

برآورد مدل‌های آمیخته خطی تعمیم‌یافته فضایی با متغیرهای پنهان چوله نرمال بسته

فاطمه حسینی: دانشگاه سمنان، گروه آمار
* محسن محمدزاده: دانشگاه تربیت مدرس، گروه آمار

چکیده

مدل‌های آمیخته خطی تعمیم‌یافته فضایی معمولاً برای مدل‌بندی پاسخ‌های فضایی گسسته به‌کار می‌روند، که در آن‌ها ساختار همبستگی فضایی داده‌ها از طریق متغیرهای پنهان در نظر گرفته می‌شود. مسئله مهم در این مدل‌ها، برآورد متغیرهای پنهان فضایی در موقعیت‌های دارای مشاهده پاسخ و پارامترهای مدل و در نهایت پیش‌گویی متغیرهای پنهان در موقعیت‌های فاقد مشاهده است. در این راستا اغلب کاربران برای سهولت، توزیع نرمال را برای متغیرهای پنهان در نظر می‌گیرند. اگرچه این فرض باعث سهولت محاسبات می‌شود، اما گاهی در عمل واقع‌گرایانه نیست، یا به‌دلیل پنهان بودن بررسی آن میسر نیست. لذا در این مقاله استفاده از توزیع چوله‌نرمال بسته که در حالت خاص شامل توزیع نرمال است و تحت حاشیه‌سازی، شرطی کردن و تبدیلات خطی بسته است، برای متغیرهای پنهان پیشنهاد می‌شود. در این مدل‌ها تابع درست‌نمایی فرم بسته‌ای ندارد و به‌دست آوردن برآوردهای ماکسیمم درست‌نمایی پارامترها به راحتی امکان‌پذیر نیست، لذا الگوریتمی تقریبی برای برآورد ماکسیمم درست‌نمایی پارامترهای مدل و پیش‌گویی تقریبی متغیرهای پنهان ارائه می‌شود، که در مقایسه با روش‌های موجود بسیار سریع‌تر است. اعتبار مدل و الگوریتم پیشنهادی در یک مطالعه شبیه‌سازی بررسی می‌شود.

مقدمه

مدل آمیخته خطی تعمیم‌یافته فضایی^۱ (SGLM) معمولاً برای مدل‌بندی داده‌های فضایی گسسته روی یک ناحیه پیوسته فضایی به‌کار گرفته می‌شود. اولین بار برسلو^۲ و کلایتون^۳ از این مدل در بررسی‌های پزشکی استفاده کردند. مسئله‌ای مهم در مدل‌های پیش‌گویی متغیرهای پنهان در موقعیت‌های فاقد مشاهده است. این امر مستلزم برآورد پارامترهای مدل و متغیرهای پنهان است. حضور متغیرهای پنهان در این مدل‌ها و عدم وجود شکل بسته برای توزیع شرطی متغیرهای پنهان به‌شرط متغیرهای پاسخ، برآورد حداکثر درست‌نمایی پارامترهای مدل را

واژه‌های کلیدی: مدل آمیخته خطی تعمیم‌یافته فضایی، متغیر پنهان، توزیع چوله‌نرمال بسته

دریافت ۹۰/۵/۲۳

پذیرش ۹۱/۴/۶

mohsen_m@modares.ac.ir

*نویسنده مسئول

۱. Spatial Generalized Linear Mixed Model

۲. Breslow

۳. Clayton

دشوار و گاهی غیرممکن می‌سازد. دیگل^۱ و همکاران [۴] با رهیافت بیزی و مینیم کردن میانگین توان دوم خطاها، پیش‌گوی بهینه را برای متغیرهای پنهان به‌دست آوردند. زانگ [۸] با ترکیب روش مونت‌کارلو و الگوریتم گرادیانت بیشینه‌سازی امید ریاضی (EMG)، الگوریتم جدید گرادیانت بیشینه‌سازی امید ریاضی مونت کارلویی^۲ (MCEMG) را برای برآورد ماکسیمم درست‌نمایی پارامترهای همبستگی و پیش‌گویی متغیرهای پنهان نرمال در مدل SGLM ارائه کرد. کریستینسن [۲] با روش ماکسیمم درست‌نمایی و الگوریتم مونت کارلو، پارامترها و پیش‌گوی بهینه را در این مدل‌های همراه با فرض نرمال بودن متغیرهای پنهان به‌دست آورد. اما در عمل، به‌دلیل غیرقابل مشاهده بودن متغیرهای پنهان در مدل‌های SGLM، بررسی نرمال بودن آن‌ها مقدور نیست و پذیرش ناصحیح این فرض می‌تواند بر دقت برآورد پارامترها و پیش‌گوها تأثیر سو داشته باشد. دامین‌گوئز توزیع چوله نرمال بسته^۳ (CSN) معرفی کرد [۳]. این خانواده از توزیع‌ها از خانواده توزیع نرمال انعطاف‌پذیرتر و تحت حاشیه‌سازی، شرطی کردن و تبدیلات خطی بسته هستند. محمدزاده و حسینی الگوریتم MCEMG را برای مدل‌های با متغیرهای پنهان فضایی چوله نرمال بسته تعمیم دادند [۷]. حسینی و همکاران توزیع متغیرهای پنهان به شرط متغیرهای گسسته پاسخ را به‌طور تقریبی به دست آوردند [۵]. در این مقاله با استفاده از این توزیع تقریبی و الگوریتم EMG، الگوریتم جدید گرادیانت بیشینه‌سازی امید ریاضی تقریبی^۴ (AEMG) برای برآورد ماکسیمم درست‌نمایی پارامترهای مدل با متغیرهای پنهان چوله نرمال بسته معرفی می‌شود. سپس پیش‌گوی تقریبی مینیم میانگین توان دوم خطا^۵ (MMSE) برای متغیرهای پنهان ارائه می‌گردد. در انتها اعتبار روش‌ها و الگوریتم ارائه شده در یک بررسی شبیه‌سازی ارزیابی می‌شود.

مدل SGLM با متغیرهای پنهان چوله نرمال بسته

برای تعریف مدل SGLM ابتدا باید یک مدل درست‌نمایی برای مشاهدات فضایی گسسته (y) و سپس یک توزیع برای متغیرهای پنهان فضایی (x) در نظر گرفته شوند. توزیع y متعلق به خانواده نمایی [۶] و توزیع x را چوله نرمال بسته در نظر می‌گیریم. یعنی بردار تصادفی n بعدی x دارای تابع چگالی

$$f_{n,q}(x, \mu, \Sigma, D, \nu, \Delta) = \Phi_q^{-1}(0; \nu, \Delta + D\Sigma D') \phi_n(x; \mu, \Sigma) \Phi_q(D(x - \mu); \nu, \Delta) \quad (1)$$

است، که در آن μ بردار پارامتر مکان، Σ ماتریس $n \times n$ معین مثبت مقیاس، D ماتریس $q \times n$ چولگی، (μ, Σ) تابع چگالی نرمال n متغیره با میانگین μ و ماتریس کواریانس Σ و $\Phi_q(\cdot; \nu, \Delta)$ تابع توزیع تجمعی نرمال q متغیره با میانگین μ و ماتریس کواریانس Δ است. متغیر x که دارای تابع چگالی (1) است، به‌صورت $x \sim CSN_{n,q}(\mu, \Sigma, D, \nu, \Delta)$ نمایش داده می‌شود. بدیهی است اگر D ماتریس صفر باشد، تابع چگالی فوق تبدیل به توزیع نرمال n متغیره خواهد شد. همچنین اگر $q=1$ ، $\nu=0$ ، $\Delta=1$ و $D = \lambda \Sigma^{-\frac{1}{2}}$ آن‌گاه چگالی فوق،

۱. Diggle ۲. Monte Carlo Expectation Maximization Gradient ۳. Closed skew Normal

۴. Approximate Expectation Maximization Gradient ۵. Minimum Mean Square Error

فوق، چگالی توزیع چوله نرمال می‌شود. میانگین توزیع چوله نرمال بسته به صورت

$$E(X) = \mu + \Sigma D' \psi \quad (2)$$

است [3]، که در آن $\psi = \frac{\Phi_q^*(\mathbf{0}; \mathbf{v}, \Delta + D\Sigma D')}{\Phi_q(\mathbf{0}; \mathbf{v}, \Delta + D\Sigma D')}$ و برای هر ماتریس معین مثبت Ω ،

$$\nabla_r = \left(\frac{\partial}{\partial r_1}, \dots, \frac{\partial}{\partial r_q} \right)' \text{ و } \Phi_q^*(\mathbf{r}; \mathbf{v}, \Omega) = \left[\nabla_r \Phi_q(\mathbf{r}; \mathbf{v}, \Omega) \right]'$$

فرض کنید $\mathbf{x} = (x_1, \dots, x_n)'$ بردار متغیرهای پنهان فضایی در n موقعیت $\{s_1, \dots, s_n\}$ با چگالی

$$\mathbf{x} | \boldsymbol{\eta} \sim \text{CSN}_{n,1}(H\boldsymbol{\beta}, \Sigma_\theta, \lambda \Sigma_\theta^{-1/2}, 0, 1)$$

باشد، که در آن $\boldsymbol{\eta} = (\boldsymbol{\beta}', \boldsymbol{\theta}, \lambda)'$ پارامترهای مدل، H ماتریس

$n \times (p+1)$ متغیرهای تبیینی، $\boldsymbol{\beta} = (\beta_0, \dots, \beta_p)'$ بردار پارامترهای رگرسیونی، λ بردار پارامترهای

چولگی و $\boldsymbol{\theta}$ بردار پارامترهای همبستگی فضایی مدل هستند. بنا بر این با توجه به (1) چگالی متغیرهای پنهان

به صورت

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{2}{(2\pi)^{n/2} |\Sigma_\theta|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - H\boldsymbol{\beta})' \Sigma_\theta^{-1} (\mathbf{x} - H\boldsymbol{\beta}) \right\} \cdot \Phi_q \left(\lambda \Sigma_\theta^{-1/2} (\mathbf{x} - H\boldsymbol{\beta}) \right), \quad (3)$$

می‌شود. اکنون فرض کنید مشاهدات در موقعیت‌های فضایی $\{s_1, \dots, s_k\}$ در اختیار باشند و هدف پیش‌گویی

متغیرهای پنهان در موقعیت‌های فاقد مشاهده $\{s_{k+1}, \dots, s_n\}$ باشد. متغیرهای پنهان در k موقعیت مشاهده شده به

صورت $\mathbf{x}^{obs} = A\mathbf{x}$ نمایش داده می‌شود که در آن $A = [I_{k \times k} \mid 0_{k \times n-k}]$ است. در این صورت بردار \mathbf{x} را

می‌توان به صورت $\mathbf{x} = (\mathbf{x}^{obs'}, \mathbf{x}^{pred'})'$ تجزیه کرد، که در آن بردار متغیرهای پنهان در $n-k$ موقعیت

انتخاب شده برای پیش‌گویی است. همچنین $\mathbf{y} = (y_1, \dots, y_k)'$ بردار متغیرهای پاسخ فضایی گسسته در

موقعیت‌های دارای مشاهده $\{s_1, \dots, s_k\}$ است، با فرض استقلال شرطی این متغیرها روی متغیرهای پنهان،

$\pi(\mathbf{y} | \mathbf{x}^{obs})$ عضو خانواده نمایی با تابع چگالی $\pi(y_i | x_i) = \exp\{y_i x_i - b(x_i) + c(y_i)\}$ ، $i = 1, \dots, k$ است

است ([6]). بنا بر این مدل به صورت $E(y_i | x_i) = g^{-1}(x_i)$ است و به طور خلاصه مؤلفه‌های مدل

SGLM بدین صورت خواهد شد:

$$\pi(\mathbf{y}, \mathbf{x} | \boldsymbol{\eta}) = \pi(\mathbf{y} | \mathbf{x}) \pi(\mathbf{x} | \boldsymbol{\eta}) \propto |\Sigma_\theta|^{-1/2} \exp \left\{ \sum_{i=1}^n [y_i x_i - b(x_i) + c(y_i)] - \frac{1}{2} (\mathbf{x} - H\boldsymbol{\beta})' \Sigma_\theta^{-1} (\mathbf{x} - H\boldsymbol{\beta}) \right\} \\ \times \Phi \left(\lambda \Sigma_\theta^{-1/2} (\mathbf{x} - H\boldsymbol{\beta}) \right) \quad (4)$$

پیش‌گویی تقریبی در مدل‌های SGLM

با فرض معلوم بودن پارامترهای مدل پیش‌گویی مینیمم میانگین مربع خطای (MMSE) متغیرهای پنهان در

موقعیت زام به صورت $E(x_j | \mathbf{y})$ به دست می‌آید که برای محاسبه آن نیاز به توزیع کناری $\pi(x_j | \mathbf{y})$ است.

چون توزیع $\pi(\mathbf{x} | \mathbf{y}) \propto \pi(\mathbf{y} | \mathbf{x}) \pi(\mathbf{x})$ ، فرم بسته‌ای ندارد، نمی‌توان $E(x_j | \mathbf{y})$ را مستقیماً محاسبه کرد و به

طور تقریبی به‌دست آورده می‌شود. حسینی و همکاران [۵] نشان دادند اگر \mathbf{x} دارای توزیع چوله نرمال بسته و توزیع \mathbf{y} عضو خانواده‌ی نمایی باشد، با خطی کردن $\pi(\mathbf{y}|\mathbf{x})$ حول یک مقدار ثابت \mathbf{x}^0 ، توزیع $\pi(\mathbf{x}|\mathbf{y})$ را می‌توان با توزیع چوله نرمال بسته به‌صورت

$$\hat{\pi}(\mathbf{x}|\mathbf{y}) \approx CSN_{n,1}(\hat{\boldsymbol{\mu}}_{x|y,\eta}(\mathbf{y}, \mathbf{x}^0), \hat{\boldsymbol{\Sigma}}_{x|y,\eta}(\mathbf{x}^0), \hat{D}_{x|y,\eta}, \hat{\mathbf{v}}_{x|y,\eta}, 1) \quad (۵)$$

تقریب زد، که در آن $\hat{\boldsymbol{\mu}}_{x|y,\eta}(\mathbf{y}, \mathbf{x}^0) = H\boldsymbol{\beta} + \boldsymbol{\Sigma}_\theta A'R^{-1}(z(\mathbf{y}, \mathbf{x}^{obs}) - AH\boldsymbol{\beta})$ ، $P, R = A\boldsymbol{\Sigma}_\theta A' + P$ ،

ماتریسی قطری با درایه‌های $P(i,i) = 1/b''(x_i)$ ، $i = 1, \dots, k$

$$\ln \pi(\mathbf{x}|\boldsymbol{\eta}) = -\frac{1}{2} \ln |\boldsymbol{\Sigma}_\theta| - \frac{1}{2} (\mathbf{x} - H\boldsymbol{\beta})' \boldsymbol{\Sigma}_\theta^{-1} (\mathbf{x} - H\boldsymbol{\beta}) + \ln \{2\Phi(\boldsymbol{\lambda}' \boldsymbol{\Sigma}_\theta^{-\frac{1}{2}} (\mathbf{x} - H\boldsymbol{\beta}))\},$$

$$z_i(y_i, x_i^0) = [y_i - b'(x_i^0) + x_i b''(x_i^0)] / b''(x_i^0)$$

و $\hat{\mathbf{v}}_{x|y,\eta}(\mathbf{y}, \mathbf{x}^0) = \boldsymbol{\lambda}' \boldsymbol{\Sigma}_\theta^{-\frac{1}{2}} (H\boldsymbol{\beta} - \boldsymbol{\mu}_{x|y,\eta})$ و $\hat{D}_{x|y,\eta} = \boldsymbol{\lambda}' \boldsymbol{\Sigma}_\theta^{-\frac{1}{2}}$ ، $\hat{\boldsymbol{\Sigma}}_{x|y,\eta}(\mathbf{x}^0) = \boldsymbol{\Sigma}_\theta - \boldsymbol{\Sigma}_\theta A'R^{-1} A\boldsymbol{\Sigma}_\theta$ هستند.

با در نظر گرفتن $\mathbf{x}^* = (x_j \quad \mathbf{x}_{-j})' = A_j \mathbf{x}$ ، $j = 1, \dots, n$ که در آن $\mathbf{x}_{-j} = (x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n)$ و ماتریس واحدی است که زامین سطر آن به سطر اول منتقل شده است، بنا بر خاصیت بسته بودن توزیع چوله نرمال بسته نسبت به تبدیلات خطی، $[\mathbf{x}^* | \mathbf{y}, \boldsymbol{\eta}]$ دارای توزیع تقریبی

$$CSN(\hat{\boldsymbol{\mu}}_{x^*|y,\eta}, \hat{\boldsymbol{\Sigma}}_{x^*|y,\eta}, \hat{D}_{x^*|y,\eta}, \hat{\mathbf{v}}_{x^*|y,\eta}, \hat{\Delta}_{x^*|y,\eta})$$

است، که در آن

$$\hat{\boldsymbol{\mu}}_{x^*|y,\eta} = A_j \hat{\boldsymbol{\mu}}_{x|y,\eta}(\mathbf{y}, \mathbf{x}^0) = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \hat{\boldsymbol{\Sigma}}_{x^*|y,\eta} = A_j \hat{\boldsymbol{\Sigma}}_{x|y,\eta}(\mathbf{x}^0) = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix},$$

$$\hat{D}_{x^*|y,\eta} = \hat{D}_{x|y,\eta} A_j' = [D_1 \quad D_2], \quad \hat{\Delta}_{x^*|y,\eta} = 1.$$

طبق خاصیت بسته بودن توزیع چوله نرمال بسته نسبت به حاشیه‌سازی، $[x_j | \mathbf{y}]$ دارای توزیع تقریبی

$$CSN(\mu_1, \Sigma_{11}, D^*, \hat{\mathbf{v}}_{x|y,\eta}, \Delta^*)$$

است، که در آن $\Delta^* = 1 + D_2 \Sigma_{22,1} D_2'$ ، $D^* = D_1 + D_2 \Sigma_{21} \Sigma_{11}^{-1}$ و $\Sigma_{22,1} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$

بنا بر این از (۲)

$$E(x_j | \mathbf{y}) = \mu_1 + \Sigma_{11} D^* \boldsymbol{\psi} \quad (۶)$$

که پیش‌گوی MMSE تقریبی متغیرهای پنهان فضایی در موقعیت زام است، که در آن

$$\boldsymbol{\psi} = \frac{\Phi_q^*(r; \hat{\mathbf{v}}_{x|y}, \Delta^* + D^* \Sigma_{11} D^*)}{\Phi_q(\mathbf{0}; \hat{\mathbf{v}}_{x|y}, \Delta + D_2 \Sigma_{11} D_2')} \Big|_{r=0}.$$

و با فرض معلوم بودن بردار پارامترها قابل محاسبه است. اما در عمل اغلب پارامترهای مدل نامعلوم هستند و به‌دلیل وجود متغیرهای پنهان فضایی، برآورد پارامترهای مدل به‌راحتی قابل محاسبه نیستند. محمدزاده و حسینی [۷] الگوریتم MCEMG زانگ [۸] را برای مدل‌های SGLM با متغیرهای پنهان فضایی چوله تعمیم

دادند. اما اجرای این الگوریتم نیازمند تولید نمونه‌های مونت‌کارلویی و صرف وقت زیاد برای هم‌گرایی است. در این‌جا با استفاده از الگوریتم EMG و تقریب ارائه شده به‌وسیله محمدزاده و حسینی [۵] برای $\pi(\mathbf{x}|\mathbf{y})$ ، الگوریتم جدید AEMG پیشنهاد می‌شود که به نمونه‌های مونت‌کارلویی نیاز ندارد و از سرعت محاسبات بیش‌تری نسبت به الگوریتم MCEMG برخوردار است. فرض کنید (\mathbf{y}, \mathbf{x}) شامل بردار متغیرهای پاسخ گسسته فضایی و متغیرهای پنهان باشد، آنگاه از رابطه (۴) تابع درست‌نمایی کامل به‌صورت

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

است. برای اجرای الگوریتم AEMG، مقادیر اولیه پارامتر به‌صورت $\boldsymbol{\eta}^{(0)}$ را در نظر گرفته و با قرار دادن $m = 0$ ، با الگوریتم EMG بردار

$$\begin{aligned} \boldsymbol{\eta}^{(m+1)} &= \boldsymbol{\eta}^{(m)} - \left[E \left\{ \frac{\partial^2 \ln L_c}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}'} \middle| \mathbf{y} \right\} \right]_{\boldsymbol{\eta}=\boldsymbol{\eta}^{(m)}}^{-1} \left[E \left(\frac{\partial \ln L_c}{\partial \boldsymbol{\eta}} \middle| \mathbf{y} \right) \right]_{\boldsymbol{\eta}=\boldsymbol{\eta}^{(m)}} \\ &= \boldsymbol{\eta}^{(m)} - \left[E \left\{ \frac{\partial^2 \ln \pi(\mathbf{x}|\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}'} \middle| \mathbf{y} \right\} \right]_{\boldsymbol{\eta}=\boldsymbol{\eta}^{(m)}}^{-1} \left[E \left\{ \frac{\partial \ln \pi(\mathbf{x}|\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} \middle| \mathbf{y} \right\} \right]_{\boldsymbol{\eta}=\boldsymbol{\eta}^{(m)}} \quad (7) \end{aligned}$$

محاسبه می‌شود، که در آن

$$\ln \pi(\mathbf{x}|\boldsymbol{\eta}) = -\frac{1}{2} \ln |\Sigma_\theta| - \frac{1}{2} (\mathbf{x} - H\boldsymbol{\beta})' \Sigma_\theta^{-1} (\mathbf{x} - H\boldsymbol{\beta}) + \ln \{ 2\Phi(\lambda' \Sigma_\theta^{-\frac{1}{2}} (\mathbf{x} - H\boldsymbol{\beta})) \}$$

و امید ریاضی‌ها در رابطه (۷) به‌طور تقریبی با به‌کار بردن توزیع تقریبی (۵) قابل محاسبه هستند. در صورت هم‌گرایی، $\boldsymbol{\eta}^{(m+1)}$ برآورد به‌صورت پارامتر خواهد بود، در غیر این صورت الگوریتم تکرار می‌شود.

بررسی شبیه‌سازی

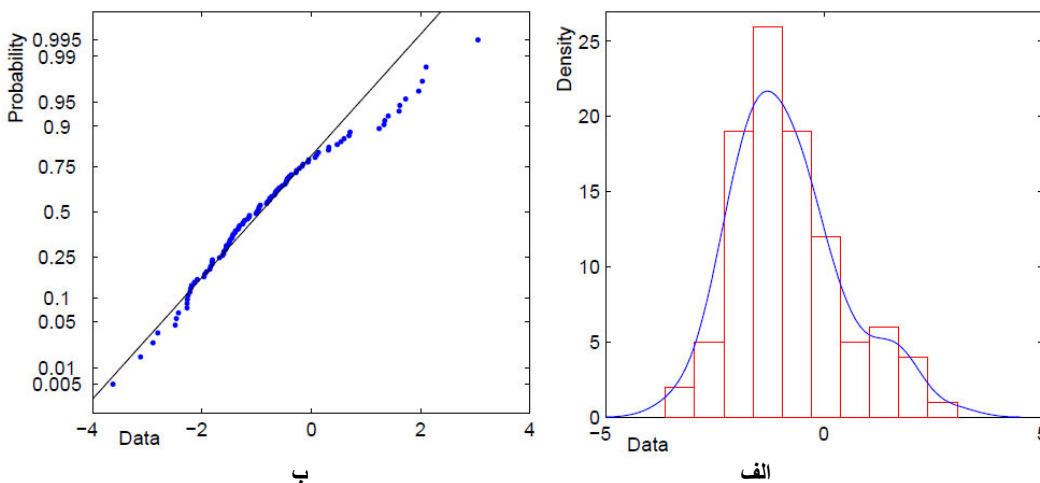
در این بخش عمل‌کرد مدل و روش‌های پیشنهاد شده در یک بررسی شبیه‌سازی براساس داده‌های تولید شده در شبکه‌ای منظم 10×10 ، به‌صورت $\{(\ell, k), \ell, k = 1, \dots, 10\}$ ارزیابی می‌شود. برای مقادیر معین پارامترهای توزیع CSN، متغیرهای پنهان فضایی \mathbf{x} از توزیع $CSN_{100,1}(\beta_1 z, \Sigma_\theta, \lambda_0 \mathbf{1}', \Sigma_\theta^{-\frac{1}{2}}, 0, 1)$ تولید شده‌اند. مقادیر $\beta_1 = 0.5$ ، $\sigma^2 = 2$ ، $\varphi = 4$ و $\lambda_0 = 2$ در نظر گرفته شده‌اند. تابع کواریانس هم‌سان‌گرد به‌صورت گرفته شده است. متغیرهای پاسخ $y_{\ell k}$ نیز با شرطی کردن روی متغیرهای پنهان فضایی از $C(h) = \sigma^2 \exp(-h/\varphi)$ ، $h > 0$ و متغیر تبیینی در موقعیت (ℓ, k) ام به‌صورت $z_{\ell k} = \log(1 + \ell)$ در نظر گرفته شده است. متغیرهای پاسخ $y_{\ell k} \sim Bin(100, \exp(x_{\ell k}) / [1 + \exp(x_{\ell k})])$ تولید شده‌اند. بافت‌نگار و نمودار Q-Q مقادیر تولید شده برای متغیرهای پنهان در شکل ۱ نشان دهنده چوله بودن توزیع داده‌ها است.

دو مدل SGLM، یکی با توزیع چوله‌نرمال بسته و دیگری با توزیع نرمال برای متغیرهای پنهان فضایی

شبیه‌سازی و براساس معیار میانگین توان دوم خطا، دقت برآوردها و براساس معیار میانگین توان دوم خطای پیش‌گویی دقت پیش‌گویی‌های حاصل مقایسه شده‌اند. برآورد پارامترها از دو الگوریتم AEMG و MCEMG محاسبه شده‌اند.

برای محاسبه مقادیر MSE تحت شرایط بالا، ۱۰۰ مجموعه داده تولید شده است. متوسط مقادیر برآورد شده پارامترها و MSE در جدول ۱ ارائه شده‌اند. چنان‌که ملاحظه می‌شود، برآورد پارامترهای مدل SGLM با توزیع چوله نرمال بسته برای متغیرهای پنهان فضایی از اریبی و MSE کوچک‌تری نسبت به مدل SGLM با متغیرهای پنهان فضایی نرمال برخوردار است. برای بررسی معنی‌داری تفاوت بین MSE‌های به دست آمده دو مدل نرمال و چوله وقتی پارامترها با الگوریتم AEMG برآورد شده‌اند، از آزمون تی زوجی استفاده شده است و مقادیر آماره آزمون، مقادیر احتمال و بازه‌های اطمینان ۹۵٪ تفاوت‌ها در جدول ۱ آورده شده است. چنان‌که ملاحظه می‌شود این آزمون‌ها تفاوت معنی‌دار بین مقادیر MSE برای دو مدل نرمال و چوله را تایید می‌کند.

همچنین چنان‌که در جدول ۱ ملاحظه می‌شود، نتایج به دست آمده از الگوریتم‌های AEMG و MCEMG مشابه هم است با این تفاوت که اجرای الگوریتم پیشنهاد شده بسیار سریع است، در صورتی که اجرای الگوریتم MCEMG ساعت‌ها به طول می‌انجامد. برای مقایسه دقت پیش‌گویی دو مدل، در موقعیت (۵/۵ و ۵/۵) مقدار پیش‌گو برای هر ۱۰۰ مجموعه داده محاسبه شده است. سپس مقدار میانگین توان دوم خطای پیش‌گویی برای دو مدل SGLM با توزیع‌های چوله نرمال بسته و نرمال محاسبه و به ترتیب مقادیر ۱/۶۴۳۲ و ۱/۸۸۴۲ حاصل شده است. این نتایج بیان‌گر آن است که توزیع چوله نرمال بسته برای متغیرهای پنهان فضایی موجب کاهش مقدار میانگین توان دوم خطای پیش‌گویی می‌شود.



نمودار ۱. الف. نمودار Q-Q، ب. بافت‌نگار متغیرهای پنهان تولید شده

جدول ۱. متوسط برآورد و میانگین مربع خطای پارامترها براساس ۱۰۰ مجموعه داده

الگوریتم	پارامتر	مدل			
		N SGLM		CSN SGLM	
		MSE	اریبی	MSE	اریبی
AEMG	β_1	۰/۳۴۵۵	-۰/۰۹۱۸	۰/۲۶۳۱	-۰/۰۵۷۳
	σ	۰/۰۶۱۳	-۰/۱۳۲۶	۰/۰۵۲۳	-۰/۰۳۷۳
	φ	۱/۸۲۸۹	۰/۱۸۱۹	۱/۴۹۳۱	-۰/۱۶۸۶
	λ	-	-	۱/۳۱۲۴	-۰/۴۷۲۲
		-	-	-	-
MCEMG	β_1	۰/۳۳۲۱	-۰/۰۸۸۷	۰/۲۱۴۶	-۰/۰۷۶۳
	σ	۰/۰۶۲۱	-۰/۱۲۳۶	۰/۰۵۱۷	-۰/۰۳۰۵
	φ	۱/۹۱۲۵	۰/۳۹۱۹	۱/۴۸۳۴	۰/۱۵۲۳
	λ	-	-	۱/۲۸۴۱	-۰/۴۵۷۴
		-	-	-	-

بحث و نتیجه‌گیری

وجود متغیرهای پنهان فضایی در مدل‌های SGLMM و نامعلوم بودن توزیع واقعی آن‌ها، روی برآورد پارامترهای مدل و دقت پیش‌گویی تأثیرگذار است. لذا در عمل فرض نرمال بودن متغیرهای پنهان فضایی گاهی فرض نادرست و گمراه‌کننده‌ای است و استنباط‌های آماری براساس این فرض، از جمله برآورد پارامترها و پیش‌گوها را بی‌اعتبار می‌سازد. برای افزایش دقت برآورد پارامترها و پیش‌گوها، استفاده از توزیع چوله‌نرمال بسته که کلاس بزرگتر و انعطاف‌پذیرتری از کلاس توزیع نرمال است، برای متغیرهای پنهان پیشنهاد گردید و الگوریتمی برای برآورد تقریبی پارامترهای مدل و پیش‌گویی تقریبی بهینه برای متغیرهای پنهان معرفی شدند. در بررسی شبیه‌سازی نشان داده شد که استفاده از توزیع چوله نرمال بسته به‌عنوان توزیع متغیرهای پنهان فضایی در مدل‌های SGLM به‌جای توزیع نرمال، موجب افزایش دقت برآورد پارامترها و پیش‌گوها می‌گردد. همچنین الگوریتم AEMG برای به‌دست آوردن پارامترهای مدل معرفی شد به‌طوری که استفاده از الگوریتم AEMG در مقایسه با الگوریتم MCEMG باعث ارائه نتایج مشابه و کاهش زمان محاسبات می‌گردد. زمان اجرای الگوریتم MCEMG با رایانه شخصی با مشخصات Cpu=2.5 Ghz و نرم‌افزار Matlab برای هر مجموعه داده حدود چهار دقیقه و برای ۱۰۰ مجموعه داده حدود سه روز و اجرای الگوریتم AEMG برای هر مجموعه داده حدود بیست ثانیه و برای ۱۰۰ مجموعه داده حدود نیم ساعت به طول انجامید. لذا تفاوت زیادی در زمان اجرای الگوریتم‌ها وجود دارد، در صورتی‌که در نتایج تفاوت زیادی مشاهده نمی‌شود.

قدردانی و تشکر

نویسندگان از داوران محترم مجله که نظرات ارزنده آن‌ها موجب بهبود مقاله گردید و از حمایت قطب علمی داده‌های ترتیبی و فضایی دانشگاه فردوسی مشهد قدردانی می‌کنند.

منابع

1. NE Breslow, D. G. Clayton, "Approximate Inference in Generalized Linear Mixed Models", *Journal of the American Statistical Association*, 88 (1993) 9-25.
2. O. F. Christensen, "Monte Carlo Maximum Likelihood in Model-Based Geostatistics", *Journal of Computational and Graphical Statistics*, 13 (2004) 702-718.
3. J. Dominguez-Molina, G. Gonzalez-Farias, A. Gupta, "The Multivariate Closed Skew Normal Distribution", *Technical Report 03-12*, Department of Mathematics and Statistics, Bowling Green State University (2003).
4. P. Diggle, J. A. Tawn, R. A. Moyeed, "Model-Based Geostatistics (with Discussion)", *Journal of the Royal Statistical Society. Series C. Applied Statistics* 47 (1998) 299-350.
5. F. Hosseini, J. Eidsvik, M. Mohammadzadeh, "Approximate Bayesian Inference in Spatial GLMM with Skew Normal Latent Variables", *Computational Statistics and Data Analysis*, 55 (2011) 1791-1806.
6. P. McCullagh, J. A. Nelder, *Generalized Linear Model*, Chapman and Hall, London (1989).
7. M. Mohammadzadeh, F. Hosseini, "Maximum-Likelihood Estimation for Spatial GLM Models with Closed-Skew Normal Latent Variables", *Procedia Environmental Science*, 3 (2011) 63-68.
8. H. Zhang, "On Estimation and Prediction for Spatial Generalized Linear Mixed Models", *Biometrics*, 58 (2002) 129-136.