

مدل رگرسیون آمیخته خطی تعمیم یافته بتا فضایی- زمانی و کاربرد آن در آمارگیری هزینه و درآمد

لیدا کلهری
پژوهشکده آمار

چکیده: در سال‌های اخیر مدل‌بندی داده‌های محدود در بازه $(0, 1)$ مورد توجه قرار گرفته است و مدل رگرسیون بتا برای مدل‌بندی این‌گونه مشاهدات معرفی شده است. در این مقاله مدل‌بندی داده‌ها در بازه $(0, 1)$ با در نظر گرفتن ساختار همبستگی فضایی- زمانی تفکیک‌پذیر بررسی شده است و برآورد پارامترها با رهیافت بی‌زی به دست آمده‌اند. کاربرد این مدل برای تحلیل داده‌های آمارگیری هزینه و درآمد در شهر تهران ارائه شده است.

واژه‌های کلیدی: رگرسیون بتا، همبستگی فضایی- زمانی، رهیافت بی‌زی، آمارگیری هزینه و درآمد
کد موضوع‌بندی ریاضی (۲۰۱۰): 62M30، 62F15.

۱ مقدمه

در سال‌های اخیر تمرکز زیادی بر مدل‌بندی داده‌هایی مانند نرخ‌ها و نسبت‌ها که در بازه $(0, 1)$ تغییر می‌کنند، شده است. اهمیت این موضوع به دلیل محدودیت‌های به‌کارگیری مدل‌های موجود برای تحلیل این‌گونه داده‌ها است. به عنوان مثال، استفاده از مدل‌های رگرسیونی یا رگرسیون تعمیم‌یافته برای مدل‌بندی داده‌های محدود در بازه $(0, 1)$ می‌تواند به پیشگویی خارج از دامنه منجر شود. **فراری و کریباری (۲۰۰۴)** برای اجتناب از اعمال تبدیل به داده‌های اولیه و با در نظر گرفتن خواص انعطاف‌پذیری توزیع بتا، استفاده از مدل رگرسیونی بتا را پیشنهاد کردند که در آن متغیر پاسخ دارای توزیع بتای بازپارامتریده است. **فیگورا و همکاران (۲۰۱۳)** پارامترهای میانگین و دقت توزیع بتای بازپارامتریده را با در نظر گرفتن اثرات تصادفی مدل‌بندی نموده و برآورد پارامترهای مدل با رهیافت بی‌زی را مورد بررسی قرار داده‌اند.

کپدا و همکاران (۲۰۱۲)، **لاگوس الوارز و همکاران (۲۰۱۷)** و **کلهری و محمدزاده (۲۰۱۸)** مدل‌بندی داده‌های فضایی با استفاده از مدل رگرسیون بتای تعمیم‌یافته را مورد مطالعه قرار داده‌اند. بسیاری از انواع داده‌ها در حوزه آمار رسمی

علاوه بر همبستگی فضایی در طول زمان نیز وابسته هستند، که به عنوان داده‌های فضایی-زمانی شناخته می‌شوند. تاکنون مطالعات متعددی برای مدل‌بندی این‌گونه داده‌ها صورت گرفته است. با این وجود مدل‌بندی داده‌های فضایی-زمانی با دامنه تغییرات محدود در یک بازه، کمتر مورد توجه قرار گرفته است. در این مقاله مدل‌بندی داده‌های همبسته فضایی محدود در بازه $(0, 1)$ ، در طول زمان مورد مطالعه قرار خواهد گرفت و برآورد پارامترها با رهیافت بیزی انجام می‌شود. در بخش ۲ مدل آماری معرفی می‌شود. بخش ۳ به مطالعه شبیه‌سازی اختصاص یافته است و در بخش ۴ کاربست مدل بر داده‌های آمارگیری هزینه و درآمد ارائه شده است.

۲ مدل آماری

در این بخش ابتدا مدل رگرسیون بتا به اختصار معرفی می‌شود و پس از آن مدل مورد استفاده برای تحلیل داده‌های فضایی-زمانی در بازه $(0, 1)$ بیان خواهد شد. تاکنون طیف گسترده‌ای از مدل‌های آماری برای مدل‌بندی داده‌های فضایی-زمانی معرفی شده و قابل استفاده هستند. از آن‌جا که همبستگی فضایی و زمانی به صورت جداگانه قابل مدل‌بندی هستند، می‌توان از همه فنون موجود در سری‌های زمانی و آمار فضایی برای مدل‌بندی ساختار همبستگی فضایی-زمانی تفکیک‌پذیر بهره گرفت (محمدزاده، ۱۳۹۴). هدف این مطالعه تحلیل داده‌های فضایی-زمانی با استفاده از یک مدل تفکیک‌پذیر است و اثرات متقابل فضایی-زمانی در دامنه این تحقیق نیست.

توزیع بتا از انعطاف‌پذیری بالایی برخوردار است و با تغییر مقادیر پارامترهایش به فرم‌های متمایزی ظاهر می‌شود. با توجه به این که در مدل‌های رگرسیونی، الگوی رفتار میانگین متغیر پاسخ مشروط بر متغیرهای تبیینی مورد بررسی قرار می‌گیرد، فراری و کریباری (۲۰۰۴) توزیع بتای بازپازامتریده را برای مطالعه نسبت‌ها پیشنهاد کردند. آن‌ها پارامترهای توزیع بتا را به‌گونه‌ای بازنویسی کردند که مدل رگرسیونی بر اساس میانگین متغیر پاسخ معین شود. بنابراین توزیع بتا با پارامترهای (a, b) را با فرض $\mu = \frac{a}{a+b}$ و $\phi = a + b$ به صورت

(۱.۲)

$$\pi(y, \mu, \phi) = \frac{\Gamma(\phi)}{\Gamma(\mu\phi)\Gamma((1-\mu)\phi)} y^{\mu\phi-1} (1-y)^{(1-\mu)\phi-1} \quad 0 < y < 1,$$

در نظر گرفتند که در آن $0 < \mu < 1$ و $\phi > 0$ ، در نتیجه $E(Y) = \mu$ ، $V(Y) = \frac{\mu(1-\mu)}{1+\phi}$ ، که در آن ϕ پارامتر دقت نامیده می‌شود. با فرض ثابت بودن پارامتر دقت، مدل رگرسیون بتا توسط فراری و کریباری (۲۰۰۴) به صورت زیر

$$g(\mu_i) = \mathbf{x}_i \boldsymbol{\beta} \quad i = 1, \dots, n, \quad (2.2)$$

ارائه شد که در آن بردار متغیرهای تبیینی، $\mathbf{x}_i = (x_{i0}, \dots, x_{ip})^T$ بردار ضرایب رگرسیونی، و $g(\cdot)$ تابع پیوند مناسب است. در این مقاله، برای مطالعه اثرات فضایی-زمانی یک مدل آمیخته خطی تعمیم یافته از رگرسیون بتا با اثرات تفکیک‌پذیر معرفی می‌شود.

فرض کنید $\mathbf{y}(s, t) = (y(s_1, t_1), \dots, y(s_n, t_T))^T$ مشاهدات متغیر $Y(s, t) = (Y(s_1, t_1), \dots, Y(s_n, t_T))^T$ در n موقعیت مکانی s_1, \dots, s_n و T زمان t_1, \dots, t_T باشند. به منظور ساده‌سازی در نوشتار از این پس فرض می‌شود y_{ij} و Y_{ij} برای $i = 1, \dots, n$ و $j = 1, \dots, T$ به ترتیب نمایشگر $y(s_i, t_j)$ و $Y(s_i, t_j)$ باشند. هدف مدل‌بندی متغیرهای تصادفی همبسته فضایی-زمانی Y_{nT}, \dots, Y_{11} از توزیع بتا است که همبستگی فضایی آن‌ها از طریق یک میدان تصادفی گاوسی^۱ (GRF) و همبستگی زمانی آن‌ها از طریق یک مدل سری زمانی خودهمبسته در نظر گرفته شود.

¹Gaussian Random Field

فرض کنید $\tau = (\tau(s_1), \dots, \tau(s_n))^T \equiv (\tau_1, \dots, \tau_n)^T$ نمایانگر یک میدان تصادفی مانای مرتبه دوم گاوسی با میانگین صفر باشد، یعنی $E(\tau) = 0$. همچنین فرض کنید $\nu = (\nu(t_1), \dots, \nu(t_T))^T \equiv (\nu_1, \dots, \nu_T)^T$ نمایانگر مؤلفه‌های یک سری زمانی خودهمبسته باشند. متغیرهای تصادفی Y به شرط مؤلفه تصادفی τ و ν از یکدیگر مستقل هستند و از توزیع بتا پیروی می‌کنند. به عبارت دیگر

$$Y | \tau, \nu \sim \text{Beta}(\mu(s, t)\phi(s, t), (1 - \mu(s, t))\phi(s, t)),$$

که توزیع آن از جایگذاری $\mu(s, t)$ و $\phi(s, t)$ در (۱.۲) به دست می‌آید. بنابراین مدل آمیخته خطی تعمیم‌یافته بتا با اثرات تفکیک‌پذیر فضایی-زمانی به صورت

$$g(E(Y_{ij} | \tau_i, \nu_j)) = \mathbf{x}_{ij}\beta + \tau_i + \nu_j \quad i = s_1, \dots, s_n, \quad j = t_1, \dots, t_T, \quad (3.2)$$

تعریف می‌شود که یک مدل خطی در اثرات ثابت و تصادفی است. در این مدل $g(\cdot)$ تابع پیوند یک به یک و مشتق‌پذیر با دامنه اعداد حقیقی، $\mathbf{x}_{ij} = (x_{ij1}, \dots, x_{ijm})^T$ و $\beta = (\beta_1, \dots, \beta_m)^T$ به ترتیب بردار متغیرهای تبیینی و ضرایب رگرسیونی، τ_i نمایشگر اثر تصادفی است که ساختار همبستگی فضایی داده‌ها را در مدل وارد می‌کند و ν_j نمایانگر روند داده‌ها در طول زمان بر اساس یک مدل سری زمانی خودهمبسته مرتبه p به صورت $\nu_j = c + \sum_{k=1}^p \varphi_k \nu_{j-k}$ است، که در آن c ثابت مدل و $\varphi_1, \dots, \varphi_p$ پارامترهای مدل هستند. فرایند تصادفی $\tau = (\tau_1, \dots, \tau_n)^T$ از توزیع نرمال چند متغیره به صورت $\tau \sim N_n(0, \Sigma_\tau)$ پیروی می‌کند، که در آن

$$(\Sigma_\tau)_{i\ell} = \text{Cov}(\tau_i, \tau_\ell) = \sigma_\tau^2 \text{Corr}(\tau_i, \tau_\ell), \quad i, \ell = s_1, \dots, s_n$$

مؤلفه‌های Σ_τ تابعی از فاصله بین نقاط هستند به گونه‌ای که تشریح کننده ساختار کوواریانس فضایی باشند. بر اساس مطالعات گذشته در مدل‌های رگرسیون بتا، تابع پیوند لوجیت برای مدل‌بندی میانگین در نظر گرفته می‌شود. بنابراین (۳.۲) به صورت

$$\text{logit}(\mu_{ij}) = \mathbf{x}_{ij}\beta + \tau_i + \nu_j, \quad i = s_1, \dots, s_n, \quad j = t_1, \dots, t_T, \quad (4.2)$$

در نظر گرفته می‌شود. در نتیجه با فرض ثابت بودن پارامتر دقت توزیع بتا، مدل سلسه مراتبی به صورت زیر خواهد بود:

$$Y_{ij} | \beta, \phi, \tau_i, \nu_j \sim \text{Beta}(\mu(s_i, t_j)\phi, (1 - \mu(s_i, t_j))\phi),$$

$$\tau | \eta \sim N_n(0, \Sigma_\tau),$$

$$\nu_j | \varphi \sim N(\nu_{j-1}, \sigma_\nu^2), \quad j = 2, \dots, t_T,$$

که در آن $\eta = (\sigma_\tau^2, \psi^{-1})$ و φ به ترتیب پارامترهای مربوط به ساختار کوواریانس فضایی و مدل سری زمانی هستند. با در نظر گرفتن توزیع‌های پیشین دلخواه $\pi(\beta)$ ، $\pi(\phi)$ ، $\pi(\sigma_\tau^2)$ ، $\pi(\psi^{-1})$ و $\pi(\sigma_\nu^2)$ توزیع پسین توام پارامترها به صورت

$$f(\beta, \sigma_\tau^2, \psi^{-1}, \sigma_\nu^2, \phi, \tau, \nu | y) \propto \prod_{i=1}^{s_n} \prod_{j=1}^{t_T} f(y_{ij} | \tau_i, \nu_j, \phi) \prod_{i=1}^{s_n} f(\tau_i | \beta, \sigma_\tau^2, \psi^{-1}) \prod_{j=1}^{t_T} f(\nu_j | \sigma_\nu^2) \quad (5.2)$$

$$\times \pi(\beta) \pi(\sigma_\tau^2) \pi(\psi^{-1}) \pi(\sigma_\nu^2) \pi(\phi),$$

حاصل می‌شود که داری فرم بسته نیست و بنابراین برآورد پارامترها با استفاده از نمونه‌گیری گیبز از توزیع‌های شرطی کامل حاصل می‌شود.

۳ مطالعه شبیه سازی

در این بخش با استفاده مطالعه شبیه سازی برآورد پارامترهای مدل با رهیافت بیزی به دست می آیند. تولید داده ها با در نظر گرفتن مقادیر $\beta = (\beta_0, \beta_1)^T = (-1, 1/5)^T$ ، $\varphi = (\varphi_1, \varphi_2, \varphi_3) = (0/45, -0/2, 0/15)$ و با به کار بردن تابع کوواریانس همسانگرد نمایی $(\Sigma_\tau)_{il} = \sigma^2 \exp(-|s_i - s_l| \psi^{-1})$ انجام شده است. متغیرهای کمکی از توزیع یکنواخت تولید شده اند. پاسخ ها به شرط متغیر پنهان و مؤلفه زمان از توزیع بتا با پارامترهای $(\mu_{ij} \phi, (1 - \mu_{ij}) \phi)$ تولید شده اند که در آن، $\phi = 50$ و $\mu_{ij} = (\exp(x_{ij} \beta + \tau_i + \nu_j)) / (1 + \exp(x_{ij} \beta + \tau_i + \nu_j))$ برای برآورد بیزی پارامترهای مدل، توزیع های پیشین های متداول به صورت $\beta_i \sim N(0, 100)$ ، $\sigma^2 \sim IG(0/01, 0/01)$ ، $\psi^{-1} \sim U(0/01, 20)$ ، $\nu_j \sim N(\nu_j, \sigma_{\nu_j}^2)$ ، $\nu_1 \sim N(0, \sigma_{\nu_1}^2)$ ، به ازای $j > 1$ و $\sigma_{\nu_j}^2 \sim IG(0/01, 0/01)$ و $\phi \sim G(0/01, 0/01)$ در نظر گرفته می شوند. مقادیر پارامترهای توزیع های پیشین به گونه ای انتخاب شده اند که توزیع های پیشین ناآگاهی بخش ۲ باشند. نتایج شبیه سازی بر روی یک شبکه 10×10 با $h = 1000$ تکرار برای مدل های $AR(1)$ ، $AR(2)$ و $AR(3)$ در جدول ۱ آمده است. ارزیابی نسبی RelBias و مجذور میانگین مربعات خطا (RMSE) برای هر پارامتر θ از روابط

جدول ۱: برآورد پارامترهای مدل رگرسیون خطی تعمیم یافته بتا فضایی-زمانی

پارامترها									
ϕ	σ^2	ψ^{-1}	φ_3	φ_2	φ_1	β_1	β_0		
50	0/5	10	0/15	-0/2	0/45	1/5	-1	مقدار واقعی	مدل
45/2	0/65	11/01	*	*	0/52	1/68	-1/17	برآورد	
-0/10	0/16	0/10	*	*	0/16	0/13	0/17	اریبی نسبی	AR(1)
5/57	0/17	1/88	*	*	0/27	0/39	0/30	مجذور میانگین مربعات خطا	
45/44	0/59	9/21	*	-0/24	0/44	1/53	-1/14	برآورد	
-0/09	0/18	-0/08	*	0/23	-0/03	0/02	0/14	اریبی نسبی	AR(2)
5/61	0/18	1/72	*	0/15	0/21	0/22	0/31	مجذور میانگین مربعات خطا	
47/36	0/45	10/30	0/12	-0/17	0/43	1/38	-1/04	برآورد	
-0/05	-0/10	0/03	-0/18	-0/17	0/04-	-0/08	0/04	اریبی نسبی	AR(3)
7/26	0/11	1/08	0/16	0/07	0/19	0/29	0/37	مجذور میانگین مربعات خطا	

$$RelBias(\theta) = \frac{1}{h} \sum_{a=1}^h (\hat{\theta}_a / -1), \quad RMSE(\theta) = \left\{ \frac{1}{h} \sum_{a=1}^h (\hat{\theta}_a - \theta)^2 \right\}^{1/2} \quad (1.3)$$

محاسبه شده است.

۴ تحلیل سهم هزینه خوراک از کل هزینه های خانوار در شهر تهران

سهم هزینه خوراک خانوار از کل هزینه ها می تواند به عنوان معیاری برای تشخیص وضعیت رفاه خانوارها به کار گرفته شود. از آنجا که مقدار این سهم در بازه (۱، ۰) تغییر می کند، از مدل معرفی شده در بخش قبل برای تحلیل سهم هزینه خوراک خانوارها از کل هزینه های خانوارهای شهر تهران در فصل های متوالی یک سال استفاده می شود. بعد خانوار، درآمد، و سطح زیر بنای محل سکونت متغیرهای کمکی هستند که در مدل بندی داده ها مورد بررسی قرار گرفته اند. برآورد پارامترها پس

² Noninformative

³ Deviance Information Criterion

از داغیدن ۵۰۰۰۰۰ نمونه از ۷۵۰۰۰۰ نمونه توزیع پسین و انتخاب یک نمونه از هر ۵۰ نمونه باقی مانده انجام می‌شود. مقایسه مقدار معیار انحراف اطلاع^۳ مدل با اثرات فضایی-زمانی (۱۲۹/۳۰-)، مدل با اثر فضایی (۷۲/۹۱-)، مدل با مؤلفه اثر زمان (۱۲۴/۰۹-) و مدل رگرسیون بتا (۶۷/۰۵-) نشان می‌دهد که به‌کارگیری اثرات فضایی-زمانی باعث بهبود برازش مدل شده است. نتایج برازش مدل با اثرات فضایی-زمانی در جدول ۲ آمده است.

جدول ۲: برآورد پارامترها برای مدل سهم هزینه خوراک خانوار از کل هزینه‌ها در شهر تهران

منبع	پارامتر	برآورد	خطای زنجیره مارکوف	بازه باور ۹۵٪
عرض از مبدأ	β_0	-۱/۷۸	۰/۰۳	(-۳/۳۲, -۰/۵۱)
درآمد	β_1	-۰/۲۱	۰/۰۱	(-۰/۴۱, -۰/۰۰۵)
بعد خانوار	β_2	۰/۱۹	۰/۰۱	(۰/۰۰۵, ۰/۳۹)
پارامتر دقت	ϕ	۳۴/۸۷	۰/۱۴	(۲۰/۲۸, ۵۵/۵۲)
واریانس فضایی	σ^2	۳/۴۶	۰/۷۵	(۰/۹۱, ۹/۶۹)
دامنه فضایی	ψ^{-1}	۷/۳۴	۱/۴۷	(۱/۷۵, ۹/۸۳)
پارامتر سری زمانی	φ	۰/۶۶	۰/۰۸	(۰/۱۸, ۰/۹۳)

همان‌طور که در جدول ۲ ملاحظه می‌شود درآمد و بعد خانوار بر سهم هزینه خوراک از کل هزینه‌های خانوار مؤثر بوده‌اند و با افزایش درآمد سهم هزینه خوراک کاهش می‌یابد و با افزایش بعد خانوار این سهم افزایش می‌یابد. پارامترهای مربوط به اثر فضایی و سری زمانی نیز معنی دارند.

بحث و نتیجه‌گیری

مدل‌بندی داده‌های محدود در بازه (۰, ۱) از اهمیت ویژه‌ای برخوردار است. به عنوان مثال در حوزه آمار رسمی نرخ بیکاری، سهم هزینه خوراک از هزینه‌های خانوار، میزان خانوارهای بهره‌مند از خدمات سلامت و بسیاری موارد دیگر در این حوزه می‌گنجد. ارائه یک مدل مناسب با لحاظ کردن همبستگی فضایی-زمانی آن‌ها از اهمیت ویژه‌ای برخوردار است. در این مقاله یک مدل رگرسیون بتا با در نظر گرفتن اثرات تفکیک‌پذیر فضایی-زمانی معرفی شد. تحلیل سهم هزینه خوراک خانوارها از کل هزینه‌ها در شهر تهران نشان داد که علاوه بر درآمد و بعد خانوار، اثرات فضایی و زمانی نیز معنی‌دار است و لحاظ کردن آن‌ها در مدل می‌تواند در بهبود برازش مدل بر داده‌ها مؤثر باشد. در حوزه آمار رسمی، توسعه این مدل با الگوهای پیچیده‌تر همبستگی فضایی-زمانی، می‌تواند کمک مؤثری در راستای تحلیل داده‌ها و پیشگویی مناسب در جهت برنامه‌ریزی و سیاست‌گذاری‌های آتی داشته باشد.

مراجع

محمدزاده، م.، (۱۳۹۸)، آمار فضایی و کاربردهای آن، چاپ سوم، مرکز نشر آثار علمی دانشگاه تربیت مدرس، تهران.

Cepeda-Cuervo, E., Urdinola, B. P., and Rodríguez, D. (2012), Double Generalized Spatial Econometric Models, *Communications in Statistics-Simulation and Computation*, **41**, 671-685.

Ferrari, S., Cribari-Nieto, F. (2004), Beta Regression for Modelling Rates and Proportions, *Journal*

of Applied Statistics, **31**, 799-815.

Figuroa-Zúñiga, J. I., Arellano-Valle, R. B., and Ferrari, S. L. (2013), Mixed Beta Regression: A Bayesian Perspective, *Computational Statistics and Data Analysis* , **61**, 137– 147.

Kalhari Nadrabadi, L., Mohhammadzadeh, M. (2018), Bayesian Inference for Spatial Beta Generalized Linear Mixed Models, *Journal of Sciences, Islamic Republic of Iran*, **29**, 173-185.

Lagos-Alvarez, B. M., Fustos-Toribio, R., Figuroa-Zuniga, J., and Mateu, J. (2017), Geostatistical Mixed Beta regression: a Bayesian Approach, *Stochastic Environmental Research and Risk Assessment*, **31**, 571-584.