

مدل‌های رگرسیون لجستیک و پروبیت فضایی برای تحلیل داده‌های یخ‌زدگی گیاهان در استان مازندران

وحید رضائی‌تبار، محسن محمدزاده^۱

گروه آمار، دانشگاه تربیت مدرس

تاریخ دریافت: ۱۳۸۸/۳/۵ تاریخ پذیرش: ۱۳۹۱/۳/۱۶

چکیده: تحلیل رگرسیون‌های لجستیک و پروبیت که در مدل‌بندی متغیرهای پاسخ دودویی به کار می‌روند با فرض استقلال خطاها صورت می‌گیرد. اما در عمل با موارد زیادی مانند داده‌های فضایی مواجه می‌شویم که مشاهدات دودویی از لحاظ موقعیت قرار گرفتن در فضای مورد مطالعه به یکدیگر وابسته‌اند و لازم است همبستگی آن‌ها در تحلیل رگرسیون لجستیک و پروبیت منظور گردد. معمولاً در آمار فضایی، پیش‌گویی فضایی برای داده‌های دودویی با استفاده از کریگیدن نشانگر صورت می‌پذیرد. در این روش لازم است ساختار همبستگی داده‌ها از طریق تغییرنگار نشانگر تعیین شود که به نوبه خود دشوار و موثر در کارایی نتایج است. در این مقاله پیش‌گویی داده‌های دودویی فضایی با استفاده از رگرسیون لجستیک و پروبیت فضایی با فرض این‌که مشاهدات روی یک مستطیل توری واقع شده و خطاهای مدل خودهمبسته کامل باشند با دو رهیافت بسامدی و بی‌زی انجام پذیرفته است. سپس کارایی پیش‌گویی توسط مدل‌های ارائه شده و روش کریگیدن نشانگر در یک مطالعه شبیه‌سازی مورد استفاده قرار گرفته است. در انتها نحوه کاربست آن‌ها برای داده‌های دما در ۲۴ مرکز هواشناسی استان مازندران بررسی و مدل‌های مناسب ارائه گردیده است.

واژه‌های کلیدی: داده‌های فضایی دودویی، رگرسیون لجستیک و پروبیت فضایی، شبکه مستطیلی

رده‌بندی ریاضی: ۶۲J۱۲ و ۶۲J۰۲

۱- مقدمه

در تحلیل رگرسیون با این فرض که خطاها ناپسته و دارای توزیع نرمال با میانگین صفر و واریانس ثابت هستند، ارتباط بین متغیرهای پاسخ و تبیینی مورد بررسی قرار می‌گیرد. در مسائلی که متغیر پاسخ، دودویی با دو مقدار ۰ و ۱ است می‌توان از مدل‌های خطی تعمیم یافته، استفاده نمود [۱]. در این مدل، تابعی از میانگین متغیر پاسخ با متغیرهای تبیینی رابطه خطی دارد، که تابع پیوند^۱ نامیده می‌شود. وقتی متغیر پاسخ دارای توزیع برنولی با احتمال موفقیت p باشد و تابع پیوند، لوجیت اختیار شود، مدل لجستیک حاصل می‌شود و چنانچه تابع پیوند، تابع توزیع تجمعی نرمال استاندارد باشد، مدل خطی تعمیم یافته را مدل پروبیت نامند. در عمل با موارد زیادی مواجه می‌شویم که مشاهدات دودویی از لحاظ موقعیت قرار گرفتن در فضای مورد مطالعه به یکدیگر وابسته‌اند، بنابراین خطاهای مدل نیز همبسته خواهند بود و لازم است این همبستگی فضایی در مدل‌بندی داده‌ها لحاظ شود. معمولا داده‌های دودویی فضایی با استفاده از کریگیدن نشانگر^۲ تحلیل می‌شوند، که اولین بار توسط مرجع [۲] به عنوان یک روش ناپارامتری برای پیش‌گویی غیرخطی از مشاهدات مورد استفاده قرار گرفت. باسو و رینسل [۳] با فرض آن‌که داده‌ها بر روی یک مستطیل توری مشاهده شده باشند و خطاهای مدل در هر موقعیت فقط وابسته به موقعیت‌های قبلی خود باشند، با روش ماکسیمم درست‌نمایی پارامترهای مدل را برآورد کردند. این مدل‌ها توسط رینسل و چینگ [۴] برای داده‌های فضایی بسط و توسعه داده شدند. همچنین رهیافت بیزی این مدل‌ها اولین بار در مطالعه انجام شده از سوی بست [۵] مورد توجه قرار گرفت و توسط بانرجی و همکاران [۶] توسعه داده شد.

یکی از مراحل مهم پیش‌گویی داده‌های فضایی دودویی، تعیین ساختار همبستگی داده‌ها از طریق تغییرنگار نشانگر است [۷]. در مقابل در مدل‌های رگرسیون لجستیک و پروبیت داده‌ها به صورت خودهمبسته در نظر گرفته می‌شوند و نیاز به تعیین ساختار همبستگی داده‌ها از طریق تغییرنگار ندارند. لذا هدف این مقاله معرفی مدل‌های رگرسیون لجستیک و پروبیت فضایی و مقایسه آن‌ها با کریگیدن نشانگر به لحاظ میزان دقت پیش‌گوهی حاصل از آن‌ها است. این مهم توسط مطالعه‌ای شبیه‌سازی صورت گرفته و نشان داده شده است برآزش مدل‌های رگرسیون لجستیک و پروبیت نه تنها ساده‌تر از کریگیدن نشانگر هستند بلکه برای داده‌های فضایی دودویی که بر شبکه‌ای منظم مشاهده شده باشند از دقت بالاتری نیز برخوردارند.

در این مقاله روش باسو و رینسل برای مدل‌بندی داده‌های دودویی فضایی با رگرسیون‌های لجستیک و پروبیت فضایی برای حالتی که خطاهای مدل خودهمبسته کامل باشند با دو

رهیافت بسامدی و بیزی توسعه داده شده است. در نهایت، کارایی این مدل‌ها و کریگیدن نشانگر در یک مطالعه شبیه‌سازی مورد مقایسه قرار گرفته و نحوه کاربست آنها درخصوص داده‌های دما و یخ‌زدگی گیاهان در ۲۴ مرکز هواشناسی استان مازندران بررسی شده است. بخش ۴ فقط به منظور ارائه نحوه تحلیل مدل‌های رگرسیون لجستیک و پروبیت با رهیافت بیزی ارائه شده و بدیهی است در شرایطی که فرض‌های اولیه مسئله برقرار باشند همواره این رهیافت نتایج بهتری از روش بسامدی ارائه می‌کند. در عین حال دقت پیش‌گوهای حاصل از روش‌های بسامدی و بیزی با روش اعتبارسنجی متقابل مورد مقایسه قرار گرفته و مدل برتر معرفی شده است.

۲- مدل‌های آماری

فرض کنید متغیر پاسخ دودویی y_{ij} ، در مکان $s = (i, j)$ روی شبکه مستطیلی $A_{mn} = \{(i, j) : i = 1, \dots, m, j = 1, \dots, n\}$ دارای توزیع برنولی با پارامتر p_{ij} مشاهده شده باشد به‌طوری که

$$y_{ij} = x'_{ij}\beta + z_{ij} \quad (i, j) \in A_{mn} \quad (1)$$

که در آن، x_{ij} بردار k بعدی متغیرهای تبیینی، z_{ij} عبارت خطا و β بردار k بعدی ضرایب رگرسیون هستند. باسو و رینسل [۳] با فرض آن که خطاها در مدل (۱) خودهمبسته ناقص به-صورت $z_{ij} = \alpha_1 z_{i-1,j} + \alpha_2 z_{i,j-1} + \alpha_3 z_{i-1,j-1}$ باشند، پارامترهای مدل را به‌روش ماکسیمم درست‌نمایی برآورد کردند. در این مقاله روش باسو و رینسل برای خطاهای خودهمبسته فضایی کامل به فرم‌های

$$z_{ij}^{(1)} = \alpha_1 z_{i-1,j} + \alpha_2 z_{i,j-1} + \alpha_3 z_{i+1,j} + \alpha_4 z_{i,j+1} + \varepsilon_{ij} \quad (2)$$

$$z_{ij}^{(2)} = \alpha_1 z_{i-1,j} + \alpha_2 z_{i,j-1} + \alpha_3 z_{i+1,j} + \alpha_4 z_{i,j+1} + \alpha_5 z_{i-1,j-1} + \alpha_6 z_{i+1,j-1} + \alpha_7 z_{i-1,j+1} + \alpha_8 z_{i+1,j+1} + \varepsilon_{ij}, \quad (3)$$

$$z_{ij}^{(r)} = \alpha_1 z_{i-1,j} + \alpha_2 z_{i,j-1} + \alpha_3 z_{i+1,j} + \alpha_4 z_{i,j+1} - \alpha_1 \alpha_2 z_{i-1,j-1} - \alpha_1 \alpha_3 z_{i+1,j-1} - \alpha_1 \alpha_4 z_{i+1,j+1} + \varepsilon_{ij}, \quad (4)$$

توسعه داده می‌شود. برای این کار مدل (۱) با خطاهای خودهمبسته فضایی کامل را می‌توان به‌صورت کلی

$$y_{ij} = x'_{ij}\beta + z_{ij}^{(w)}, \quad w = 1, 2, 3 \quad (5)$$

با امید ریاضی و واریانس زیر نوشت.

$$p_{ij}^{(w)} = E(y_{ij}) = x'_{ij}\beta + z_{ij}^{(w)} - \varepsilon_{ij}, \quad w = 1, 2, 3 \quad (6)$$

$$Var(y_{ij}) = p_{ij}^{(w)}(1 - p_{ij}^{(w)}).$$

مدل لجستیک فضایی

با انتخاب تابع پیوند لوجیت، مدل رگرسیون لجستیک فضایی به صورت

$$\text{logit}(p_{ij}^{(w)}) = \ln \frac{p_{ij}^{(w)}}{1 - p_{ij}^{(w)}} = x'_{ij}\beta + z_{ij}^{(w)} - \varepsilon_{ij}, \quad w = 1, 2, 3 \quad (7)$$

است. برای نمونه تصادفی $y = (y_{11}, \dots, y_{mn})$ ، که در آن y_{ij} از توزیع برنولی با پارامتر $P(y_{ij} = 1) = p_{ij}$ پیروی می‌کند، با قرار دادن $p = (p_{11}, \dots, p_{mn})$ تابع درست‌نمایی به صورت

$$\begin{aligned} L(p; y) &= p(y_{11} | y_{12}, \dots, y_{mn}; p) p(y_{12} | y_{12}, \dots, y_{mn}; p) \dots p(y_{mn}; p) \\ &= p_{11}^{y_{11}} (1 - p_{11})^{1 - y_{11}} \sum_{y_{11}} p(y_{11} | y_{12}, \dots, y_{mn}; p) \dots \sum_{y_{11}, \dots, y_{mn-1}} p(y_{11}, \dots, y_{mn}; p) \end{aligned}$$

خواهد بود که به دلیل پیچیدگی، ماکسیمم کردن آن دشوار است. لذا از لگاریتم تابع شبه-درست‌نمایی یعنی ضرب توابع چگالی کناری شرطی به صورت

$$\ell(p, y) = \sum_{(i,j) \in A_{mn}} y_{ij} \ln \frac{p_{ij}}{1 - p_{ij}} + \sum_{(i,j) \in A_{mn}} \ln(1 - p_{ij}) \quad (8)$$

استفاده می‌شود [۸]. اکنون با جایگذاری (۷) در رابطه (۸) لگاریتم تابع شبه‌درست‌نمایی برای β به صورت

$$\ell(\beta, \alpha, y) = \sum_{(i,j) \in A_{mn}} y_{ij} c_{ij}^{(w)} - \sum_{(i,j) \in A_{mn}} \ln(1 + e^{c_{ij}^{(w)}}) \quad (9)$$

حاصل می‌شود، که در آن

$$c_{ij}^{(w)} = x'_{ij}\beta + z_{ij}^{(w)}, \quad w = 1, 2, 3.$$

اکنون با مشتق گرفتن از رابطه (۹) نسبت به پارامترهای α و β و حل معادلات به روش‌های عددی برآورد ضرایب مدل، امید ریاضی و واریانس آن‌ها قابل محاسبه‌اند.

مدل پروبیت فضایی:

چنانچه تابع پیوند پروبیت اختیار شود، مدل پروبیت فضایی به صورت

$$\text{probit}(p_{ij}^{(w)}) = \Phi^{-1}(p_{ij}^{(w)}) = x'_{ij}\beta + z_{ij}^{(w)} - \varepsilon_{ij} \quad (10)$$

خواهد بود، که در آن $\Phi^{-1}(\cdot)$ معکوس تابع توزیع تجمعی نرمال استاندارد است. در این صورت با جایگذاری

$$p_{ij}^{(w)} = \int_{-\infty}^{x'_{ij}\beta + z_{ij}^{(w)} - \varepsilon_{ij}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx \quad (i, j) \in A_{mn}$$

در رابطه (۸)، لگاریتم تابع شبه‌درست‌نمایی براساس مدل پروبیت به دست می‌آید.

کریگیدن نشانگر: فرض کنیم $N = nm$ و $I = \{I(s_1, z), \dots, I(s_N, z)\}$ مشاهدات فضایی دودویی از میدان تصادفی $\{I(s, z) : s \in A_{mn}\}$ با نیم‌تغییرنگار نشانگر $\gamma_z(h)$ باشند، که در آن آستانه z به صورت

$$I(s, z) = \begin{cases} 1 & Z(s) \leq z \\ 0 & Z(s) > z \end{cases}$$

تعریف می‌شود. در این صورت کریگیدن نشانگر این میدان تصادفی در یک موقعیت جدید s_0 براساس مشاهدات I عبارتست از

$$\hat{I}(s_0, z) = \sum_{i=1}^N \lambda_i(z) I(s_i, z)$$

به طوری که بردار ضرایب آن با شرط $\sum_{i=1}^N \lambda_i(z) = 1$ از رابطه

$$\lambda'_z = \left(\gamma_z + J \frac{1 - J' \Gamma_z^{-1} \gamma_z}{J' \Gamma_z^{-1} J} \right) \Gamma_z^{-1} \quad (11)$$

به دست می آید، که در آن $\gamma_z = (\gamma_z(s_1 - s_0), \dots, \gamma_z(s_N - s_0))'$ ماتریسی $N \times N$ با (i, j) امین درایه $\gamma_z(s_i - s_j)$ و J یک بردار N تایی با درایه های ۱ است. میانگین توان دوم خطای پیش گو نیز به صورت زیر تعیین می شود [۷]:

$$\sigma^z(s_0, z) = \gamma_z' \Gamma_z^{-1} \gamma_z - \frac{(J' \Gamma_z^{-1} \gamma_z - 1)^2}{J' \Gamma_z^{-1} \gamma_z J}$$

۳- تحلیل مدل های رهیافت بیزی

پاسکوتو [۹] و بیلی [۱۰] تحلیل داده های دودویی فضایی با رهیافت بیزی را مورد مطالعه قرار داده و بانرجی و همکاران [۶] از آن در اپیدمیولوژی استفاده کردند، که در این بخش با الهام از روش آن ها رگرسیون های لجستیک و پروبیت فضایی بیزی ارائه می شوند. فرض کنیم متغیر پاسخ y_{ij} در موقعیت $(i, j) \in A_{mn}$ دارای توزیع برنولی با احتمال موفقیت

$$P_{ij}^{(w)} = \text{logit}^{-1}(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij}) = \frac{\exp(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij})}{1 + \exp(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij})}$$

است. با فرض آن که بردار β دارای توزیع پیشین نرمال k متغیره $N_k(\theta, \Sigma)$ باشد، توزیع پسین آن به صورت پیچیده ی

$$\begin{aligned} \pi(\beta | y) &\propto \phi_k(\beta; \theta, \Sigma) p(y_{11}, \dots, y_{mn}) \\ &= \phi_k(\beta; \theta, \Sigma) p(y_{11} | y_{12}, \dots, y_{mn}) p(y_{12} | y_{13}, \dots, y_{mn}) \dots p(y_{mn}) \\ &= \phi_k(\beta; \theta, \Sigma) p_{11}^{y_{11}} (1 - p_{11})^{1 - y_{11}} \sum_{y_{12}} (y_{12} | y_{13}, \dots, y_{mn}) \dots \sum_{y_{11}, \dots, y_{mn-1}} p(y_{11}, \dots, y_{mn}) \end{aligned}$$

خواهد بود، که در آن $\phi_k(\beta; \theta, \Sigma)$ چگالی توزیع نرمال k متغیره با میانگین θ و ماتریس کواریانس Σ است. با استفاده از تابع شبه درست نمایی

$$\begin{aligned} p(y_{11}, \dots, y_{mn}; p) &\approx p(y_{11} | p_{12}, \dots, y_{mn}; p) \\ &\quad \times p(y_{12} | p_{11}, y_{13}, \dots, y_{mn}; p) \\ &\quad \vdots \\ &\quad \times p(y_{mn} | y_{11}, \dots, y_{m-1}; p) \end{aligned}$$

$$= p_{11}^{y_{11}} (1-p_{11})^{1-y_{11}} \dots p_{mn}^{y_{mn}} (1-p_{mn})^{1-y_{11}}$$

$$\prod_{ij} \text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij} \right)^{y_{ij}} \left[1 - \text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij} \right) \right]^{(1-y_{ij})}$$

توزیع پسینی به صورت

$$\pi(\beta | y) \propto \prod_{ij} \text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij} \right)^{y_{ij}} \left[1 - \text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij} \right) \right]^{(1-y_{ij})} \quad (12)$$

خواهد شد. برای نمونه‌گیری از توزیع پسین (۱۲) با الگوریتم متروپولیس-هستینگز می‌توان توزیع t -استیودنت چند متغیره با k درجه آزادی به صورت

$$q(\beta | y) \propto \left[1 + k^{-1} (\beta - \beta_0)' \Sigma^{-1} (\beta - \beta_0) \right]^{-\frac{k+p}{2}}$$

را به عنوان تابع چگالی پیشنهادی انتخاب کرد [۱۱]. برای کم کردن حساسیت توزیع پسینی به ابرپارامترها، از توزیع‌های پیشینی زیر استفاده شده است:

$$y_{ij} \sim \text{Ber} \left[\text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij} \right) \right]$$

$$\beta \sim N_k(\theta, \Sigma)$$

$$Z_{ij} | \sigma^2 \sim N(0, 100 \cdot \sigma^2)$$

$$\sigma^2 \sim f(\sigma^2) = (1 + \sigma^2)^{-\tau}.$$

چون اطلاعاتی در مورد پارامترهای توزیع پیشینی در اختیار نیست و از طرفی اگر پیشینی ناسره برای پارامترها در نظر گرفته شود، سره بودن توزیع پسینی باید اثبات شود، برای $Z_{ij} | \sigma^2$ توزیع ناآگاهی بخش در نظر گرفته شده است که واریانس آن بسیار زیاد و معادل $100 \cdot \sigma^2$ است. در این صورت با توجه به فرض استقلال β و σ^2 توزیع پسینی به صورت

$$\pi(\beta, Z_{ij}, \sigma^2 | y_{ij}) \propto p(y_{ij} | \beta, Z_{ij}, \sigma^2) \pi(\beta, Z_{ij}, \sigma^2)$$

$$= p(y_{ij} | \beta, Z_{ij}, \sigma^2) \pi(Z_{ij} | \sigma^2) \pi(\beta) \pi(\sigma^2)$$

$$\propto \prod_{ij} C \text{Ber} \left[\text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} \right) \right] \phi_k(\beta; \theta, \Sigma) \phi(Z_{ij}; 0, 100 \cdot \sigma^2) (1 + \sigma^2)^{-\tau}$$

خواهد بود. برای نمونه‌گیری از توزیع پسینی به منظور برآورد پارامترها از الگوریتم نمونه‌گیری گیبز استفاده می‌شود. برای این منظور توزیع‌های شرطی کامل به صورت زیر به دست آمده‌اند:

$$\pi(\beta | y_{ij}, Z_{ij}, \sigma^r) \propto \prod_{ij} C_{\gamma} \text{Ber} \left[\text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} \right) \right] \phi_k(\beta; \theta, \Sigma)$$

$$\pi(Z_{ij} | y_{ij}, \beta, \sigma^r) \propto \prod_{ij} C_{\gamma} \text{Ber} \left[\text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} \right) \right] \phi(Z_{ij}; \sigma, \mathbf{1}, \sigma^r)$$

$$\pi(\sigma^r | y_{ij}, \beta, Z_{ij}) \propto \prod_{ij} C_{\gamma} \text{Ber} \left[\text{logit}^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} \right) \right] (1 + \sigma^r)^{-r}$$

برای مدل پروبیت فضایی با فرض‌های

$$y_{ij} \sim \text{Ber} \left[\Phi^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij} \right) \right], w = 1, 2, 3$$

$$\beta \sim N_k(\theta, \Sigma)$$

توزیع پسینی به صورت

$$\pi(\beta | y) \propto \phi_k(\beta; \theta, \Sigma) \prod_{ij} \Phi^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij} \right)^{y_{ij}} \times \left[1 - \Phi^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij} \right) \right]^{(1-y_{ij})} \quad (13)$$

خواهد شد، که در این حالت نیز برای نمونه‌گیری از توزیع پسینی از الگوریتم متروپلیس-هستینگر با تابع چگالی کاندید t -استیودنت چند متغیره با درجه آزادی k استفاده می‌شود. برای برآورد β با تابع زیان درجه دوم می‌توان از توزیع‌های شرطی کامل نمونه‌گیری کرده، میانگین آن‌ها را به عنوان برآورد پارامترهای مورد نظر استفاده نمود.

۴- ارزیابی مدل‌ها

در این مقاله برای ارزیابی پیش‌گویی فضایی با مدل‌های رگرسیون لجستیک، پروبیت و کریگیدن نشانگر از ملاک ریشه دوم میانگین توان دوم خطاهای اعتبارسنجی متقابل و ملاک آکائیک^۱ (AIC) و برای ارزیابی مدل‌های رگرسیون لجستیک و پروبیت حاصل از توزیع‌های پیشینی مختلف، از ملاک اطلاع شوارتز یا ملاک اطلاع بیزی^۲ (AIC) استفاده شده است.

در روش اعتبارسنجی متقابل یک مشاهده را از مجموعه مشاهدات حذف نموده و مقدار آن براساس $N - 1$ مشاهده باقیمانده پیش‌گویی می‌شود. این عمل را برای تمام N موقعیت تکرار کرده و ملاک ریشه دوم میانگین توان دوم خطاهای اعتبارسنجی متقابل به صورت

1- Akaike Information Criterion

2- Bayes Information Criterion

$$RMSE = \left[\frac{1}{N} \sum_{i=1}^N e_{(-i)}^2 \right]^{\frac{1}{2}} \quad (14)$$

محاسبه می‌شود، که در آن $e_{(-i)}$ خطای پیش‌گوی محاسبه شده براساس مشاهدات به جز مشاهده i ام است. بدیهی است هر چقدر مقدار این ملاک به صفر نزدیکتر باشد، پیش‌گو از دقت بیشتری برخوردار است.

معیارهای AIC و BIC برای داده‌های فضایی به ترتیب به صورت

$$AIC = -2\ell p(\hat{p}, y) + 2k$$

$$BIC = -2\ell p(\hat{p}, y) + k \log N$$

می‌باشند، که در آن‌ها $\ell p(\hat{p}, y)$ ماکسیمم لگاریتم تابع شبه‌درست‌نمایی و k تعداد پارامترها است [۱۲]. هر چه مقدار AIC و BIC برای یک مدل کمتر باشد، برازش مدل به داده‌ها مناسب‌تر و بر مدل‌های رقیب ترجیح دارد.

۵- شبیه‌سازی

در این بخش مدل‌های رگرسیون لجستیک، پروبیت فضایی و کریگیدن نشانگر براساس معیارهای RMSE و BIC مورد مقایسه قرار می‌گیرند. برای این منظور x_1 از توزیع نرمال با میانگین صفر و واریانس‌های 0.5 ، 1 و 2 و x_2 از توزیع χ^2 دو با 5 درجه آزادی، به عنوان متغیرهای تبیینی در نظر گرفته شده‌اند. با توجه به رابطه (۷) و اختیار کردن $w = 1$ مقادیر $\{z_{i-1,j}, z_{i,j-1}, z_{i+1,j}, z_{i,j+1}\}$ از توزیع نرمال استاندارد و مقادیر y_{ij} نیز از توزیع برنولی با احتمال موفقیت $P_{ij}^{(1)}$ تولید شده‌اند. با استفاده از الگوریتم نیوتن رافسون و مقادیر اولیه $\alpha_0 = 0.05$ ، $\beta_0 = -0.05$ ، $\beta_1 = 0.025$ ، $\beta_2 = 0.01$ ، $\alpha_1 = 0.02$ ، $\alpha_2 = -0.001$ و $\alpha_3 = 0.004$ پارامترهای $\alpha = (\alpha_0, \alpha_1, \alpha_2, \alpha_3)$ و $\beta = (\beta_0, \beta_1, \beta_2)$ برآورد شده‌اند. برای پنج حجم نمونه 25 ، 50 ، 100 ، 200 و 500 شبیه‌سازی 10000 بار تکرار شده و مقادیر RMSE محاسبه شده‌اند. از آنجا که مقادیر چارک‌های RMSE بیان‌گر متوسط و کران‌های پایین و بالای خطا هستند، مقادیر آن‌ها برای واریانس‌های 0.5 ، 1 و 2 محاسبه و در جداول ۱ تا ۳ ارائه شده‌اند. همان‌طور که ملاحظه می‌شود مقادیر RMSE با افزایش حجم نمونه برای هر سه مدل رگرسیون لجستیک، پروبیت فضایی و کریگیدن نشانگر کاهش می‌یابند. ولی مدل‌های رگرسیون لجستیک، پروبیت فضایی نسبت به کریگیدن نشانگر همواره از خطای کمتری برخوردار هستند. این نتایج دور از انتظار نیستند؛ زیرا با توجه به توری بودن داده‌ها

انتظار می‌رود مدل‌های رگرسیون لجستیک، پروبیت فضایی با خطاهای خودهمبسته، برازش مناسب‌تری نسبت به کریگیدن نشانگر داشته باشند.

جدول ۱: مقدار RMSE لجستیک، پروبیت فضایی و کریگیدن نشانگر برای $\sigma^2 = 0.5$

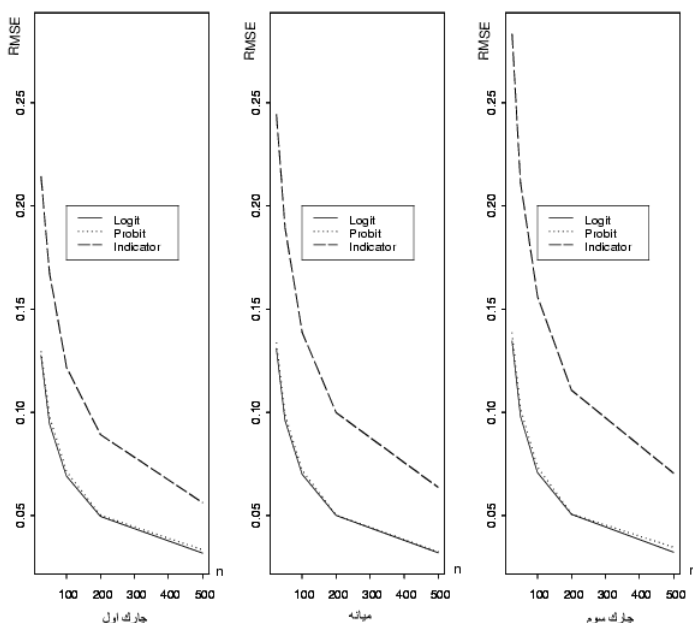
حجم نمونه	لجستیک فضایی		پروبیت فضایی		کریگیدن نشانگر	
	چارک اول	میانه	چارک اول	میانه	چارک اول	میانه
25	0.123	0.158	0.198	0.125	0.189	0.218
50	0.097	0.124	0.157	0.099	0.159	0.193
100	0.073	0.099	0.123	0.079	0.110	0.168
200	0.053	0.066	0.081	0.057	0.071	0.149
500	0.033	0.041	0.051	0.036	0.045	0.075

جدول ۲: مقدار RMSE لجستیک، پروبیت فضایی و کریگیدن نشانگر برای $\sigma^2 = 1$

حجم نمونه	لجستیک فضایی		پروبیت فضایی		کریگیدن نشانگر	
	چارک اول	میانه	چارک اول	میانه	چارک اول	میانه
25	0.123	0.162	0.194	0.138	0.201	0.218
50	0.110	0.124	0.158	0.111	0.159	0.195
100	0.071	0.090	0.132	0.075	0.099	0.173
200	0.055	0.069	0.082	0.056	0.071	0.155
500	0.037	0.045	0.057	0.038	0.046	0.135

جدول ۳: مقدار RMSE لجستیک، پروبیت فضایی و کریگیدن نشانگر برای $\sigma^2 = 2$

حجم نمونه	لجستیک فضایی		پروبیت فضایی		کریگیدن نشانگر	
	چارک اول	میانه	چارک اول	میانه	چارک اول	میانه
25	0.135	0.136	0.138	0.141	0.214	0.246
50	0.120	0.123	0.130	0.124	0.197	0.199
100	0.088	0.091	0.135	0.090	0.171	0.183
200	0.057	0.075	0.089	0.068	0.078	0.161
500	0.040	0.047	0.059	0.051	0.069	0.142



شکل ۱: مقادیر RMSE لجستیک، پروبیت فضایی و کریگیدن نشانگر

همچنین با مقایسه مقادیر جداول ۱ تا ۳ ملاحظه می‌شود که با افزایش واریانس متغیرهای تبیینی، مقادیر RMSE افزایش پیدا می‌کند، که علت آن افزایش تغییرات مقادیر متغیرهای پاسخ و در نتیجه افزایش مقدار خطا به دلیل پراکندگی مقادیر متغیرهای تبیینی است. شکل ۱ که نمودار مقادیر RMSE را برای خطای خودهمبسته $z_{ij}^{(2)}$ نشان می‌دهد، نشانگر دقت بیشتر مدل‌های لجستیک و پروبیت فضایی نسبت به کریگیدن نشانگر هستند. زیرا این مدل‌ها نسبت به مدل با خطای $z_{ij}^{(1)}$ کامل‌تر هستند. بنابراین وقتی موقعیت داده‌ها به صورت توری باشند، دو مدل پروبیت و لجستیک فضایی با خطاهای خودهمبسته کامل عمل کرد بهتری دارد و مقادیر پیش‌گویی دقیق‌تری نسبت به کریگیدن نشانگر ارائه می‌کند.

۶- مثال کاربردی

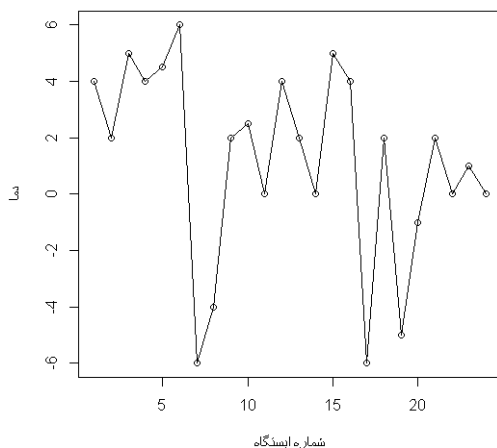
محمدزاده، و کیاپور [۱۳] داده‌های دما در ۲۴ ایستگاه هواشناسی استان مازندران را به عنوان مجموعه‌ای از داده‌های دودویی فضایی مورد مطالعه قرار داده و با استفاده از کریگیدن نشانگر، احتمال یخ‌زدگی در ایستگاه‌ها را به‌دست آورده‌اند. در این بخش با استفاده از مدل‌های

رگرسیون لجستیک و پروبیت فضایی نقشه پیش‌گویی فضایی احتمال یخ‌زدگی در استان مازندران تعیین و دقت آن با نتایج به‌دست آمده توسط مرجع [۱۳] مورد مقایسه قرار می‌گیرد. از آنجا که احتمال یخ‌زدگی نواحی نزدیک در هر منطقه به یکدیگر مرتبط هستند و این ارتباط در نواحی دورتر از هم کاهش می‌یابد، می‌توان از مدل‌های رگرسیون لجستیک و پروبیت فضایی با خطاهای خودهمبسته کامل نیز برای پیش‌گویی احتمال یخ‌زدگی در هر موقعیت استفاده کرد. شکل ۲ موقعیت قرار گرفتن مشاهدات دمای روز پانزدهم بهمن ماه ۱۳۸۴ را در ۲۴ ایستگاه هواشناسی استان مازندران برحسب سانتی‌گراد نشان می‌دهد. در داده‌های نشانگر، کد یک به ناحیه یخ‌زده (دمای منفی) و کد صفر به ناحیه یخ‌نزده اختصاص داده شده است.



شکل ۲: موقعیت ایستگاه‌های هواشناسی استان مازندران

با فرض آن‌که دما از مدل خودهمبسته کامل پیروی می‌کند، موقعیت داده‌ها در یک مستطیل توری به صورت $A_{mn} = \{(i, j) : i = 1, \dots, 4, j = 1, \dots, 4\}$ در نظر گرفته شده است. میزان دما بدون در نظر گرفتن محل قرار گرفتن آن‌ها نیز در شکل ۳ نشان داده شده است. در مدل‌های رگرسیون لجستیک و پروبیت فضایی به ازای $i, j = 1$ ، مقادیر اولیه‌ای برای دما مورد نیاز است. بنابراین هر چه مرتبه خطا در مدل‌های رگرسیون لجستیک و پروبیت فضایی بالاتر باشد، مقادیر اولیه زیادتری مورد نیاز است. لذا امکان بروز خطا افزایش می‌یابد. بنابراین لزوماً افزایش مرتبه یک مدل، به معنی بهتر بودن مدل نیست، بلکه باید به همبستگی آن از لحاظ مکانی با سایر مشاهدات توجه شود. با توجه به شکل ۳ ملاحظه می‌شود که مشاهدات دما در ۲۴ شهر استان مازندران به‌صورت تصادفی در بازه $[-6, 6]$ توزیع شده‌اند. بنابراین یک ایده مناسب برای تولید مقادیر اولیه می‌تواند اختیار یک نمونه تصادفی در این بازه باشد.



شکل ۳: نمودار میزان دما در ۲۴ شهر استان مازندران

جدول ۴: معیارهای ارزیابی مقایسه مدل

مدل	نوع خطا	RMSE	AIC
پروبیت فضایی	۱	۰/۴۰	۲۰/۷۷۸
	۲	۰/۱۴	۱۵/۱۱۸
	۳	۰/۱۸	۱۲/۶۵۴
لجستیک فضایی	۱	۰/۴۳	۲۲/۲۹۱
	۲	۰/۲۳	۱۷/۲۲۰
	۳	۰/۲۶	۱۹/۷۱۷
کریگیدن نشانگر	-	۰/۴۸	۲۵/۳۷۱

مقادیر ریشه دوم میانگین توان دوم خطاها و ملاک آکائیک برای مدل‌های رگرسیون لجستیک و پروبیت فضایی در جدول ۴ ارائه شده است. با توجه به این جدول، مدل مناسب برآزنده شده به داده‌های دما، بر اساس RMSE، مدل پروبیت فضایی با خطاهای خودهمبسته $z_{ij}^{(r)}$ و براساس AIC، مدل پروبیت فضایی با خطاهای خودهمبسته $z_{ij}^{(r)}$ است. چون مدل با خطای $z_{ij}^{(r)}$ از پارامترهای کمتری نسبت به مدل با خطای $z_{ij}^{(r)}$ برخوردار است، برای تحلیل داده‌های دما مناسب‌تر و به صورت زیر است:

$$\begin{aligned} \text{Probit}(P_{ij}^{(r)}) = & -0.128 - 0.934x_{ij} + 0.561z_{i-1,j} - 0.1012z_{i,j-1} + 0.767z_{i+1,j} \\ & - 0.0008z_{i,j+1} + 0.0056z_{i-1,j-1} + 0.0077z_{i+1,j-1} + 0.0004z_{i-1,j+1} \\ & - 0.0908z_{i+1,j+1}, \quad i = 1, \dots, 4, j = 1, \dots, 6 \end{aligned} \quad (15)$$

با توجه به مدل (۱۵)، آستانه تحمل گیاهان به ازای هر متر افزایش ارتفاع، به صورت ضریبی از ۰/۰۹۳۴ کاهش می‌یابد.

برای تحلیل بیزی داده‌های دما با توزیع پیشینی $\beta \sim N_k(0, 10^4 I)$ برای تولید نمونه از توزیع پسینی

$$\begin{aligned} \pi(\beta | y) = & \frac{1}{p(y)} \prod_{ij} \log it^{-1} \left[x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij}^{(w)} \right]^{y_{ij}} \\ & \times \left[1 - \log it^{-1} \left(x'_{ij} \beta + z_{ij}^{(w)} - \varepsilon_{ij}^{(w)} \right) \right]^{(1-y_{ij})} \phi_k(\beta; \alpha, 0, 10^4 I), w = 1, 2, 3 \end{aligned}$$

از الگوریتم متروپلیس- هستینگز استفاده شده است و توزیع t -استیودنت چند متغیره با درجه آزادی $k = 5$ ، تعداد متغیرهای $p = 6$ و $\Sigma = 10^4 k$ به صورت

$$q(\beta | y) \propto \left[1 + k^{-1} (\beta - \beta_0)' \Sigma^{-1} (\beta - \beta_0) \right]^{-\frac{k+p}{2}}$$

به‌عنوان تابع چگالی کاندید انتخاب شده‌اند، که در آن β_0 بردار مقادیر اولیه است. علت این انتخاب را می‌توان کلفت‌تر بودن دم‌های تابع چگالی کاندید نسبت به توزیع نرمال دانست، که با

انتخاب درجه آزادی ۵ این امر میسر شده است. همچنین $\hat{p}(y) = \left[\frac{1}{k} \sum_{i=1}^k \frac{1}{p(y | \beta_i)} \right]^{-1}$ اختیار شده است. مقادیر BIC_p برای مدل‌های لجستیک، پروبیت فضایی و کریگیدن نشانگر در جدول ۵ بیان‌گر آن است که مدل لجستیک فضایی با خطای $z_{ij}^{(l)}$ به صورت

$$\begin{aligned} \log it(P_{ij}^{(l)}) = & 1/0.04 + 0.2541x_{ij} - 0.56z_{i-1,j} + 0.67z_{i,j-1} \\ & + 0.923z_{i+1,j} + 0.1035z_{i,j+1} \end{aligned} \quad (16)$$

مدل مناسبی برای برازش به داده‌های دما است. برای مقایسه دقت پیش‌گوهای حاصل از دو مدل بسامدی و بیزی، ملاک ریشه دوم میانگین توان دوم خطای اعتبار سنجی متقابل (RMSE) محاسبه و به ترتیب مقادیر ۰/۱۸ و ۰/۱۳ حاصل شدند، که بیان‌گر برتری پیش‌گوی بیزی بر پیش‌گوی بسامدی است.

جدول ۵: مقادیر معیار ارزیابی BIC_p

BIC_p	نوع خطا	مدل
-۴/۶۴۴	۱	پروبیت فضایی
-۲/۴۴۱	۲	
-۱/۸۹۶	۳	
-۱۱/۴۱۱	۱	لجستیک فضایی
-۶/۸۸۹	۲	
-۴/۳۰۷	۳	
۸/۸۸۸	-	کریگیدن نشانگر

۷- بحث و نتیجه‌گیری

برای مدل‌بندی داده‌های دودویی فضایی می‌توان همانند کریگیدن نشانگر، از مدل‌های رگرسیون لجستیک یا پروبیت فضایی با خطاهای خودهمبسته فضایی نیز استفاده کرد. در این مقاله با فرض توری بودن ناحیه شامل موقعیت مشاهدات، مدل‌های مذکور برای پیش‌گویی فضایی داده‌های دودویی پیشنهاد گردید. در مطالعه شبیه‌سازی نشان داده شد که برای داده‌های دودویی فضایی با موقعیت‌های توری، استفاده از مدل‌های رگرسیون لجستیک و پروبیت فضایی نسبت به کریگیدن نشانگر از خطای کمتری برخوردار است.

ویژگی تحلیل داده‌های واقع بر یک مستطیل توری بی‌نیازی از برآورد ساختار همبستگی داده‌ها در قالب توابع تغییرنگار یا هم‌تغییرنگار است. در این صورت ساختار همبستگی به صورت ساده خودهمبسته در نظر گرفته می‌شود که برآزش آن‌ها ساده‌تر از برآورد توابع تغییرنگار یا هم‌تغییرنگار است. یکی از مسائل اساسی، نحوه توری کردن ناحیه مورد مطالعه است، که باید به گونه‌ای انجام شود تا حتی الامکان گره‌ها فاقد مشاهده یا دارای چند مشاهده نباشند. تعیین توزیع پیشین مناسب در رهیافت بیزی، خود مساله‌ای است که باید مورد توجه قرار گرفته و حساسیت پیش‌گوها با انواع توزیع‌های پیشین مورد بررسی قرار گیرد.

تقدیر و تشکر

نویسندگان از داوران محترم به خاطر پیشنهادات ارزنده‌ای که موجب بهبود مقاله گردید، کمال تشکر را دارند. از حمایت قطب علمی داده‌های ترتیبی و فضایی دانشگاه فردوسی مشهد نیز قدردانی می‌شود.

مراجع

- [1] Nelder, J. and Wedderburn, R. W. M. (1972), Generalized Linear Models, *Journal of the Royal Statistical Society, Series A*, **135**, 370-384.
- [2] Journel, A. G. (1983), Nonparametric Estimation of Spatial Distribution, *Mathematical Geology*, **15**, 445-468.
- [3] Basu, S. and Reinsel, G. C. (1994), Regression Models with Spatially Correlated Errors, *Journal of the American Statistical Association*, **89**, 88-99.
- [4] Reinsel, G. C. and Cheng, W. K. (2003), Approximate ML and REML Estimation for Regression Models with Spatial or Time Series AR (1) Noise, *Statistics and Probability Letters*, **62**, 123-135.
- [5] Best, N. G. (1999), Bayesian Models for Spatially Correlated Disease and Exposure Data, *Bayesian Statistics 6*, edited by J. M. Bernardo, B. O. Berger, A. P. Dawid and A. F. M. Smith, Oxford University Press, Oxford, UK. 131-156.
- [6] Banerjee, S., Carlin, B. P. and Gelfand, A. E. (2004), Hierarchical Modeling and Analysis for Spatial Data, *Chapman and Hall*, London.
- [7] Cressie, N. (1993), Statistical for Spatial Data, *Wiley, New York*.
- [8] Besage, J. E. (1974), Spatial Interaction and Statistical Analysis of Lattice System, *Journal of the Royal Statistical Society, Series B*, **36**, 192-236.
- [9] Pascutto, C. (2000), Statistical Issues in the Analysis of Disease Mapping Data, *Statistics in Medicine*, **19**, 2493-2519.
- [10] Baily, T. (2001), Spatial Statistical Analysis in Health, *Cadernos de Saude Publica*, **17**, 1083-1098.
- [11] Marjerison, W. M. (2006), Bayesian Logistic Regression with Spatial Correlation, *MSc Thesis, Worcester Polytechnic Institute, Worcester, Massachusetts, USA*.
- [12] Molenberghs, G. (2005), Models for Discrete Longitudinal Data. *Springer Science, New York*.

[۱۳] محمدزاده، م. و کیاپور، آ. (۱۳۸۷)، پیش‌گویی فضایی داده‌های سخت و نرم برای تهیه نقشه احتمال یخ‌زدگی، مجله امیرکبیر، سال نوزدهم، شماره ۵-۶۸، ۴۴-۳۷.

Spatial Logistic and Probit Regression Models for Analysis of Frosting Data in Mazandaran Province

Vahid Rezaeitabar and Mohsen Mohammadzadeh

Department of Statistics, Tarbiat Modares University, Tehran, Iran

Abstract

Logistic and Probit regression models are usually used in binary response variable analysis based on independence assumption of the observations. But, in practice, there are many situations in which, due to their different locations in the underlying space of study, this assumption does not satisfy. In spatial statistics, it is generally supposed that the binary observations are analyzed with indicator Kriging. In this paper, we considered the spatial Logistic and Probit regression models with auto correlated errors on a rectangular grid. Also, in a simulation study, the prediction accuracy of the models has been compared. Finally, the implementation of the models for a temperature data set, reported by weather stations in Mazandaran province of Iran, is shown.

Keywords: Spatial binary data, Logistic and Probit regression, Rectangular grids.

Mathematics Subject Classification (2000): 62J12, 62J02