردن برآورد ماکسیمم درست‌نمایی (MLE) برای β به صورت $\hat{\beta} = \tilde{\beta}(\hat{\sigma}^2, \hat{\Sigma})$

و در نهایت مسئله برآورد ضرایب اسپلاین، یعنی u ، مطرح است. چون u یک بردار تصادفی است، نمی‌توان از برآوردهایی مانند ماکسیمم درست‌نمایی استفاده کرد و معمولاً روش بهترین پیش‌بینی مورد استفاده قرار می‌گیرد. بنا بر مدل (۱۰) این برآورد به صورت

$$\tilde{u}(\sigma^2, \Sigma) \equiv E(Y|u) \equiv G_{\sigma^2} Z^T V(\sigma^2, \Sigma)^{-1} (y - X \tilde{\beta}(\sigma^2, \Sigma))$$

است، به طوری که $G_{\sigma^2} = \text{blockdiag}(\sigma_r^2 I_{k_r})_{1 \leq r \leq d}$ از آنجایی که مقادیر β و σ^2 و Σ نامعلوم هستند و باید مورد برآورد قرار گیرند، یک برآورد مناسب برای \tilde{u} ، بهترین پیش‌بینی تجربی به صورت

$$\hat{u} = G_{\hat{\sigma}^2} Z^T V(\hat{\sigma}^2, \hat{\Sigma})^{-1} (y - X \tilde{\beta}(\hat{\sigma}^2, \hat{\Sigma}))$$

است. حال با استفاده از این برآوردها می‌توان توابع رگرسیونی f را به‌سادگی به دست آورد و برازش رگرسیونی را انجام داد.

ممکن است این تصور پیش آید که شاید بتوان مدل (۱۰) را به‌عنوان یک مدل آمیخته در نظر گرفت و آن را با توابعی همانند $\text{lme}()$ در نرم‌افزار R برآورد کرد که در جواب باید گفت این موضوع با محدودیت‌های تابع $\text{lme}()$ قابل انجام نیست. بنابراین با تمامی محدودیت‌های موجود در نظر گرفتن استنباط بیزی و اجرای آن توسط فرایند گیبز شاید اجتناب‌ناپذیر باشد. اگرچه برآورد نقطه‌ای دارای مفهومی ساده‌تر و فرایند انجام آن سریع‌تر است، اما باید به این نکته توجه کرد که این برآورد به مدل‌های آمیخته گاوسی وابستگی شدید داشته و در مواردی که با مدل‌های ناگوسی روبرو هستیم، از لحاظ محاسباتی بسیار مشکل خواهد بود. همچنین برای به دست آوردن توزیع برآوردهای REML، جهت ساختن بازه اطمینان با توجه به قانون اعداد بزرگ نیازمند مشاهدات زیادی هستیم. این محدودیت سبب پایین آمدن دقت برآورد و نامعتبر بودن بازه اطمینان می‌گردد. از طرف دیگر روش MCMC به توزیع گاوسی وابسته نیست و با افزایش تعداد تکرار الگوریتم، دقت برآورد توزیع پسین و استنباط‌های مورد نظر بالاتر رفته و بازه اطمینان نیز قابل‌اعتمادتر می‌شوند (وست و همکاران، ۲۰۰۶)

۶ استنباط بیزی

در این قسمت سعی بر آن است که با استفاده از برآورد بیزی ماتریس کوواریانس مدل حاشیه‌ای را به دست آورده و استنباط‌های موردنظر را انجام دهیم. بر این اساس برای پارامترهای مدل آمیخته (۱۰) توزیع‌های پیشین

$$\beta \sim N(\mathbf{0}, \mathbf{F}), \quad \sigma_r^2 \sim IG(A_r, B_r), \quad \Sigma \sim IW(a, \mathbf{B})$$

پیشنهاد می‌شود، که در آن A_r و B_r ($1 \leq r \leq d$) مقادیر ثابت و مثبت، \mathbf{F} و \mathbf{B} ماتریس‌های معین مثبت ابرپارامتر و IG و IW بترتیب نمادهای توزیع‌های گامای وارون و ویشارت وارون با توابع چگالی

$$[\sigma_r^2] = \frac{B_r^{A_r}}{\Gamma(A_r)} (\sigma_r^2)^{-A_r-1} e^{-B_r/\sigma_r^2}$$

$$[\Sigma] = C_{n,a}^{-1} |B|^{a/2} |\Sigma|^{-(a+n+1)/2} \exp\{-\frac{1}{2}tr(B\Sigma^{-1})\}, \quad a > 0$$

هستند، به طوری که $C_{n,a} = 2^{an/2} \pi^{n(n-1)/4} \prod_{i=1}^n \Gamma(\frac{a+1-i}{2})$ برای انجام استنباط بیزی نیاز به تعیین توزیع‌های پسین

$$[\beta|y], \quad [u|y], \quad [\Sigma|y], \quad [\sigma_r^2|y]$$

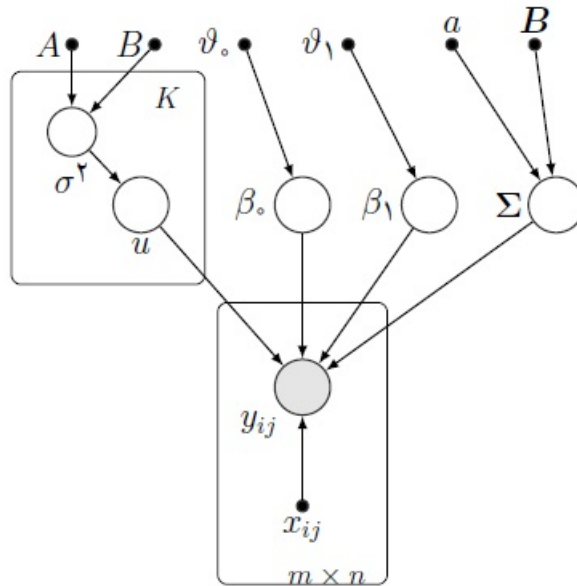
است که به دست آوردن آن‌ها مشکل و حتی شاید غیرممکن است. معمولاً این توزیع‌ها را می‌توان با روش‌های محاسباتی تصادفی و گاهی غیر تصادفی تخمین زد (مینستاسا و واند، ۲۰۱۳). در این بین روش MCMC یکی از روش‌های معمول تصادفی محاسباتی است که در ادامه با نرم‌افزارهای OpenBUGS و R انجام داده خواهد شد. برای انجام محاسبات بیزی انتخاب ابر پارامترها که تأثیر زیادی بر توزیع پسین خواهد گذاشت، بسیار مهم است. اگر در مورد توزیع پیشین اطلاعاتی در دسترس باشد، می‌توان ابر پارامترها را بر مبنای این اطلاعات انتخاب نمود. اما اغلب این اطلاعات موجود نیست، بنابراین انتخاب ابر پارامترها به گونه‌ای است که منجر به توزیع‌های پیشین ناآگاهی‌بخش شوند. بر این اساس ابر پارامترها به صورت زیر پیشنهاد می‌شوند:

$$F = 1 \times I, \quad A_r = B_r = 0 \times 1, \quad a = n, \quad B = 0 \times 1 I_n$$

حال برای استفاده از الگوریتم نمونه‌گیری گیبز به توزیع‌های تمام شرطی پارامترها نیاز است که با توجه به توزیع‌های پیشین، به صورت زیر به دست می‌آیند:

$$\begin{aligned} \Sigma|rest &\sim IW(a+m, B+(y-X\beta-Zu)(y-X\beta-Zu)^T) \\ \begin{bmatrix} \beta \\ u \end{bmatrix}|rest &\sim N((C(I_m \otimes \Sigma^{-1})C + G_{\sigma^2})^{-1}C^T \Sigma^{-1}y, (C(I_m \otimes \Sigma^{-1})C + G_{\sigma^2})^{-1}) \\ \sigma_r^2|rest &\sim IG(A_r + \frac{1}{2}K_r, B_r + \frac{1}{2}\|u_r\|^2) \end{aligned}$$

که در آن $C = [X \ Z]$



شکل ۳. مدل گرافیکی برای مدل بیزی سلسله مراتبی (۱۱)

۷ مدل ناپارامتری برای داده‌های شبیه‌سازی شده

در این بخش با استفاده از اطلاعات بخش‌های گذشته و داده‌های شبیه‌سازی شده در بخش ۳ دقت برآورد بیزی مؤلفه‌های واریانس مدل حاشیه‌ای، سنجیده می‌شود. برای بیان مدل ناپارامتری (۲) به صورت یک مدل بیزی سلسله مراتبی، داریم

$$\begin{aligned}
 [y_i | \beta, \sigma^2, u, \Sigma] &\sim N(f(x_i), \Sigma) \\
 [\beta_0] &\sim N(\circ, \vartheta_0) , \quad [\beta_1] \sim N(\circ, \vartheta_1) , \quad \Sigma \sim IW(a, B) \\
 u_1, \dots, u_K &\stackrel{iid}{\sim} N(\circ, \sigma^2) , \quad [\sigma^2] \sim IG(A, B)
 \end{aligned} \tag{۱۱}$$

همان‌طور که در بخش ۲ بیان شد، هدف استنباط بیزی تخمین توزیع پسین برای پارامترهای مدل است. در اینجا چون با مدل‌های حاشیه‌ای روبرو هستیم، تمرکز ما بر روی برآورد Σ است. در اینجا از داده‌های شبیه‌سازی شده‌ی شکل ۲، که در آن مقادیر x_{ij} با فواصل یکسان تولید شدند با این ویژگی که

اولین مقدار برای هر آزمودنی، برای مثال x_{i1} ، از توزیع یکنواخت در فاصله‌ی (۸, ۱۲) به‌طور تصادفی انتخاب شده است، استفاده می‌کنیم. حال الگوریتم نمونه‌گیری گیبز را برای ۵۰۰۰ تکرار، با میزان دورریز برابر ۵۰۰ و با میزان معیار تقلیل برابر ۵، برای این مدل ناپارامتری انجام می‌دهیم. برای این کار از بسته BRugs در نرم‌افزار R، که پل ارتباطی بین نرم‌افزارهای OpenBUGS و R است، استفاده شده است. خروجی‌های حاصل در شکل‌های ۴ و ۶ سطح مطلوبی از برآورد بیزی را نشان می‌دهد. همان‌طور که ملاحظه می‌شود این خروجی‌ها شامل ۵ قسمت اصلی به شرح زیر هستند:

۱- *trace*: این قسمت نمودار تکرار الگوریتم در مقابل مقادیر حاصل از فرایند گیبز را نشان می‌دهد. نبودن جهش در داده‌ها در طول تکرار الگوریتم می‌تواند به‌منزله همگرایی داده‌ها به یک توزیع خاص باشد.

۲- *lag1*: این نمودار به بررسی تصادفی بودن داده‌ها می‌پردازد. برای داده‌های تصادفی نمودار پراکندگی داده‌ها در مقابل تکرار نباید از روند خاصی پیروی کند.

۳- *acf*: مقدار خودهمبستگی بین داده‌ها را نشان می‌دهد. در عمل خودهمبستگی کمتر از ۰/۱ برای اولین *lag* قانع‌کننده است (اولین *lag* منظور داده‌هایی که برای مبنای پارامتر تقلیل مشخص شده، به‌دست آمده‌اند). خط‌چین‌ها در نمودار نشان دهنده خط‌های ۰/۱ و ۰/۱- است.

۴- *density*: این قسمت نمودار نشان دهنده شکل تابع توزیع پسین است که با استفاده از نمونه تولید شده از توزیع پسین در فرایند گیبز حاصل شده است.

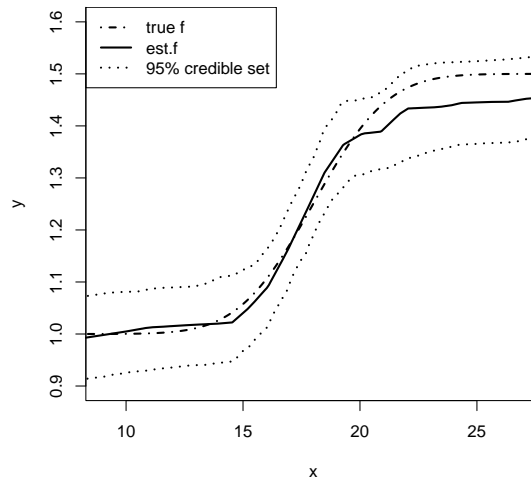
۵- *summary*: این قسمت از دو بخش میانگین پسین و بازه مؤثق^۶ تشکیل شده است. میانگین پسین، میانگین نمونه‌ای از مقادیر تولید شده توسط فرایند گیبز است. بازه مؤثق نیز یک بازه اطمینان برآورد شده برای پارامتر موردبررسی است که به‌صورت نمونه‌ای حاصل شده است. برای مثال حد پایین یک بازه مؤثق ۹۵ درصدی، چندک ۰/۲۵ داده‌ها و حد بالای آن چندک ۰/۹۷۵ داده‌ها است.

نمودار *trace* همگرایی داده‌ها به یک توزیع مانا را تأیید می‌کند و نمودارهای *lag1* و *acf* عدم وجود همبستگی بین داده‌ها را نشان می‌دهند. نمودار *density* نشان‌دهنده این موضوع است که بیشترین احتمال حول مقدار واقعی پارامتر موردنظر قرار گرفته است (خط بریده عمودی نشان‌دهنده مقدار دقیق

^۶Credible interval

parameter	trace	lag 1	acf	density	summary
Σ_{11}					p.m: 0.12 95% CI: (0.0886,0.158)
Σ_{22}					p.m: 0.116 95% CI: (0.0876,0.153)
Σ_{33}					p.m: 0.117 95% CI: (0.0899,0.156)
Σ_{44}					p.m: 0.114 95% CI: (0.0864,0.15)
Σ_{55}					p.m: 0.127 95% CI: (0.098,0.165)

شکل ۴. جزئیات استنباط بیزی از مقادیر قطر اصلی ماتریس Σ برای مدل ناپارامتری (۱۱) (مقادیر درست پارامترها با خط‌های بریده افقی نمایش داده شده است)



شکل ۵. نمودار f (خط ممتد)، برآورد تابع با الگوریتم گیبز (خط چین) و بازه مؤثق (نقطه چین)

پارامتر است). همچنین در قسمت *summary* می‌بینیم که میانگین نمونه‌ای به مقدار اصلی پارامتر نزدیک بود و بازه مؤثقت نیز این مقدار را در برمی‌گیرد. شکل ۵ نیز نشان دهنده مقایسه برآورد تابع f (از روش MCMC با مقدار واقعی آن است. همان‌طور که ملاحظه می‌شود، فرایند گیبز تا حد قابل قبولی این برآورد را درست انجام داده است. همچنین در این شکل بازه مؤثقت ۹۵ درصد برای این برآورد نیز به نمایش درآمده است. به‌طور کلی می‌توان گفت که استنباط بیزی عملکرد قابل قبولی را نشان می‌دهد (لون و همکاران، ۲۰۰۰).

parameter	trace	lag 1	acf	density	summary
Σ_{12}					p.m: 0.0941 95% CI: (0.0668,0.127)
Σ_{13}					p.m: 0.0668 95% CI: (0.0426,0.0983)
Σ_{14}					p.m: 0.0562 95% CI: (0.0323,0.0862)
Σ_{15}					p.m: 0.0422 95% CI: (0.0193,0.0702)
Σ_{23}					p.m: 0.0915 95% CI: (0.0671,0.128)
Σ_{24}					p.m: 0.0736 95% CI: (0.0505,0.106)
Σ_{25}					p.m: 0.0585 95% CI: (0.0342,0.0895)
Σ_{34}					p.m: 0.0945 95% CI: (0.0699,0.131)
Σ_{35}					p.m: 0.0788 95% CI: (0.0552,0.111)
Σ_{45}					p.m: 0.0958 95% CI: (0.0688,0.13)

شکل ۶. جزئیات استنباط بیزی از مقادیر کوواریانس ماتریس Σ برای مدل ناپارامتری (۱۱) (مقادیر درست پارامترها با خط‌های بریده افقی نمایش داده شده است)

۸ مدل ضرایب متغیر برای داده‌های همه گیرشناسی

در این بخش رگرسیون ضرایب متغیر (۸) نمونه‌ای از داده‌های اپیدمیولوژی تغذیه برازش داده می‌شود. این داده‌ها بر اساس گزارش ۲۴ ساعته بر روی نشانگرهای انرژی و پروتئین (نیتروژن ادرار)، همراه با پرسشنامه مصرف مواد غذایی در ۲۴ ساعت انجام شده است. متغیرهای موجود در این مدل به صورت زیر تعریف می‌شوند:

y : لگاریتم مصرف پروتئین اندازه‌گیری شده با نشانگر نیتروژن ادرار

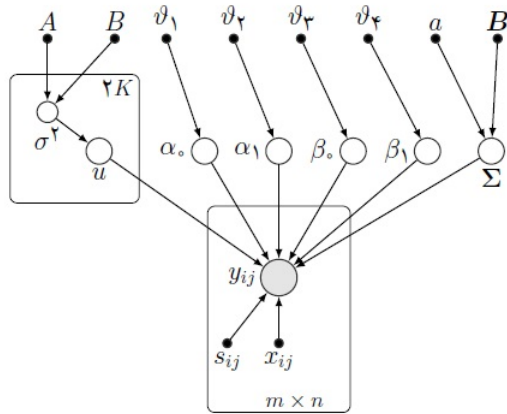
s : شاخص توده‌ی بدن

x : لگاریتم مصرف پروتئین اندازه‌گیری شده در گزارش ۲۴ ساعته

در اینجا $m = 294$ ، $n = 2$ و تعداد گره‌ها برابر با ۳۰ اختیار شده است. در این قسمت نیز با استفاده از روش نمونه‌گیری گیبز برآوردی برای Σ به دست آورده می‌شود. از این رو مدل بیز سلسله مراتبی برای مدل (۸) را می‌توان به صورت

$$\begin{aligned} [y_i | \beta, \gamma, \sigma^2, u, \Sigma] &\sim N(f_0(s_i) + f_1(s_i)x_i, \Sigma) \\ [\beta_0] &\sim N(0, \vartheta_2) \quad , \quad [\beta_1] \sim N(0, \vartheta_3) \\ [\alpha_0] &\sim N(0, \vartheta_1) \quad , \quad [\alpha_1] \sim N(0, \vartheta_4) \\ u_{11}, \dots, u_{1K}, u_{2K} &\overset{iid}{\sim} N(0, \sigma^2) \\ \Sigma &\sim IW(a, B) \quad , \quad [\sigma^2] \sim IG(A, B) \end{aligned} \quad (12)$$

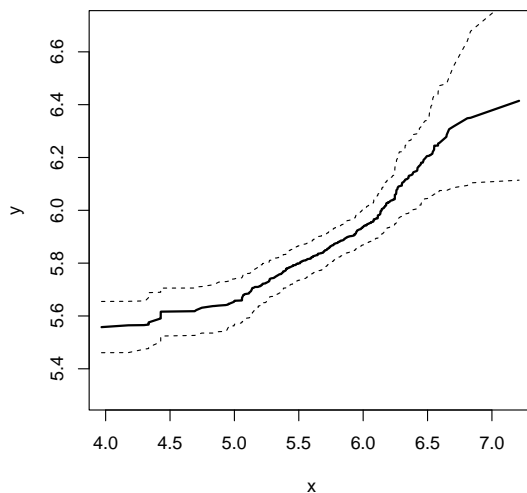
در نظر گرفت. همچنین مدل گرافیکی این روابط در شکل ۷ آمده است. باید دقت شود که در این مدل f_0 و f_1 مانند معادله (۹) تعریف شده است. همچنین تعداد گره‌ها به صورت $K_1 = K_2 = 30$ در نظر گرفته شده و مکان گره‌ها نیز در فواصل یکسان از هم انتخاب شده‌اند. برای این مدل نیز الگوریتم نمونه‌گیری گیبز برای ۵۰۰۰ تکرار، با میزان دورریز ۵۰۰۰ و با میزان معیار تقلیل ۵ انجام شده است. در شکل ۸ جزئیات استنباط بیزی انجام شده توسط فرایند گیبز نمایش داده شده است. همان‌طور که ملاحظه می‌شود نمودار $trace$ همگرایی مطلوبی را نشان می‌دهد و بر مبنای نمودارهای acf و $lag1$ ، مقدار همبستگی بین نمونه تولید شده توسط فرایند گیبز، ناچیز است. همچنین توزیع مؤلفه‌های واریانس، توزیع تقریباً متقارن برآورد شده است. لازم به ذکر است از آنجایی که داده‌ها در دو زمان تکرار شده‌اند ماتریس Σ ، یک ماتریس 2×2 است.



شکل ۷. مدل گرافیکی برای مدل بیزی سلسله مراتبی (۱۲)

parameter	trace	lag 1	acf	density	summary
Σ_{11}					p.m: 0.0729 95% CI: (0.0616,0.0859)
Σ_{22}					p.m: 0.0828 95% CI: (0.0703,0.0971)
Σ_{12}					p.m: 0.045 95% CI: (0.0352,0.0566)

شکل ۸. استنباط بیزی برای مؤلفه‌های ماتریس Σ ، مدل ضرایب متغیر (۱۲)



شکل ۹. برآورد بیزی تغییرات متغیر پاسخ y در مقابل متغیر پیش‌بین x

همچنین در شکل ۹ برآورد بیزی تغییرات متغیر پاسخ y در مقابل متغیر پیش‌بین x (حاصل از فرایند گیبز) به نمایش درآمده است. همچنین خط‌چین‌ها در این شکل نشان‌دهنده بازه مؤثق ۹۵ درصدی برای این برآورد است.

۱۰.۸ بررسی همگرایی مدل با آماره گلمن-روبین

از مسائل مهم در برآورد یابی به روش MCMC، بررسی همگرایی فرایند گیبز به یک توزیع مانا یا توزیع پسین پارامتر موردعلاقه است. روش‌های گوناگونی برای بررسی همگرایی وجود دارد که همگی در مجموعه روش‌های تعقیبی^۷ قرار دارند. یکی از راه‌های تشخیص همگرایی، روش گلمن-روبین^۸ است، که آماره آن طی گام‌های زیر به دست می‌آید:

۱- اجرا کردن فرایند گیبز برای $m \geq 2$ زنجیره با طول $2n$ و با مقادیر اولیه متفاوت.

۲- کنار گذاشتن n نمونه اول از هر زنجیره.

^۷Post hoc

^۸Gelman-Rubin

۳- به دست آوردن واریانس‌های بین و درون زنجیر. واریانس درون زنجیر برای پارامتر θ به صورت زیر حاصل می‌شود:

$$W = \frac{1}{m} \sum_{j=1}^m s_j^2 \quad \text{به طوری که} \quad s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\theta_{ij} - \bar{\theta}_j)^2$$

اگر داده‌های به دست آمده توسط فرایند گیبز، به یک توزیع مانا همگرا نشده باشند، مقدار W ، واریانس پارامتر را درست برآورد نمی‌کند. همچنین واریانس بین گروهی برای زنجیره‌های متفاوت به صورت زیر حاصل می‌شود:

$$B = \frac{n}{m-1} \sum_{j=1}^m (\bar{\theta}_j - \bar{\theta})^2 \quad \text{به طوری که} \quad \bar{\theta} = \frac{1}{m} \sum_{j=1}^m \bar{\theta}_j \quad (13)$$

۴- برآورد واریانس پارامتر مورد نظر به وسیله واریانس بین و درون زنجیره به صورت زیر به دست می‌آید:

$$\widehat{Var}(\theta) = \left(1 - \frac{1}{n}\right)W + \frac{1}{n}B \quad (14)$$

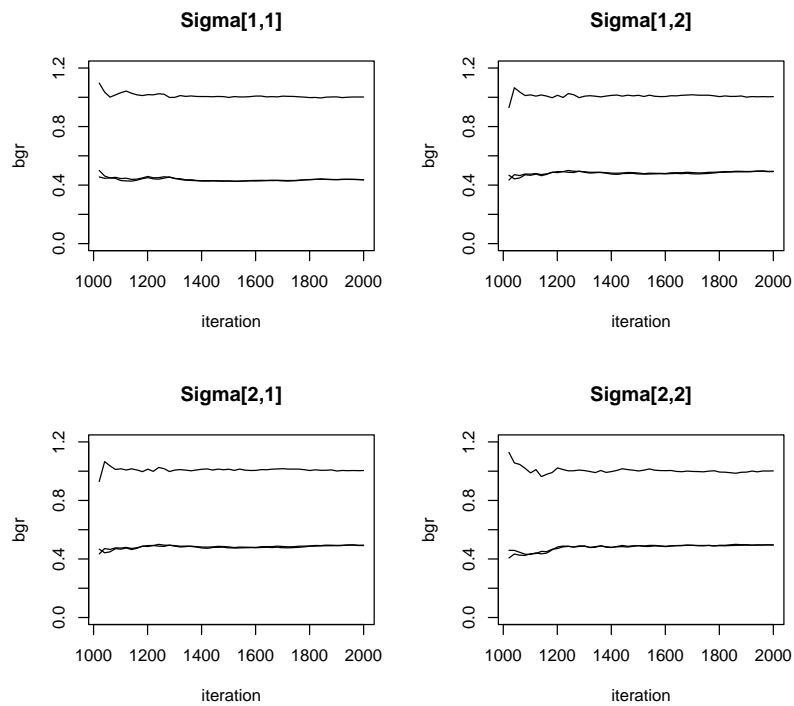
۵- به دست آوردن ضریب کاهش بالقوه^۹ به صورت زیر:

$$R = \frac{\widehat{Var}(\theta)}{W} \quad (15)$$

حال اگر تمامی زنجیره‌های تولید شده به توزیع پسین همگرا شده باشند، تمام زنجیره‌ها هم‌پوشانی داشته و $B = W$ است و $R = 1$ به دست می‌آید. ولی اگر این مقدار بیشتر از ۱/۱ یا ۱/۲ باشد، باید تعداد تکرارها را تا زمانی که مقدار R به سطح مطلوب (حول ۱، با اختلاف کمتر از ۱/۱۰ تا حداکثر ۲/۱۰) برسد افزایش می‌دهیم. در شکل ۱۰ نمودار حاصل از این آماره، برای مؤلفه‌های ماتریس Σ به تصویر کشیده شده است. این نمودار برحسب سه زنجیره که هر کدام با میزان دورریز ۵۰۰۰ و تکرار ۵۰۰۰ و با تقلیل ۵ اجرا شده‌اند، حاصل شده است. در هر یک از این نمودارها دو خط افقی پایین نشان دهنده صورت و مخرج

^۹Potential scale reduction factor

معادله (۱۵) هستند. دو خط پایین به علت قرار گرفتن آن‌ها روی هم در اکثر نقاط، شاید به خوبی قابل مشاهده نباشند). مقدار R یا آماره گلمن-روبین نیز در خط بالای نمودار مشخص شده است که از همگرایی مدل به یک توزیع مانا حمایت می‌کند (کراسچک، ۲۰۱۴).



شکل ۱۰. بررسی آماره گلمن-روبین برای مؤلفه‌های ماتریس Σ ، دو خط پایین صورت و مخرج معادله (۱۵) هستند.

همچنین در جدول ۱ نیز مقادیر حاصل شده برای آماره گلمن-روبین گزارش شده است. همان‌طور که ملاحظه می‌شود، تکرارها به 5° بلوک تقسیم شده و برای هر بلوک واریانس‌های بین و درون زنجیر به دست آمده و در نهایت مقدار R طبق معادله (۱۵) حاصل شده است.

جدول ۰۱. مقادیر آماره گلمن-روبین برای مؤلفه‌های ماتریس Σ

بلوک	تکرار	$R_{\Sigma_{11}}$	$R_{\Sigma_{12}}$	$R_{\Sigma_{22}}$
۱	۱۰۲۰-۱۰۰۱	۰۹۷۱/۱	۹۸۹۸/۰	۱۲۸۹/۱
۲	۱۰۴۰-۱۰۲۱	۰۳۵۰/۱	۰۶۵۲/۱	۰۵۶۶/۱
۳	۱۰۶۰-۱۰۴۱	۰۰۱۲/۱	۰۳۶۲/۱	۰۴۶۰/۱
۴	۱۰۸۰-۱۰۶۱	۰۱۶۷/۱	۰۱۱۷/۱	۰۱۹۴/۱
۵	۱۲۰۰-۱۰۸۱	۰۴۲۴/۱	۰۱۶۵/۱	۹۸۸۱/۰
۶	۱۲۲۰-۱۲۰۱	۰۱۶۵/۱	۰۰۸۱/۱	۹۶۳۰/۰
⋮	⋮	⋮	⋮	⋮
۴۴	۱۸۸۰-۱۸۶۱	۰۰۰۸/۱	۰۰۸۸/۱	۹۹۲۵/۰
۴۵	۱۹۰۰-۱۸۸۱	۰۰۲۴/۱	۰۰۱۲/۱	۹۹۳۶/۰
۴۶	۱۹۲۰-۱۹۰۱	۹۹۷۸/۰	۰۰۴۹/۱	۰۰۰۷/۱
۴۷	۱۹۴۰-۱۹۲۱	۰۰۰۶/۱	۰۰۳۴/۱	۹۹۵۵/۰
۴۸	۱۹۶۰-۱۹۴۱	۰۰۲۰/۱	۰۰۴۶/۱	۰۰۱۲/۱
۴۹	۱۹۸۰-۱۹۶۱	۰۰۱۹/۱	۰۰۳۵/۱	۰۰۱۱/۱
۵۰	۲۰۰۰-۱۹۸۱	۰۰۱۹/۱	۰۰۴۵/۱	۰۰۱۹/۱

بحث و نتیجه‌گیری

مدل‌های طولی حاشیه‌ای مدل‌های مناسبی برای تحلیل داده‌های طولی هستند. در این مقاله مدل را بر اساس روش مدل‌های آمیخته بر پایه اسپلاین‌های تاوانیده که کمتر به صورت عمیق به آن پرداخته شده است را مورد بررسی قرار داده‌ایم. همچنین با استفاده از BUGS و بسته نرم‌افزاری *BRugs* استنباط‌های بیزی مورد نظر انجام شده است. علاوه بر آن با استفاده از شبیه‌سازی نشان داده شده است که استنباط‌های مورد نظر از دقت نسبتاً بالایی برخوردارند. لازم به ذکر است که از دیگر مزیت‌های این روش می‌توان به کاربرد آن در مدل‌های پیچیده‌تر دیگر مانند مدل‌های لگ خطی، رگرسیون چندکی و یا حتی مدل‌ها با داده‌های گم‌شده نیز اشاره کرد که در این مقاله به آن پرداخته نشده است و می‌تواند پیشنهادی برای پژوهش‌های بعدی در نظر گرفته شوند.

تقدیر و تشکر

نویسندگان از داوران گرامی برای ارائه پیشنهادهای سازنده در راستای بهبود این پژوهش و از دقت نظر ویراستار محترم مجله بسیار سپاس‌گزارند.

مراجع

- Al Kadiri, M., Carroll, R. J. and Wand, M. P., (2010), Marginal Longitudinal Semiparametric Regression via Penalized Splines, *Statistics and Probability Letters*, **80**, 1242–1252.
- Carroll, R. J., Maity, A., Mammen, E., Yu, K., (2009), Efficient Semiparametric Marginal Estimation for the Partially Linear Additive Model for Longitudinal/Clustered Data, *Statistics in Biosciences*, **1**, 10-31.
- Fan, J., Huang, T. and Li, R., (2007), Analysis of Longitudinal Data with Semiparametric Estimation of Covariance Function, *Journal of the American Statistical Association*, **102**, 632-641.
- Kruschke, K. J., (2014), *Donig Bayesian Data Analysis , a Tutorial with R and BUGS* , Academic Press ,Elsevier.
- Lunn, D. J., Thomas, A., Best, N., Spiegelhalter, D., (2000), WinBUGS-a Bayesian Modelling framework: Concepts, Structure, and Extensibility, *Journal of Statistics and Computing*, Stat. **10**, 325-337.
- Menictasa, M., and Wand, M. P., (2013), Variational Inference for Marginal Longitudinal Semiparametric Regression, *Dissemination of Statistics Research*, **2**, 61–71.
- Ruppert, D., Wand, P. and Carroll, R. J., (2003), *Semiparametric Regression*, Cambridge, New York.
- West, B. T., Galecki, A. T. and Welch, K. (2006), *Linear Mixed Models A Practical Guide Using Statistical Software*, Taylor & Francis Group New York.
- Wand, M. P., and Ormerod, J., (2008), On Semiparametric Regression With O’Sullivan Penalized Splines, *Australian Statistical Publishing Association*, Stat. **50**, 179–198.
- Wand, M. P., (2008), *Semiparametric Regression and Graphical Models*, Centre for Statistical and Survey Methodology, The University of Wollongong.

- Wang, N., Carroll, R. J. and Lin, X., (2005), Efficient Semiparametric Marginal Estimation for Longitudinal/Clustered Data, *Journal of the American Statistical Association*, **100**, 147-157.
- Welham, S. J., Cullis, B. R., Kenward, M. G. and Thompson, R., (2007), A Comparison of Mixed Model Splines for Curve Fitting, *Australian and New Zealand Journal of Statistics*, **49**, 1-23.
- Zeger, S. and Diggle, P. J., (1994), Semiparametric Models for Longitudinal Data with Application to CD4 Cell Numbers in HIV Seroconverters, *Biometrics*, **50**, 689-699.
- Zhao, Y., Staudenmayer, J., Coull, B. A., and Wand, M. P., (2006), General Design Bayesian Generalized Linear Mixed Models, *Statistical Science*, **21**, 35-51.