

تحلیل دو سطحی با اثرات تصادفی چوله نرمال و مدل‌بندی داده‌های طولی

محمد غلامی فشارکی^۱، انوشیروان کاظم‌نژاد^۱، فرید زایری^۲

^۱ گروه آمار زیستی، دانشگاه تربیت مدرس

^۲ گروه آمار زیستی، دانشگاه علوم پزشکی شهید بهشتی

تاریخ دریافت: ۱۳۹۲/۲/۱۷ تاریخ آخرین بازنگری: ۱۳۹۲/۱۰/۲۶

چکیده: مدل‌سازی داده‌های دوسطحی با فرض نرمال بودن مولفه تصادفی و خطا انجام می‌شود. عدم برقراری این فرض باعث استنباط غلط در مورد پارامترهای مدل می‌گردد. در این مقاله، استفاده از خانواده توزیع چوله نرمال که خانواده‌ای انعطاف‌پذیرتر از توزیع نرمال است مطرح می‌شود. سپس در یک مطالعه شبیه‌سازی نشان داده می‌شود عدم در نظر گرفتن چولگی مثبت (منفی) در مدل باعث بیش برآوردی (کم برآوردی) عرض از مبدا و کم برآوردی (بیش برآوردی) شیب مدل می‌گردد سپس با مدل به‌دست آمده رابطه نوبت کاری و کلسترول خون تعیین می‌شود.

واژه‌های کلیدی: تحلیل دوسطحی، توزیع چوله نرمال، رهیافت بیزی.

آدرس الکترونیک مسئول مقاله: انوشیروان کاظم‌نژاد، kazem_an@modares.ac.ir
کد موضوع‌بندی ریاضی (۲۰۱۰): ۶۰E۰۵, ۶۲F۱۵

بسیاری از داده‌های حوزه علوم بهداشتی دارای ساختار خوشه‌ای یا سلسله مراتبی هستند (گلدستین، ۲۰۰۸). از آنجا که استفاده از روش‌های متداول مدل‌بندی رگرسیونی در تحلیل داده‌های خوشه‌ای به دلیل عدم در نظر گرفتن همبستگی بین مشاهدات منجر به کم برآوردی خطای برآورد ضرایب رگرسیونی و معنی‌داری نادرست ضرایب و به تناسب آن افزایش خطای نوع اول می‌گردد (پینهایرو و بتس، ۲۰۰۰) از این رو محققان برای تحلیل داده‌های سلسله مراتبی از روش تحلیل چندسطحی^۱ استفاده می‌نمایند (گلدستین، ۲۰۰۸؛ لوک، ۲۰۰۴).

تحلیل چندسطحی بسط مدل‌های خطی تعمیم‌یافته است که در آن علاوه بر مدل‌بندی متغیر پاسخ ضرایب رگرسیونی نیز مدل‌بندی می‌گردد. این روش کارا در مدل‌سازی داده‌هایی با ساختار آشیانه‌ای بوده و هدف آن مدل‌بندی متغیر وابسته براساس تابعی از متغیرهای پیشگو در بیشتر از یک سطح است (هاکس، ۲۰۱۰). معمولاً برای سادگی فرض می‌شود مولفه‌های تصادفی و خطای مدل دارای توزیع نرمال هستند. از آنجا که عدم برقراری این فرض منجر به برآورد نادرست پارامترهای مدل و استنباط نامعتبر در مورد پارامترهای مدل می‌شود، محققان به روش‌هایی چون تبدیل داده‌ها (به‌ویژه استفاده از تبدیل باکس-کاکس^۲)، استفاده از روش‌های ناپارامتری و توزیع‌های غیر نرمال رو می‌آورند. استفاده از روش‌های ناپارامتری علی‌رغم کاهش اطلاعات از لحاظ عملی دارای محدودیت‌های زیادی برای محققان بوده و از طرف دیگر استفاده از روش تبدیل متغیرها نیز باعث مشکلاتی چون تفسیر نامناسب پارامترهای تحت بررسی، عدم احراز فرض نرمال توام متغیرها در حالات چند متغیره (تی سانگ و لی، ۲۰۰۳؛ لی و همکاران، ۲۰۰۵؛ جارا و همکاران، ۲۰۰۸)، پنهان شدن فرآیند تولید داده‌ها و کاهش اطلاعات (جارا و همکاران، ۲۰۰۸)، یکتا نبودن تبدیل در مطالعات مختلف می‌گردد (بندیویدهای و همکاران، ۲۰۱۰). همچنین استفاده از توزیع‌های غیر نرمال، محقق را از دسترسی به خصوصیات مناسب توزیع‌های بیضوی مانند نرمال دور می‌سازد. راهکار دیگر در

^۱ Multilevel Modeling

^۲ Box-Cox

حل این معضل استفاده از توزیع‌هایی است که به لحاظ نظری قادر به تبیین تغییرات مشاهده بدون استفاده از فرض نرمال و با داشتن خصوصیات توزیع‌های بیضوی می‌باشد که از آن جمله می‌توان به مطالعات پینهایرو و همکاران (۲۰۰۱)، ژو و همی (۲۰۰۸)، روسا و همکاران (۲۰۰۳)، لین و لی (۲۰۰۶، ۲۰۰۷ و ۲۰۰۸)، لانگ و سینشیمیر (۱۹۹۳)، ما و همکاران (۲۰۰۴)، لاجوس و همکاران (۲۰۰۹)، جارا و همکاران (۲۰۰۸) و بندیوپدهایی و همکاران (۲۰۱۰) در برآزش تابع درست‌نمایی و رهیافت بیزی و همچنین مقایسه روش‌های انتخاب مدل اشاره نمود. اما علی‌رغم مطالعات متعدد، وجود مطالعه‌ای که در آن به مشاهده تاثیر مولفه تصادفی و خطای غیرنرمال بر پارامترهای یک مدل دوسطحی می‌پردازد احساس می‌گردد. از این رو این موضوع با انجام یک مطالعه شبیه‌سازی بررسی و نهایتاً این مدل آماری در بررسی طولی رابطه نوبت کاری با کلسترول خون به‌کار گرفته شده است.

۲ مدل آماری و تحلیل بیزی

۱.۲ تابع چگالی توام مدل دو سطحی

فرض کنید محقق می‌خواهد تاثیر کار در نوبت کاری را بر میزان کلسترول خون، وقتی که مشاهدات کارگران بصورت طولی جمع‌آوری شده است را اندازه‌گیری نماید. از آنجا که کلسترول خون افراد در زمان‌های مختلفی اندازه‌گیری شده‌اند، داده‌ها مطالعه طولی و به یکدیگر وابسته هستند. در اینجا مجموعه اندازه‌های گرفته شده در زمان‌های مختلف برای هر یک از اعضای نمونه، به‌عنوان واحدهای سطح اول و هر یک از افراد انتخاب شده به‌عنوان سطح دوم در نظر گرفته می‌شود. این ساختار در مدل‌های چندسطحی به‌صورت

$$\begin{aligned} y_{i(j)} &= \beta_{0(j)} + \beta_1 x_{i(j)} + \epsilon_{i(j)} \\ \beta_{0(j)} &= \beta_0 + u_{(j)} \end{aligned} \quad (1)$$

به ترتیب در سطح اول و دوم قابل ارائه است، که با ادغام آن‌ها به‌صورت

$$y_{i(j)} = \beta_0 + \beta_1 x_{i(j)} + u_{(j)} + \epsilon_{i(j)} \quad i = 1, \dots, M_j, \quad j = 1, \dots, M \quad (2)$$

۲۳۶ تحلیل دوسطحی با مولفه تصادفی چوله نرمال

تبدیل می شود. فرض های متداول در برازش مدل های چندسطحی، استقلال و نرمال بودن مولفه های تصادفی و خطا است، یعنی

$$\begin{aligned} Cov(\epsilon_{i(j)}, u_{(j)}) &= 0 \\ u_{(j)} &\sim N(0, \sigma_u^2), \quad \epsilon_{i(j)} \sim N(0, \sigma_\epsilon^2) \end{aligned} \quad (3)$$

چون فرض نرمال بودن مولفه تصادفی و خطا در بسیاری از مواقع به دلیل مشاهدات دور افتاده و یا فرم نامتقارن توزیعی برقرار نیستند می توان به جای آن از توزیع چوله نرمال استفاده نمود.

$$u_{(j)} \sim SN(0, \lambda_u, \sigma_u^2), \quad \epsilon_{i(j)} \sim SN(0, \lambda_\epsilon, \sigma_\epsilon^2) \quad (4)$$

که در آن، $(Z \sim SN(0, \sigma^2, \lambda))$ نماد توزیع چوله نرمال با میانگین صفر، واریانس σ^2 و پارامتر چولگی λ با تابع چگالی (آزالینی، ۱۹۸۵)

$$f(z) = \frac{1}{\sigma} \varphi\left(\frac{z-\mu}{\sigma}\right) \Phi\left(\lambda \frac{z-\mu}{\sigma}\right) \quad (5)$$

است به طوری که $\varphi(0)$ و $\Phi(0)$ به ترتیب توابع چگالی و توزیع احتمال نرمال استاندارد هستند. آزالینی (۱۹۸۵) نشان داد اگر X_1 و X_2 هر دو دارای توزیع نرمال استاندارد باشند، آنگاه توزیع متغیر تصادفی Z

$$Z = \lambda |X_1| + X_2 \quad (6)$$

چوله نرمال استاندارد با پارامتر λ است. بنابراین دو بردار تصادفی u و ϵ را می توان به صورت

$$\begin{aligned} u_{(j)} &= \lambda_u |t_{(j)}^u| + t_{(j)}^{u2}, & t_{(j)}^u, t_{(j)}^{u2} &\sim N(0, \sigma_u^2) \\ \epsilon_{i(j)} &= \lambda_\epsilon |t_{i(j)}^\epsilon| + t_{i(j)}^{\epsilon2}, & t_{i(j)}^\epsilon, t_{i(j)}^{\epsilon2} &\sim N(0, \sigma_\epsilon^2) \end{aligned} \quad (7)$$

نوشت. با در نظر گرفتن

$$X = \begin{pmatrix} 1 & x_{11} \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ 1 & x_{MM} \end{pmatrix}, \quad \beta' = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \quad Y = \begin{pmatrix} y_{11} \\ \cdot \\ \cdot \\ \cdot \\ y_{M1} \end{pmatrix}, \quad T^\epsilon = \begin{pmatrix} t_{11}^\epsilon \\ \cdot \\ \cdot \\ \cdot \\ t_{M1}^\epsilon \end{pmatrix},$$

$$U = \begin{pmatrix} u_1 \\ \vdots \\ u_M \end{pmatrix}, \quad T^{u_1} = \begin{pmatrix} t_1^{u_1} \\ \vdots \\ t_M^{u_1} \end{pmatrix},$$

و با توجه به ساختار آشیانی داده‌ها، تابع چگالی توام داده‌ها براساس مولفه‌های تصادفی را می‌توان به صورت

$$f(Y, T^{\epsilon_1}, U, T^{u_1}) = f(Y|T^{\epsilon_1}, U, T^{u_1})f(T^{\epsilon_1}|U, T^{u_1})f(U|T^{u_1})f(T^{u_1}) \quad (8)$$

نوشت. با در نظر گرفتن $n = \sum_{j=1}^M M_j$ داریم

$$Y|T^{\epsilon_1}, U, T^{u_1} \sim N_n(X\beta + U + \lambda_{\epsilon} T^{\epsilon_1}, \sigma_{\epsilon}^2 I_n), \quad T^{u_1} \sim N_M(0, \sigma_u^2 I_M), \\ U|T^{u_1} \sim N_M(\lambda_u T^{u_1}, \sigma_u^2 I_M), \quad T^{\epsilon_1}|U, T^{u_1} \sim N_n(0, \sigma_{\epsilon}^2 I_n) \quad (9)$$

بنابراین تابع چگالی توام به صورت زیر حاصل می‌شود.

$$f(Y, T^{\epsilon_1}, U, T^{u_1}) = \\ \prod_{j=1}^M \prod_{i=1}^{M_j} (\pi \sigma_{\epsilon}^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma_{\epsilon}^2} (y_{i(j)} - \beta_0 - \beta_1 x_{i(j)} - u_{(j)} - \lambda_{\epsilon} |t_{i(j)}^{\epsilon_1}|)^2\right\} \\ \times \prod_{j=1}^M \prod_{i=1}^{M_j} \frac{1}{\sqrt{2\pi\sigma_{\epsilon}^2}} \exp\left\{-\frac{t_{i(j)}^{\epsilon_1}}{2\sigma_{\epsilon}^2}\right\} \\ \times \prod_{j=1}^M \sqrt{\frac{1}{2\pi\sigma_u^2}} \exp\left\{-\frac{(u_j - \lambda_u |t_{(j)}^{u_1}|)^2}{2\sigma_u^2}\right\} \\ \times \prod_{j=1}^M \frac{1}{\sqrt{2\pi\sigma_u^2}} \exp\left\{-\frac{t_{(j)}^{u_1}}{2\sigma_u^2}\right\} \quad (10)$$

۲.۲ توزیع پیشین و توزیع پسین

در این مطالعه توزیع‌های توابع پیشین با استفاده از کارهای هوبرت و کسلا (۱۹۹۶) و ژاو و همکاران (۲۰۰۶) به صورت جدول ۱ انتخاب گردید. با در نظر گرفتن

جدول ۱: توزیع‌های پیشین برای پارامترهای مدل

پارامتر پراکندگی	پارامتر چولگی	ضرایب بتا
$\sigma_u^2 \sim IG(\alpha_u, \beta_u)$	$\lambda_u \sim N(\mu_{\lambda_u}, \sigma_{\lambda_u}^2)$	$\beta_o \sim N(\mu_{\beta_o}, \sigma_{\beta_o}^2)$
$\sigma_\epsilon^2 \sim IG(\alpha_\epsilon, \beta_\epsilon)$	$\lambda_\epsilon \sim N(\mu_{\lambda_\epsilon}, \sigma_{\lambda_\epsilon}^2)$	$\beta_1 \sim N(\mu_{\beta_1}, \sigma_{\beta_1}^2)$

$\theta = (\beta_o, \beta_1, \sigma_u^2, \sigma_\epsilon^2, \lambda_u, \lambda_\epsilon)$ و مستقل فرض کردن پارامترهای مدل، توزیع توام پیشین θ به صورت

$$\pi(\theta) = \pi(\beta_o)\pi(\beta_1)\pi(\sigma_u^2)\pi(\sigma_\epsilon^2)\pi(\lambda_u)\pi(\lambda_\epsilon) \quad (11)$$

است. با ضرب توزیع پیشین θ در تابع چگالی توام، توزیع‌های شرطی کامل پارامترها با انجام محاسبات تحلیلی به صورت زیر حاصل می‌شوند.

$$\beta_o | others \sim N\left(\frac{\frac{\mu_{\beta_o}}{\sigma_{\beta_o}^2} + \frac{1}{\sigma_\epsilon^2} \sum_{j=1}^M \sum_{i=1}^{M_k} (y_{i(j)} - \beta_1 x_{i(j)} - u_{j(k)} - \lambda_\epsilon |t_{i(j)}^{\epsilon}|)}{\frac{1}{\sigma_{\beta_o}^2} + \frac{n}{\sigma_\epsilon^2}}, \left(\frac{1}{\sigma_{\beta_o}^2} + \frac{n}{\sigma_\epsilon^2}\right)^{-1}\right)$$

$$\beta_1 | others \sim N\left(\frac{\frac{\mu_{\beta_1}}{\sigma_{\beta_1}^2} + \frac{1}{\sigma_\epsilon^2} \sum_{j=1}^M \sum_{i=1}^{M_k} (y_{i(j)} - \beta_o - u_j - \lambda_\epsilon |t_{i(j)}^{\epsilon}|) x_{i(jk)}}{\frac{1}{\sigma_{\beta_1}^2} + \frac{\sum_{j=1}^M \sum_{i=1}^{M_k} x_{i(jk)}^2}{\sigma_\epsilon^2}}, \left(\frac{1}{\sigma_{\beta_1}^2} + \frac{\sum_{j=1}^M \sum_{i=1}^{M_k} x_{i(jk)}^2}{\sigma_\epsilon^2}\right)^{-1}\right)$$

$$\lambda_\epsilon | others \sim N\left(\frac{\frac{\mu_{\lambda_\epsilon}}{\sigma_{\lambda_\epsilon}^2} + \frac{1}{\sigma_\epsilon^2} \sum_{j=1}^M \sum_{i=1}^{M_k} (y_{i(j)} - \beta_o - \beta_1 x_{i(j)} - u_{(j)}) |t_{i(j)}^{\epsilon}|}{\frac{1}{\sigma_{\lambda_\epsilon}^2} + \frac{1}{\sigma_\epsilon^2} \sum_{j=1}^M \sum_{i=1}^{M_k} |t_{i(j)}^{\epsilon}|^2}, \left(\frac{1}{\sigma_{\lambda_\epsilon}^2} + \frac{1}{\sigma_\epsilon^2} \sum_{j=1}^M \sum_{i=1}^{M_k} |t_{i(j)}^{\epsilon}|^2\right)^{-1}\right)$$

$$\lambda_u | others \sim N\left(\frac{\frac{\mu_{\lambda_u}}{\sigma_{\lambda_u}^2} + \frac{1}{\sigma_u^2} \sum_{j=1}^M \sum_{i=1}^{M_k} |t_{i(j)}^u|}{\frac{1}{\sigma_{\lambda_u}^2} + \frac{1}{\sigma_u^2} \sum_{j=1}^M \sum_{i=1}^{M_k} |t_{j(i)}^u|^2}, \left(\frac{1}{\sigma_{\lambda_u}^2} + \frac{1}{\sigma_u^2} \sum_{j=1}^M \sum_{i=1}^{M_k} |t_{j(i)}^u|^2\right)^{-1}\right)$$

$$\sigma_e^2 | \text{others} \sim IG(n + \alpha_\epsilon, \circ / \Delta \sum_{j=1}^M \sum_{i=1}^{M_k} (t_{i(j)}^\epsilon)^2 + (y_{i(j)} - \beta_0 - \beta_1 x_{i(j)} - u(j) - \lambda |t_{i(j)}^\epsilon|)^2) + \beta_\epsilon)$$

$$\sigma_u^2 | \text{others} \sim IG(M + \alpha_u, \circ / \Delta \sum_{j=1}^M (t_{(j)}^u)^2 + (u(j) - \lambda_u |t_{(j)}^u|)^2) + \beta_u)$$

۳ مطالعه شبیه‌سازی

برای شبیه‌سازی، ابتدا مدل دوسطحی (۱۲) را با پارامترهای $\beta_0 = 2$ ، $\beta_1 = -3$ ، استراتژی، ۱- مولفه خطا نرمال، مولفه تصادفی نرمال، ۲- مولفه خطا چوله نرمال، مولفه تصادفی نرمال، ۳- مولفه خطا نرمال، مولفه تصادفی چوله نرمال و ۴- مولفه خطا چوله نرمال، مولفه تصادفی چوله نرمال و با استفاده از مقادیر جدول ۲ و به کمک دستور sn در بسته نرم افزاری R شبیه‌سازی شده است. سپس در هر شبیه‌سازی هر چهار مدل با بسته نرم‌افزاری R2WinBUGS و نرم افزار WinBUGS به داده‌ها برازش داده شده است. این عمل ۱۰۰ بار تکرار و خلاصه نتایج در جدول ۳ و ۴ ارائه شده است.

در برازش بیزی بعد از بررسی نمودارهای لازم، زمان داغیدن ۱۰۰۰ تعیین شده سپس ۱۰۰۰۰ نمونه از توزیع پسین استخراج و از هر ۱۰۰ نمونه یک نمونه برای تحلیل نهایی در نظر گرفته شده است. همچنین در این مطالعه برای β_0 ، β_1 و λ توزیع پیشین $N(0, 100)$ و برای σ_e^2 و σ_u^2 توزیع پیشین $IG(0/001, 0/001)$ در نظر گرفته شده است.

$$y_{i(j)} = \beta_0 + \beta_1 x_{i(j)} + u(j) + \epsilon_{i(j)}, \quad i = 1, \dots, 5, \quad j = 1, \dots, 50 \quad (12)$$

برای مقایسه مدل‌های برازش داده شده از ملاک متوسط قدرمطلق خطا به صورت

$$MAE(\theta_k) = \frac{\sum_{i=1}^n |\hat{\theta}_{ik} - \theta_k|}{n} \quad (13)$$

استفاده شده است، که در آن θ_k پارامتر مورد بررسی و $\hat{\theta}_k$ برآورد بیزی آن است.

جدول ۲: چهار استراتژی برای مولفه تصادفی $\epsilon_{i(j)}$ و $u_{(j)}$

توزیع مولفه تصادفی		نماد	استراتژی
$\epsilon_{i(j)}$	$u_{(j)}$		
$N(0, 9)$	$N(0, 4)$	$N - N$	۱
$SN(0, 9, 3)$	$N(0, 4)$	$SN - N$	۲
$N(0, 9)$	$SN(0, 4, 2)$	$N - SN$	۳
$SN(0, 9, 3)$	$SN(0, 4, 2)$	$SN - SN$	۴

نتایج شبیه‌سازی برآورد پارامتر، میزان خطای برآورد براساس ملاک MAE و درصد بیش برآوردی در جداول ۳ و ۴ خلاصه شده‌اند. میزان خطای واقعی برابر با میزان خطای محاسبه شده براساس شاخص MAE است. اما میزان خطای مربوط به مدل میزان انحراف برآورده‌ها در ۱۰۰ بار شبیه‌سازی است. همان‌طور که ملاحظه می‌شود در مدل دوسطحی وقتی مولفه تصادفی و مولفه خطا هر دو نرمال هستند، برازش مدل‌های چوله در مولفه تصادفی و خطای مدل باعث کم برآوردی واریانس مولفه چوله می‌شود. این در حالی است که هنگامی که واقعا مولفه خطا چوله است، و ما از مدل دوسطحی معمولی (توزیع خطا نرمال و توزیع مولفه تصادفی نرمال استفاده نماییم) در صورتی که پارامتر چولگی مثبت باشد مقدار عرض از مبدا با بیش برآوردی و مقدار ضریب متغیر کمکی با کم برآوردی مواجه و در صورت چولگی منفی این رفتار با کم برآوردی عرض از مبدا و بیش برآوردی ضریب متغیر تبیینی مواجه خواهد داد. از طرف دیگر نتایج این شبیه‌سازی نشان داد که در سطح دوم حتی وقتی داده‌ها واقعا چوله هستند، پارامتر چولگی معنی دار نمی‌شود.

۴ مثال کاربردی

داده‌های این پژوهش از نوع مطالعات طولی گذشته‌نگر بوده و جامعه مورد مطالعه آن را کلیه کارکنان شاغل در کارخانه فولاد مبارکه اصفهان طی سالهای ۱۳۷۰ تا ۱۳۹۰ تشکیل می‌داد. تعداد هم‌گروه مورد نظر در این مطالعه ۵۷۴ نفر بوده که از میان ۶۷۱۳ نفر پرسنل این کارخانه با روش نمونه‌گیری طبقه‌ای خوشه‌ای از نواحی کاری و سپس مشاهدات یک فرد با مراجعه به پرونده‌های پزشکی کارکنان انجام

پذیرفت. هدف از این مطالعه بررسی رابطه کار در نوبت کاری و شاخص چاقی (BMI) با کلسترول خون بود. در شکل ۱ بافت‌نگار میزان کلسترول و متوسط کلسترول هر فرد نمایش داده شده است. میزان چولگی برای کلیه مشاهدات ۴/۱۶ و برای میانگین هر نفر در زمانهای متفاوت ۳/۳۷ محاسبه گردید. براساس مقادیر چولگی و بافت‌نگار به خوبی می‌توان به اختلاف اساسی فرض نرمال بودن این دو مولفه در مدلسازی دوسطحی اطمینان نمود.

جدول ۳: برآورد خطای واقعی برآورد پارامترها

برآورد خطای واقعی پارامترهای مدل						مدل در نظر گرفته شده			مدل واقعی
λ_u	λ_e	σ_u^2	σ_e^2	β_1	β_0	u	e		
---	---	۰/۸۵	۰/۷۲	۰/۰۹	۰/۳۸	E	N	N	N - N
---	---	۰/۴۲	۰/۵۷	۰/۵۸	۰/۴۸	P			
---	۲/۰۴	۰/۸۶	۲/۲	۰/۱	۱/۰۹	E	N	SN	
---	۰/۰۱	۰/۴۱	۰/۰۲	۰/۵۵	۰/۴۶	P			
۳/۳	---	۲/۰۸	۰/۷۲	۰/۰۹	۰/۹۴	E	SN	N	
۰/۰۳	---	صفر	۰/۵۳	۰/۵۸	۰/۵۱	P			
۳/۲۴	۲/۰۶	۲/۰۶	۲/۲۳	۰/۰۹	۱/۱۶	E	SN	SN	
۰/۰۳	۲/۰۶	صفر	۰/۰۱	۰/۵۷	۰/۵۴	P			
---	---	۰/۷۴	۱۵/۲۲	۰/۵	۶/۸۶	E	N	N	
---	---	صفر	۱	صفر	۱	P			
---	۰/۵۷	۰/۲	۱/۲	۰/۴۹	۱/۵۵	E	N	SN	
---	۱	صفر	صفر	صفر	۱	P			
۴/۲۷	---	۲/۴۸	۱۴/۸۹	۰/۵	۷/۳۹	E	SN	N	
صفر	---	صفر	۱	صفر	۱	P			
۴/۲۷	۰/۳۶	۲/۵۷	۰/۴۶	۰/۴۹	۲/۵۳	E	SN	SN	
صفر	۰/۳۶	صفر	صفر	صفر	۱	P			
---	---	۱۳/۲۷	۰/۷۸	۰/۱	۴/۷۱	E	N	N	
---	---	۱	۰/۴۷	۰/۶	۱	P			
---	۱/۷۴	۱۳/۳۳	۲/۲۷	۰/۱	۴/۲۶	E	N	SN	
---	صفر	۱	۰/۰۳	۰/۵۸	۱	P			
۱/۱۲	---	۲/۹	۰/۷۸	۰/۱	۱/۶۴	E	SN	N	
۰/۴۰	---	۰/۸۵	۰/۴۵	۰/۵۹	۰/۷۶	P			
۱/۱۲	۱/۷۸	۲/۸۲	۲/۳۵	۰/۱	۱/۷۵	E	SN	SN	
۰/۴۸	۱/۷۸	۰/۸۱	۰/۰۲	۰/۵۴	۰/۶۲	P			
---	---	۱۴/۶۷	۱۳/۳۳	۰/۱۵	۹/۴۸	E	N	N	
---	---	۱	۱	۰/۵۷	۱	P			
---	۰/۷۹	۱۴/۶۴	۳/۰۸	۰/۱۶	۵/۲۷	E	N	SN	
---	۰/۵۷	۱	۰/۵۴	۰/۵۵	۱	P			
۱/۰۸	---	۳/۳۱	۱۳/۱۸	۰/۱۵	۶/۰۵	E	SN	N	
۰/۴۱	---	۰/۸۴	۱	۰/۵۴	۱	P			
/۲۵	۰/۸۳	۳/۳۸	۳/۲۳	۰/۱۶	۲/۵۳	E	SN	SN	
۰/۴۱	۰/۸۳	۰/۸۹	۰/۵۶	۰/۵۲	۰/۸	P			

E: متوسط میزان خطا براساس شاخص MSE، P: درصد بیش برآوردی

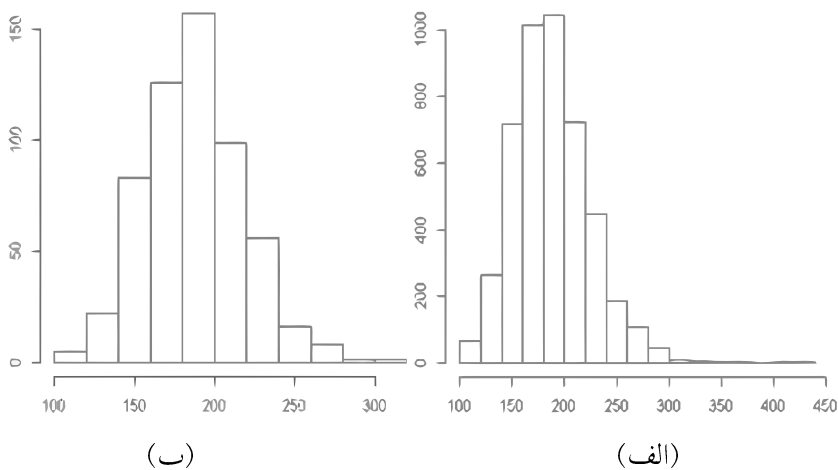
جدول ۴: برآورد میانگین و انحراف معیار برآورد پارامترها

برآورد پارامترهای مدل						مدل در نظر گرفته شده			مدل واقعی
λ_{ii}	λ_e	σ_{ii}^2	σ_e^2	β_1	β_0	u	e		
---	۳/۹۲	---	۹/۲۱	-۲/۹۸	۱/۹۰	C	N	N	
---	۱/۲۱	---	.۹۳	۰/۱۴	۰/۵۳	D			
---	۳/۹۵	-۰/۰۴	۶/۸۱	-۲/۹۹	۱/۹۷	C	N	SN	
---	۱/۲۵	۰/۸۸	۱/۸۱	۰/۱۳	۱/۶۱	D			N - N
-۰/۲۶	۱/۹۲	---	۹/۱۸	-۲/۹۹	۲/۰۵	C	SN	N	
۲/۳۷	۱/۳۶	---	۰/۹۴	۰/۱۳	۱/۶۱	D			
۱/۹۴	۱/۹۴	-۰/۰۴	۶/۷۸	-۲/۹۹	۲/۱۲	C	SN	SN	
۱/۳۹	۱/۳۹	۰/۹۲	۱/۸۶	۰/۱۳	۲/۱۷	D			
---	۳/۲۶	---	۲۴/۲۲	-۳/۵	۲/۸۶	C	N	N	
---	۱/۹۵	---	۲/۵۵	۰/۲۲	۰/۷۷	D			
---	۳/۸	۲/۵۷	۷/۸	-۳/۴۹	۳/۵۵	C	N	SN	
---	۱/۶۸	۰/۷۷	۳/۲۴	۰/۲۱	۰/۹۷	D			SN - N
-۱/۲۷	۱/۵۲	---	۲۳/۸۹	-۳/۵	۹/۳۹	C	SN	N	
۳/۱۹	۱/۵۶	---	۲/۴۳	۰/۲۲	۱/۸۷	D			
۱/۴۳	۱/۱۴۳	۲/۳۶	۸/۵۴	-۳/۴۹	۴/۵۳	C	SN	SN	
۱/۳۷	۱/۳۷	۰/۶۳	۲/۸۸	۰/۲۱	۱/۸۵	D			
---	۱۷/۲۷	---	۹/۱	-۲/۹۹	۶/۷۱	C	N	N	
---	۳/۹۸	---	۰/۹	۰/۱۳	۰/۷۴	D			
---	۱۷/۳۳	۰/۲۶	۶/۷۴	-۲/۹۹	۶/۲۶	C	N	SN	
---	۴/۰۴	۰/۹۴	۱/۸۸	۰/۱۳	۱/۷۵	D			N - SN
۲/۵۲	۶/۶۴	---	۹/۰۵	-۲/۹۹	۳/۳۶	C	SN	N	
۱/۸۲	۴/۴	---	۰/۹۲	۰/۱۳	۲	D			
۶/۵۲	۶/۵۲	۰/۲۳	۶/۶۷	-۳	۲/۹	C	SN	SN	
۴/۳۸	۴/۳۸	۰/۹۴	۱/۸۵	۰/۱۳	۲/۵	D			
---	۱۸/۶۷	---	۲۲/۳۳	-۲/۹۷	۱۱/۴۸	C	N	N	
---	۴/۸۴	---	۲/۲۲	۰/۲۱	۰/۹۳	D			
---	۱۸/۶۵	۲/۱۶	۹/۷۱	-۲/۹۸	۷/۲۷	C	N	SN	
---	۴/۸۸	۰/۹۶	۳/۸۷	۰/۲	۱/۶۱	D			SN - SN
۲/۵۸	۷/۰۱	---	۲۲/۱۸	-۲/۹۸	۸/۰۵	C	SN	N	
۱/۹۹	۵/۱۳	---	۲/۲۶	۰/۲۱	۲/۲۹	D			
۷/۱۸	۷/۱۸	۲/۱۲	۹/۸۴	-۲/۹۹	۴/۰۴	C	SN	SN	
۵/۰۴	۵/۰۴	۰/۹۵	۳/۸۷	۰/۲	۲/۶۲	D			

در مدل برازش داده شده (۱۴) متغیر پاسخ کلسترول خون (Cholesterol)، متغیر نشانگر کار در نوبت کار چرخشی (Shift۱)، متغیر نشانگر کار در نوبت کار هفتگی (Shift۱)، شاخص توده بدنی (BMI)، مولفه تصادفی تکرار (u_(j)) و مولفه تصادفی خطا (ε_{i(j)}) می باشد.

$$Cholesterol_{i(j)} = \beta_0 + \beta_1 Shift_{i(j)} + \beta_2 Shift_{i(j)} + \beta_3 BMI_{i(j)} + u_{(j)} + \epsilon_{i(j)}, \quad i = 1, \dots, n_j \quad j = 1, \dots, 574 \quad (14)$$

پس از بررسی توصیفی داده‌های مطالعه، جهت انجام تحلیل بیزی بعد از بررسی



شکل ۱: (الف) بافت‌نگار میزان کلمسترویل کلیه مشاهدات، (ب) بافت‌نگار میانگین میزان کلمسترویل هر فرد.

نمودارهای لازم، زمان داغیدن ۱۰۰۰ تعیین گردید. سپس ۱۰۰۰۰ نمونه از توزیع پسین استخراج گردید. از هر ۱۰۰ نمونه به دست آمده یکی به عنوان نمونه نهایی در نظر گرفته شد. به عبارتی دیگر در پایان از ۱۰۰ نمونه برای استنباطها استفاده گردید. در این تحلیل برای مقادیر β و λ از توزیع پیشین $N(0, 100)$ و برای پارامترهای واریانس مدل از توزیع پیشین $IG(0/001, 0.001)$ استفاده شده است.

برای مقایسه مدل‌های برازش داده شده از ملاک اطلاع انحراف $(DIC = \overline{D(\theta)} + P_D)$ استفاده شد، که در آن میانگین پسین انحراف بوده و میزان برازش را نشان می‌دهد، P_D برابر با تعداد پارامترهای موثر است و میزان پیچیدگی مدل را نشان می‌دهد. بر اساس این معیار مدلی که دارای کمترین مقدار DIC باشد به عنوان بهترین مدل انتخاب می‌شود. این معیار برای هر اندازه نمونه قابل استفاده بوده و به آسانی توسط روش‌های مونت کار لوی زنجیره مارکوفی قابل محاسبه است (اسپیگرهالتر و همکاران، ۲۰۰۲)

در جدول ۵ خلاصه اطلاعات مربوط به مقادیر DIC محاسبه شده برای ۴ مدل

۲۴۴ تحلیل دوسطحی با مولفه تصادفی چوله نرمال

ارائه شده‌اند. با توجه به مقادیر DIC مدل شماره ۲ (توزیع نرمال چوله برای مولفه خطا و توزیع نرمال برای مولفه های تصادفی فرد) به عنوان بهترین مدل برازش داده شده انتخاب گردید. نتایج حاصل از برآورد و فواصل اطمینان ناپارامتری بیزی ضرایب در جدول ۶ نشان‌دهنده آن است که متغیر BMI بر کلمستروول خون تاثیر مثبت معنی‌دار دارد و متغیر نوبت کاری بر میزان کلمستروول تاثیر معنی‌داری ندارد. همچنین پارامتر چولگی، واریانس سطح اول و خطا معنی‌دار است.

جدول ۵: خلاصه اطلاعات DIC برای مدل های برازش داده شده

شماره مدل	مولفه تصادفی		pD	D(θ)	DIC	زمان برازش
	$u_{(j)}$	$\epsilon_{i(j)}$				
۱	نرمال	نرمال	۵۳۵	۴۱۲۵۵	۴۱۷۶۰	۱۵
۲	نرمال	چوله نرمال	۳۲۱۰	۳۷۱۴۰	۴۰۳۵۰	۱۸۲
۳	چوله نرمال	نرمال	۵۳۱	۴۰۱۵۹	۴۰۶۹۰	۲۹
۴	چوله نرمال	چوله نرمال	۲۵۰۹	۳۷۹۳۱	۴۰۴۴۰	۱۴۷

جدول ۶: نتایج حاصل از برآورد و بازه اطمینان ناپارامتری بیزی ضرایب

	λ_e	$\sigma_{u_i}^2$	σ_e^2	β_3	β_2	β_1	β_0	
برآورد	۱/۸۵	۷۴۱/۶	۲۱۰/۴	۳/۷۷	-۲/۶۸	-۳/۹۰	۷۲/۹۲	
انحراف معیار	۰/۴۱	۵۷/۷۹	۲۹/۷۱	۰/۲۰	۴/۷۱	۲/۷۳	۵/۵۵	
حد پایین	۱/۵۵	۶۵۰/۸	۸۳/۳۵	۲/۳۰	-۱۱/۸۴	-۸/۷۱	۶۱/۲۷	
حد بالا	۳/۸۱	۸۴۷/۵	۲۴۵/۱	۴/۲۱	۶/۶۲	۳/۱۴	۸۵/۲۶	

۵ بحث و نتیجه‌گیری

در این مقاله با توجه به اینکه کلمستروول خون متغیری ذاتاً چوله است (لاچوس و همکاران، ۲۰۰۹) بنابراین انتظار ما پیش از برازش مدل این بود که انتخاب مدل دوسطحی معمول برازش ضعیفی را نسبت به مدل‌های چوله ارائه نماید. این مطلب با برازش ۴ مدل پیشنهادی مورد بررسی قرار گرفته و همانگونه که مقادیر DIC ارائه شده در جدول ۵ نمایش می‌دهد از بین تمام مدل‌های برازش داده شده مدل شماره ۲ (توزیع نرمال چوله برای مولفه خطا و توزیع نرمال برای مولفه‌های تصادفی فرد)

دارای برآزش بهتری نسبت به مدل دوسطحی معمول با توزیع نرمال بوده و توزیع های چوله بخوبی توانسته بود تا اثر تخطی داده‌ها از توزیع نرمال را کنترل نماید. همچنین نتایج حاصل از مطالعه شبیه‌سازی نشان داد که در حضور چولگی مثبت مولفه خطای مدل، مقدار عرض از مبدا بیشتر از مقدار واقعی برآورد و در صورتی که پارامتر چولگی منفی باشد این میزان کمتر از مقدار واقعی برآورد می‌شود. بیش برآوردی عرض از مبدا باعث کم برآوردی ضریب متغیر تبیینی و همچنین کم برآوردی عرض از مبدا باعث بیش برآوردی ضرایب متغیرهای تبیینی می‌گردد (نتر و همکاران، ۱۹۸۷) نتایج بیشتر نشان دهنده عدم معنی داری پارامتر چولگی متغیر تصادفی سطح دوم در داده‌هایی که سطح دوم آن از توزیع چوله شبیه‌سازی شده بود می‌باشد. از آنجایی که برآورد تقریبی مولفه تصادفی خوشه σ_i^2 برابر با $\bar{y}_i - \bar{y}$ است، که در آن \bar{y}_i متوسط متغیر پاسخ برای همه مشاهدات خوشه i ام می‌باشد و بخاطر متوسط گیری و طبق قضیه حد مرکزی که میانگین نمونه به توزیع نرمال میل می‌نماید می‌توان گفت که چولگی مولفه تصادفی سطح دو با متوسط گیری به نوعی از داده‌ها حذف می‌شود.

بنابراین استفاده از خانواده توزیع‌های نامتقارن به جای در نظر گرفتن توزیع نرمال برای مولفه‌های تصادفی و خطای مدل روشی کارا در تحلیل داده‌های دوسطحی با متغیر چوله نرمال بوده و به عنوان جایگزینی مناسب برای روش‌هایی مانند روش تبدیل متغیر و استفاده از روش‌های ناپارامتری پیشنهاد می‌گردد.

تقدیر و تشکر

نویسندگان از پیشنهادات ارزنده داوران گرامی مجله که باعث ارائه بهتر و بهبود مقاله شده است، کمال تشکر را دارند.

- Arellano-Valle, R. B., Bolfarine, H. and Lachos V. H. (2005), Skew-Normal Linear Mixed Models, *Journal of Data Science*, **3**, 415-438.
- Azzalini A. (1985), A Class of Distributions which Includes the Normal Ones, *Scandinavian Journal of Statistics Theory and Applications*, **12**, 171-178.
- Bandyopadhyay, D., Lachos, V. H., Abanto-Valle, C. A. and Ghosh, P. (2010), Linear Mixed Models for Skew-Normal/Independent Bivariate Responses with an Application to Periodontal Disease, *Statistics in Medicine*, **29**, 2643-2655.
- Goldstein, H. (2005), *Handbook of Multilevel Analysis*, Springer, New York.
- Hobert, J. and Casellas, G. (1996), The Effect of Improper Priors on Gibbs Sampling in Hierarchical Linear Mixed Models, *Journal of the American Statistical Association*, **91**, 1461-1473.
- Hox, J. J. (2010), *Multilevel Analysis: Techniques and Applications (Quantitative Methodology Series) Great Britain*, Routledge.
- Jara, A., Quintana, F. and San Martin, E. (2008), Linear Mixed Models with Skew-Elliptical Distributions: a Bayesian Approach, *Computational Statistics and Data Analysis*, **52**, 5033-5045.
- Lachos, V. H., Bolfarine, H. and Arellano-Valle, B. (2009), Bayesian Inference for Skew-Normal Linear Mixed Models, *Journal of Applied Statistics*, **34**, 663-682.

- Lange, K. and Sinsheimer J. S. (1993), Normal/Independent Distributions and Their Applications in Robust Regression, *Journal of Computational and Graphical Statistics*, **2**, 175-198.
- Lee, J. C., Lin, T. I., Lee, K. J. and Hsu, Y. L. (2005), Bayesian Analysis of Box-Cox Transformed Linear Mixed Models with (p,q)ARMA Dependence, *Journal of Statistical Planning and Inference*, **133**, 435-451.
- Lin, T. I. and Lee, J. C. (2006), A Robust Approach to t Linear Mixed Models Applied to Multiple Sclerosis Data, *Statistics in Medicine*, **25**, 1397-1412.
- Lin, T. I. and Lee, J. C. (2007), Bayesian Analysis of hierarchical Linear Mixed Modeling Using Multivariate t Distributions, *Journal of Statistical Planning and Inference*, **137**, 484-495.
- Lin, T. I. and Lee, J. C. (2008), Estimation and Prediction in Linear Mixed Models with Skew-Normal Random Effects for Longitudinal Data, *Statistics in Medicine*, **27**, 1490-1507.
- Luke, A. D. (2004), *Multilevel Modeling (Quantitative Applications in the Social Sciences)*, A Sage University Paper Series.
- Ma, Y. M., Genton, G. and Davidian, M. (2004), *Skew-Elliptical Distributions and their Applications: A Journey Beyond Normality, Linear Mixed Models with Flexible Generalized Skew-Elliptical Random Effect*, Chapman & Hall, CRC Press, Boca Raton, FL.
- Neter, J., Wasserman, W. and Whitmore, G. A. (1987), *Applied Statistics*, Boston Allyn and Bacon Publisher.

تحليل دوسطحی با مولفه تصادفی چوله نرمال ۲۴۸

Pinheiro, J. C. and D. M. Bates (2000), *Mixed-Effects Models in S and S-PLUS*, New York, Springer-Verlag.

Pinheiro, J. C., Liu, C. H. and Wu, Y. N. (2001), Efficient Algorithms for Robust Estimation in Linear Mixed-Effects Models Using a Multivariate t-Distribution, *Journal of Computational and Graphical Statistics*, **10**, 249-276.

Rosa, G. J. M., Padovani, C. R. and Gianola, D. (2003), Robust Linear Mixed Models with Normal/Independent Distributions and Bayesian MCMC Implementation, *Biometrical Journal*, **45**, 573-590.

Spiegelhalter, D. J., Best, N .G., Carlin, B. P. and Van Der Linde, A. (2002), Bayesian Measures of Model Complexity and Fit (with discussion), *Journal of Royal Statistical Society, B*, **64**, 583-640.

Tsung, I. L. and Lee, J. C. (2003), On Modelling Data from Degradation Sample Paths over Time, *Australian and New Zealand Journal of statistics*, **45**, 257-270.

Zhao, Y., Staudenmayer, J., Coull, B. and Wand, M. (2006), General Design Bayesian Generalized Linear Mixed Models, *Statistical Science*, **21**, 35-51.

Zhou, T. and He, X. (2008), Three-Step Estimation in Linear Mixed Models with Skew-t Distributions, *Journal of Statistical Planning and Inference*, **138**, 1542-1555.