

تحلیل مدل‌های آمیخته خطی تعمیم‌یافته فضایی با استفاده از تقریب لاپلاس آشیانی جمع‌بسته

فاطمه حسینی^۱، امید کریمی، منور محمدکریمی

گروه آمار، دانشگاه سمنان

تاریخ پذیرش: ۱۳۹۳/۶/۲۳

تاریخ دریافت: ۱۳۹۳/۳/۲۵

چکیده: برای مدل‌بندی پاسخ‌های گسسته فضایی زمین‌آمار از مدل‌های آمیخته خطی تعمیم‌یافته فضایی استفاده می‌شود و ساختار همبستگی فضایی داده‌ها از طریق متغیرهای پنهان در نظر گرفته می‌شود. از مهم‌ترین اهداف در بررسی این مدل‌ها پیش‌گویی متغیرهای پنهان و برآورد پارامترهای مدل است. در این مقاله برای تحلیل این مدل‌ها، ابتدا یک روش پیش‌گویی ارائه و سپس به بیان رهیافت بیزی و الگوریتم‌های مونت‌کارلویی پرداخته می‌شود. به دلیل پیچیدگی این مدل‌ها و استفاده از نمونه‌های مونت‌کارلویی در تحلیل بیزی، زمان محاسبات بسیار طولانی است. برای رفع این مشکل روش بیزی تقریبی با استفاده از تقریب لاپلاس آشیانی جمع‌بسته بررسی می‌شود. در نهایت یک مجموعه داده واقعی مربوط به تعداد روزهای دارای بارندگی استان سمنان در سال ۱۳۹۱، مشاهده شده در ایستگاه‌های هواشناسی این استان با مدل و روش‌های معرفی شده مورد مطالعه قرار می‌گیرد.

واژه‌های کلیدی: مدل آمیخته خطی تعمیم‌یافته فضایی، الگوریتم‌های مونت‌کارلویی، رهیافت بیزی

رده‌بندی ریاضی (۲۰۱۰): ۶۲F۱۵، ۹۱B۷۲

۱- مقدمه

مدل‌های خطی تعمیم‌یافته شامل دامنه گسترده‌ای از مدل‌های آماری هستند که می‌توان آن‌ها را به مدل‌هایی با متغیرهای پاسخ گسسته تعمیم داد. در این مدل‌ها با فرض استقلال مشاهدات، با استفاده از یک تابع پیوند بین میانگین مشاهدات و متغیرهای کمکی ارتباط برقرار می‌شود. نلدر و ودربرن [۱] اولین کسانی بودند که چارچوب واحدی برای این مدل‌ها ارائه کردند. مک‌کلا و نلدر [۲] به تعمیم مدل‌های خطی برای مدل‌بندی متغیرهای پاسخ گسسته

پرداختند و به‌طور مفصل پیرامون مدل‌های خطی تعمیم‌یافته مطالعه نمودند. در حالتی که بین مشاهدات همبستگی وجود دارد، در این مدل فرض استقلال مشاهدات به استقلال شرطی تعدیل و همبستگی بین آن‌ها با اضافه کردن اثرات تصادفی از طریق متغیرهای پنهان به مدل در نظر گرفته می‌شود و در این صورت مدل خطی تعمیم‌یافته تبدیل به مدل آمیخته خطی تعمیم‌یافته می‌شود. در صورتی که همبستگی مذکور از نوع فضایی باشد، مدل آمیخته خطی تعمیم‌یافته به مدل آمیخته خطی تعمیم‌یافته فضایی^۱ (SGLM) تبدیل می‌شود. برسلو و کلایتون [۳] از این مدل در مطالعات پزشکی استفاده کردند. دیگل و همکاران [۴] این مدل‌ها را برای تحلیل متغیرهای فضایی گسسته در یک ناحیه پیوسته به کار بردند.

یک مسئله مهم در مدل SGLM پیش‌گویی متغیرهای پنهان در موقعیت‌های فاقد مشاهده است، که مستلزم برآورد پارامترهای مدل و متغیرهای پنهان در موقعیت‌های دارای مشاهده‌ی پاسخ است. لازم به ذکر است که چون در این مدل توزیع پاسخ‌های فضایی گاوسی نیست و علاوه بر این به دلیل وجود متغیرهای پنهان، برخلاف مدل‌های خطی، تابع درستنمایی شکل بسته‌ای ندارد و بنابراین برآورد پارامترها و پیش‌گویی‌ها به راحتی و از راه‌های مستقیم امکان‌پذیر نیستند، لذا در اکثر مقالات با پذیرش فرض نرمال بودن متغیرهای پنهان به ارائه راه‌حلی برای برآورد پارامترهای مدل و متغیرهای پنهان با ماکسیمم کردن توابع درستنمایی، شبه درستنمایی تاوانیده یا درستنمایی سلسله مراتبی به روش‌های عددی پرداخته شده است. از جمله مک‌کلا [۵]، با به کار بردن روش‌های عددی مثل ماکسیمم‌سازی امید ریاضی مونت کارلویی، نیوتن رافسون مونت کارلویی و ماکسیمم درستنمایی شبیه‌سازی شده، برآورد ماکسیمم درستنمایی پارامترها را در مدل‌های آمیخته خطی تعمیم‌یافته با اثرات تصادفی غیر فضایی به دست آورد. همچنین پن و تامسون [۶]، برآورد شبه مونت کارلو و لین [۷]، برآورد پارامترهای مدل آمیخته خطی تعمیم‌یافته پواسونی را با ترکیب درستنمایی نما و شبه درستنمایی تاوانیده مورد مطالعه قراردادند. ناتارجان و کاس [۸]، با بسط روش‌های بیزی در تحلیل مدل‌های آمیخته خطی تعمیم‌یافته نحوه انتخاب پیشین در این مدل‌ها را مطالعه کردند.

ژانگ [۹] یک الگوریتم مونت کارلویی با عنوان الگوریتم گرادیانت ماکسیمم‌سازی امید ریاضی برای به دست آوردن برآورد ماکسیمم درستنمایی پارامترهای مدل ارائه نمود. کریستن سن و همکارانش مقالات زیادی در راستای حل مشکل برآورد پارامترها و پیش‌گویی در مدل‌های SGLM معرفی نمودند که اکثر مقالات ایشان بر پایه استنباط بیزی و روش‌های مونت کارلویی می‌باشد. کریستن سن و همکاران [۱۰] و کریستن سن و وگپترسون [۱۱] با استفاده از رهیافت بیزی و روش‌های مونت کارلویی به تحلیل مدل‌های SGLM پرداختند و نشان دادند که

استفاده از الگوریتم لانژون-هستینگس، برای شبیه‌سازی از توزیع پسین متغیرهای پنهان گاوسی مفید می‌باشد. کریستن‌سن [۱۲] با روش ماکسیمم درست‌نمایی و الگوریتم مونت کارلو، پارامترها و پیش‌گویی بهینه را در مدل‌های SGLM با فرض نرمال بودن متغیرهای پنهان به دست آورد. کریستن‌سن و همکاران [۱۳] با معرفی روش‌های مونت کارلویی استوار به تحلیل مدل‌های SGLM با متغیرهای پنهان نرمال پرداختند و اینسورث و دین [۱۴] برآوردهای بیزی را با برآوردهای شبه درست‌نمایی توانیده در این مدل‌ها مقایسه کردند.

باغیشنی و محمدزاده [۱۵]، الگوریتم عددی همسانه‌سازی داده‌ها را برای محاسبه برآورد درست‌نمایی معرفی کردند. حسینی [۱۶]، حسینی و همکاران [۱۷] و حسینی [۱۸] الگوریتم‌های مفیدی برای پیش‌گویی و به‌دست آوردن برآوردهای درست‌نمایی معرفی نمودند.

روش‌های فوق، روش‌هایی مفید برای به دست آوردن برآوردها و تحلیل بیزی پارامترهاست اما مشکل بزرگی مانند طولانی بودن زمان محاسبات را دارد. برای حل این مشکل رو و همکاران [۱۹] روش تقریب لاپلاس آشیانی جمع‌بسته^۱ (INLA) را برای تحلیل مدل‌های پیچیده فضایی معرفی کردند و نشان دادند این روش ضمن حفظ دقت برآورد پارامترها سرعت محاسبات را نیز افزایش می‌دهد. ایدسویک و همکاران [۲۰] از تقریب معرفی‌شده توسط رو و مارتینو [۲۱] استفاده کردند و به تحلیل بیزی تقریبی مدل‌های SGLM پرداختند. حسینی و همکاران [۲۲] و حسینی و محمدزاده [۲۳] با به‌کار بردن توزیع چوله نرمال بسته برای متغیرهای پنهان تحلیل بیزی و بیزی تقریبی این مدل‌ها را بیان نمودند. قلی‌زاده و همکاران [۲۴] از تقریب لاپلاس آشیانی جمع‌بسته استفاده نمودند و به تحلیل بیزی تقریبی مدل‌های رگرسیون جمعی ساختاری پرداختند و داده‌های جرم شهر تهران را با این مدل و روش مورد تحلیل قراردادند.

در این مقاله ابتدا مدل SGLM معرفی می‌شود، سپس برآورد پارامترهای مدل و متغیرهای پنهان در موقعیت‌های دارای مشاهده پاسخ و پیش‌گویی متغیرهای پنهان در موقعیت‌های فاقد مشاهده پاسخ، با دو رهیافت بیزی معمولی و تقریبی ارائه می‌شود. درنهایت به مطالعه مجموعه داده‌های تعداد روزهای بارندگی در استان سمنان پرداخته می‌شود.

۲- مدل

فرض کنید $\mathbf{x} = (x_1, \dots, x_n)'$ بردار متغیرهای پنهان فضایی در n موقعیت $\{s_1, \dots, s_n\}$ با چگالی $\pi(\mathbf{x} | \boldsymbol{\eta}) = N_n(H\boldsymbol{\beta}, \Sigma_0)$ از یک میدان تصادفی گاوسی باشد، که در آن

مدل هستند. همچنین فرض کنید بردار پارامترهای مدل، $\boldsymbol{\eta} = (\boldsymbol{\beta}', \boldsymbol{\theta}')$ بردار پارامترهای رگرسیونی و $\boldsymbol{\theta}$ بردار پارامترهای همبستگی فضایی $n \times (p+1)$ ماتریس H متغیرهای تبیینی، میدان تصادفی فضایی ناگوسی در موقعیت‌های $\{s_1, \dots, s_k\}$ باشد. با فرض استقلال شرطی روی متغیرهای پنهان، $\pi(\mathbf{y} | \mathbf{x})$ متعلق به خانواده نمایی با تابع چگالی $\pi(\mathbf{y}_i | x_i) = \exp\{y_i x_i - b(x_i) + c(y_i)\}$ است، که در آن $b(\cdot)$ و $c(\cdot)$ توابع معلوم‌اند. اکنون با توجه به خانواده توزیع متغیر پاسخ و تعریف یک تابع پیوند مشخص مانند $g(\cdot)$ مدل به صورت $E(y_i | x_i) = g^{-1}(x_i)$ تعریف می‌شود و مؤلفه‌های مدل به فرم

$$\pi(\mathbf{x} | \mathbf{y}, \boldsymbol{\eta}) \propto \pi(\mathbf{y} | \mathbf{x}^{obs}) \pi(\mathbf{x} | \boldsymbol{\eta})$$

$$\propto |\Sigma_{\theta}|^{-1/2} \exp \left\{ \sum_{i=1}^k [y_i x_i - b(x_i) + c(y_i)] - \frac{1}{2} (\mathbf{x} - H\boldsymbol{\beta})' \Sigma_{\theta}^{-1} (\mathbf{x} - H\boldsymbol{\beta}) \right\}, \quad (1)$$

می‌باشد.

۳- پیش‌گویی مدل

در یک مدل SGLM هدف پیش‌گویی متغیرهای پنهان در موقعیت‌های فاقد مشاهده پاسخ $\{s_{k+1}, \dots, s_n\}$ و برآورد متغیرهای پنهان در موقعیت‌های دارای مشاهده پاسخ $\{s_1, \dots, s_k\}$ است. برای این منظور بردار متغیرهای پنهان به صورت $\mathbf{x} = (\mathbf{x}^{obs'}, \mathbf{x}^{pred'})'$ تجزیه می‌شود. که $\mathbf{x}^{obs} = A\mathbf{x}$ ، متعلق به k موقعیت دارای مشاهده پاسخ است، که در آن $A = [I_{k \times k} | 0_{k \times n-k}]$ و بردار متغیرهای پنهان در $n-k$ موقعیت انتخاب‌شده برای پیش‌گویی است. حسینی [۱۸] قضیه زیر را با فرض معلوم بودن پارامترهای مدل برای محاسبه پیش‌گویی مینیمم میانگین مربع خطا متغیرهای پنهان در $n-k$ موقعیت ارائه نمود.

قضیه ۱: فرض کنید توزیع متغیرهای پنهان فضایی $\mathbf{x} = (\mathbf{x}^{obs'}, \mathbf{x}^{pred'})'$ متعلق به خانواده نرمال $N_n(H\boldsymbol{\beta}, \Sigma_{\theta})$ باشد، به طوری که $\Sigma_{\theta} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$ و $H\boldsymbol{\beta} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$ با شرط روی متغیرهای پنهان فضایی، توزیع متغیرهای مستقل شرطی پاسخ متعلق به خانواده نمایی باشد، در این صورت $E(\mathbf{x}^{pred} | \mathbf{y}) = \boldsymbol{\mu}_2 + \Sigma_{21} \Sigma_{11}^{-1} (E(\mathbf{x}^{obs} | \mathbf{y}) - \boldsymbol{\mu}_1)$ پیش‌گویی مینیمم میانگین مربع خطا متغیرهای پنهان در $n-k$ موقعیت $j = k+1, \dots, n$ ام می‌باشد.

اما در عمل اغلب پارامترهای مدل نامعلوم هستند و به دلیل وجود متغیرهای پنهان فضایی و ناگوسی بودن متغیرهای پاسخ، در این مدل‌ها محاسبه برآورد پارامترهای مدل به راحتی

قابل محاسبه نیستند. یک روش معمول برای برآورد پارامترها، رهیافت ماکسیمم درست‌نمایی است.

فرض کنید (\mathbf{y}, \mathbf{x}) شامل بردار متغیرهای پاسخ گسسته فضایی و متغیرهای پنهان باشد، آنگاه

تابع درست‌نمایی به صورت $L(\boldsymbol{\eta} | \mathbf{y}) = \int \prod_{i=1}^k \pi(y_i | x_i) \pi(\mathbf{x} | \boldsymbol{\eta}) d\mathbf{x}$ است و تابع درست‌نمایی کامل به صورت

$$L_c(\boldsymbol{\eta} | \mathbf{y}, \mathbf{x}) = \prod_{i=1}^k \pi(y_i | x_i) \pi(\mathbf{x} | \boldsymbol{\eta}), \quad (2)$$

خواهد شد که شکل بسته‌ای ندارد و برای به‌دست آوردن برآورد ماکسیمم درست‌نمایی باید از الگوریتم‌های عددی استفاده کرد. علاوه بر این چون در عمل گاهی اوقات گردآوری نمونه با حجم زیاد امکان‌پذیر نیست، در این صورت نتایج تحلیل کلاسیک از دقت کمی برخوردار است و همچنین در تحلیل‌های کلاسیک از اطلاعات پیشینی در خصوص پارامترها استفاده نمی‌شود و لذا نتایجی با دقت کم‌تری نسبت به تحلیل‌های بیزی ارائه می‌کند در بخش بعد تحلیل مدل-SGLM با رهیافت بیزی ارائه می‌شود.

۴- تحلیل بیزی مدل

در این قسمت تحلیل بیزی مدل‌های SGLM روی یک ناحیه پیوسته فضایی مورد بررسی قرار می‌گیرد. یک مدل پارامتری به صورت $C_\theta(s_i, s_j) = \sigma^2 \exp\{-\|s_i - s_j\| / \varphi\}$ برای کوواریانس فضایی در نظر بگیرید. که در آن $\|\cdot\|$ نرم اقلیدسی و σ و φ به ترتیب پارامترهای مقیاس و همبستگی فضایی هستند. برای سره بودن توزیع پسین، برای پارامترهای مدل پیشین‌های سره در نظر گرفته می‌شود. پیشین متداول برای β نرمال و برای σ و φ به ترتیب گامای معکوس و گاما در نظر گرفته می‌شود. اکنون با فرض استقلال پیشین‌ها، توزیع پیشین توأم $\pi(\boldsymbol{\beta}, \sigma, \varphi) = \pi(\boldsymbol{\beta})\pi(\sigma)\pi(\varphi)$ حاصل می‌شود. لذا توزیع پسین

$$\pi(\mathbf{x}, \boldsymbol{\beta}, \sigma, \varphi | \mathbf{y}) \propto \pi(\mathbf{y} | \mathbf{x})\pi(\mathbf{x} | \boldsymbol{\beta}, \sigma, \varphi)\pi(\boldsymbol{\beta})\pi(\sigma)\pi(\varphi)$$

خواهد شد که دارای شکل پیچیده‌ای است. بنابراین برای تولید نمونه و شبیه‌سازی از توزیع پسین از روش‌های MCMC استفاده می‌شود. فرض کنید $\boldsymbol{\beta} \sim N(a, B)$ آنگاه توزیع شرطی $\boldsymbol{\beta}$ کامل

$$\begin{aligned} \pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{x}, \sigma, \varphi) &\propto \pi(x | \boldsymbol{\beta}, \sigma, \varphi) \pi(\boldsymbol{\beta}) \\ &\propto \exp\left\{-\frac{1}{2}[\boldsymbol{\beta}'(H'\Sigma_{\theta}^{-1}H + B^{-1})\boldsymbol{\beta} - 2\boldsymbol{\beta}'(B^{-1}a + H'\Sigma_{\theta}^{-1}\mathbf{x})]\right\} \quad (۳) \\ &\propto \exp\left\{-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_{\beta})'\Sigma_{\beta}^{-1}(\boldsymbol{\beta} - \boldsymbol{\mu}_{\beta})\right\}, \end{aligned}$$

است که در آن $\Sigma_{\beta} = (H'\Sigma_{\theta}^{-1}H + B^{-1})$ و $\boldsymbol{\mu}_{\beta} = \Sigma_{\beta}(B^{-1}a + H'\Sigma_{\theta}^{-1}\mathbf{x})$ بنا بر این تولید نمونه از توزیع شرطی کامل $\boldsymbol{\beta}$ به راحتی امکان پذیر است. یک مزیت دیگر در نظر گرفتن پیشین نرمال برای $\boldsymbol{\beta}$ کمک به کاهش تعداد پارامترهای مدل است به طوری که اگر فرض کنیم $\boldsymbol{\beta} \sim N(a, B)$ و $\mathbf{x} | \boldsymbol{\beta}, \sigma, \varphi \sim N_n(H\boldsymbol{\beta}, \Sigma_{\theta})$ می توان نشان داد

$$\mathbf{x} | \sigma, \varphi \sim N_n(\boldsymbol{\mu}_x, \Sigma_x),$$

که در آن $\boldsymbol{\mu}_x = Ha$ و $\Sigma_x = (\Sigma_{\theta} + HBH')$. توزیع شرطی کامل برای سایر پارامترها

$$\begin{aligned} \pi(\sigma | \mathbf{y}, \mathbf{x}, \boldsymbol{\beta}, \varphi) &\propto \pi(x | \boldsymbol{\beta}, \sigma, \varphi) \pi(\sigma) \\ &= IG(\alpha, \tau) \phi_n(\mathbf{x}; H\boldsymbol{\beta}, \Sigma_{\theta}), \\ \pi(\varphi | \mathbf{y}, \mathbf{x}, \boldsymbol{\beta}, \sigma) &\propto \pi(x | \boldsymbol{\beta}, \sigma, \varphi) \pi(\varphi) \\ &= \Gamma(\gamma, \omega) \phi_n(\mathbf{x}; H\boldsymbol{\beta}, \Sigma_{\theta}), \end{aligned} \quad (۴)$$

که دارای شکل خاصی از توزیع های شناخته شده نیستند و برای تولید نمونه از الگوریتم متروپلیس- هستینگس استفاده می شود. توزیع شرطی کامل برای هر مؤلفه بردار متغیرهای پنهان نیز به صورت

$$\pi(x_k | \mathbf{x}_{-k}, \mathbf{y}, \boldsymbol{\beta}, \sigma, \varphi) \propto \pi(x_k | \mathbf{x}_{-k}, \boldsymbol{\beta}, \sigma, \varphi) \pi(\mathbf{y} | \mathbf{x})$$

می باشد. برای پیش گویی بیزی متغیرهای پنهان، توزیع پیش گو به صورت

$$\pi(x_0 | \mathbf{y}) = \int \pi(x_0 | \mathbf{x}, \boldsymbol{\beta}, \sigma, \varphi) \pi(\mathbf{x}, \boldsymbol{\beta}, \sigma, \varphi | \mathbf{y}) dx d\boldsymbol{\beta} d\sigma d\varphi$$

است که در آن $\pi(x_0 | \mathbf{x}, \boldsymbol{\beta}, \sigma, \varphi)$ طبق خواص توزیع نرمال دارای توزیع نرمال است. توزیع پیش گوی بیزی y_0 نیز به صورت $\pi(y_0 | \mathbf{y}) = \int \pi(y_0 | x_0) \pi(x_0 | \mathbf{y}) dx_0$ است.

۴-۱- تحلیل بیزی مدل با تقریب لاپلاس آشیانی جمع بسته (INLA)

استفاده از الگوریتم های مونت کارلوی زنجیر مارکوفی به دلیل پیچیده بودن مدل های SGLM زمان بر هستند. ایدسویک و همکاران [۲۰] روش بیزی تقریبی را به جای روش بیزی معمولی و الگوریتم های نمونه گیری مونت کارلویی با استفاده از INLA معرفی کرده و نشان دادند

محاسبات این روش سریع‌تر و نتایجی معادل روش‌های بیزی معمولی ارائه می‌کند. مسئله اصلی در مدل‌های بیزی SGLM برای برآورد پارامترها و پیش‌گویی‌ها، محاسبه دو توزیع پسین $\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\eta})$ و $\pi(\boldsymbol{\eta}|\mathbf{y})$ است. اگر \mathbf{x} دارای توزیع نرمال و \mathbf{y} متعلق به خانواده نمایی باشد، با خطی کردن $\pi(\mathbf{y}|\mathbf{x}^{obs})$ حول یک مقدار ثابت \mathbf{x}° ، ایدسویک و همکاران [۲۰] نشان دادند توزیع $\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\eta})$ را می‌توان با توزیع نرمال به صورت

$$\hat{\pi}(\mathbf{x}|\mathbf{y}, \boldsymbol{\eta}) \approx N_n(\hat{\mu}_{\mathbf{x}|\mathbf{y}, \boldsymbol{\eta}}(\mathbf{x}^{\circ}), \hat{\Sigma}_{\mathbf{x}|\mathbf{y}, \boldsymbol{\eta}}(\mathbf{x}^{\circ})), \quad (۵)$$

تقریب زد، که در آن $R = A\Sigma_{\theta}A' + P$ ، $\hat{\mu}_{\mathbf{x}|\mathbf{y}, \boldsymbol{\eta}}(\mathbf{x}^{\circ}) = H\boldsymbol{\beta} + \Sigma_{\theta}A'R^{-1}(z(\mathbf{y}, \mathbf{x}^{\circ}) - AH\boldsymbol{\beta})$ ، $\hat{\Sigma}_{\mathbf{x}|\mathbf{y}, \boldsymbol{\eta}}(\mathbf{x}^{\circ}) = \Sigma_{\theta} - \Sigma_{\theta}A'R^{-1}A\Sigma_{\theta}$ همچنین P یک ماتریس قطری $k \times k$ با عناصر $P(i, i) = 1/b''(x_i)$ است و داریم

$$z_i(y_i, x_i^{\circ}) = [y_i - b'(x_i^{\circ}) + x_i b''(x_i^{\circ})] / b''(x_i^{\circ}).$$

یک تقریب لاپلاس اصلاح‌شده نیز برای توزیع‌های حاشیه‌ای پسین متغیرهای پنهان توسط ایدسویک و همکاران [۲۰] معرفی شده است که به صورت

$$\tilde{\pi}_{LA}(x_j|\mathbf{y}, \boldsymbol{\eta}) \propto \frac{\pi(\mathbf{y}|\mathbf{x})\pi(\mathbf{x}|\boldsymbol{\eta})}{\hat{\pi}(\mathbf{x}_{-j}|x_j, \mathbf{y}, \boldsymbol{\eta})} \Big|_{x=\hat{E}(\mathbf{x}_{-j}|x_j, \mathbf{y}, \boldsymbol{\eta})}, \quad (۶)$$

می‌باشد، بنا به خواص توزیع نرمال $\hat{\pi}(\mathbf{x}_{-j}|x_j, \mathbf{y}, \boldsymbol{\eta})$ یک توزیع تقریبی نرمال است. پس از تعیین این توزیع پسین تقریبی می‌توان تقریب لاپلاس توزیع پسین پارامترها را از رابطه

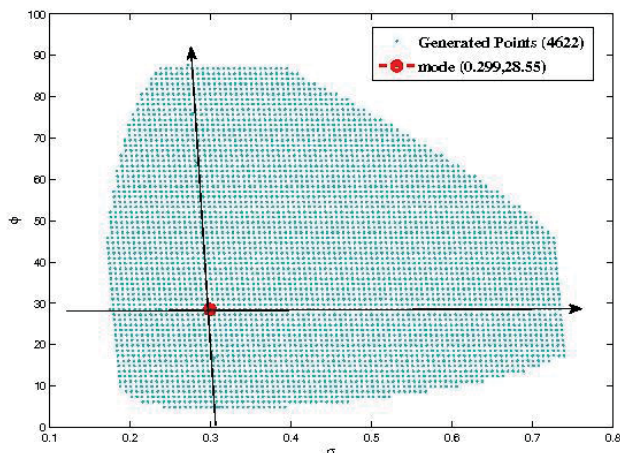
$$\hat{\pi}_{LA}(\boldsymbol{\eta}|\mathbf{y}) \propto \frac{\pi(\mathbf{y}|\mathbf{x})\pi(\mathbf{x}|\boldsymbol{\eta})\pi(\boldsymbol{\eta})}{\hat{\pi}(\mathbf{x}|\mathbf{y}, \boldsymbol{\eta})} \Big|_{x=\hat{E}(\mathbf{x}|\mathbf{y}, \boldsymbol{\eta})}, \quad (۷)$$

تعیین نمود. برای تعیین تقریب (۷) مد تابع $\log \hat{\pi}(\boldsymbol{\eta}|\mathbf{y})$ با روش‌های بهینه‌سازی به صورت

$$\boldsymbol{\eta}^* \text{ تعیین می‌شود. ماتریس هسین } H = \frac{\partial^2 \log \hat{\pi}(\boldsymbol{\eta}|\mathbf{y})}{\partial \eta^i \partial \eta^j} \text{ محاسبه و عکس آن به صورت}$$

$H^{-1} = V\Lambda V'$ تجزیه می‌شود، که در آن V ماتریس بردارهای ویژه و Λ ماریس قطری مقادیر ویژه آن است. مبدأ مختصات را به مد $\boldsymbol{\eta}^*$ انتقال داده، فرمول مختصات در مبدأ $\boldsymbol{\eta}^*$ به صورت $\boldsymbol{\eta}(t) = \boldsymbol{\eta}^* + V\Lambda^{-1}t$ تعریف می‌شود که در آن t مقادیر استاندارد شده هستند. با شروع از مبدأ مختصات جدید روی هر یک از محورهای نقطه‌ای به فاصله مقادیر صحیح δ_i به گونه‌ای اختیار می‌شوند که شرط $\log \hat{\pi}(\boldsymbol{\eta}(0)|\mathbf{y}) - \log \hat{\pi}(\boldsymbol{\eta}(t)|\mathbf{y}) < \delta_i$ برقرار باشد. سپس به‌طور مشابه نقاط درون صفحات نیز تعیین می‌شوند. با به‌کار بردن این الگوریتم η_i ‌هایی از توزیع $\hat{\pi}(\boldsymbol{\eta}|\mathbf{y})$ تولید می‌شوند و می‌توان تقریبی از توزیع را به دست آورد، (برای جزئیات

بیشتر به رو و همکاران [۱۹] مراجعه شود). به عنوان مثال برای مجموعه داده استفاده شده در این مقاله از روش توضیح داده شده تعداد ۴۶۲۲ نقطه تولید و در شکل ۱ رسم شده است.



شکل ۱: نقاط تولیدشده برای توزیع پسین تقریبی پارامترها

اکنون با داشتن نمونه‌های تولیدشده برای توزیع پسین تقریبی $\hat{\pi}_{LA}(\boldsymbol{\eta} | \mathbf{y})$ و مشخص شدن توزیع پسین تقریبی $\tilde{\pi}_{LA}(\mathbf{x}_j | \mathbf{y}, \boldsymbol{\eta})$ توزیع تقریبی پیش‌گو برای متغیرهای پنهان به صورت

$$\begin{aligned} \hat{\pi}(x_j | \mathbf{y}) &= \int \hat{\pi}(x_j | \mathbf{y}, \boldsymbol{\eta}_\ell) \hat{\pi}(\boldsymbol{\eta} | \mathbf{y}) \\ &= \sum_{\ell} \tilde{\pi}_{LA}(x_j | \mathbf{y}, \boldsymbol{\eta}_\ell) \hat{\pi}_{LA}(\boldsymbol{\eta}_\ell | \mathbf{y}), \quad j=1, \dots, n, \end{aligned} \quad (8)$$

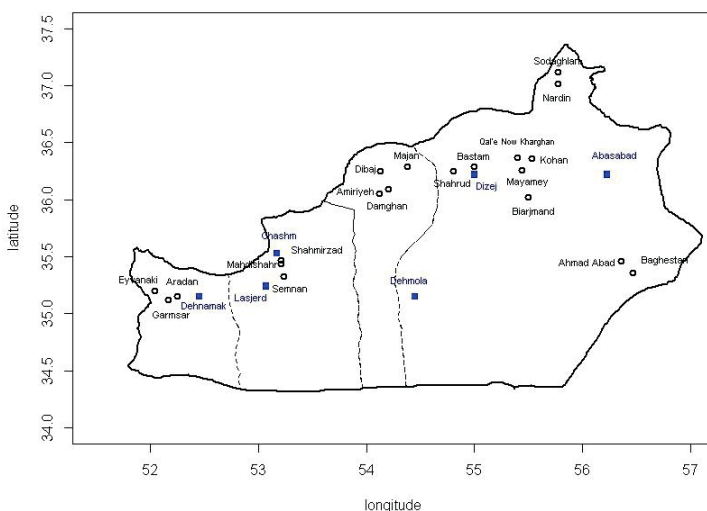
حاصل می‌شود، که در آن ℓ تعداد نقاط تولیدشده از الگوریتم تعیین توزیع پسین لاپلاس تقریبی پارامترها می‌باشد و بنابراین با توجه به رابطه (۸) این روش محاسبه توزیع تقریبی پیش‌گو به روش INLA معروف شده است. در عمل پیش‌گویی در موقعیت‌های $j = k+1, \dots, n$ مدنظر می‌باشد. همچنین می‌توان توزیع پیشگویی y_j در موقعیت فاقد مشاهده را به صورت

$$\pi(y_j | \mathbf{y}) = \int \pi(y_j | x_j) \pi(x_j | \mathbf{y}) dx_j, \quad j = k+1, \dots, n$$

به دست آورد که در آن به جای $\pi(x_j | \mathbf{y})$ از تقریب مستقیم یا اصلاح شده آن استفاده می‌شود.

۵- بررسی داده‌های بارندگی استان سمنان

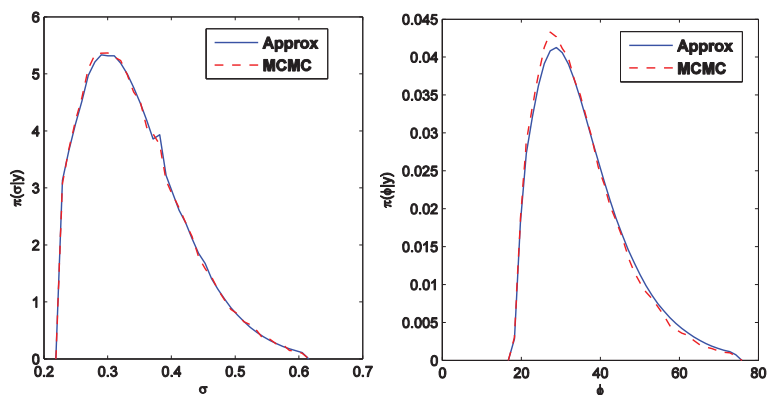
در این بخش مدل SGLM به داده‌های تعداد روزهای دارای بارندگی در استان سمنان برازش داده می‌شود و مورد تحلیل آماری قرار می‌گیرد. متغیر پاسخ تعداد روزهای دارای بارندگی است، که طی پنج ماه از اول آبان تا پایان اسفند ۱۳۹۱ در ۲۰ ایستگاه هواشناسی استان سمنان مشاهده شده‌اند. در شکل ۲ موقعیت ۲۰ ایستگاهی که تعداد روزهای دارای بارندگی در آن‌ها ثبت شده است و موقعیت‌هایی که به منظور پیش‌گویی انتخاب شده‌اند بر روی نقشه استان سمنان مشخص شده است. فرض می‌شود متغیرهای پاسخ دوجمله‌ای شرطی مستقل هستند و متغیرهای پنهان X نیز دارای توزیع نرمال به شکل $N(\beta_0, \Sigma_\theta)$ باشند. برای ساختار همبستگی فضایی تابع کوواریانس نمایی همسان‌گرد در نظر گرفته شده است. برای پارامتر مقیاس σ ، پیشین گامای معکوس با پارامترهای $(\frac{0}{5}, 5)$ ، برای پارامتر همبستگی پیشین گاما با پارامترهای ۸ و ۵ منظور شده است.



شکل ۲: موقعیت جغرافیایی ایستگاه‌های هواشناسی روی نقشه استان سمنان

برای این مجموعه داده‌ها متغیر کمکی وجود ندارد و برای جمله ثابت β_0 پیشین نرمال با میانگین ۱ و واریانس ۵ در نظر گرفته شد که به دلیل عدم اهمیت این پارامتر و با به کار بردن رابطه (۴)، از مدل حذف می‌شود. شکل ۳ نشان‌دهنده توزیع‌های پسین حاشیه‌ای پارامترهاست، که با ۱۰۰۰۰ تکرار از روش‌های MCMC و به کار بردن الگوریتم متروپولیس-هستینگس و از روش بیزی تقریبی با به کار بردن ۴۶۲۲ نقطه از توزیع تقریبی پسین و الگوریتم توضیح داده شده حاصل شده‌اند. خط ممتد روش بیزی تقریبی و خط چین روش‌های

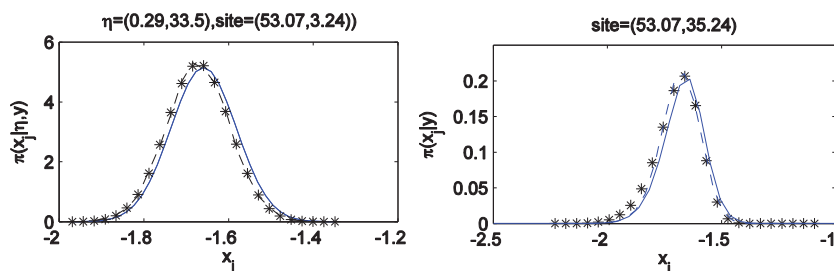
MCMC هستند، مشاهده می‌شود که دو روش بسیار شبیه به هم هستند. برآورد بیزی پارامترها به صورت $(\hat{\sigma}, \hat{\phi}) = (0/29, 33/5)$ به دست آمده است. اکنون با جایگذاری برآوردهای بیزی پارامترها می‌توان در نقاط جدید پیش‌گویی انجام داد.



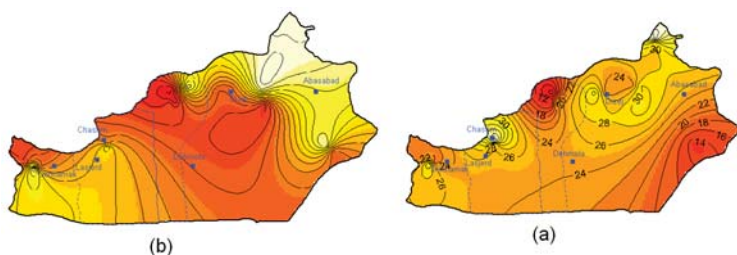
شکل ۳: نمودار توزیع‌های حاشیه‌ای پسین پارامترها با روش بیزی تقریبی (خطوط ممتد) و روش MCMC (خط چین)

شکل ۴ (چپ) نمودار توزیع پیش‌گو برای یک موقعیت فاقد مشاهده دلخواه با طول و عرض جغرافیایی $(53/07, 35/24)$ رسم شده است. برای مدل SGLM با متغیرهای پنهان نرمال، نمودار توزیع شرطی پیش‌گو با جایگذاری مد پسین تقریبی به جای بردار پارامتر رسم شده‌اند. نمودار مشخص شده با خطوط ممتد از روش بیزی تقریبی، خطوط مشخص شده با * از روش بیزی تقریبی اصلاح شده و خطوط خط چین از روش MCMC به دست آمده‌اند. همچنین شکل ۴ (راست) نشان‌دهنده چگالی حاشیه‌ای پیش‌گو $\pi(x_j | \mathbf{y})$ است که حاصل از سه روش بیزی تقریبی، بیزی تقریبی اصلاح شده و روش MCMC، می‌باشد و تشابه هر سه روش در این شکل مشخص است. همچنین برای بررسی بیشتر معادل بودن سه روش، مقادیر پیش‌گویی در شش موقعیت مشخص شده در شکل ۲ محاسبه و مقدار میانگین مربع انحرافات و میانگین قدرمطلق انحرافات مقادیر پیش‌گویی روش MCMC از روش تقریبی اصلاح شده به ترتیب $0/0000452$ و $0/00557$ و میانگین مربع انحرافات و میانگین قدرمطلق انحرافات مقادیر پیش‌گویی روش MCMC از روش تقریبی مستقیم به ترتیب $0/00246$ و $0/003918$ به دست آمد. این مقادیر بیان‌گر نزدیک بودن روش تقریبی اصلاح شده به روش MCMC نسبت به روش تقریبی مستقیم است. شکل ۵ (چپ) نقشه پیش‌گویی متغیر پنهان به صورت لایه‌هایی در ۱۵ سطح، روی نقشه سمنان و پیش‌گویی تعداد روزهای دارای بارندگی، در ۶ ایستگاه مذکور و کل نقشه می‌باشد. بیشترین بارندگی در این استان در مناطق شمالی استان می‌باشد که به دلیل

همسایگی با استان‌های شمالی است. پیش‌گویی تعداد روزهای دارای بارندگی در ۶ ایستگاه هواشناسی موردنظر، لاسجرد، ده نمک، ده‌ملا، چاشم، دیزج و عباس‌آباد نیز در نمودار سمت راست مشخص شده است. زمان اجرای برنامه‌ها برای روش بیزی تقریبی در حدود ۵۰ ثانیه، برای بیزی تقریبی اصلاح‌شده حدود ۴ دقیقه و برای روش‌های MCMC حدود یک روز به طول انجامیده است.



شکل ۴: (چپ) چگالی شرطی $\hat{\pi}(x_j | \mathbf{y}, \boldsymbol{\eta})$ حاصل از بیزی تقریبی مستقیم (خطوط ممتد)، روش تقریبی اصلاح‌شده (*) و روش MCMC (خط‌چین). (راست) چگالی حاشیه‌ای پیش‌گو $\pi(x_j | \mathbf{y})$ حاصل از بیزی تقریبی مستقیم (خطوط ممتد)، روش تقریبی اصلاح‌شده (*) و روش MCMC (خط‌چین).



شکل ۵: نقشه پیش‌گویی متغیرهای پنهان (چپ)، نقشه پیش‌گویی بیزی تعداد روزهای دارای بارندگی (راست).

۶- بحث و نتیجه‌گیری

در این مقاله دو رهیافت بیزی معمولی و تقریبی برای تحلیل مدل‌های آمیخته خطی تعمیم‌یافته فضایی ارائه شد و از آن‌جا که اجرای الگوریتم‌های مونت‌کارلویی برای مدل‌ها بسیار طولانی و زمان‌بر هستند، روش بیزی تقریبی بررسی شد. در یک مثال کاربردی بارندگی در استان سمنان نشان داده شد که روش بیزی تقریبی معادل روش بیزی معمولی عمل می‌کند با

این تفاوت که روش تقریبی ارائه شده نیاز به چند ثانیه و روش تقریبی اصلاح شده نیاز به چند دقیقه اجرای کامپیوتری دارند، در صورتی که روش های MCMC نیاز به صرف زمان طولانی تری برای محاسبات کامپیوتری دارند. حجم نمونه داده های مورد بررسی در این مقاله کوچک است در صورتی که در بسیاری از مسائل با مجموعه داده ها با حجم بالا سروکار داریم که استفاده از استنباط های معمول کلاسیک و بیزی بسیار پیچیده و زمان بر است و روش بیزی تقریبی لاپلاس آشیانی جمع بسته به دلیل استفاده از تقریب های گاوسی و لاپلاس بسیار مفید به نظر می رسد. از ایرادات اصلی روش INLA، کاهش دقت نتایج با افزایش پارامترها می باشد. یکی از محدودیت های مدل های SGLM این است که متغیرهای پنهان به صورت میدان تصادفی گاوسی پنهان مدل بندی می شوند و به دلیل پنهان بودن، توزیع دقیق آن ها مشخص نیست و فرض نرمال بودن بر روی نتایج تأثیرگذار است. تاکنون مطالعاتی در خصوص به کار بردن کلاس چوله نرمال به عنوان توزیع متغیرهای پنهان انجام شده است که تا حدودی نتایج بهتری در بعضی موارد حاصل شده است. به عنوان پیشنهاد می توان متغیرهای پنهان را با کلاس بزرگ تری از توزیع ها مثل کلاس توزیع چوله t مدل بندی نمود.

مراجع

- [1] Nelder, J.A., and Wedderburn, R.W.M. (1972). Generalized Linear Mixed Models, *Journal of the Royal Statistical Association, Ser. A*, **135**, 370-384.
- [2] McCulloch, P. and Nelder, J.A. (1989), *Generalized Linear Models*, London, Chapman and Hall.
- [3] Breslow, N.E. and Clayton, D.G. (1993). Approximate Inference in Generalized Linear Mixed Models, *Journal of the American Statistical Association*, **88**, 9-25.
- [4] Diggle, P., Tawn, J.A. and Moyeed, R.A. (1998). Model-Based Geostatistic, *Journal of the Royal Statistical Society, Series C. Applied Statistics*, **47**, 299-350.
- [5] McCulloch, C. (1997). Maximum Likelihood Algorithms for Generalized Linear Mixed Models, *Journal of the American Statistical Association*, **92**, 162-170.
- [6] Pan, J. and Thompson, R. (2007). Quasi-Monte Carlo Estimation in Generalized Linear Mixed Models, *Computational Statistics and Data Analysis*, **51**, 5765-5775.

- [7] Lin, X. (2007). Estimation Using Penalized Quasi-Likelihood and Quasi-Pseudo-Likelihood in Poisson Mixed Models, *Lifetime Data Anal*, **13**, 533-544.
- [8] Natarajan, R. and Kass, R.E. (2000). Reference Bayesian methods for Generalized Linear Mixed Models, *Journal of the American Statistical Association*, **95**, 227-237.
- [9] Zhang, H. (2002). On Estimation and Prediction for Spatial Generalized Linear Mixed Models, *Biometrics*, **58**, 129-136.
- [10] Christensen, O.F., Moller, J. and Waagepetersen R.P. (2000). Analysis of Spatial Data Using Generalized Linear Mixed Models and Langevin-Type Markov Chain Monte Carlo, Research Report R-00-2009, *Department of Mathematical Sciences, Aalborg University*.
- [11] Christensen, O.F. and Waagepetersen R.P. (2002). Bayesian Prediction of Spatial Count Data Using Generalized Linear Mixed models, *Biometrics*, **58**, 280-286.
- [12] Christensen, O.F. (2004). Monte Carlo Maximum Likelihood in Model-Based Geostatistics, *Journal of Computational and Graphical Statistics*, **13**, 702-718.
- [13] Christensen, O.F., Roberts, G.O. and Skold, M. (2006). Robust MCMC Methods for Spatial Generalized Linear Mixed Models *Journal of Computational and Graphical Statistics*, **15**, 1-17.
- [14] Ainsworth, L.M. and Dean, C.B. (2006). Approximate Inference for Disease Mapping, *Computational Statistics and Data Analysis*, **50**, 2552-2570.
- [15] Baghishani, H. and Mohammadzade, M. (2011). A Data Cloning Algorithm for Computing Maximum Likelihood Estimates in Spatial Generalized Linear Mixed Models, *Computational Statistics and Data Analysis*, **55**, 1748-1759.
- [16] Hosseini, F. (2014). An efficient Algorithm to Analysis of Categorical Spatial Data, *proceedings of 12th Iranian Statistical Conference, Razi university, Iran*, 83-89.
- [17] Hosseini, F. and Mohammadzadeh, M. (2013). Estimation of Spatial Generalized Linear Mixed Models with Closed Skew Normal Latent Variables, *Journal of Science Kharazmi University*, **12**, 305-312.

- [18] Hosseini, F. and Mohammadzadeh, M. and Karimi, O. (2014). Pseudo-likelihood Inference for Discrete Spatial Response (A Case Study of the Semnan rainfall data), *Journal of Science Kharazmi University*, **13**, 797-808.
- [19] Rue, H. and Martino, S. (2007). Approximate Bayesian Inference for Hierarchical Gaussian Markov Random fields models, *Journal of Statistical Planning and Inference*, **137**, 3177-3199.
- [20] Eidsvik, J., Martino, S. and Rue, H. (2009). Approximate Bayesian Inference in Spatial Generalized Linear Mixed Models, *Scandinavian Journal of Statistics*, **36**, 1-22.
- [21] Rue, H., Martino, S. and Chopin, N. (2009). Approximate Bayesian Inference for models, *Journal of Statistical Planning and Inference*, **137**, 3177-3199.
- [22].Hosseini, J., Eidsvik, J and Mohammadzadeh, M. (2011). Approximate Bayesian Inference in Spatial GLMM with Skew Normal Latent Variables, *Computational Statistics and Data Analysis*, **55**, 1791-1806.
- [23] Hosseini, F. and Mohammadzadeh, M. (2012). Bayesian Prediction for Spatial GLMM's with Closed Skew Normal Latent Variables, *Australian & New Zealand Journal of Statistics*, **54**, 43-62
- [24] Gholizadeh, K., Mohammadzadeh, M. and Ghayyomi, Z., (2013). Spatial Analysis of Structured Additive Regression and Modeling of Crime Data in Tehran City using Integrated Nested Laplace Approximation, *Journal of Statistical Science*, **7**, 103-124.

Inference of Spatial Generalized Linear Mixed Models using Integrated Laplace Nested Approximation

Fatemeh Hosseini, Omid Karimi, Monavar Mohammad Karimi

Department of Statistics, University of Semnan, Semnan, Iran.

Abstract

Spatial generalized linear mixed models are used for modeling geostatistical discrete spatial responses and spatial correlation of the data is considered via latent variables. The most important interest in these models is estimation of the parameters and prediction of the latent variables. In this paper, first, a prediction method is presented. Then a Bayesian approach and MCMC algorithms are proposed. Since these models are complicated and Monte Carlo sampling is used in the Bayesian inference of these models, computation time is long. In order to resolve this problem, the Approximate Bayesian methods are considered. Finally, the proposed methods are applied to a case study on rainfall data observed in the weather stations of Semnan in 1391.

Keywords: Spatial generalized linear mixed model, Monte Carlo algorithm, Approximation Bayesian approach.

Mathematics Subject Classification (2000): 91B72, 62F15

