

رهیافت بیزی نیمه پارامتری و کاربرد آن در مدل‌های با اثرات آمیخته

کیوان عسلی صاف^۱

چکیده:

متداول ترین مدل‌ها برای تحلیل داده‌های وابسته در علوم مختلف، مدل‌های خطی تعمیم یافته با اثرهای آمیخته است. یک مسأله اساسی در استفاده از این مدل‌ها بکارگیری فرآیند مناسبی برای تولید توزیع پیشین اثرهای آمیخته است. با فرض این که توزیع این اثرها از فرآیند دیریکله با توزیع پایه و پارامتر دقت مشخص پیروی می‌کنند، مدل‌های بیزی سلسله مراتبی نیمه پارامتری برای تحلیل داده‌های وابسته در رگرسیون خطی تعمیم یافته مورد بررسی قرار می‌گیرند. بدین منظور، ابتدا فرآیند دیریکله معرفی و سپس روش بکارگیری آن به عنوان فرآیند تولید توزیع‌های پیشین توضیح داده خواهد شد. در ادامه، توزیع‌های پسین شرطی کامل پارامترهای مدل‌های رگرسیونی با اثرهای آمیخته را یافته و توسط رهیافت بیز، استنباط آماری را با الگوریتم نمونه‌گیری گیبز انجام خواهیم داد. در نهایت در قالب مثالی، مدل مذکور به داده‌های واقعی برازش داده خواهد شد تا اهمیت کاربرد روش‌های بیز نیمه پارامتری در مقایسه با روش‌های متداول مشخص می‌شود.

واژه‌های کلیدی: مدل‌های رگرسیون خطی تعمیم یافته، مدل‌های بیز سلسله مراتبی، اثرهای آمیخته، فرآیند دیریکله، نمونه‌گیری گیبز، توزیع‌های پسین شرطی کامل.

۱ مقدمه

دیریکله با توزیع پایه و پارامتر دقت مشخص پیروی می‌کنند، ایده‌ی اصلی این روش بر طبق فرگوسن^۲ [۴] و آنتونیاک^۳ [۱] می‌باشد. در نهایت برای مقایسه‌ی روش‌های ارائه شده با روش‌های متداول بیز، به برآورد پارامترهای مجهول و اثرات آمیخته در مدل‌های خطی می‌پردازیم که این امر با محاسبه توزیع‌های شرطی کامل پارامترهای

واقف بودن محقق به پاره‌ای از اطلاعات قبلی در مورد مدل‌های با اثرات آمیخته و پارامترهای آن اهمیت ویژه‌ای در حصول نتایج برآوردیابی دارد. بنابراین یک مسأله‌ی اساسی در استفاده از این مدل‌ها بکارگیری فرآیند مناسبی برای تولید توزیع پیشین اثرهای آمیخته است. در این مقاله، فرض می‌کنیم که توزیع این اثرها از فرآیند

^۱ کارشناس ارشد آمار، مدرس دانشگاه محقق اردبیلی
^۲ Ferguson
^۳ Antoniak

۲ تعاریف اولیه

قبل از اینکه روش‌های مدل‌سازی در جوامع آمیخته را شرح دهیم و به چگونگی استفاده از این روش‌ها بپردازیم، برخی از مفاهیمی را که در برارزش مدل‌های آمیخته مفید هستند شرح می‌دهیم.

۱.۲ توزیع آمیخته

فرض کنید که متغیر تصادفی Y بتواند از یکی از k جامعه (مؤلفه) مفروض به دست آید به طوری که هر کدام از این مؤلفه‌ها دارای یک تابع چگالی احتمال متفاوت با بقیه باشند. اگر Y از i امین مؤلفه با احتمال w_i و تابع احتمال $f(Y|\theta_i)$ آمده باشد، آنگاه تابع چگالی احتمال یک مشاهده Y را می‌توان به شکل آمیخته زیر معرفی کرد

$$f(Y|\phi) = \sum_{i=1}^k w_i f(Y|\theta_i) \quad (1)$$

که در آن $\phi = (k, w_1, \dots, w_k, \theta_1, \dots, \theta_k)$ تعداد مؤلفه‌ها و معلوم و w_i ها وزن‌های آمیختگی هستند که نامنفی بوده و مجموعشان برابر یک است و θ_i ها پارامترهای نامعلوم مربوط به مؤلفه i ام را معرفی می‌کند.

۲.۲ توزیع دیریکله

توزیع دیریکله ما را به سوی درک فرآیند دیریکله آمیخته که به اختصار با DPM^{۱۰} نشان خواهیم

مدل با اثرات آمیخته بر طبق ایشوارن و جیمز^۴ میسر می‌باشد. سپس با توجه به معیارهای اعتبار مدل‌ها از قبیل معیار اطلاع پراکندگی (DIC^۵)، پراکندگی^۶، AIC^۷ و غیره، به مقایسه و سنجش روش ارائه شده با روش‌های معمول بیز خواهیم پرداخت که معیارهای اعتبار ذکر شده برای اولین بار توسط اسپینگل هالتر و همکاران^۸ [۹] مطرح گردید.

در بحث مدل‌سازی برخی از جوامع مورد بررسی را می‌توان به صورت افزایی از چند جامعه کوچکتر در نظر گرفت و می‌توان فرض کرد که داده‌های مورد بررسی با یک سری احتمالات خاص متعلق به این زیر جوامع هستند. لذا با توجه به افزای داده‌هایی که مورد بررسی محقق می‌باشند، برآورد پارامترهای این نوع جوامع نیازمند مدل‌سازی‌های مناسب می‌باشند که بر اساس کاتر و همکاران^۹ [۷] به این امر پرداخته خواهد شد. علاوه بر این چگونگی اتخاذ توزیع‌های پیشین برای پارامترهای این نوع جوامع و مدل‌های آمیخته در آمار بیز حائز اهمیت فراوان می‌باشد. از این رو با اتخاذ پیشین‌های مناسب، روش مطلوبی را برای برآورد پارامترهای مجهول مدل‌های برارزش داده شده به جوامع ذکر شده پیشنهاد می‌دهیم.

^۴ Ishwaran and James

^۵ Deviance Information Criterion

^۶ Deviance

^۷ Akaike Information Criterion

^۸ Spiegelhalter and et al.

^۹ Kutner and et al.

^{۱۰} Dirichlet Process Mixture

۱.۳.۲ خواص مهم فرآیند دیریکله

فرض کنید که G یک فرآیند دیریکله روی فضای (X, A) با پارامترهای ثابت G و α باشد، بنابراین فرگوسن [۴] و آنتویاک [۱] نتایج زیر را داریم: اگر π_1, \dots, π_n یک نمونه تصادفی n تایی از G باشد، داریم:

۱- اگر فرض کنیم $\pi \sim G$ آنگاه:

$$(\pi_2 | \pi_1) \sim \frac{\alpha}{\alpha+1} G_0 + \frac{1}{\alpha+1} \delta_{\pi_1}(\pi_2)$$

⋮

$$(\pi_n | \pi_1, \dots, \pi_{n-1}) \sim \frac{\alpha}{\alpha+n-1} G_0$$

$$+ \frac{1}{\alpha+n-1} \sum_{j=1}^{n-1} \delta_{\pi_j}(\pi_n),$$

که در آن $\delta_{\pi_j}(\pi_n)$ مشخص کننده جرم نقطه‌ای واحد در $\pi_n = \pi_j$ می‌باشد، یعنی

$$\delta_{\pi_j}(\pi_n) = \begin{cases} 1 & , \quad \pi_j = \pi_n \\ 0 & , \quad \pi_j \neq \pi_n \end{cases}$$

۲- توزیع شرطی کامل

$(\pi_i | \pi_1, \dots, \pi_{i-1}, \pi_{i+1}, \dots, \pi_n)$ برای $i = 1, 2, \dots, n$ به صورت زیر است:

$$f(\pi_i | \pi_1, \dots, \pi_{i-1}, \pi_{i+1}, \dots, \pi_n) = \frac{\alpha}{\alpha+n-1} G_0 + \frac{1}{\alpha+n-1} \sum_{j \neq i} \delta_{\pi_j}(\pi_i) \quad (3)$$

۳- اگر $\pi = (\pi_1, \dots, \pi_n)$ یک نمونه تصادفی n تایی از G باشد، آنگاه نمونه متوالی $(n+1)$ ام دارای توزیع زیر است:

$$(\pi_{n+1} | \pi) = \frac{\alpha}{\alpha+n} G_0(\pi_{n+1}) + \frac{1}{\alpha+n} \sum_{i=1}^n \delta_{\pi_i}(\pi_{n+1}). \quad (4)$$

داد، سوق می‌دهد. توزیع دیریکله یک توزیع چند پارامتری تعمیم یافته از توزیع بتا است. فرض کنید $\Theta = \{\theta_1, \theta_2, \dots, \theta_n\}$ یک توزیع احتمال روی فضای گسسته $X = \{X_1, X_2, \dots, X_n\}$ باشد، به طوری که $P(X = X_i) = \theta_i$ که در آن X یک متغیر تصادفی در فضای X می‌باشد. توزیع دیریکله روی Θ به صورت زیر است:

$$P(\Theta | \alpha, M) = \frac{\Gamma(\alpha)}{\prod_{i=1}^n \Gamma(\alpha m_i)} \prod \theta_i^{\alpha m_i - 1} \quad (2)$$

که در آن $i = 1, 2, \dots, n$ و $M = \{m_1, m_2, \dots, m_n\}$ اندازه پایه تعریف شده روی X و بردار میانگین Θ و پارامتر دقتی است که چگونگی تمرکز توزیع حول M را نشان می‌دهد.

۳.۲ فرآیند دیریکله

اگر فرض کنیم که (X, A) یک فضای اندازه با یک اندازه احتمال G_0 و α یک عدد حقیقی مثبت باشد، یک فرآیند دیریکله، توزیع یک اندازه تصادفی G روی فضای (X, A) می‌باشد، به طوری که اگر برای هر افراز (A_1, \dots, A_k) از X ، بردار تصادفی $(G(A_1), \dots, G(A_k))$ دارای یک توزیع دیریکله با بعد متناهی به صورت زیر باشد:

$$(G(A_1), \dots, G(A_k)) \sim$$

$$Dir(\alpha \circ G_0(A_1), \dots, \alpha \circ G_0(A_k))$$

در این صورت G را دارای فرآیند دیریکله با پارامتر تمرکز α و اندازه پایه G_0 گویند که با نماد زیر نشان داده می‌شود:

$$G \sim DP(\alpha \circ G_0)$$

پیشین معلوم برای پارامترهای مدل یعنی θ_i ها است. برای ضرایب آمیختگی مفروض، متغیرهای نشانگر c_i ها دارای توزیع چند جمله‌ای می‌باشند.

توجه کنید که متغیرهای نشانگر پنهان یعنی c_i ها تنها به دلیل سهولت محاسبات بکار برده شده‌اند. حال اگر تعداد اجزا یا موارد آمیخته از قبل معلوم باشد، پارامترها برای هر مورد می‌توانند از G استخراج گردند و می‌توان توزیع احتمال دیریکله را به عنوان تابع چگالی احتمال ضرایب آمیختگی انتخاب کرد.

در اینجا حالت حدی این مدل را وقتی که $k \rightarrow \infty$ بررسی می‌کنیم. فرگوسن [۴] نشان داد که در حالت حدی، فرآیند دیریکله جایگزین توزیع دیریکله می‌گردد، به طوریکه برای هر نشانگر c_i به شرط نشانگرهای قبلی از توزیع چندجمله‌ای، θ_i متنظاری وجود دارد که از توزیع G استخراج می‌شود. در حالت حدی $k \rightarrow \infty$ ، این برچسب‌ها دارای توزیعی پیوسته فرض می‌شوند و می‌توان از استفاده از این برچسب‌ها در مدل صرف نظر و در عوض فرض کرد که پارامترها از یک فرآیند دیریکله با اندازه پایه G استخراج می‌شوند. از این رو مدل آمیخته فرآیند دیریکله (DPM) به صورت زیر معرفی می‌گردد:

$$\begin{aligned} Y_i | \theta_i &\sim f(Y | \theta_i) \\ \theta_i | G &\sim G(\theta) \\ G &\sim DP(\alpha G_0(\theta)), \end{aligned} \quad (6)$$

بنابراین، نمونه بعدی π_{n+1} به شرط π ، مشخص کننده یک مقدار مجزای جدید با احتمال $\frac{\alpha}{\alpha+n}$ می‌باشد و به عبارت دیگر نمونه بعدی به طور یکنواخت از بین n مقدار اولیه π_1, \dots, π_n استخراج می‌شود که مقدار جدید می‌تواند از توزیع G نیز استخراج گردد.

۳ مدل آمیخته فرآیند دیریکله (DPM)

یک مدل آمیخته به صورت $Y_i \sim \sum_{i=1}^k \pi_i f(Y | \theta_i)$ را در نظر بگیرید، که در آن Y به صورت آمیخته‌ای از توزیع‌های دارای شکل پارامتری یکسان f توزیع شده‌است که این توزیع‌های پارامتری تنها در پارامترهایشان با هم متفاوت می‌باشند. همچنین فرض بر این است که برای $i = 1, \dots, k$ ، پارامترهای θ_i از توزیع معلوم G استخراج شده‌اند. این مدل آمیخته را می‌توان به صورت سلسله مراتبی به صورت زیر بیان کرد:

$$\begin{aligned} Y_i | c_i, \Theta &\sim f(Y | \theta_{c_i}) \\ c &= ((c_1, \dots, c_k) | \pi_{1:k}) \\ &\sim Discrete(\pi_1, \dots, \pi_k) \\ \theta_i &\sim G_0(\theta), \\ (\pi_1, \dots, \pi_k) &\sim Dir(\alpha, M) \end{aligned} \quad (5)$$

که در آن c_i ها نشانگرها یا برچسب‌هایی هستند که مقادیر Y_i را به مقدار پارامتر θ_{c_i} نسبت می‌دهند و π_i ها ضرایب آمیختگی می‌باشند که از یک توزیع دیریکله استخراج شده‌اند و G توزیع

بنابراین کاربرد DPM برای یک مسأله در حالت کلی به این صورت مطرح می‌گردد که هر مدل پارامتری برای مقادیر Y_i به صورت سلسله مراتبی به حالت

$$\begin{aligned} Y_i | \theta_i &\sim f(Y | \theta_i) \\ \theta_i | \Psi &\sim h(\theta | \Psi) \end{aligned} \quad (۸)$$

می‌تواند با یک مدل DPM از رابطه (۶) جایگزین گردد. البته باید در نظر داشت که در مرحله دوم، می‌توان از پیشین پارامتری $h(\theta | \Psi)$ صرف نظر کرد و در عوض آن را با یک توزیع عمومی G که دارای پیشین فرآیند دیریکله است، جایگزین می‌کنیم.

حال به منظور نشان دادن نحوه بکارگیری فرآیند دیریکله به عنوان توزیع پیشین در آمار بیز و چگونگی استفاده از آن به منظور استخراج نمونه از جوامع آمیخته، در این بخش به تشریح نمونه‌گیری با استفاده از مدل‌های آمیخته فرآیند دیریکله می‌پردازیم.

۴ نمونه‌گیری پسین تحت

DPM

مدل (۶) را در نظر بگیرید که برای $i = 1, \dots, n$ توزیع پسین حاشیه‌ای θ_i به شرط دیگر پارامترها در آن مطلوب است. بردار $\theta^{(i-)}$ را به مفهوم همه

که در آن $DP(\alpha G_0(\theta))$ فرآیند دیریکله با اندازه پایه G_0 و پارامتر دقت α را نشان می‌دهد و G یک توزیع تصادفی استخراج شده از DP می‌باشد. نکته قابل ذکر در اینجا این است که از آنجایی که در بخش اول مدل (۶)، $f(Y | \theta_i)$ یک توزیع معلوم و در بخش دوم $G(\theta)$ یک توزیع عمومی در نظر گرفته می‌شود، در اغلب موارد از مدل (۵) به عنوان یک مدل نیمه پارامتری سلسله مراتبی بیزی^{۱۱} نام برده می‌شود.

توزیع‌های حاشیه‌ای پارامترهای مجهول مدل (۶) از آن جهت حائز اهمیت می‌باشد که با در دست داشتن این توزیع‌ها، می‌توانیم با استفاده از روش نمونه‌گیری گیبز به برآورد این پارامترها می‌پردازیم. بنابراین در این قسمت به یافتن توزیع‌های حاشیه‌ای پارامترهای مجهول مدل می‌پردازیم.

برای یافتن توزیع‌های حاشیه‌ای پارامترهای مجهول مدل (۶)، با توجه به پیشین $\theta_i | G \sim G(\theta)$ و $G \sim DP(\alpha G_0(\theta))$ و با استفاده از خاصیت ۳ از خواص دیریکله، می‌توانیم برای $i = 1, \dots, k$ توزیع حاشیه‌ای $\theta_i | (\theta_1, \dots, \theta_{i-1})$ را به صورت

$$\begin{aligned} P(\theta_i = \theta | \theta_1, \dots, \theta_{i-1}) &= \frac{\alpha}{\alpha + i - 1} G_0(\theta) \quad (۷) \\ &+ \frac{1}{\alpha + i - 1} \sum_{j=1}^{i-1} \delta_{\theta_j}(\theta) \end{aligned}$$

به دست آوریم که این فرمول توسط بلکول و مک کوئین^{۱۲} [۲] مطرح شده است.

^{۱۱} Bayesian Hierarchical Semi - Parametric Models
^{۱۲} Blackwell and McQueen

توزیع پسین حاشیه‌ای پارامتر θ_i به صورت زیر به دست می‌آید:

$$P(\theta_i | \theta^{(i-)}, Y_i) = b \alpha G_{\circ}(\theta_i) P(Y_i; \theta_i) + \sum_{\substack{j=1 \\ j \neq i}}^n P(Y_i; \theta_j) \delta_{\theta_j}(\theta_i) \quad (13)$$

که در آن به b ثابت نرمال ساز می‌گوییم و $\delta_{\theta_j}(\theta_i)$ تابع جرم نقطه‌ای جرم احتمال در نقطه θ_i می‌باشد و در واقع q_{\circ} نیز توزیع حاشیه‌ای Y_i است. بنابراین با توجه به توزیع پسین θ_i به صورت

$$h(\theta_i | Y_i) = \frac{G_{\circ}(\theta_i) P(Y_i; \theta_i)}{\int_{\theta} G_{\circ}(\theta_i) P(Y_i; \theta_i)}$$

توزیع پسین حاشیه‌ای پارامتر θ_i را در (۱۳) می‌توان به صورت زیر نوشت:

$$P(\theta_i | \theta^{(i-)}, Y_i) = b \alpha q_{\circ} h(\theta_i | Y_i) + b \sum_{\substack{j=1 \\ j \neq i}}^n P(Y_i; \theta_j) \delta_{\theta_j}(\theta_i) \quad (14)$$

این عبارت را می‌توان به شکل آمیخته‌ای از توزیع پسین حاشیه‌ای روی θ_i به صورت زیر نوشت:

$$P(\theta_i | \theta^{(i-)}, Y_i) = \begin{cases} \theta_j & \text{با احتمال } P(Y_i; \theta_j) \\ \sim h(\theta | Y_i) & \text{با احتمال } \alpha q_{\circ} \end{cases} \quad (15)$$

بنابراین با در دست داشتن توزیع پسین حاشیه‌ای پارامترهای مدل، به سادگی می‌توان الگوریتم نمونه گیر گیبز را به کار برد.

θ_i ها، به جز θ_i ، معرفی می‌کنیم. هدف کلی پیدا کردن توزیع پسین θ_i به شرط داده‌ها می‌باشد که با استفاده از قانون بیز به صورت زیر به دست می‌آید:

$$P(\theta_i | \theta^{(i-)}, Y_i) \propto P(Y_i | \theta_i) P(\theta_i | \theta^{(i-)}) \quad (9)$$

که در آن احتمالات به طور ضمنی روی پارامترهای فرآیند دیریکله، شرطی شده‌اند. تابع $P(\theta_i | \theta^{(i-)})$ به عنوان توزیع پیشین، طبق (۶) عبارت است از

$$P(\theta_i | \theta^{(i-)}) = \frac{\alpha}{\alpha + n - 1} G_{\circ}(\theta_i) + \frac{1}{\alpha + n - 1} \sum_{\substack{j=1 \\ j \neq i}}^n \delta_{\theta_j}(\theta_i) \quad (10)$$

توجه شود که $P(Y_i | \theta_i)$ به عنوان تابع درست‌نمایی و معادل $f(Y | \theta) = f(Y; \theta)$ در مدل (۵) می‌باشد. بنابراین توزیع پسین θ_i را می‌توان با استفاده از قانون بیز بدست آورد. در نتیجه داریم:

$$P(\theta_i = \theta | \theta^{(i-)}, Y_i) = \frac{P(Y | \theta_i) P(\theta_i | \theta^{(i-)})}{\int P(Y | \theta_i) P(\theta_i | \theta^{(i-)}) d\theta_i} \quad (11)$$

حال با جایگذاری (۱۰) در (۱۱) به دست می‌آید:

$$P(\theta_i | \theta^{(i-)}, Y_i) = \frac{\alpha P(Y_i | \theta_i) G_{\circ}(\theta_i) + P(Y_i | \theta_i) \sum_{\substack{j=1 \\ j \neq i}}^n \delta_{\theta_j}(\theta_i)}{\int \alpha P(Y_i | \theta_i) G_{\circ}(\theta_i) d\theta_i + \int P(Y_i | \theta_i) \sum_{\substack{j=1 \\ j \neq i}}^n \delta_{\theta_j}(\theta_i) d\theta_i}$$

که با در نظر گرفتن

$$b = (\alpha q_{\circ} + \sum_{\substack{j=1 \\ j \neq i}}^n P(Y_i; \theta_j))^{-1} \quad (12)$$

$$q_{\circ} = \int_{\theta} G_{\circ}(\theta) P(Y_i | \theta) d\theta$$

۵ DPM در مدل رگرسیون خطی چندگانه با اثرهای آمیخته

مدل رگرسیون خطی با اثرهای تصادفی را به صورت زیر در نظر بگیرید:

$$Y_i = X_i\beta + Z_i b_i + e_i,$$

که در آن $i = 1, \dots, N$. برای یک i منحصر به فرد با n_i تکرار، Y_i یک بردار $(n_i \times 1)$ و X_i یک ماتریس با بعد $(n_i \times p)$ ، به عنوان متغیرهای کمکی ثابت و β بردار پارامتر $(p \times 1)$ از ضرایب رگرسیون می باشد که معمولاً در این مدل ها به عنوان بردار اثرهای ثابت از آن یاد می شود. Z_i هم یک ماتریس با بعد $(n_i \times v)$ از متغیرهای کمکی برای بردار b_i است که b_i نیز یک بردار $(v \times 1)$ از اثرهای تصادفی می باشد. e_i هم بردار $(n_i \times 1)$ خطاها در نظر گرفته می شود. در بررسی و تحلیل این مدل ها، مرسوم است که فرض کنیم b_i و e_i مستقل هستند و هر دو دارای توزیع نرمال می باشند، به طوری که

$$e_i \sim N_{n_i}(0, \sigma^2 I_{n_i}), \quad \sim N_v(0, V)$$

که در آن I_{n_i} ماتریس همانی با بعد $(n_i \times n_i)$ می باشد. تحت این فرض ها داریم:

$$Y_i | \beta, b_i \sim N_{n_i}(X_i\beta + Z_i b_i, \sigma^2 I_{n_i}). \quad (16)$$

مدل (۱۶) را مدل اثرهای تصادفی نرمال می نامیم. از آنجایی که در این بخش هدف کاربرد فرآیند دیریکله در این مدل هاست، پس فرض نرمال بودن

b_i را در نظر نگرفته و در عوض فرض می کنیم $b_i \sim G$ ، که در آن توزیع G از فرآیند دیریکله با اندازه پایه G_0 پیروی می کند. اگر فرض شود که توزیع نمونه ای Y_i ها نرمال است، آنگاه یک اندازه پایه نرمال برای اثرهای تصادفی، یک مدل DPM مزدوج را مشخص می کنند. حال فرض کنید که توزیع بردار حاصل Y_i برای i امین واحد به صورت زیر باشد:

$$Y_i | \beta, b_i \sim N_{n_i}(X_i\beta + Z_i b_i, \sigma^2 I_{n_i}).$$

با فرض $\tau = \sigma^{-2}$ ، توزیع پیشین برای τ را به صورت زیر در نظر می گیریم:

$$\tau \sim \Gamma\left(\frac{\alpha_0}{\nu}, \frac{\lambda_0}{\nu}\right),$$

که در آن Γ نشانگر توزیع گاما می باشد. بنابراین مدل فرآیند دیریکله برای رگرسیون خطی با اثرهای تصادفی خطی نرمال به صورت

$$\beta \sim N_p(\mu_0, \Sigma_0),$$

$$b_i \sim G, \quad (17)$$

$$G \sim DP(M, N_v(0, V)),$$

معرفی می گردد، که در آن μ_0 ، Σ_0 ، M و V پارامترهای معلوم فرض می شوند.

وقتی که G یک پیشین کاملاً پارامتری باشد، آنگاه بعد از یافتن توزیع پسین کامل پارامترهای مدل، با استفاده از قانون بیز و انجام چند محاسبه جبری و نیز حذف عباراتی که به β بستگی ندارند،

پارامترهای موجود در مدل و از جمله سایر b_i ها، یعنی $b^{(-i)}$ هستیم. وست و همکاران^{۱۶} [۱۰] توزیع پسین هر کدام از b_i ها را براساس (۷) به صورت زیر به دست آوردند:

$$\begin{aligned} & \pi(b_i | \beta, \tau, Y, b^{(-i)}) \\ & \propto \sum_{i \neq j} \phi_{n_i}(Y_i | X_i \beta + Z_i b_j, \tau^{-1} I_{n_i}) \delta_{b_j} \\ & + \{M \int \phi_{n_i}(Y_i | X_i \beta + Z_i b_i^*, \tau^{-1} I_{n_i}) \\ & \quad \times \phi_v(b_i^* | \circ, V) db_i^*\} \\ & \times \phi_{n_i}(Y_i | X_i \beta + Z_i b_i, \tau^{-1} I_{n_i}) \phi_v(b_i | \circ, V), \end{aligned}$$

که در آن $b^{(-i)}$ به بردار اثرهای تصادفی به استثنای جز i اشاره دارد. همچنین چگالی نرمال چندمتغیره با $\phi_{n_i}(Y_i | X_i \beta + Z_i b_i, \tau^{-1} I_{n_i})$ و نیز با $\phi_v(b_i | \circ, V)$ نشان داده شده اند. پس از انجام چند محاسبه جبری داریم:

$$\begin{aligned} & \pi(b_i | \beta, \tau, Y, b^{(-i)}) \\ & \propto \left(\sum_{i \neq j} \tau^{\frac{n_i}{\nu}} \exp \left[-\frac{\tau}{\nu} (Y_i - X_i \beta + Z_i b_j)' \right. \right. \\ & \quad \left. \left. (Y_i - X_i \beta + Z_i b_j) \right] \delta_{b_j} \right) \\ & + M |Q_i|^{\frac{1}{\nu}} |V|^{\frac{1}{\nu}} \tau^{\frac{n_i}{\nu}} \exp \left\{ \frac{\tau}{\nu} [(Y_i - X_i \beta)' \right. \\ & \quad \left. U_i (Y_i - X_i \beta)] \right\} \\ & \times \phi_{n_i}(Y_i | X_i \beta + Z_i b_i, \tau^{-1} I_{n_i}) \phi_v(b_i | \circ, V), \end{aligned}$$

که در آن

$$Q_i = (V^{-1} + \tau Z_i' Z_i)^{-1}$$

توزیع پسین شرطی کامل β به صورت زیر به دست می آید:

$$\beta | b, \tau, Y \sim N_{n_i}(\hat{\beta}, T), \quad (18)$$

که در آن

$$b = (b_1, \dots, b_N)$$

$$Y = (Y_1, \dots, Y_N)$$

$$T = (\tau \sum_{i=1}^N X_i' X_i + \Sigma_0^{-1})^{-1}$$

$$\hat{\beta} = T \sum_{i=1}^N X_i' (Y_i - Z_i b_i) + \Sigma_0^{-1} \mu_0.$$

همچنین توزیع پسین شرطی کامل τ به صورت زیر محاسبه می شود:

$$\tau | b, \beta, Y \sim \Gamma\left(\frac{n + \alpha_0}{\nu}, \frac{\sum_{i=1}^N r_i' r_i + \lambda_0}{\nu}\right) \quad (19)$$

که در آن

$$n = \sum_{i=1}^N n_i \quad \text{و} \quad r_i = Y_i - X_i \beta - Z_i b_i.$$

برای مطالعه جزئیات بیشتر می توانید به بوش و مک ایچرن^{۱۳} [۳]، کلینمن و ابراهیم^{۱۴} [۶] و کاتر و همکاران^{۱۵} [۵] مراجعه کنید.

۱.۵ توزیع پسین شرطی کامل اثرات تصادفی در رگرسیون خطی چندگانه با اثرات آمیخته

حال به منظور تحلیل روی هر کدام از b_i ها به دنبال توزیع پسین شرطی کامل b_i به شرط تمام

Bush and McEachern^{۱۳}
Kleinman and Ibrahim^{۱۴}
Kutner and et al.^{۱۵}
West, Muller and Escobar^{۱۶}

$$\propto |V^{-1}|^{(d_0+k-v-1)/2} \quad \text{و} \quad (20)$$

$$\times \exp \left\{ -\frac{1}{\varphi} \text{tr}((c_0 R_0)^{-1} V^{-1} U_i = (\tau Z_i Q_i Z_i' - I_n) \right.$$

$$\left. -\frac{1}{\varphi} \sum_{l=1}^k \gamma_l V^{-1} \gamma_l \right\}.$$

بنابراین داریم

$$b_i | \beta, \tau, Y \sim N_v(\tau Q_i Z_i' (Y_i - X_i \beta), Q_i).$$

بنابراین

$$(V^{-1} | b, \beta, \tau, Y) \sim W_v(d_0 + k, (\frac{1}{c_0 R_0} + \sum_{l=1}^k \gamma_l \gamma_l')^{-1}). \quad (21)$$

حال در این قسمت به منظور بررسی بیشتر بر روی مدل‌های رگرسیون با اثرات تصادفی و چگونگی بکارگیری فرآیند دیریکله در این مدل‌ها و تحلیل تأثیر این فرآیند در مدل به بیان مثال می‌پردازیم.

این نتایج به همانند استخراج نمونه از یک توزیع آمیخته می‌باشد که یک بخش از آن، توزیع نرمال و قسمت‌های دیگرش جرم احتمال هستند.

اغلب ماتریس کواریانس V به عنوان اندازه پایه در مدل احتمال فرآیند دیریکله (۱۶) نامعلوم فرض می‌شود. بنابراین باید یک توزیع پیشین مناسب برای آن تعیین شود تا بتوانیم به برآورد آن بپردازیم. یک فرض آن است که

$$V^{-1} \sim W_v(d_0, c_0 R_0),$$

که در آن $W_v(d_0, c_0 R_0)$ نشانگر توزیع ویشارت می‌باشد که $d_0 > v$ ، $c_0 > 0$ و R_0 ماتریس معین مثبت با بعد $(v \times v)$ می‌باشد. لذا یک توزیع پیشین را به صورت زیر داریم

$$\pi(V^{-1} | d_0, c_0 R_0) \propto |V^{-1}|^{(d_0-v-1)/2} \exp \left\{ -\frac{1}{\varphi} \text{tr}((c_0 R_0)^{-1} V^{-1}) \right\}.$$

توزیع پسین شرطی کامل V^{-1} به صورت زیر بدست می‌آید:

$$\pi(V^{-1} | b, \beta, \tau, Y) = \pi(V^{-1} | \gamma, \beta, \tau, b)$$

مثال ۱ در این مثال یک مدل اثرهای تصادفی نرمال را مورد بررسی قرار می‌دهیم. تحقیق مورد بحث در یکی از ایالات امریکا به انجام رسیده است. نمرات پیشرفت ریاضی ۱۷ دانش آموزان ۱۶۰ دبیرستان را که متغیر مستقل مدل می‌باشد، مورد تحقیق قرار می‌دهیم و این متغیر را متأثر از چهار متغیر مستقل جنسیت^{۱۸}، منطقه سکونت^{۱۹}، اقلیت مذهبی^{۲۰} و وضعیت اقتصادی^{۲۱} دانش آموزان فرض می‌کنیم. با فرض توزیع نرمال برای نمرات ریاضی دانش آموزان به صورت $mathach_{ij} \sim N(\mu_i, \sigma_i^2)$ برای $i = 1, \dots, 160$

Math Achievement^{۱۷}
Gender^{۱۸}
Sector^{۱۹}
Religion Minority^{۲۰}
Socio-economic^{۲۱}

و $n_i = 1, \dots, n_i$ یک مدل رگرسیونی با اثرات آمیخته را به صورت زیر برآزش می‌دهیم:

$$mathach_{ij} = \alpha_i + \beta_1 Gender + \beta_2 Sector + \beta_3 Religion + \beta_4 Socio + e_{ij}$$

جز پارامترهای نامعلوم مدل محسوب می‌شوند و با فرض $\tau_e = \sigma_e^{-2}$ و $\tau_\alpha = \sigma_\alpha^{-2}$ توزیع‌های پیشین یکسان

$$\tau_\alpha \sim \Gamma(0.1, 0.1) \quad \text{و} \quad \tau_e \sim \Gamma(0.1, 0.1)$$

که در آن α_i اثر تصادفی برای دبیرستان i ام و n_i تعداد دانش آموزان در مدرسه i ام می‌باشد و فرض بر این است که برای $i = 1, \dots, 160$ داریم:

$$\alpha_i \sim N(0, \sigma_\alpha^2)$$

را برای آنها در نظر می‌گیریم که در آن نشان دهنده توزیع گاما می‌باشد. توسط شبیه‌سازی کامپیوتری توسط نرم افزار *WinBUGS* به برآورد پارامترهای مجهول مدل می‌پردازیم. با تولید ۶۰۰۰ نمونه و دورریز ۱۰۰۰ و بعد از اینکه توزیع پارامترهای مجهول مدل به همگرایی مطلوب رسیدند، برآوردهای پارامترهای مجهول مدل را طبق جدول (۱) داریم.

لازم به ذکر است که برای $j = 1, 2, 3, 4$ ، توزیع پیشین نرمال با پارامترهای معلوم برای β_j ها فرض شده است و مؤلفه‌های واریانس مدل یعنی σ_α^2 و

جدول (۱). برآورد پارامترهای مجهول در مثال ۱

	مقدار برآورد	انحراف استاندارد	کران پایین فاصله طمینان	کران بالای فاصله اطمینان
β_1	-۲/۸۹۱	۰/۲۲۰۱	-۳/۳۲۹	-۲/۴۷۹
β_2	-۱/۱۲۸	۰/۱۶۵۸	-۱/۴۳۳	-۰/۷۹۷۳
β_3	۱/۹۱۳	۰/۱۱۰۶	۱/۶۹۳	۲/۱۲۴
β_4	-۱۲۲/۰	۳/۰۷۷	-۱۲۶/۷	-۱۱۵/۶
σ_α^2	۸۵۳۱/۰	۱۰۲۷/۰	۶۷۱۹/۰	۱۰۷۱۰/۰
σ_e^2	۳۵/۹۳	۰/۶۱۳۷	۳۴/۷۴	۳۷/۱۵

جدول (۲). برآورد اثرهای تصادفی مدارس در مثال ۱

	مقدار برآورد	انحراف استاندارد	کران پایین فاصله طمینان	کران بالای فاصله اطمینان
α_1	۱۰/۶۵	۰/۸۸۰۴	۸/۹۴۱	۱۲/۳۷
α_2	۱۴/۳۶	۱/۲۱	۱۱/۹۸	۱۶/۷۷
⋮	⋮	⋮	⋮	⋮
α_{38}	۱۱۵/۳	۲۰/۲۳	۷۴/۱۶	۱۴۵/۵
⋮	⋮	⋮	⋮	⋮
α_{142}	۵/۴۱۹	۱/۰۶۱	۳/۳۳۳	۱۲/۳۷
⋮	⋮	⋮	⋮	⋮
α_{159}	۱۲/۴۲	۱/۱۲۵	۱۰/۲۱	۱۴/۶۳
α_{160}	۱۱۱/۱	۲۰/۲۲	۷۰/۱	۱۴۱/۲

برای اثرهای تصادفی در مدل مورد استفاده قرار می‌دهیم، معیار AIC را به کار می‌بریم. از این رو با توجه به رابطه

$$AIC = Deviance + 2p_{\psi}$$

بین معیار انحراف و AIC وجود دارد و پراکندگی به صورت جدول (۳) برآورد شده است، بنابراین AIC برابر با ۴۶۳۳۲ به دست می‌آید. لازم به ذکر است که p_{ψ} تعداد پارامترهای مجهول مدل می‌باشد.

تعدادی از اثرات تصادفی برآورد شده مربوط به مدارس نیز در جدول (۲) قابل مشاهده است. با توجه به برآورد ۱۶۰ اثر تصادفی مدارس، بیشترین اثر تصادفی برابر ۱۱۵/۳ و مربوط به مدرسه ۳۸ام و کمترین اثر تصادفی نیز برابر ۵/۴۱۹ و مربوط به مدرسه ۱۴۳ام می‌باشد. سایر اثرهای تصادفی بین این دو حد ماکزیمم و مینیمم می‌باشند. حال برای مقایسه این مدل با مدل‌های دیگر از جمله زمانی که فرآیند دیریکله را به عنوان پیشین

جدول (۳). برآورد $Deviance$ در مثال ۱

مقدار برآورد	انحراف استاندارد	کران پایین فاصله طمینان	کران بالای فاصله اطمینان
$463320/0$	$16/84$	$46302/0$	$46340/0$
$Deviance$			

می‌گیریم. بنابراین نشانگرهای خوشه‌ای Z_i را برای خوشه i ام به صورت

$$Z_i \sim \text{Categorical}(p)$$

تعریف می‌کنیم که در آن $p = (p_1, \dots, p_k)$ به عنوان بردار ضرایب آمیختگی می‌باشد که برای $j = 1, \dots, k$ داریم $0 \leq p_j \leq 1$ و $\sum_{j=1}^k p_j = 1$. با فرض $kT\infty$ با توجه به مفهوم فرآیند دیریکله می‌توان فرض کرد که α_j ها از جامعه آمیخته

$$q(\cdot) = \sum_{j=1}^{\infty} p_j h(\cdot | \alpha_j)$$

استخراج می‌شوند که در آن $h(\cdot | \alpha_j)$ تابع چگالی احتمال خوشه i ام می‌باشد که α_j از آن استخراج

مثال ۲ در ادامه بررسی مثال قبل، در این قسمت نحوه برازش مدل به روش بیز نیمه پارامتری و چگونگی استفاده از این روش را مورد تحلیل قرار می‌دهیم. بدین منظور فرآیند دیریکله را به عنوان پیشین برای اثرات آمیخته استفاده می‌کنیم. از آنجایی که جامعه مورد بررسی شامل نمرات ریاضی دانش آموزان ۱۶۰ دبیرستان هستند که این دبیرستان‌ها خوشه‌های مورد بررسی در طرح فوق می‌باشند، بنابراین یک جامعه آمیخته را برای مطالعه پیش رو داریم لذا برای $\alpha_i, i = 1, \dots, 160$ ها را به عنوان اثرات تصادفی مدل و جزو پارامترهای مجهول در نظر می‌گیریم که بایستی برآورد گردند. یک توزیع پیشین دیریکله را با پارامترهای معلوم برای α_i ها در نظر

می‌گردد. لازم به ذکر است که اندازه پایه مربوط به فرآیند دیریکله را نیز یک توزیع نرمال با میانگین صفر و واریانس نامعلوم و پارامتر دقت مرتبط با آن را یک اسکالر معلوم فرض می‌کنیم. اثرات تصادفی مربوط به ۱۶۰ مدرسه به صورت بالقوه در ۲ خوشه، دسته بندی می‌شوند. بنابراین تعداد خوشه‌ها را ۲ در نظر می‌گیریم. بعد از اجرای برنامه کامپیوتری پارامترهای مجهول به صورت جدول (۴) برآورد می‌گردند.

پراکندگی و AIC به صورت

$$AIC = Deviance + 2p_{\psi}$$

مقدار AIC برابر ۴۶۱۲۲ به دست می‌آید که با مقایسه این معیار با مدل مثال یک که برابر ۴۶۳۳۲ به دست آمده بود، به این نتیجه می‌رسیم که مدل دوم با در نظر گرفتن فرآیند دیریکله به عنوان پیشین برای اثرات تصادفی نسبت به مدل اول برتری دارد. برای مطالعه جزئیات بیشتر در مورد معیارهای اعتبار مدل‌ها می‌توانید به اسپیکل هالتر و همکاران [۹] مراجعه کنید.

می‌توانیم با توجه به حدسی که در مورد تعداد خوشه‌ها زدیم، اثرات تصادفی را، در ۲ خوشه دسته‌بندی کنیم که خوشه اول اثرات تصادفی را بین صفر و ۲۰ و خوشه دوم را نیز بین ۱۰۰ و ۲۰۰ در نظر می‌گیریم. بنابراین p_1 و p_2 به ترتیب احتمال متعلق بودن اثرات تصادفی به خوشه اول و دوم را نشان می‌دهند که به صورت جدول (۵) برآورد می‌گردند. تعدادی از اثرات تصادفی مربوط به مدارس نیز به صورت جدول (۶) برآورد شده‌اند. با توجه به جدول (۶) همانند مثال یک

جدول (۴). برآورد پارامترهای مجهول در مثال ۲

	مقدار برآورد	انحراف استاندارد	کران پایین فاصله طمینان	کران بالای فاصله اطمینان
β_1	-۳/۳۷۴	۰/۲۱۶۱	-۳/۸۰۶	-۲/۹۴۸ - ۲/۴۷۹
β_2	-۱/۴۵۹	۰/۱۵۸۳	-۱/۴۶	-۱/۱۴۵
β_3	۱/۸۸۴	۰/۱۱۱۳	۱/۸۸۳	۲/۰۹۹
β_4	۳/۱۱	۰/۰۰۱۸	۳/۱۱۲	۳/۵۸۳
σ_e^2	۳۶/۹۲	۰/۰۱۱۴	۳۶/۹۱	۳۸/۱۷

جدول (۵). برآورد احتمالات متعلق بودن اثرات تصادفی به خوشه ۱ و ۲ در مثال ۲

	مقدار برآورد	انحراف استاندارد	کران پایین فاصله طمینان	کران بالای فاصله اطمینان
p_1	۰/۴۹۳۸	۰/۰۵۲۹۱	۰/۳۹۰۶	۰/۵۹۶۲
p_2	۰/۵۰۶۲	۰/۰۵۳۴۲	۰/۴۰۳۰	۰/۶۰۹۴

جدول (۶). برآورد اثرات تصادفی مدارس در مثال ۲

	مقدار برآورد	انحراف استاندارد	کران پایین فاصله طمینان	کران بالای فاصله اطمینان
α_1	۱۱/۳۵	۰/۲۹۸۱	۱۰/۷۲	۱۱/۹
α_2	۱۴/۴۵	۰/۷۵۰۱	۱۱/۴۹	۱۵/۱۵
⋮	⋮	⋮	⋮	⋮
α_{38}	۱۴/۶۲	۰/۲۷۴	۱۴/۱۳	۱۵/۱۶
⋮	⋮	⋮	⋮	⋮
α_{142}	۱۱/۳۵	۰/۲۹۸۸	۱۰/۷۲	۱۱/۹
⋮	⋮	⋮	⋮	⋮
α_{159}	۱۲/۸۲	۱/۵۸۳	۱۰/۹۴	۱۴/۹۸
α_{160}	۱۴/۰۳	۱/۲۷۱	۱۱/۰۶	۱۵/۱۴

جدول (۷). برآورد $Deviance$ در مثال ۲

	مقدار برآورد	انحراف استاندارد	کران پایین فاصله طمینان	کران بالای فاصله اطمینان
$Deviance$	۴۶۱۱۰/۰	۱۸/۴	۴۶۰۸۰/۰	

برای سنجش و مقایسه دو مدل فوق علاوه بر معیار

$$f(Y_{it}|\theta_{it}, \tau) = \exp\{\tau(Y_{it}\theta_{it} - q(\theta_{it})) + c(Y_{it}, \tau)\}$$

AIC می‌توانیم تک تک اثرات تصادفی را به همراه

انحراف استاندارد برآورد شده برای هر کدام را نیز

که در آن

با هم مقایسه کنیم.

$$\begin{aligned} E(Y_{it}|\theta_{it}, \tau) &= \mu_{it} \\ &= \frac{dq(\theta_{it})}{d\theta_{it}} \end{aligned}$$

۶ MDP در مدل رگرسیون

خطی تعمیم یافته با اثرهای

آمیخته

و τ یک اسکالر به عنوان پارامتر پراکنندگی می

باشد. در مدل آمیخته خطی تعمیم یافته فوق،

پارامتر کانونی θ_{it} توسط تابع پیوند $h(\cdot)$ با

فرض کنید برای $i = 1, \dots, N$ و $t = 1, \dots, n_i$

توزیع متغیرهای تصادفی Y_{ij} متعلق به خانواده

برازش مدل می‌پردازیم. لذا برای $i = 1, \dots, N$ فرض می‌کنیم که b_i ها دارای یک توزیع آزاد مثل G باشند. به عبارت دیگر G یک توزیع عمومی بوده و فرض می‌شود که توزیع b_i ها نرمال است. بنابراین یک مدل خطی تعمیم یافته با اثرات تصادفی خواهیم داشت که توزیع آن اثرات از فرآیند دیریکله پیروی می‌کند. برای روشن شدن بیشتر موضوع فرض کنید که توزیع Y_{it} برای جز i ام در زمان t ، بر طبق (۲۲) برابر $f(Y_{it}|\beta, b_i, \tau)$ باشد. در ادامه، توزیع‌های پیشین برای پارامترهای مدل آمیخته فرآیند دیریکله در رگرسیون خطی تعمیم یافته را به صورت زیر در نظر می‌گیریم:

$$\tau \sim \Gamma(\alpha_0, \lambda_0)$$

و

$$\beta \sim N_p(\mu_0, \Sigma_0)$$

و با فرض توزیع G برای b_i ها، داریم:

$$G \sim DP(M.N_v(0, V)) \quad (24)$$

که در آن برای ماتریس V^{-1} توزیع پیشین ویشارت با پارامترهای معلوم را فرض می‌کنیم. بنابراین با توجه به پیشین‌های فوق به منظور برآورد پارامترهای مجهول مدل توزیع پسین توأم پارامترها را به صورت زیر به دست می‌آوریم:

$$\begin{aligned} \pi(\beta, b, \tau) &\propto \exp \text{Ln} f(Y|\beta, b, \tau) \pi(\beta, b, \tau) \\ &\propto \exp \left\{ \sum_{i=1}^N \sum_{t=1}^{n_i} \text{Ln} f(Y_{it}|\beta, b_i, \tau) \right. \\ &\quad \left. - \frac{1}{\rho} (\beta - \mu_0)' \Sigma_0^{-1} (\beta - \mu_0) \right\} \quad (25) \end{aligned}$$

متغیرهای توضیحی X و متغیرهای کمکی Z به صورت $\theta_{it} = h(\eta_{it})$ مرتبط می‌شود که در آن پیش بینی کننده خطی به صورت $\eta_{it} = x'_{it} + z'_{it}b_i$ است که در آن x'_{it} و z'_{it} به ترتب سطرهای ماتریس X_i و Z_i می‌باشند و تابع $h(\cdot)$ را یک تابع مشتق پذیر یکنوا فرض می‌کنیم. لازم به ذکر است که η_{it} به پیش بینی کننده خطی معروف است. بنابراین داریم:

$$f(Y_{it}|\theta_{it}, \tau) \equiv f(Y_{it}|\beta, b_i, \tau)$$

که در آن

$$f(Y_{it}|\beta, b_i, \tau) = \quad (22)$$

$$\exp\{\tau(Y_{it}h(\eta_{it}) - q(h(\eta_{it}))) + c(Y_{it}, \tau)\}$$

می‌دانیم که اگر $\theta_{it} = h(\eta_{it}) = \eta_{it}$ آنگاه $h(\cdot)$ تابع پیوند کانونی گویند.

توجه کنید که در مدل‌های خطی تعمیم یافته، فرض بر آن است که اثرهای تصادفی معمولاً از توزیع نرمال پیروی می‌کنند و مشاهدات در واحد i ام از هم مستقل می‌باشند. بنابراین تابع درستنمایی به صورت زیر است:

$$f(Y|\beta, b, \tau) \propto \prod_{1 \leq i \leq N} \left(\prod_{1 \leq t \leq n_i} f(Y_{it}|\beta, b_i, \tau) \right) \quad (23)$$

که در آن

$$b = (b_1, \dots, b_N)', Y = (Y_{11}, \dots, Y_{Nn_N})$$

حال به منظور برآورد اثرات تصادفی به روش نیمه پارامتری بیزی، با استفاده از فرآیند دیریکله به عنوان توزیع پیشین برای اثرات تصادفی، به

بیشتر در زمینه کاربرد MDP در رگرسیون خطی تعمیم یافته با اثرهای آمیخته می‌توانید به اولسن و همکاران ۲۲ [۸] مراجعه کنید.

۷ نتیجه‌گیری

در این مقاله بعد از آشنایی با اهمیت کاربرد فرآیند دیریکله به عنوان پیشین برای اثرات تصادفی در مدل‌های رگرسیون بیزی با اثرات آمیخته به تعمیم این موضوع پرداختیم و کاربرد مدل‌های دیریکله آمیخته (DPM) را برای مدل‌های رگرسیون خطی تعمیم یافته با اثرات آمیخته در حالت بیزی مورد بررسی قرار دادیم و در ادامه با ارائه مثالی کاربردی و شبیه‌سازی کامپیوتری با نرم افزار $WinBUGS$ نحوه برازش (DPM) برای این مدل‌ها را تشریح کردیم.

در نهایت به مقایسه مدل‌های برازش داده شده پرداختیم. نتیجه ای که می‌توان بی شک به آن اشاره کرد، بهبود برآوردگرهایی می‌باشد که برای پارامترها و اثرات آمیخته یافتیم و دلیل این مدعا معیارهای سنجش اعتبار مدل از قبیل AIC و $Deviance$ می‌باشند.

$$\times \left\{ -\tau\lambda_0 - \frac{1}{\nu} \sum_{i=1}^N b_i' V^{-1} b_i \right\} \times \tau^{\alpha_0 - 1}$$

که در آن $\pi(\beta, b, \tau)$ نشان‌دهنده چگالی پیشین توام β و b و τ می‌باشد. حال به منظور استنباط برای این پارامترها، ابتدا چگالی پسین شرطی کامل آنها را به دست می‌آوریم. از این رو در عبارت (۲۵) با حذف عباراتی که به β بستگی ندارند، چگالی پسین حاشیه‌ای شرطی کامل β را به صورت

$$\pi(\beta|Y, b, \tau) \propto \exp \left\{ \sum_{i=1}^N \sum_{t=1}^{n_i} \ln f(Y_{it}|\beta, b_i, \tau) - \frac{1}{\nu} (\beta - \mu_0)' \Sigma_0^{-1} (\beta - \mu_0) \right\}$$

به دست می‌آوریم و از آنجایی که این توزیع شکل صریحی ندارد، با کمک روش‌های $MCMC$ و به ویژه نمونه‌گیر گیبز می‌توان به استخراج نمونه از این توزیع پرداخت. توزیع پسین حاشیه‌ای شرطی کامل نیز به صورت

$$\begin{aligned} \pi(\tau|\beta, b, Y) & \propto \exp \left\{ \sum_{i=1}^N \sum_{t=1}^{n_i} c(Y_{it}, \tau) \right\} \\ & \times \exp \left\{ \sum_{i=1}^N \sum_{t=1}^{n_i} (Y_{it} \theta_{it} - q(\theta_{it})) - \tau\lambda_0 \right\} \end{aligned}$$

می‌باشد که با نمونه‌گیری گیبز از این توزیع می‌توان به برآورد τ پرداخت. برای مطالعه مثال‌های

مراجع

- [1] Antoniak, C.E. (1974), Mixture of Dirichlet processes with applications to non-parametric problems, *Annals of Statistics*, 2, 1152-1174.

- [2] Blackwell, D. and McQueen, J.B. (1973), Ferguson distribution via Polya Urn schemes, *Annals of Statistics*, 1, 353-355.
- [3] Bush, C. A. and McEachern, S. N. (1996), A semi-parametric Bayesian model for randomized block designs, *Biometrika*, 83, 275-286 .
- [4] Ferguson, T.S. (1973), A Bayesian analysis of some nonparametric problems, *Annals of Statistics* 1, 209-230.
- [5] Ishwaran, H. and James, L.F. (2001), Gibbs sampling methods for stick-breaking priors, *Journal of the American Statistical Association*, 96, 161-173.
- [6] Kleinman, K.P. and Ibrahim, J.G. (1998), A semiparametric Bayesian approach to the random effects model, *Biometrics*, 54, 921-938.
- [7] Kutner, M.H., Nachtsheim, C.J., Neter, J. and Li, W. (2005), *Applied Linear Statistical Models*, 5th Ed, New York: McGraw-Hill/Irwin.
- [8] Ohlssen, D.I., Sharples, L.D. and Spiegelhalter, D.J. (2006), Flexible random-effects models using Bayesian semiparametric models, *Statistics in Medicine* 26(9), 2088-2112.
- [9] Spiegelhalter, D.J., Best, N.G., Carlin, B.P. and Linde, A.V. (2002), Bayesian measures of model comolexity and fit, *Journal of the Royal Statistical Society B*. 64(4), 583-639
- [10] West, M., Muller, P. and Escobar, M.D. (1994), *Hierarical Priors and Mixture Models with Application in Regression and Density Estination*, Aspects of Uncertainly, eds. P.R. Freeman and A. FM. Smith, New York: Willey, 363-386.