



مدل‌بندی همزمان میانگین و دقت در مدل‌های آمیخته رگرسیون بتای افزوده

زهرة فلاح محسن‌خانی^۱، پروین اژدری^{۲*}

^(۱) استادیار، پژوهشکده آمار، تهران، ایران

^(۲) استادیار، دانشکده علوم پایه، دانشگاه آزاد اسلامی واحد تهران شمال، تهران، ایران

دبیر مسئول: محمد رضا زادکرمی

تاریخ پذیرش: ۱۴۰۰/۸/۹

تاریخ دریافت: ۱۳۹۹/۵/۱۹

چکیده: مدل رگرسیون بتای افزوده برای مدل‌بندی داده‌هایی از جنس نرخ، نسبت یا درصد استفاده می‌شود. این مدل از آمیختن توزیع بتا روی بازه (۰, ۱) و دو توزیع تباهیده در صفر و یک ایجاد می‌شود. با بازپارامتریدن توزیع بتا، پارامترهای میانگین و دقت این مدل با ساختاری شامل اثرات ثابت و تصادفی مدل‌بندی می‌شود. عموماً برای راحتی در مطالعات، پارامتر دقت ثابت در نظر گرفته می‌شود و مدل‌بندی فقط بر اساس پارامتر میانگین انجام می‌شود. در این مقاله مدل‌بندی همزمان میانگین و دقت در مدل‌های آمیخته رگرسیون بتای افزوده ارائه و کارایی مدل در مطالعات شبیه سازی با رهیافت بیزی مورد بررسی قرار می‌گیرد. سپس نحوه کاربست این مدل برای تحلیل مدل‌بندی سهم شاغلین در خانوار بر اساس نتایج آمارگیری نیروی کار مرکز آمار نشان داده می‌شود و در انتها بحث و نتیجه‌گیری ارائه خواهد شد.

واژه‌های کلیدی: رگرسیون بتای افزوده، پارامتر دقت، مدل آمیخته، تحلیل بیزی، آمارگیری نیروی کار.

رده‌بندی ریاضی: 62F15, 62J12

۱ مقدمه

برخی از مطالعات در حوزه‌های مختلف شامل داده‌هایی هستند که به صورت نرخ یا نسبت در بازه‌ی (۰ و ۱) اندازه‌گیری می‌شوند. مدل رگرسیون بتا یک انتخاب مناسب برای تحلیل متغیرهای پاسخ پیوسته روی بازه واحد است. پائولینو (۲۰۰۱) برای اولین بار به منظور مدل‌بندی متغیرهای پاسخ از جنس نسبت، فرض کرد متغیر پاسخ از توزیع بتا $B(a, b)$ تابع چگالی

$$\pi(y; a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^{a-1} (1-y)^{b-1} \quad 0 < y < 1, \quad a, b > 0$$

پیروی می کند، که میانگین و واریانس آن به ترتیب عبارتند از $E(Y) = \frac{a}{a+b}$ و $Var(Y) = \frac{ab}{(a+b)^2(a+b+1)}$. وی پارامترهای توزیع بتا، را نیز مدل بندی نموده و برآورد ماکسیمم درستنمایی آن ها را به دست آورد این روش از لحاظ محاسباتی با دشواری بسیاری همراه بود. فراری و کریباری (۲۰۰۴) توزیع بتای بازپارامتریده را برای مطالعه این نوع داده ها پیشنهاد کردند. آن ها پارامترهای توزیع بتا را به گونه ای بازنویسی کردند که مدل رگرسیونی براساس میانگین متغیر پاسخ بیان شود. برای این منظور با قرار دادن $\mu = \frac{a}{a+b}$ و $\phi = a + b$ توزیع بتا به صورت:

$$\pi(y; \mu, \phi) = \frac{\Gamma(\phi)}{\Gamma(\mu\phi) \Gamma((1-\mu)\phi)} y^{\mu\phi-1} (1-y)^{(1-\mu)\phi-1} \quad 0 < y < 1 \quad (1.1)$$

بازپارامتریده می شود، که در آن $\Gamma(\cdot)$ تابع گاما و μ میانگین و ϕ پارامتر دقت است ($0 < \mu < 1$ و $\phi > 0$). سپس با در نظر گرفتن الگوی خطی برای متغیرهای تبیینی به برازش مدل پرداختند. در این مدل ها متغیرهای تبیینی و میانگین متغیر پاسخ از طریق یک تابع پیوند مناسب به صورت $g(\mu_i) = \sum_{j=1}^k x_{ij}\beta_j$ به هم مربوط می شوند، که در آن β بردار پارامترهای رگرسیونی است. چپیدا و گامرن (۰۵۲) و اسمیتسون و وركولین (۰۶۲)، مدل رگرسیون بتا را با متغیر در نظر گرفتن پارامتر دقت بررسی کردند. آن ها به طور همزمان لوجیت میانگین و لگاریتم پارامتر دقت را به صورت الگوی خطی از متغیرهای تبیینی لحاظ نمودند. سپس برانساگام و همکاران (۲۰۰۷) یک مدل رگرسیون بتا بیزی را ارائه کردند و کاربرد آن را در داده های ژنتیک مطالعه نمودند. مدل آمیخته رگرسیون بتا اولین بار توسط زیمیریچ (۲۰۱۰) در مطالعه داده های طولی از طریق وارد کردن اثر تصادفی در مدل میانگین برای بررسی تغییر سرعت واکنش افراد نسبت به محرک ها با افزایش سن، مطرح شد، وركولین و اسمیتسون (۲۰۱۲) با در نظر گرفتن توزیع بتای آمیخته برای متغیر پاسخ، برآورد پارامترهای مدل را با رهیافت بیزی و روش ماکسیمم درستنمایی ارائه نمودند. فیگیورا-زونیکا و همکاران (۱۳۰۲)، پارامترهای میانگین و دقت مدل آمیخته خطی تعمیم یافته بتا را با رهیافت بیزی ارائه کردند. همچنین بونات و همکاران (۲۰۱۳) مدل آمیخته رگرسیون بتا را برای داده های واقعی استفاده و با فرض ثابت بودن پارامتر دقت، برآوردهای ماکسیمم درستنمایی پارامترهای مدل را ارائه کردند. چپیدا و همکاران (۲۰۱۴) مدل رگرسیون بتا را با در نظر گرفتن یک اثر تصادفی در مدل بندی پارامتر دقت مطالعه کرده و برآورد بیزی پارامترهای مدل را به دست آوردند.

در مسائل کاربردی ممکن است با مواردی مواجه شویم که نرخ ها هر یک از مقادیر مجموعه $\{0$ و $1\}$ را نیز اختیار کنند، یک راه حل استفاده از آمیختن گسسته-پیوسته توزیع است که توسط اوسپینا و فراری (۲۰۱۰) پیشنهاد شد که بسته به شرایط، جرم احتمال های صفر، یک یا هر دو را نیز شامل می شود، همچنین کاربرد این مدل برای داده های واقعی در ۲۰۱۲ توسط وركولین و سمیتسون ارائه گردید. گالویس و همکاران (۲۰۱۴) مدل رگرسیون بتای افزوده را برای داده های مشاهده شده در بازه $[0, 1]$ پیشنهاد کردند که از آمیختن دو توزیع تباهیده در نقاط صفر و یک و توزیع بتا با تکیه گاه $(0, 1)$ به دست می آید. پارکر و همکاران (۲۰۱۴)، اثرات فضایی را در مدل بندی رگرسیون بتای افزوده در برآورد پارامتر میانگین بررسی و در داده های دندان پزشکی بکار گرفتند. فلاح و همکاران (۲۰۱۹) اثرات طولی را در مدل بندی پارامتر میانگین این مدل بررسی کردند. نوگارتو و همکاران (۲۰۲۰) برآوردهای روش بیزی و روش ماکسیمم درستنمایی را در مدل رگرسیون بتای افزوده با یکدیگر مقایسه کردند. نتایج تحقیق آنها نشان داده است که رهیافت بیزی به اندازه روش ماکسیمم درستنمایی دقیق است. در سال های اخیر استفاده از مدل رگرسیون بتای افزوده برای داده های نسبتی و لحاظ اثرات ثابت و تصادفی با استفاده از روش های گوناگون و با لحاظ رهیافت های متفاوت در مدل بندی پارامتر میانگین، بحث و بررسی شده است، در این مقاله مدل آمیخته رگرسیون بتای افزوده برای حالتی که پارامتر دقت توزیع بتا نیز متغیر باشد بررسی می شود. به عبارتی به طور همزمان یک الگوی خطی برای پارامتر میانگین و یک الگوی خطی برای پارامتر دقت در نظر گرفته، و با استفاده از رهیافت بیزی برآورد پارامترها ارائه می شود. سپس مدل معرفی شده برای این نوع داده ها در مطالعات شبیه سازی مورد بررسی قرار می گیرد. همچنین پیشین های مناسب برای پارامترهای این مدل معرفی می شوند و عملکرد آن ها در حالت های مختلف بررسی می شوند. توزیع های پیشین انتخاب شده برای مدل بندی سهم شاغلین در خانوار براساس نتایج آمارگیری نیروی کار مرکز آمار که داده های نسبتی هستند، بکار گرفته و در پایان بحث و نتیجه گیری ارائه می شود.

۲ مدل بندی داده های نسبتی در دامنه بسته $[0$ و $1]$

توزیع بتا اغلب انتخابی مناسب برای برازش داده های پیوسته در بازه $(0, 1)$ است. تابع چگالی احتمال بتای بازپارامتریده شده برای متغیر تصادفی Y برحسب پارامترهای میانگین و دقت به صورت رابطه (۱.۱) است،

$$E(Y) = \mu \text{ و } Var(Y) = (\mu(1-\mu))/(1+\phi)$$

علی رغم سازگاری مدل رگرسیون بتا با این گونه داده ها، گاهی وجود تعداد قابل ملاحظه ای صفر یا یک در مشاهدات، تطبیق پذیری این مدل را به چالش می کشاند، زیرا تکیه گاه توزیع بتا بازه $(0, 1)$ است. راه حل مواجهه با این مشکل افزودن احتمال های صفر و یک به تابع

چگالی بتا و ایجاد توزیع آمیخته است. گالویس و همکاران (۲۰۱۴) مدل بتای افزوده‌صفر و یک را پیشنهاد کردند که شامل یک توزیع سه قسمتی است به‌گونه‌ای که در نقاط صفر و یک تباهیده و در بازه (۰، ۱) دارای چگالی بتا به صورت زیر است:

$$f(y|p_0, p_1, \mu, \phi) = \begin{cases} p_0 & y = 0 \\ p_1 & y = 1 \\ (\lambda - p_0 - p_1)\pi(y|\mu, \phi) & 0 < y < 1 \end{cases}$$

که در آن تابع چگالی بتا و $0 \leq p_0 + p_1 \leq 1$ و $p_0, p_1 \geq 0$ میانگین و واریانس این توزیع به ترتیب عبارتند از:

$$E(Y) = (\lambda - p_0 - p_1)\mu + p_1,$$

$$Var(Y) = p_1(\lambda - p_1) + (\lambda - p_0 - p_1) \left[\frac{\mu(\lambda - \mu)}{(\lambda + \phi)} + (p_0 + p_1)\mu^2 - 2\mu p_1 \right]$$

۱.۲ مدل رگرسیون آمیخته بتای افزوده:

فرض کنید $y_i = (y_{i1}, \dots, y_{in_i})$ برداری به طول n_i برای واحد نمونه i ام است که از توزیعی به صورت:

$$f(y_{ij}|p_0, p_1, \mu_{ij}, \phi_{ij}) = \begin{cases} p_0 & y_{ij} = 0 \\ p_1 & y_{ij} = 1 \\ (\lambda - p_0 - p_1)\pi(y_{ij}|\mu_{ij}, \phi_{ij}) & 0 < y_{ij} < 1 \end{cases} \quad (1.2)$$

پیروی می‌کند. که در آن تابع چگالی بتای (۱.۱) است، $p_0 = P(Y_{ij} = 0)$ و $p_1 = P(Y_{ij} = 1)$ بنابراین تابع چگالی y_{ij} برحسب تابع نشانگر برابر است با:

$$\pi(y_{ij}; p_0, p_1, \mu_{ij}, \phi_{ij}) = p_0^{I(y_{ij}=0)} p_1^{I(y_{ij}=1)} \{(\lambda - p_0 - p_1)\pi(y_{ij}; \mu_{ij}, \phi_{ij})\}^{(\lambda - I(y_{ij}=0))(\lambda - I(y_{ij}=1))}$$

در مدل رگرسیون بتا وقتی پارامتر دقت ثابت باشد، شرط همگنی واریانس که از جمله شرط‌های مدل‌های خطی تعمیم یافته است، برقرار می‌باشد؛ کرباری و زلیس (۲۰۱۰). بنابراین متغیرهای تبیینی و اثرات تصادفی را می‌توان براساس یک تبدیل مناسب پارامتر متغیر میانگین به صورت:

$$g(E(\mathbf{Y}_i | \mathbf{b}_i)) = g(\boldsymbol{\mu}_i) = \mathbf{X}_i^T \boldsymbol{\beta} + \mathbf{Z}_i^T \mathbf{b}_i$$

مدل‌بندی کرد. که در آن ماتریس طرح \mathbf{X}_i طرح $p \times n_i$ $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$ بردار اثرات ثابت، \mathbf{Z}_i ماتریس طرح $q \times n_i$ و $\mathbf{b}_i = (b_{i1}, \dots, b_{iq})^T$ بردار اثرات تصادفی است. برای تابع پیوند $g(\cdot)$ انتخاب‌های متفاوتی از جمله لجیت را می‌توان اختیار نمود.

۲.۲ مدل بندی همزمان میانگین و دقت

برای مدل بندی همزمان پارامترهای میانگین و دقت، لازم است تابع پیوند مناسبی برای پارامتر دقت نیز در نظر گرفته شود. فرض کنید $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})$ برداری به طول n_i برای واحد نمونه i ام است که از توزیعی براساس رابطه (۱.۲) با ثابت در نظر گرفتن پارامترهای p_0 و p_1 تبعیت می‌کند، مدل بندی پارامترهای میانگین و دقت را می‌توان به صورت زیر در نظر گرفت:

$$g(\boldsymbol{\mu}_i) = \mathbf{X}_i^T \boldsymbol{\beta} + \mathbf{Z}_i^T \mathbf{b}_i$$

$$h(\boldsymbol{\phi}_i) = \mathbf{W}_i^T \boldsymbol{\delta} + \mathbf{U}_i^T \mathbf{d}_i$$

که در آن \mathbf{U}_i ماتریس طرح $q^* \times n_i$ بعدی، $\mathbf{d}_i = (d_{i1}, \dots, d_{iq^*})^T$ بردار اثرات تصادفی می‌باشند. فیوگورا-زونیکا و همکاران (۲۰۱۳) برآوردهای پارامترهای مدل آمیخته خطی تعمیم‌یافته بتا را در وضعیت متغیر بودن پارامتر دقت با رهیافت بیزی مطالعه نمودند و

برای مدل بندی پارامتر میانگین، از پیوند لوجیت و برای مدل بندی پارامتر دقت، از پیوند لگاریتم استفاده کردند. با در نظر گرفتن پیوند لوجیت برای مدل بندی پارامتر میانگین و پیوند لگاریتم برای پارامتر دقت و صرف نظر کردن از اثرات تصادفی در مدل بندی پارامتر دقت داریم:

$$\begin{aligned} \log it(\mu_i) &= \mathbf{X}_i^T \beta + \mathbf{Z}_i^T \mathbf{b}_i \\ \ln(\phi_i) &= \mathbf{W}_i^T \delta \end{aligned}$$

بعباری پارامترها را می توان به صورت زیر در نظر گرفت:

$$\begin{aligned} \mu_{ij} &= \frac{\exp(\mathbf{X}_i^T \beta + \mathbf{Z}_i^T \mathbf{b}_i)}{1 - \exp(\mathbf{X}_i^T \beta + \mathbf{Z}_i^T \mathbf{b}_i)} = \frac{\exp(\eta_{ij})}{1 - \exp(\eta_{ij})} \\ \phi_{ij} &= \exp(\mathbf{W}_i^T \delta) = \exp(\tau_{ij}) \end{aligned}$$

که در آن β پارامترهای رگرسیونی و \mathbf{b}_i بردار اثرات تصادفی در مدل بندی پارامتر میانگین و δ پارامترهای رگرسیونی در مدل بندی پارامتر دقت رگرسیون آمیخته بتای افزوده است. همچنین i نشان دهنده فرد مورد نظر و اندیس j معرف تعداد تکرار فرد i ام است. در تحلیل بیزی برای اثرات تصادفی توزیع نرمال با میانگین صفر و واریانس σ_b^2 و برای پارامترها توزیع های پیشین ناآگاهی بخش یا کم آگاهی بخش در نظر گرفته می شود، بدین منظور در مواردی از بیز سلسله مراتبی نیز استفاده شده است. برای پارامترهای رگرسیونی در مدل بندی پارامتر دقت رگرسیون آمیخته بتای افزوده، گروه های متفاوتی از توزیع های پیشین را در نظر می گیریم. گروه اول، توزیع نرمال $\delta \sim N(0, \Sigma_\delta)$ است، که در آن Σ_δ ماتریس مثبت معین کوواریانس است. گروه دوم پیشین سلسله مراتبی معرفی شده توسط فیوگیورا-زونیکا و همکاران (۱۳۰۲) در مدل های آمیخته رگرسیون بتا می باشد، آنها برای δ توزیع $t - t$ استودنت بصورت $\delta \sim t(v_d, 0, \Sigma_\delta)$ و همچنین برای v توزیع نمایی $v \sim \exp(a)$ به ازای $a = 0.1$ در نظر گرفتند. گروه سوم مفروضات گروه دوم را داراست علاوه بر آن یک بیز سلسله مراتبی برای پارامتر Σ_δ نیز در نظر گرفته می شود، فونگ و همکاران (۲۰۱۰) برای مدل های آمیخته خطی تعمیم یافته واریانس توزیع $t - t$ استودنت، Σ_δ را دارای توزیع ویشارت معکوس $\Sigma_\delta \sim IWishart(\xi, \nu)$ در نظر گرفتند.

در سه مدل معرفی شده، برای بعضی از پارامترها، توزیع های پیشین یکسان است که به صورت زیر فرض می شوند:

$$\begin{aligned} \beta &\sim N_p(0, \Sigma_\beta) \quad \Sigma_\beta = \text{diag}(100, 100, 100) \\ p_0 &\sim U(0, 1) \quad p_1 \sim U(0, 1 - p_0) \\ \sigma_b^2 &\sim I\Gamma(0.1, 0.1) \end{aligned} \quad (2.2)$$

در مدل اول توزیع پیشین پارامتر δ به صورت زیر در نظر گرفته شده است:

$$\delta \sim N_q(0, \Sigma_\delta) \quad \Sigma_\delta = \text{diag}(100, 100, 100) \quad (3.2)$$

برای مدل دوم توزیع پیشین پارامتر δ به صورت زیر تعریف می شود:

$$\begin{aligned} \delta &\sim t_q(0, \Sigma_\delta) \quad \Sigma_\delta = \text{diag}(100, 100, 100) \\ v &\sim \exp(a), \quad a = 0, 1 \end{aligned} \quad (4.2)$$

در مدل سوم توزیع پیشین به صورت زیر به کار گرفته شده است:

$$\begin{aligned} \delta &\sim t_q(v, 0, \Sigma_\delta) \\ v &\sim \exp(a), \quad a = 0, 1 \\ \Sigma_\delta &\sim IWishart(\xi, \nu), \quad \xi = \text{diag}(0.439, 0.591), \quad \nu = 5 \end{aligned} \quad (5.2)$$

مقادیر پیشنهادی فونگ و همکاران (۲۰۱۰) بر اساس شبیه سازی های انجام شده $\xi = \text{diag}(0.439, 0.591)$ ، $\nu = 5$ می باشد، که در این مقاله نیز همان مقادیر در نظر گرفته می شود. با فرض استقلال پارامترها، چگالی پسین توام عبارت است از:

$$f(\boldsymbol{\eta}, \boldsymbol{\tau}, \beta, \delta, p_0, p_1; \mathbf{y}) \propto \prod_{i=1}^{n_i} \pi(\mathbf{y}_i | \boldsymbol{\tau}_i, \boldsymbol{\eta}_i) \prod_{i=1}^{n_i} \pi(\boldsymbol{\eta}_i | \beta) \prod_{i=1}^{n_i} \pi(\boldsymbol{\tau}_i | \delta) \pi(\beta) \pi(\delta) \pi(p_0) \pi(p_1)$$

از آن جا که به دست آوردن توزیع های پسین حاشیه ای بسیار پیچیده است از الگوریتم MCMC و نمونه گیر گیبز استفاده می شود. مراحل MCMC برای برآورد پارامترها با استفاده از بسته R2WinBUGS می باشد که دو نرم افزار R و WinBUGS را به هم متصل می کند.

۳.۲ انتخاب مدل بیزی

معیارهای مختلفی برای انتخاب مدل مناسب در استنباط‌های بیزی تعریف شده است. یکی از معیارهای معمول، آمار ترتیب پیش‌بینی شرطی (CPO) است. فرض کنید Y مجموعه داده‌ها و Y^{-i} به تمام داده‌ها بجز مشاهده i ام اشاره داشته باشد، پارامترهای نامعلوم به صورت Θ تعریف شوند و $\pi(\Theta|Y^{(-i)})$ توزیع پسین بردار پارامتر مدل، یعنی Θ به شرط $Y^{(-i)}$ باشد. برای هر مشاهده i ، معیار

$$CPO_i = \int_{\Theta} f(y_i|\theta) \pi(\theta|Y^{(-i)}) d\theta$$

می‌تواند با

$$CPO_i = \left(\frac{1}{m} \sum_{k=1}^m \frac{1}{f(y_i|\theta^{(k)}, Y)} \right)^{-1}$$

به‌دست آید که در آن m تعداد تکرارها بعد از دوره دورریز است (گلفن و دی، ۱۹۹۴). یک آماره CPO_i با لگاریتم درست‌نمایی حاشیه‌ای‌نما (LPML) تعریف می‌شود به‌طوری‌که:

$$LPML = \sum_{i=1}^n \log(CPO_i).$$

مقادیر بزرگ‌تر LPML برازش بهتر به مدل را نشان می‌دهند. برای مقایسه مدل‌ها معیارهای تشخیصی دیگری مانند EAIC و EBIC نیز می‌توان در نظر گرفت (کارلین و لوئیس، ۲۰۰۸). فرض کنید $\theta^{(k)}$ نمونه پسین MCMC تولید شده در تکرار k الگوریتم، η تعداد پارامترها، n تعداد مشاهدات و

$$\bar{D} = \frac{1}{m} \sum_{k=1}^m \sum_{i=1}^n \log f(y_i|\theta^{(k)}, Y), \quad k = 1, \dots, m$$

باشد آنگاه داریم:

$$EAIC = -2\bar{D} + 2\eta$$

$$EBIC = -2\bar{D} + \eta \log n$$

برعکس LPML مقادیر کوچک‌تر EAIC و EBIC بیانگر مناسب بودن مدل مربوط است.

۳ مطالعه شبیه‌سازی

در این قسمت برآورد بیزی پارامترهای مدل و دقت توزیع پیشین‌های معرفی شده مورد ارزیابی و مقایسه قرار می‌گیرند. در تمام حالت‌ها، ۱۰۰ مجموعه داده با اندازه نمونه $n = 100$ برای ۵ تکرار تولید شده‌است. برای پارامتر μ_{ij} ، پیوند لوجیت به‌صورت:

$$\logit(\mu_{ij}) = \beta_1 + \beta_2 x_{2ij} + \beta_3 x_{3ij} + b_i, \quad i = 1, \dots, 100, \quad j = 1, \dots, 5$$

و برای پارامتر ϕ_{ij} پیوند لگاریتم به‌صورت:

$$\log(\phi_{ij}) = \delta_1 + \delta_2 x_{2ij}, \quad i = 1, \dots, 100, \quad j = 1, \dots, 5$$

در نظر گرفته شده است، که در آن‌ها b_i دارای توزیع نرمال استاندارد، پارامترهای p و p_1 ثابت و برابر با ۰/۸ و x_{2ij} و x_{3ij} متغیرهای مستقل یکنواخت $U(0, 1)$ در نظر گرفته شده‌اند. برای پارامترهای رگرسیونی β_1 ، β_2 و β_3 به ترتیب مقادیر ۰/۵، ۰/۴ و ۰/۶ و برای پارامترهای رگرسیونی δ_1 و δ_2 به ترتیب مقادیر ۰/۳ و ۰/۷ در نظر گرفته شده است. الگوریتم تولید مقادیر y_{ij} به صورت زیر است:

گام اول - یک مقدار تصادفی d از توزیع برنولی با احتمال $\frac{1}{8}$ تولید شود.
 گام دوم - اگر $d = 1$ مقدار y_{ij} از توزیع بتای
 $B(\mu_{ij}\phi_{ij}, (1 - \mu_{ij})\phi_{ij})$ تولید شود اگر $d = 0$ مقدار y_{ij} از توزیع برنولی با احتمال موفقیت $\frac{1}{5}$ تولید شود.
 با تکرار این الگو، دنباله‌ای از y_{ij} هایی به دست می آید که با احتمال $\frac{1}{8}$ دارای توزیع بتا، با احتمال $\frac{1}{8}$ مقدار صفر و با احتمال $\frac{1}{8}$ مقدار یک را اختیار می کند. تحلیل های بیزی با لحاظ پیشین های (۲.۲) و سه مدل معرفی شده، انجام می شوند.
 برای هر مدل، دو زنجیر با مقادیر اولیه متفاوت، با $1,000,000$ تکرار برای هر زنجیر اجرا می شود. نتایج برای $500,000$ تکرار آخر ارائه می شود. به علاوه برای اجتناب از مسئله همبستگی زنجیرهای تولید شده، فاصله زنجیره ها 100 در نظر گرفته می شود.

جدول ۱: برآوردهای بیزی پارامترها پس از برازش مدل در اندازه نمونه های متفاوت

پارامتر	مقدار واقعی	مدل ۱ Est.(SD)	مدل ۲ Est.(SD)	مدل ۳ Est.(SD)
β_1	۰.۵	(۰.۱۴۵) ۰.۳۶۰	(۰.۱۳۰) ۰.۴۴۶	(۰.۱۳۵) ۰.۴۴۳
β_2	۰.۴	(۰.۲۱۸) ۰.۴۸۳	(۰.۲۳۳) ۰.۴۴۰	(۰.۲۰۵) ۰.۴۲۸
β_3	۰.۶	(۰.۱۶۷) ۰.۵۴۵	(۰.۱۷۹) ۰.۵۵۴	(۰.۱۷۳) ۰.۵۵۶
δ_1	۰.۳	(۰.۱۴۹) ۰.۳۱۵	(۰.۱۴۳) ۰.۲۸۸	(۰.۱۴۲) ۰.۳۱۲
δ_2	۰.۷	(۰.۲۲۳) ۰.۶۷۹	(۰.۲۶۷) ۰.۶۸۹	(۰.۲۵۸) ۰.۶۸۷
p_0	۰.۱	(۰.۰۱۶) ۰.۱۰۳	(۰.۰۱۳) ۰.۱۰۳	(۰.۰۱۱) ۰.۱۰۳
p_1	۰.۱	(۰.۰۱۷) ۰.۱۲۶	(۰.۰۱۵) ۰.۱۱۲	(۰.۰۱۵) ۰.۱۱۱
σ_b^2	۱	(۰.۱۴۹) ۰.۷۰۶	(۰.۱۴۲) ۰.۷۳۲	(۰.۱۳۲) ۰.۷۴۳

آزمون های همگرایی تشخیصی گلن روبین (۱۹۹۲) و همگرایی تشخیصی هیدل برگ و ولج (۱۹۸۱) دلالت بر همگرایی زنجیرهای تولید شده دارند. همچنین معیار همگرایی گلن و روبین توام نیز حدود ۱ به دست آمد که این معیار هم همگرایی زنجیرها را نشان می دهد. همان طور که در جدول ۱ ملاحظه می شود، سه مدل برآوردهای قابل قبولی برای پارامترها ارائه می نمایند. البته مدل های ۲ و ۳ نسبت به مدل ۱ از عملکرد بهتری برخوردار هستند و برآوردهای مناسب تری برای پارامترها ارائه می کنند و مدل ۳ نیز نسبت به مدل ۲ از مقبولیت بیشتری برخوردار است.

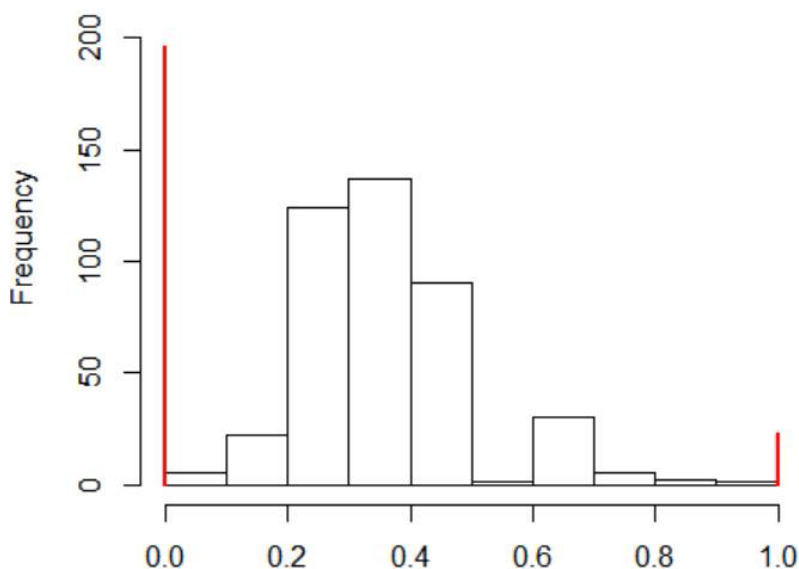
۴ تحلیل داده های نیروی کار

آمارگیری نیروی کار یکی از مهمترین نمونه گیری های خانواری دنیا است که جامعه آن همه اعضای خانوارهای معمولی ساکن هستند. در این نمونه گیری هر فصل از حدود ۶۰۰۰۰ خانوار آمارگیری می شود. برای هر فرد ۱۰ ساله و بیشتر خانوار منتخب نمونه، یک فرم وضع فعالیت تکمیل می شود. نمونه ها بر اساس الگوی چرخش ۲-۲-۲ انتخاب می شوند. یعنی هر خانوار ۲ فصل پی در پی در آمارگیری حضور دارد، در فصل های سوم و چهارم در نمونه ها قرار نمی گیرد و در دو فصل بعدی پنجم و ششم دوباره به عنوان عضو نمونه محسوب می شوند. بنابراین برای هر خانوار نمونه ۴ تکرار وجود دارد. برای برآورد پارامترهای جامعه، داده هایی که از این آمارگیری برای هر عضو خانوارهای نمونه جمع آوری شده است با سه مرحله وزن دهی بر اساس وزن پایه، وزن بی پاسخی و وزن پیش بینی های جمعیتی تعدیل می شوند. در این آمارگیری بر اساس وضع فعالیت تمام افراد ۱۰ ساله و بیشتر خانوارهای نمونه شاخص های مربوط به اشتغال و بیکاری و نسبت شاغلین در هر خانوار تعیین می شود. اگر تعداد اعضای شاغل خانوار کمتر از تعداد اعضای خانوار باشد، این نسبت، عددی بین صفر و یک، اگر تمام اعضای خانوار شاغل باشند، عدد یک و در صورتی که هیچ یک از اعضای خانوار شاغل نباشند، عدد صفر را اختیار می کند. بنابراین متغیر پاسخ یعنی نسبت شاغلین در خانوار، متغیری پیوسته با تحقق هایی در بازه بسته $[0, 1]$ است. در این مطالعه داده های مربوط به خانوارهای مشترک در فصل های اول و دوم سال های ۱۳۹۲ و ۱۳۹۳ آمارگیری نیروی کار شهر تهران در نظر گرفته شده است که شامل اطلاعات ۱۵۹ خانوار است. همان طور که در شکل ۱ ملاحظه می شود، نسبت شاغلان در خانوار بر اساس نتایج ۴ نوبت آمارگیری، متغیری پیوسته با تحقق هایی در بازه بسته $[0, 1]$ است.

مدل آمیخته رگرسیون بتای افزوده با $n = 159$ ، برای ۴ تکرار را برای این داده ها در نظر می گیریم. با فرض ثابت بودن پارامترهای p_0 و p_1 برای پارامتر μ_{ij} پیوند لوجیت به صورت

$$\log it(\mu_{ij}) = \beta_1 + \beta_2 x_{2ij} + \beta_3 x_{3ij} + b_i \quad 1, \dots, 159; \quad j = 1, \dots, 4$$

$$b_i \sim N(0, \sigma_b^2)$$



شکل ۱: هیستوگرام نسبت شاغلان در خانوارهای نمونه در فصل‌های اول و دوم سال‌های ۱۳۹۲ و ۱۳۹۳ آمارگیری نیروی کار شهر تهران.

و به‌طور همزمان برای پارامتر ϕ_{ij} نیز پیوند لگاریتم به‌صورت

$$\log(\phi_{ij}) = \delta_1 + \delta_2 x_{2ij} \quad i = 1, \dots, 159; \quad j = 1, \dots, 4$$

در نظر گرفته شده است، که در آن y_{ij} نسبت افراد شاغل در خانوار، x_{2ij} تعداد افراد ۱۰ ساله و بیشتر در خانوار و x_{3ij} بعد خانوار است. بر اساس توزیع‌های پیشین (۲.۲)، سه مدل معرفی شده به داده‌ها برازش شده‌اند. نتایج برای هر مدل، دو زنجیر مستقل با مقادیر اولیه متفاوت، با ۱,۰۰۰,۰۰۰ تکرار و با استفاده از ۵۰۰,۰۰۰ تکرار آخر و احتساب فاصله ۱۰۰ ارائه شده است.

جدول ۲: ارزیابی مدل‌های ۱ و ۲ و ۳ برای داده‌های نیروی کار

مدل	LPML	EAIC	EBIC
۳	-۲۴/۱۰	-۴۸/۱۸	-۱۸/۱۵
۲	-۲۵/۲۱	-۴۴/۱۹	-۱۵/۵۸
۱	-۲۸/۱۱	-۳۳/۷۸	-۱۱/۲۳

در جدول ۲ عملکرد مدل‌های ۱ و ۲ و ۳ برای داده‌های نیروی کار با استفاده از معیارهای انتخاب مورد ارزیابی و مقایسه قرار گرفته‌اند. همان‌طور که ملاحظه می‌شود عملکرد مدل ۳، نسبت به سایر مدل‌ها بهتر است.

جدول ۳: برآورد و خطای استاندارد پارامترهای مدل‌های ۱، ۲ و ۳ بر روی داده‌های آمارگیری نیروی کار شهر تهران برای سال‌های ۱۳۹۲ و ۱۳۹۳

پارامتر	مدل ۳		مدل ۲		مدل ۱	
	Est.	(SD)	Est.	(SD)	Est.	(SD)
β_1	۰/۶۳۹	۰/۱۵۱	۰/۶۳۸	۰/۱۹۹	۰/۴۸۹	۰/۲۰۳
β_2	۰/۱۸۸	۰/۰۵۹	۰/۶۰۳	۰/۲۲۸	۰/۵۷۳	۰/۲۵۲
β_3	-۰/۵۱۵	۰/۰۶۷	-۰/۶۰۶	۰/۲۱۷	-۰/۶۴۳	۰/۲۵۲
δ_1	۰/۳۱۰	۰/۰۱۸	۰/۴۳۰	۰/۱۵۱	۰/۲۱۲	۰/۱۵۵
δ_2	۰/۳۱۵	۰/۱۵۱	۰/۳۱۰	۰/۱۵۹	۰/۸۹۹	۰/۲۵۸
p_0	۰/۳۰۹	۰/۰۱۸	۰/۳۰۹	۰/۰۱۹	۰/۳۱۱	۰/۰۲۱
p_1	۰/۰۳۷	۰/۰۰۷	۰/۰۳۸	۰/۰۰۸	۰/۰۳۸	۰/۰۰۹
σ_b^2	۰/۱۲۶	۰/۰۲۲	۰/۱۵۶	۰/۱۰۳	۰/۲۵۵	۰/۲۱۴

جدول ۳ برآوردهای پارامترهای پسین بدست آمده از برازش مدل ها را بر داده های آمارگیری نیروی کار نشان می دهد. بازه باور ۹۵٪ تنها برای مدل برتر ارائه شده است. نتایج نشان می دهند که افراد ۱۰ ساله و بیشتر خانوار با اثر مثبت و متغیر مستقل بعد خانوار بالعکس با اثر منفی در مدل میانگین وارد می شود، همچنین افراد ۱۰ ساله و بیشتر خانوار با اثر مثبت بر مدل دقت تاثیرگذار است. پارامترهای p_1 و p_0 که بیانگر خانوارهای فاقد شاغل و خانوارهایی با تمام اعضای شاغل هستند نیز به ترتیب برابر ۰/۳۱ و ۰/۰۴ برآورد شدند. همبستگی داده ها طی ۴ فصل، با اثر تصادفی در نظر گرفته شده است و مقدار واریانس آن ۰/۱۳ برآورد شده است.

۵ ارزیابی و مقایسه مدل بندی همزمان میانگین و دقت با مدل بندی میانگین

حال برآوردهای حاصل از مدل منتخب که برای مدل بندی همزمان میانگین و دقت در مدل های آمیخته رگرسیون بتای افزوده انتخاب شد با برآوردهایی که از مدل بندی میانگین در مدل آمیخته رگرسیون بتای افزوده با اثر تصادفی نرمال به دست آمد، ارزیابی و مقایسه می شوند.

$$\begin{aligned} \text{logit}(\mu_{ij}) &= \beta_1 + \beta_2 x_{2ij} + \beta_3 x_{3ij} + b_i, \quad b_i \sim N(0, \sigma_b^2), \quad i = 1, \dots, 100, \quad j = 1, \dots, 5, \\ \log(\phi_{ij}) &= \delta_1 + \delta_2 x_{2ij} \end{aligned} \quad (1.5)$$

$$\text{logit}(\mu_{ij}) = \beta_1 + \beta_2 x_{2ij} + \beta_3 x_{3ij} + b_i, \quad b_i \sim N(0, \sigma_b^2), \quad i = 1, \dots, 100, \quad j = 1, \dots, 5 \quad (2.5)$$

در جدول ۴ عملکرد دو مدل برای برازش به داده های نیروی کار براساس معیارهای انتخاب نشان داده می شود. مقادیر بزرگتر LPML مربوط به مدل بندی همزمان میانگین و دقت با لحاظ اثر تصادفی نرمال است که برازش بهتر مدل را نشان می دهد. همین مدل مقادیر کوچکتر EAIC و EBIC که نشان دهنده برازش بهتر به داده ها است را نیز ارائه می کند. بنابراین مدل بندی همزمان میانگین و دقت برازش بهتری نسبت به مدل بندی میانگین در مدل های آمیخته رگرسیون بتای افزوده بر داده های نیروی کار ارائه می نماید.

جدول ۴: ارزیابی مدل همزمان میانگین و دقت

EBIC	EAIC	LPML	مدل بندی میانگین
-۴/۲۱	-۲۸/۷۶	-۳۰/۵۶۱	میانگین
-۱۸/۱۵	-۴۸/۱۸	-۲۴/۱۰	همزمان میانگین و دقت

۶ نتیجه گیری و پیشنهادات

در این مقاله مدل بندی همزمان میانگین و دقت در مدل های آمیخته رگرسیون بتای افزوده با پیشین های پیشنهادی معرفی شدند و در مطالعات شبیه سازی، با لحاظ اریبی و خطای استاندارد برآورد پارامترها، همچنین براساس معیارهای مدل های بیزی، بهترین مدل معرفی شد. برای تحلیل داده های آمارگیری نیروی کار مرکز آمار ایران، داده های مربوط به خانوارهای مشترک در فصل های اول و دوم سال های ۱۳۹۲ و ۱۳۹۳ شهر تهران، در نظر گرفته و نسبت شاغلین در خانوار براساس دو مدل معرفی شده مدل بندی گردید. از آن جا که مهمترین مشکل برآورد وضعیت نیروی کار در سطح کمتر از استان عدم وجود متغیرهای کمکی مناسب است یکی از موارد کاربرد این مدل ساختن متغیر تبیینی مناسب برای برآورد وضعیت نیروی کار نواحی کوچک است زیرا با ضرب سهم شاغلین در خانوار در تعداد خانوار این نواحی، تعداد شاغلین این نواحی قابل برآورد است.

استفاده از الگوریتم های MCMC در مدل های آمیخته رگرسیون بتای افزوده خصوصا وقتی تعداد پارامترها از حالت معمول بیشتر می شوند، بسیار زمان بر است. بنابراین استفاده از روش های تقریب برای کوتاه کردن زمان اجرای برنامه ها و همچنین مدل بندی همزمان میانگین و دقت در مدل های آمیخته رگرسیون بتای افزوده با اثرات تصادفی وابسته فضایی می توانند موضوعاتی برای مطالعات بعدی باشند تا بتوان مدل های جامع تری ارائه نمود.

فهرست منابع

[۱] نتایج آمارگیری نیروی کار، (۱۳۹۲)، تهران، مرکز آمار ایران.

- [2] Branscum, A.J., Johnson, W.O. and Thurmond, M. (2007), *Bayesian Beta Regression Applications to Household Expenditure Data and Genetic Distance Between Food and Mouth Diseases Viruses*, Australian & New Zealand Journal of Statistics, 49, 287-301.
- [3] Bonat, W.H., Ribeiro, P.J. and Zeviani, W.M. (2013), *Likelihood Analysis for a Class of Beta Mixed Models* Cornell University Library, arXiv Preprint arXiv: 1312.2413.
- [4] Cepeda, E. D., and Gamerman, D. (2005), *Bayesian Methodology for Modeling Parameters in the Two Parameter Exponential Family*, Revista Estadística, 57, 168-169.
- [5] Cepeda, E. D., Migon, H. S., Garrido, L. and Achcar, J. A. (2014), *Generalized Linear Models with Random Effects in the Two-Parameter Exponential Family*, Journal of Statistical Computation and Simulation, 84, 513-525.
- [6] Carlin, B. P. and Louis, T. A. (2008), *Bayesian Methods for Data Analysis*, Mineapolis, CRC Press.
- [7] Ferrari, S. and Cribari, F. (2004), *Beta Regression for Modelling Rates and Proportions*, Journal of Applied Statistics, 31, 799-815.
- [8] Figueroa-Zúñiga, J. I., Arellano-Valle, R. B. and Ferrari, S. L. (2013), *Mixed Beta Regression: A Bayesian Perspective*, Computational Statistics & Data Analysis, 61, 137– 147.
- [9] Fallah Mohsenkhani, Z., Mohammadzadeh, M. and Baghfalaki, T. (2019), *Augmented Mixed Beta Regression Models with Skew-Normal Independent Distributions: Bayesian Analysis of Labor Force Data*, Communications in Statistics-Simulation and Computation, Volume 48, Issue 7, 2147-2164.
- [10] Fong, Y., Rue, H. and Wakefield, J. (2010), *Bayesian Inference for Generalized Linear Mixed Models*, Biostatistics, 11, 397-412.
- [11] Galvis, M. D., Dipankar, B. and Victor, H. L. (2014), *Augmented Mixed Beta Regression Models for Periodontal Proportion Data*, Preprinted, (In Press), Statistics in Medicine.
- [12] Gelman, A., Rubin, D. B., (1992). *Inference from Iterative Simulation Using Multiple Sequences*, Statistical Science, 7, 457–511.
- [13] Gelfand, A. E. and Dey, D. K. (1994), *Bayesian Model Choice: Asymptotics and Exact Calculations*, Journal of the Royal Statistical Society, Series B (Methodological), 501-514.
- [14] Heidelberger, P. and Welch, P. D. (1981), *A Spectral Method for Confidence Interval Generation and Run Length Control in Simulations*, Communications of the ACM, 24, 233-245.
- [15] Nogarotto, D, Azevedo, C, Bazan, J. Bayesian (2020), *modeling and prior sensitivity analysis for zero-one augmented beta regression models with an application to psychometric data*, Brazilian Journal of Probability and Statistics, 304-322.
- [16] Ospina, R. and Ferrari, S. L. (2010), *Inflated Beta Distributions*, Statistical Papers. 51, 111- 126.
- [17] Paolino, P. (2001), *Maximum Likelihood Estimation of Models with Beta-Distributed Dependent Variables*, Political Analysis, 9, 325-346.
- [18] Parker, A, Bandyopadhyay, D and Slate, E. (2014), *A spatial augmented beta regression model for periodontal proportion data*, Statistical Modelling, vol. 14, 503-521.

- [19] Smithson, M. and Verkuilen, J. (2006), *A Better Lemon Squeezer? Maximum-Likelihood Regression with Beta-Distributed Dependent Variables*, Psychological Methods, 11, 54-71.
- [20] Verkuilen, J. and Smithson, M. (2012), *Mixed and Mixture Regression Models for Continuous Bounded Responses Using the Beta Distribution*, Journal of Educational and Behavioral Statistics, 37,82-113.
- [21] Zimprich, D. (2010), *Modeling Change in Skewed Variables Using Mixed Beta Regression Models*, Research in Human Development, 7, 9-26.



Simultaneous Modeling of Mean and Variance In Augmented Mixed Beta Regression

Zohreh Fallah Mohsenkhani¹, Parvin Azhdari² †

⁽¹⁾ Statistical Research and Training Center, Tehran, Iran

⁽²⁾ Department of Statistics, North Tehran Branch, Islamic Azad University, Tehran, Iran

Communicated by: Mohammad Reza Zadkarami

Received: 2020/8/9

Accepted: 2021/10/31

Abstract: Augmented Beta Regression models are used for modeling data such as rate, ratio or percentage. This model is made by combining the Beta distribution on the interval (0,1) and two degenerate distributions at 0 and 1. By reparameterizing the beta distribution, the mean and precision parameters are modeled with a structure including fixed and random effects. In this paper, simultaneous modeling of mean and precision the augmented mixed beta regression models is presented and the model efficiency in simulation studies by Bayesian approach is investigated. Next, the application of this model to analyze the proportions of employed persons in every household based on the results of the Statistical Centre of Iran is shown and at the end, conclusion and results are presented.

Keywords: Augmented Beta Regression, Precision Parameter, Mixed Model, Bayesian Analysis, Labour Force Survey.



©2022 Shahid Chamran University of Ahvaz, Ahvaz, Iran. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0 license) (<http://creativecommons.org/licenses/by-nc/4.0/>).

†Corresponding author.

E-mail addresses: zoherhf@yahoo.com (Z. Falah), p_azhdari@iau-tnb.ac.ir (P. Azhdari).