

مجله علوم آماری، بهار و تابستان ۱۳۹۱

جلد ۶، شماره ۱، ص ۹۹-۱۱۰

مدل مکانی عام برای پاسخ‌های همبسته پیوسته، ترتیبی و اسمی

سیده فاطمه میری، احسان بهرامی سامانی

گروه آمار، دانشگاه شهید بهشتی

تاریخ دریافت: ۱۳۹۱/۲/۳ تاریخ آخرین بازنگری: ۱۳۹۱/۶/۲۳

چکیده: هدف این مقاله معرفی یک مدل تعمیم‌یافته برای توزیع توام متغیرهای اسمی، ترتیبی و پیوسته با و بدون داده گم‌شده است. فرم‌های بسته‌ای برای تابع درست‌نمایی مربوط به مدل‌های مکانی عام ارائه می‌شود. همچنین تقریب جو، برای برآورد پارامترهای مدل مکانی عام با پاسخ‌های اسمی، پیوسته و ترتیبی با و بدون داده‌های گم‌شده به کار برده شده است. برای نشان دادن قابلیت مدل‌های پیشنهادی، مطالعه‌های شبیه‌سازی انجام شده است. همچنین مدل‌های ارائه شده بر روی داده‌های واقعی مربوط به آموزش زبان خارجی مورد تحلیل قرار گرفته است.

واژه‌های کلیدی: پاسخ‌های همبسته پیوسته، ترتیبی و اسمی، تابع درست‌نمایی، داده‌های گم‌شده، متغیر پنهان، مدل مکانی عام.

۱ مقدمه

یکی از موضوع‌های مهمی که مورد توجه محققان علوم آموزشی و علوم پزشکی قرار می‌گیرد، بررسی پیوند بین برداری از متغیرهای پیوسته و گسسته است. یکی از مدل‌هایی که برای این منظور مورد استفاده قرار می‌گیرد، مدل مکانی عام برای داده‌های آمیخته همبسته

آدرس الکترونیک مسئول مقاله: سیده فاطمه میری، sadatmiri88@yahoo.com
کد موضوع‌بندی ریاضی (۲۰۰۰): ۶۲E۹۹، ۶۲H۱۰

پیوسته و ترتیبی است. در این مدل متغیرهای اسمی یک جدول پیش‌آیندی به وسیله تقاطع سطوح‌شان ایجاد می‌کنند که در هر خانه فرض می‌شود متغیرهای پیوسته دارای توزیع نرمال چند متغیره هستند. در واقع مساله اصلی پیوند دو متغیر پیوسته و ترتیبی بر پایه متغیر پنهان و تعیین تاثیر متغیرهای تبیینی روی این متغیرها است. در این مدل یک سطر ماتریس داده‌ها ممکن است برای بعضی یا همه متغیرهای ترتیبی یا پیوسته دارای مقادیر گم‌شده باشند. با توجه به اهمیت مدل مکانی عام اولین بار اولکین و تیت (۱۹۶۱) حالت‌هایی را که تیت (۱۹۵۵) به آن‌ها اشاره کرده بود تعمیم دادند و مدلی تحت عنوان مدل مکانی عام با پاسخ‌های همبسته اسمی و پیوسته مطرح کردند. کاکس (۱۹۷۲) با استفاده از تجزیه توام متغیرهای پیوسته و گسسته به توزیع حاشیه‌ای متغیر پیوسته و توزیع شرطی متغیر گسسته به شرط متغیر پیوسته، مدلی را ارائه نمود که پیوند بین متغیرهای گسسته و پیوسته را مورد بررسی قرار می‌داد. اولین بار رابین (۱۹۷۶) مفهوم گم‌شدگی داده‌ها را معرفی کرد و سپس لیتل و شالتز (۱۹۸۵) برآورد ماکسیمم درستی‌مایی را برای مدل‌های آمیخته پیوسته که دارای پاسخ‌های گم‌شده باشند معرفی کرد. هکمن (۱۹۷۸) مدل همزمان را برای پاسخ‌های آمیخته پیوسته و گسسته ارائه کرد. کاکس و ویرموتس (۱۹۹۲) به مقایسه مدل مکانی عام و مدل کاکس (۱۹۷۲) پرداختند. کرزانوسکی (۱۹۹۳) مدل‌های مکانی را برای داده‌های آمیخته پیوسته و گسسته ارائه کرد. دیگل و همکاران (۱۹۹۴) مدل‌هایی را روی داده‌های طولی آمیخته مطرح کردند. داده‌های طولی آمیخته داده‌هایی هستند که پاسخ‌های آمیخته اسمی و پیوسته برای هر آزمودنی در طول زمان تکرار می‌شوند. فیتز موریس و لیرد (۱۹۹۷) برای حالتی که برخی از پاسخ‌ها دارای گم‌شدگی هستند به کار بردند. لیون و کریر (۲۰۰۷) مدل‌های مکانی را برای داده‌های آمیخته پیوسته و ترتیبی با و بدون داده‌های گم‌شده بر اساس مدل‌های همزمان، پرداختند. به‌طور کلی آنچه در این مقاله مورد بررسی قرار گرفته به شرح زیر است: بخش دوم، تعمیم مدل داده‌های آمیخته اسمی، پیوسته و ترتیبی با و بدون داده‌های گم‌شده معرفی شده و سپس برآورد پارامترها با روش ماکسیمم درستی‌مایی به‌دست آمده است. بخش سوم به مطالعات شبیه‌سازی برای نشان دادن قابلیت مدل ارائه شده و در انتها نتایجی که از داده‌های واقعی به‌دست می‌آید ارائه می‌شود. بخش چهارم، یک نتیجه‌گیری کلی از مدل مکانی عام برای پاسخ‌های آمیخته اسمی، پیوسته و ترتیبی بیان می‌شود.

۲. تعمیم مدل داده‌های آمیخته

در این بخش ابتدا مدل مکانی عام را برای حالتی که پاسخ‌های آمیخته پیوسته، ترتیبی و اسمی معرفی شده، ارائه خواهد شد. سپس در ادامه به معرفی تعمیم مدل مکانی عام برای داده‌های آمیخته پیوسته، ترتیبی و اسمی با داده‌های گم شده که گم‌شدگی از نوع غیر قابل چشم‌پوشی است بیان می‌شود. فرض کنید X_1, \dots, X_S نشان‌دهنده مجموعه‌ای از متغیرهای اسمی و Y_1, \dots, Y_C نشان‌دهنده مجموعه‌ای از متغیرهای پیوسته و Z_1, \dots, Z_Q نشان‌دهنده مجموعه‌ای از متغیرهای ترتیبی باشند. اگر هر یک از این متغیرها برای n فرد در نظر گرفته شوند، نتیجه یک ماتریس داده $n \times (S + C + Q)$ بعدی $M = (X, Y, Z)$ است، که در آن $X = (X_1, \dots, X_S)$ ، $Y = (Y_1, \dots, Y_C)$ و $Z = (Z_1, \dots, Z_Q)$ به ترتیب نشان‌دهنده‌ی بخش‌های اسمی، پیوسته و ترتیبی ماتریس M هستند. توزیع توام $[X, Y, Z]$ به صورت

$$[X, Y, Z] = [X][Y|X][Z|X, Y]$$

تجزیه می‌شود، که در آن $[X]$ ، $[Y|X]$ و $[Z|X, Y]$ به ترتیب مشخص کننده توزیع کناری X ، توزیع شرطی Y به شرط X و توزیع شرطی Z به شرط X و Y هستند. همچنین $[Z|X, Y]$ بر مبنای $[Y^*|X, Y]$ توزیع شرطی Y^* به شرط X و Y است.

فرض کنید U بردار $D \times 1$ از متغیرهای اسمی می‌باشد که با d امین مولفه از آن دارای s_d حالت ممکن ($d = 1, \dots, D$) است. بردار U یک جدول پیش‌بینی با $S = \prod_{d=1}^D s_d$ حالت ممکن برای مقادیر U تعریف می‌کند. یک بردار $1 \times S$ ، $x = (X_1, \dots, X_S)$ می‌تواند تعریف شود، که اگر U در حالت s قرار گیرد X_s برابر یک و در غیر این صورت برابر صفر است. در این صورت هر سطر U فقط شامل یک ۱ است. تحت مدل مکانی عام متغیر X دارای توزیع چند جمله‌ای است، یعنی $X|\pi \propto \prod_{s=1}^S \pi_s^{x_s}$ که در آن $\pi = (\pi_1, \dots, \pi_S)^T$ یک بردار از احتمال‌های خانه‌های مربوط به جدول پیش‌بینی است. برای $[y, y^*|x_{(s)}]$ ، که y برداری $1 \times C$ و y^* برداری $1 \times Q$ ، یک مدل مکانی تعمیم یافته با میانگین‌های $E(Y|x_s) = \mu_s$ و $E(Y^*|X_s) = \mu_s^*$ و ماتریس‌های کواریانس $Cov(Y|x_s) = Cov(Y) = \Sigma$ ، $Cov(Y^*|x_s) = Cov(Y^*) = \Sigma^*$ و $Cov(Y, Y^*|x_s) = Cov(Y, Y^*) = \Sigma_{Y, Y^*}$ هستند. ارتباط بین متغیر پنهان و متغیر ترتیبی به وسیله مدل آستانه‌ای تعریف می‌شود، یعنی اگر $\alpha_q^{\ell_q} < Y_q^* < \alpha_q^{\ell_q}$ ، $Z_q = \alpha_q^{\ell_q}$ که در آن $\alpha_q^{\ell_q}$ ، $q = 1, \dots, Q$ نقاط آستانه‌ای، $\alpha_q^{\ell_q} = -\infty$ و $\alpha_q^{\ell_q+1} = +\infty$ ، $\{\alpha_q^{\ell_q}, \dots, \alpha_q^{\ell_q}\}$ و $\alpha_q^{\ell_q+1}$ نقاط نامعلوم هستند. یک رابطه ترتیبی برای نقاط α به صورت $\alpha_q^{\ell_q+1} < \dots < \alpha_q^{\ell_q}$ است و

تحت مدل مکانی عام توزیع Y^* به شرط Y و x_s نرمال چند متغیره

$$Y^*|Y, x_s \sim N(\mu_s^* + \Sigma_{y,y^*}^T \Sigma^{-1}(y - \mu_s), \Sigma^* - \Sigma_{y,y^*}^T \Sigma^{-1} \Sigma_{y,y^*})$$

است. ماتریس کواریانس به فرم $DRD = \Sigma^* - \Sigma_{y,y^*}^T \Sigma^{-1} \Sigma_{y,y^*}$ است، که در آن D ماتریس قطری از انحراف معیار شرطی d_q و R ماتریس متقارن $Q \times Q$ از همبستگی چند رشته‌ای Y^* به شرط X_s و Y است. برای کاهش دادن تعداد پارامترهای مدل، S ثابت در نظر گرفته شده است، که در آن

$$\mu_s = \xi + \xi_s, \quad \mu_s^* = \xi^* + \xi_s^*, \quad s = 1, \dots, S-1$$

همچنین $\mu_S = \xi$ و $\mu_S^* = \xi^*$ به ترتیب میانگین‌های Y و Y^* برای حالت s هستند. استاندارد $[Y^*|x_s, Y]$ برای $s \neq S$ به صورت

$$D^{-1}(y^* - \mu_s^* - \Sigma_{y,y^*}^T \Sigma^{-1}(y - \mu_s)) = D^{-1}(y^* - \xi^*) - \tau_s - B(y - \xi) \quad (1)$$

انجام می‌شود، که در آن

$$\tau_s = D^{-1} \xi_s^* - B \xi_s, \quad B = D^{-1} \Sigma_{y,y^*}^T \Sigma^{-1}$$

عبارت (۱) دارای توزیع نرمال چند متغیره با میانگین صفر و ماتریس کواریانس R است. به‌طور مشابه استاندارد نقاط آستانه‌ای شرطی $\{\alpha_q^1, \dots, \alpha_q^{L_q}\}$ ($q = 1, \dots, Q$) به صورت $\gamma_q^{L_q} - \tau_{sq} - \beta_q^T y$ که در آن $\gamma_q^{L_q} = \frac{\xi_q^*}{d_q} - \beta_q^T \xi$ و $\tau_{sq} = \frac{\xi_{sq}^*}{d_q} - \beta_q \xi_{sq}$ q امین عنصر $\beta_q = d_q^{-1} \Sigma_{y,y_q^*}^T \Sigma^{-1}$ $\beta_q^T = (\beta_{q1}, \dots, \beta_{qC})$ q امین ردیف B است. توجه شود $\tau_{Sq} = 0$ و $\gamma_q^{L_q+1} = +\infty$ و $\gamma_q^0 = -\infty$ واضح است که

$$P(Z = \ell | x = x_s, y, \theta) = \int_S \phi_Q(v|R) dv,$$

که در آن $\phi_Q(\cdot|R)$ تابع چگالی توزیع نرمال با میانگین صفر و ماتریس کواریانس R ، $\ell = (\ell_1, \dots, \ell_Q)^T$ ، $S = \{(v_1, \dots, v_Q) : \nu_{sq}^{L_q-1} < v_q < \nu_{sq}^{L_q}, q = 1, \dots, Q\}$ ، که $\nu_{sq}^{L_q} = \gamma_q^{L_q} - \beta_q^T y - \tau_{sq}$ و $\theta = (\mu, \Sigma, R)$ هستند. چگالی توام $[x, y, z]$ به صورت

$$P(x = x_s, y, z = \ell | \theta) = \pi_s \phi(y - \mu_s | \Sigma) \int_S \phi_Q(v|R) dv$$

تعریف می‌شود و متغیرهای پیوسته Y^* ، $I_{R_z^*}$ و $I_{R_y^*}$ به ترتیب نشان‌دهنده متغیر پنهان برای پاسخ ترتیبی Z_i ، متغیر پنهان مربوط به مکانیسم گم‌شدگی Y_i و متغیر پنهان مربوط به

مکانیسم گم‌شدگی Z_i هستند. از طرفی دیگر متغیرهای نشانگر $I_{R_{y_i}}$ و $I_{R_{z_i}}$ به ترتیب به صورت

$$I_{R_{y_i}} = \begin{cases} 1 & R_{y_i}^* > 0 \\ 0 & o.w. \end{cases} \quad I_{R_{z_i}} = \begin{cases} 1 & R_{z_i}^* > 0 \\ 0 & o.w. \end{cases}$$

تعریف می‌شوند، که در آن‌ها $R_{y_i}^*$ و $R_{z_i}^*$ ، متغیرهای مربوط به R_{y_i} و R_{z_i} هستند. یک بودن تابع نشانگر نشان‌دهنده این است که متغیر پاسخ مشاهده شده است و صفر بودن آن نشان دهنده این است که متغیر پاسخ گم‌شده است. مدل مکانی عام مربوط به پاسخ‌های پیوسته و ترتیبی که دارای گم‌شدگی هستند به صورت

$$Y_i^* = \beta_1' T_{i1} + \varepsilon_{i1}, \quad Y_i = \beta_2' T_{i2} + \varepsilon_{i2}$$

$$R_{Y_i}^* = \alpha_1' T_{i3} + \varepsilon_{i3}, \quad R_{Z_i}^* = \alpha_2' T_{i4} + \varepsilon_{i4}$$

در نظر گرفته می‌شوند، که در آن $(\varepsilon_{11}, \varepsilon_{12}, \varepsilon_{23}, \varepsilon_{24})$ دارای توزیع نرمال چندمتغیره با میانگین صفر و ماتریس کواریانس Σ_{1234} است. بنابراین پارامترهای مدل در حالتی که بردار متغیرهای Z و Y عبارتند از $\beta_1', \beta_2', \eta_1'$ و η_2' و نقاط آستانه‌ای $\alpha_1^L < \dots < \alpha_1^q$ و ضریب همبستگی R_{1234} هستند. به خاطر نشان‌پذیر بودن پارامترهای مدل Σ_{1234} ، $\Sigma_{R_z^*, R_z^*} = Cov(R_z^*, R_z^*) = I$ ، $\Sigma_{y^*, y^*} = Cov(y^*, y^*) = I$ ، $\Sigma_{R_y^*, R_y^*} = Cov(R_y^*, R_y^*) = I$ و $\Sigma_{y^*, R_z^*} = Cov(y^*, R_z^*)$ و در نهایت $\Sigma_{y^*, R_y^*} = Cov(y^*, R_y^*)$ ، $\Sigma_{R_z^*, R_y^*} = Cov(R_z^*, R_y^*)$ لازم به ذکر است پارامترها در صورت صفر بودن ماتریس‌های Σ_{y^*, R_y^*} ، Σ_{y^*, R_z^*} و $\Sigma_{R_z^*, R_y^*}$ ، از نوع قابل چشم‌پوشی است.

۱.۲ تابع درستی‌مندی مدل مکانی عام برای داده‌های کامل

فرض کنید نمونه تصادفی به اندازه N از مدل داده آمیخته تعمیم یافته^۱ (GMDM) مشاهده شده است. لگاریتم تابع درستی‌مندی داده‌ها به صورت

$$\log L(\pi_s, \mu, \Sigma, R | x_s, y, z) = \ell(\pi_s | x_s) + \ell(\mu, \Sigma | x_s, y) + \ell(R | x_s, y, z)$$

^۱ General Mixed Data Model

به طوری که

$$\begin{aligned} \ell(\pi|X) &= \sum_{s=1}^S n_s \log(\pi_s) \\ \ell(\mu, \sigma|X, Y) &= \frac{-N}{\nu} \log |\Sigma| - \frac{1}{\nu} \sum_{s=1}^S \sum_{i(s) \in N} (y_{i(s)} - \mu_s)^T \Sigma^{-1} (y_i - \mu_s) \\ \ell(R|X, Y, Z) &= \sum_{i=1}^N \sum_{s=1}^S \log \sum_{\epsilon_q=0}^Q (-1)^{q+1} \Phi_Q(\dots, v_{i(s,\ell)}^{\epsilon_q - \epsilon_q}, \dots | R) \end{aligned}$$

که در آن $\Phi_Q(\dots, v_{i(s,\ell)}^{\epsilon_q - \epsilon_q}, \dots | R)$ تابع توزیع نرمال چند متغیره با میانگین صفر و ماتریس همبستگی R است.

۲.۲ تابع درستنمایی مدل مکانی عام برای داده‌های گم شده

برای بررسی و تحلیل پاسخ‌های همبسته برداری از متغیرهای پیوسته و برداری از متغیرهای ترتیبی بر پایه متغیر پنهان لازم است تابع درستنمایی پاسخ‌ها تعیین شود که محاسبه آن در شرایط خاص با پیچیدگی و دشواری همراه است. برای حل این مسأله، می‌توان تابع درستنمایی را تقریب زد (جو، ۱۹۹۵). برای این منظور فرض کنید $X = (X_1, \dots, X_m)$ برای $m \geq 3$ دارای توزیع نرمال چند متغیره با بردار میانگین صفر، واریانس‌های ۱ و ماتریس همبستگی Ω باشد.

فرض کنید $I_i = I(\omega_i \geq X_i \geq x_i)$ ، برای $i = 1, \dots, m$ یک تابع نشانگر باشد به طوری که $E(i) = \Phi(x_i) - \Phi(\omega_i)$ ، که در آن Φ تابع توزیع نرمال استاندارد است. اکنون با توجه به تجزیه

$$\begin{aligned} P(\omega_1 \geq X_1 \geq X_2, \dots, \omega \geq X_m \geq x_m) &= P(\omega_1 \geq X_1 \geq x_1, \omega_2 \geq X_2 \geq x_2) \\ &\times \prod_{k=3}^m P(\omega_k \geq X_k \geq x_k | \omega_1 \geq X_1 \geq x_1, \dots, \omega_{k-1} \geq X_{k-1} \geq x_{k-1}) \quad (2) \end{aligned}$$

در مرحله اول، تقریب عبارت دوم رابطه (۲) به صورت

$$E(I_k | I_1 = 1, \dots, I_{k-1} = 1) = E(I_k) + \Omega_{21} \Omega_{11}^{-1} (1 - E(I_1), \dots, 1 - E(I_{k-1}))^T$$

است، که در آن Ω_{21} یک بردار سطری شامل اعضای $Cov(I_i, I_k)$ است، $i = 1, \dots, k-1$ و Ω_{11} یک ماتریس $(k-1) \times (k-1)$ با درایه‌های

که $Cov(I_i, I_j) = E(I_i, I_j) - E(I_i)E(I_j)$ و $1 \leq i$ و $z \leq k - 1$ است، لازم به ذکر است که $E(I_i, I_j) = P(\omega_i \geq X_i \geq x_i, \omega_j \geq X_j \geq x_j)$ برای اطلاعات بیشتر به جوی (۱۹۹۵) رجوع شود.

تابع درست‌نمایی برای پاسخ‌های آمیخته پیوسته و ترتیبی به‌طوری که پاسخ‌ها، دارای گم‌شدگی هستند و متغیر اسمی همواره مشاهده شده باشند، برای افراد گوناگون وابسته به الگوی گم‌شدنشان به‌گونه متفاوتی بایستی به‌دست آید. با در نظر گرفتن یک متغیر پاسخ Y ، Z و X تابع درست‌نمایی در چهار حالت به‌دست آورده شده است که دارای ماتریس کواریانس با پارامترهای σ (واریانس متغیر پاسخ Y)، ρ_{12} (ضریب همبستگی بین Y و Y^*)، ρ_{13} (ضریب همبستگی بین Y و $R_{Y^*}^*$)، ρ_{23} (ضریب همبستگی بین Y^* و $R_{Y^*}^*$)، ρ_{24} (ضریب همبستگی بین Y^* و $R_{Y^*}^*$) و ρ_{34} (ضریب همبستگی بین $R_{Y^*}^*$ و $R_{Y^*}^*$) است. پارامترهای $\eta_1, \eta_2, \beta'_1, \beta'_2$ ضرایب رگرسیونی مدل مورد نظر هستند و همچنین نقاط آستانه‌ای مربوط به متغیر تبیینی X برآورد می‌شوند. در این جا فقط پاسخ Y برای فرد نام مشاهده شده مورد بررسی قرار می‌گیرد و می‌توان آن را برای حالت‌های دیگر (۱- فقط پاسخ Z برای فرد نام مشاهده شده است. ۲- هر دو پاسخ Y و Z برای فرد نام مشاهده نشده باشند. ۳- هر دو پاسخ Y و Z برای فرد نام مشاهده شده باشند.) تعمیم داد.

$$\begin{aligned} L_i &= f(y_i, I_{R_{Y_i}} = 1, I_{R_{Z_i}} = 0) \\ &= f(y_i | X_i = x_i) P(X_i = x_i) [P(I_{R_{Z_i}} = 0 | y_i, X_i = x_i) \\ &\quad - P(I_{R_{Z_i}} = 0, I_{R_{Y_i}} = 0 | y_i, X_i = x_i)] \\ &= f(y_i | X_i = x_i) P(X_i = x_i) [P(R_{Z_i}^* < 0 | y_i, X_i = x_i) \\ &\quad - P(R_{Z_i}^* < 0, R_{Y_i}^* < 0 | y_i, X_i = x_i)] \\ &= f(y_i | X_i = x_i) P(X_i = x_i) \\ &\quad \times \left[\Phi\left(\frac{-\alpha'_2 T_{i4} - \frac{\rho_{24}}{\sigma}(y_i - \beta'_2 T_{i2})}{\sqrt{1 - \rho_{24}^2}}\right) - \Phi_{12}\left(\frac{-\alpha'_2 T_{i4} - \frac{\rho_{24}}{\sigma}(y_i - \beta'_2 T_{i2})}{\sqrt{1 - \rho_{24}^2}}, \right. \right. \\ &\quad \left. \left. \frac{-\alpha'_1 T_{i3} - \frac{\rho_{13}}{\sigma}(y_i - \beta'_1 T_{i1})}{\sqrt{1 - \rho_{13}^2}}; \frac{\rho_{24} - \rho_{24}\rho_{23}}{\sqrt{1 - \rho_{23}^2}\sqrt{1 - \rho_{24}^2}}\right) \right] \end{aligned}$$

که در آن $\Phi_{12}(\cdot, \cdot; \cdot)$ تابع توزیع تجمعی نرمال متغیرهای پنهان مربوط به مکانیزم گم‌شدگی Y_i و مکانیزم گم‌شدگی Z_i است. تابع درست‌نمایی که از ضرب تمام درست‌نمایی‌ها به‌دست می‌آید را می‌توان در نرم‌افزار R ، با دستور "nlminb" بهینه کرد.

۳ مطالعه شبیه‌سازی

هدف از این شبیه‌سازی، بررسی مدل مکانی عام با داده‌های گم‌شده است. با فرض آنکه $s = 2$ و $C = Q = L = 1$ باشند، Y و Y^* (متغیر پنهان) دو متغیر پیوسته، Z یک متغیر ترتیبی با دو سطح و X یک بردار اسمی دوتایی در نظر گرفته شده‌اند. در این مدل فرض می‌شود متغیر ترتیبی دارای گم‌شدگی است. پارامترهای مدل عبارتند از $\theta = (\pi, \mu^T, \sigma^2, \gamma, \beta, \tau)$ ، که در آن $\mu^T = (\mu_1, \mu_2)$ امید Y به شرط $x = x_{(s)}$ به‌ازای $R_{z_i}^*$ ، $s = 1, 2$

$$\gamma = \frac{R_{z_i}^*}{\sqrt{1-\rho^2}} - \left(\frac{\mu_2^*}{\sqrt{1-\rho^2}} - \beta\mu_2 \right), \quad \beta = \frac{\rho}{\sigma\sqrt{1-\rho^2}}, \quad \tau = \frac{\xi^*}{\sqrt{1-\rho^2}} - \beta\xi$$

و ρ ضریب همبستگی بین Y و Y^* است. در این شبیه‌سازی نمونه‌هایی به حجم ۲۵۰، ۵۰۰ و ۱۰۰۰ با ۲۰۰۰ تکرار از GMDM با مقادیر واقعی $\pi = 0/5$ ، $\mu_1 = \mu_2 = 0/5$ ، $\mu_2 = \mu_1^* = 0/5$ ، $\sigma^2 = 1$ ، $\rho = 0/5$ ، $\gamma = -0/289$ و $\beta = 0/577$ تولید شده است. با توجه به نتایج ارائه شده در جدول ۱ ملاحظه می‌شود، هر اندازه حجم نمونه زیاد باشد برآورد پارامترها به مقدار واقعی نزدیک‌تر و اریبی کم‌تر است. همچنین با افزایش حجم نمونه کارایی نسبی افزایش می‌یابد و مقادیری که بزرگتر از یک هستند نشان می‌دهد مقدار برآورد پارامتر کاراتر از مقدار میانگین پارامترها است.

جدول ۱: برآوردهای ماکسیمم درست‌نمایی پارامترها با ۲۰۰۰ تکرار از

GMDM با داده‌های گم‌شده

پارامتر	مقدار پارامتر	حجم نمونه	میانگین برآورد	اریبی	کارایی
γ	0/866	250	0/8716	-1/6826	0/9588
		500	0/8432	-0/8864	1/0147
		1000	-0/8900	-0/6460	1/200
β		250	0/5900	-2/2530	0/9947
	0/577	500	0/5871	-1/7501	1/0001
		1000	0/5892	-2/1140	1/0166
τ		250	-0/2945	-1/9030	0/9220
	-0/289	500	-0/2932	-1/4530	0/9455
		1000	-0/2900	0/3460	1/0065

جدول ۲: برآورد پارامترها و خطای استاندارد مدل کامل GMDM

انحراف معیار	برآورد	پارامتر	
			<i>SEX</i> × <i>LAN</i>
۰/۰۲۳	۰/۱۴۳	π_1	فرانسوی و مرد
۰/۰۲۴	۰/۱۵۲	π_2	اسپانیایی و مرد
۰/۰۲۹	۰/۲۵۱	π_3	آلمانی و مرد
۰/۰۲۲	۰/۱۳۰	π_4	فرانسوی و زن
۰/۰۲۲	۰/۱۲۶	π_5	اسپانیایی و زن
۰/۰۲۶	۰/۱۹۹	π_6	آلمانی و زن
۲/۳۰۰	۸۴/۹۳۹	μ_{11}	فرانسوی و مرد
۰/۱۰۹	۲/۸۰۳	μ_{12}	
۲/۲۳۴	۷۸/۵۴۳	μ_{21}	اسپانیایی و مرد
۰/۱۰۶	۲/۶۱۰	μ_{22}	
۱/۷۳۵	۸۲/۱۷۲	μ_{31}	آلمانی و مرد
۰/۰۸۲	۲/۸۰۰	μ_{32}	
۲/۴۱۳	۸۲/۷۶۷	μ_{41}	فرانسوی و زن
۰/۱۱۴	۲/۷۶۹	μ_{42}	
۲/۴۵۴	۸۴/۲۴۱	μ_{51}	اسپانیایی و زن
۰/۱۱۶	۲/۸۹۱	μ_{52}	
۱/۹۴۸	۸۵/۵۰۰	μ_{61}	آلمانی و زن
۰/۰۹۲	۲/۷۴۳	μ_{62}	
۱۶/۲۴۷	۱۷۴/۶۱۳	σ_1^2	واریانس
۰/۰۳۶	۰/۳۹۲	σ_2^2	
۰/۹۶	۰/۰۵۹	ρ	ضریب همبستگی
۰/۲۳	۰-/۵۹۷	γ^1	نقاط برش
۰/۲۲۱	۰/۴۰۰	γ^2	
۰/۱۱۳	۰-/۶۸۰	β_1	اثر رگرسیونی FLAS
۰/۰۰۴	۰-/۰۱۲	β_2	اثر رگرسیونی HGPA
۰/۲۷۳	-۰/۸۵۸	τ_1	
۰/۲۴۶	-۰/۵۸۴	τ_2	
۰/۲۹	-۰/۵۷۸	τ_3	اثر <i>SEX</i> × <i>LAN</i>
۰/۲۹۷	-۰/۲۶۲	τ_4	
۰/۲۷۹	۰/۴۰۳	τ_5	

۱.۳ مثال کاربردی برای داده‌های کامل

این بخش شامل کاربرد روان‌سنجی زبان خارجی که از مطالعه شفر (۱۹۹۷)، به دست آمده ارائه می‌شود. هدف اصلی از مطالعه زبان خارجی، بررسی ویژگی روان‌سنجی^۲ (FLAS) که توسط ریموند و رابرتز (۱۹۸۳) برای پیشگویی مهارت در FLAS گزارش یافته است مورد بررسی قرار می‌گیرد. آن‌ها ۲۳۱ دانشجویی که در دوره‌های زبان خارجی در دانشگاه پنسیلوانیا در سال ۱۹۸۰ با پر کردن فرم‌هایی که شامل متغیرهای زیر هستند را در نظر گرفتند. ۱- LAN، یک متغیر اسمی با سه سطح، متناظر با رشته‌های تحصیلی دانشجویان (۱- زبان فرانسوی ۲- زبان اسپانیایی ۳- زبان آلمانی). ۲- SEX، یک متغیر اسمی با دو سطح، متناظر با جنسیت افراد (۱- مذکر ۲- مونث). ۳- FLAS، یک متغیر پاسخ پیوسته، نمرات دانشجویان. ۴- HGPA، یک متغیر پاسخ پیوسته، معدل دانشجویان در دوره دبیرستان. ۵- GRD^۳، یک متغیر پاسخ ترتیبی با سه سطح، نمرات دانشجویان زبان (۱- نمره C یا زیر C ۲- نمره B ۳- نمره A). شفر (۱۹۹۷) به تحلیل یک مجموعه بزرگ داده‌ها که شامل اطلاعات بر روی یک تعداد متغیرهای دیگر همانند نمرات آزمون استعداد مدرسه بود پرداخته است. هدف او از انجام این کار، به کار بردن مدل مکانی عام بوده است. حال هدف از مطالعه، بررسی ارتباط بین متغیرهای آمیخته با مدل‌بندی کردن بر روی توزیع توام، بر اساس GMDM است. با توجه به این که GRD نتیجه مورد انتظار است، متغیرها را با روش GMDM در حالتی که گم‌شدگی وجود ندارد به صورت

$$[GRD, FLAS, HGPA, LAN, SEX] = [GRD|FLAS, HGPA, LAN, SEX] \\ \times [FLAS, HGPA|LAN, SEX] \times [LAN, SEX]$$

گروه‌بندی می‌شود. نتیجه برازش مدل GMDM کامل به وسیله‌ی برآورد ماکسیمم درست‌نمایی در جدول ۲ نشان داده شده است. LAN و SEX به ترتیب سه سطحی و دو سطحی هستند، تعداد خانه‌های جدول پیش‌بینی ۶ = S است. با بررسی برآورد احتمال‌ها مشاهده می‌شود کم‌ترین برآورد احتمال ۰/۱۲، در رشته زبان اسپانیایی و جنس مونث و بیشترین آن، ۰/۲۵ در رشته زبان آلمانی و جنس مذکر است. بررسی برآورد میانگین‌ها، در متغیر پیوسته FLAS، نشان می‌دهد دانش‌آموزان مذکر که زبان اسپانیایی را مطالعه می‌کنند دارای کمترین نمره و دانش‌آموزان مونث که زبان آلمانی را مطالعه می‌کنند دارای بیشترین نمره هستند. ضریب همبستگی بین دو متغیر پیوسته FLAS و HGPA، (۰/۰۵۹۱ = $\hat{\rho}$) یک ارتباط ضعیف بین

^۲ Foreign Language Attitude Scale

^۳ Grad in The Foreign Language

دو متغیر را نشان می‌دهد. با توجه به جدول ۲ آماره Z برای حالت رگرسیون FLAS برابر $2/828-$ و HGPA برابر $6/02-$ هستند. این نشان می‌دهد، متغیرهای پیوسته، توانایی پیشگویی متغیر ترتیبی (GRD) را دارند.

۴ بحث و نتیجه‌گیری

یکی از مدل‌های مهم در پیوند بین برداری متغیرهای پیوسته و گسسته، مدل مکانی عام برای داده‌های آمیخته همبسته پیوسته و ترتیبی است. یکی از مسائل به‌کارگیری این مدل، وجود پاسخ‌های گم‌شده و تشخیص مکانیسم گم‌شدگی پاسخ‌ها و عوامل موثر در گم‌شدن آن‌ها است. در این مقاله این مدل برای پاسخ‌های آمیخته پیوسته، ترتیبی و اسمی با و بدون داده‌های گم‌شده بررسی شد. در این مدل متغیرهای اسمی، یک جدول پیش‌بینی به‌وسیله تقاطع سطوح‌شان تشکیل می‌دهند و متغیرهای پیوسته و پنهان متناظر با این متغیرها درون خانه‌های جدول قرار می‌گیرند. مدل مکانی عام برای پاسخ‌های آمیخته پیوسته، ترتیبی و اسمی با و بدون داده‌های گم‌شده مواردی هستند که در آینده می‌توانند مورد کاوش قرار گیرند.

مراجع

- Bahrami Samani, E., Ganjali, M. (2008), A Latent Variable Model for Mixed Continuous and Ordinal Responses, *Journal of Statistical Theory and Applications, Journal of Statistical*, **7**, 337-348.
- Cox, D. R. (1972), The Analysis of Multivariate Binary Data, *Applied Statistics*, **21**, 113-126.
- Cox, D. R, and Wermuth, N. (1992), Response Models for Mixed Binory and Quantitative Variables, *Biometrika*, **79**, 441-461.
- Diggle, P. J. and Kenward, M. G. (1994), Informative Dropout in Longitudinal Data Analysis. *Journal of Applied Statistics*, **43**, 49-93.
- Fitzmaurice, G. M. and Laird, N. M, (1995), Regresion Models for Bivariate Discrete and Continuous Outcome with Clustering, *Journal of the American Statistical Association*, **90**, 845-852.

- Fitzmaurice, G. M. And Laird, N. M, (1997), Regression Models for Mixed Discrete and Continuous Responses with Potentially Missing Values. *Biometrika*, **53**, 110-122.
- Heckman, J. J. (1978), Dummy Endogenous Variables in a Simultaneous Equation System, *Biometrika*, **88**, 551-561.
- Joe, H. (1995), Approximations Multivariate Normal Rectangle Probabilities Based on Conditional Expectation. *Journal of the American Statistical Association*, **90**, 957-967.
- Krzanowski, W. J. (1993), The Location Model for Mixtures of Categorical and Continuous Variables. *Journal of Classification*, **10**, 25-49.
- Leon. A. R. and Carrer. K. C. (2007), General Mixed-Data Model: Extension of General Location and Grouped Continuous Models, *The Canadian Journal of Statistics*, **35**, 533-548.
- Little, R. J. A. and Schluchter, M. D. (1985), Maximum Likelihood Estimation for Mixed Continuous and Categorical Data with Missing Values, *Biometrika*, **72**, 492-512.
- Olkin, I. and Tate, R. F. (1961), Multivariate Correlation Models with Mixed Discrete and Continuous Variables, *The Annals of Mathematical Statistics*, **32**, 448-465; Correction. *The Annals of Mathematical Statistics*, **36**, 343-344.
- Raymond, M. R. and Roberts, D. M. (1983), Development and Validation of a Foreign Language Attitude Scale, *Educational and Psychological Measurement*, **43**, 1239-1246.
- Rubin, D. B. (1976), Inference and Missing Data. *Biometrika*, **82**, 669-710.
- Schafer, J. L. (1997), *Analysis of Incomplete Multivariate Data*, Chapman, Hall, CRC, Boca Raton, Florida.
- Tate, F. R. (1955), The Teory of Correlation Between two Continuous Variabls when one is Dichatomized, *Biometrika*, **42**, 205-216.