

تحلیل کواریانس با فرض چوله نرمال بودن توزیع متغیر پاسخ

افشین فلاح*

دانشگاه بین المللی امام خمینی قزوین

زهرا گودرزی

دانشگاه بین المللی امام خمینی قزوین

چکیده

در این مقاله تحلیل کواریانس تحت فرض چوله نرمال بودن توزیع متغیر پاسخ مورد توجه قرار گرفته و برای این منظور مدلی پیشنهاد شده است. پارامترهای مدل پیشنهادی با استفاده از روش ماکسیمم درستنمایی، یک الگوریتم EM توسعه داده شده است. مدل پیشنهادی در قالب یک مطالعه شبیه سازی مورد ارزیابی قرار گرفته است.

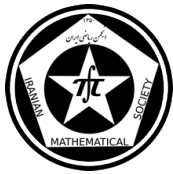
واژه‌های کلیدی: تحلیل کواریانس، توزیع چوله نرمال، تابع درستنمایی داده‌های کامل، متغیر تبیینی

Mathematics Subject Classification [2010]: 62H25, 62F12

۱ مقدمه

تعیین رابطه بین یک متغیر تصادفی پاسخ و مجموعه‌ای از متغیرهای تبیینی یا تیمارها همواره در مطالعات مختلف مورد توجه قرار گرفته است. برای این منظور، می‌توان از روش‌های مختلف تحلیل رگرسیونی و تحلیل واریانس ($ANOVA$) استفاده کرد. تحلیل کواریانس، تلفیقی از تحلیل رگرسیونی و تحلیل واریانس است، که در آن علاوه بر متغیرهای تبیینی، یک یا چند متغیر کیفی نیز در مدل وجود دارد. ایده‌ی اولیه تحلیل کواریانس توسط [۵] طرح شده است. [۴] مفاهیم و کاربردهای مختلف تحلیل کواریانس را مورد بحث قرار داده است. نظریه سنتی تحلیل کواریانس همانند تحلیل رگرسیونی و تحلیل واریانس، عمدتاً بر فرض نرمال بودن متغیر پاسخ بنا شده است. در مواردی که داده‌ها نرمال نباشند، از تبدیل‌هایی برای فراهم آوردن شرایط مورد نظر استفاده می‌شود. انتخاب تبدیل مناسب همواره مسأله‌ای چالش برانگیز است، زیرا در بسیاری موارد استفاده از تبدیل‌ها موجب پیچیده‌تر شدن مدل و تغییر ساختار تابعی آن می‌شود. راهکار دیگری که می‌توان در چنین شرایطی در پیش گرفت، برازش مدل تحلیل کواریانس تحت فرض توزیع‌های دیگری است که بتوانند به خوبی عدم تقارن مشاهدات را در مدل‌سازی لحاظ نمایند. در شرایطی که مشاهدات متقارن نیستند، جایگزین نمودن یک توزیع چوله به جای توزیع نرمال به عنوان توزیع متغیر پاسخ، می‌تواند به بهبود قابل توجه در کارایی مدل تحلیل کواریانس منجر شود. یکی از توزیع‌هایی که لزوماً متقارن نیست اما خواصی مشابه توزیع نرمال دارد، به توزیع چوله نرمال معروف است. این توزیع دارای پارامتری برای مدل‌بندی چولگی است و توزیع نرمال را به عنوان یک عضو شامل می‌شود. اگر چه شکل اولیه توزیع چوله نرمال، توسط [۶] ارائه شد، اما تعریف رایج این خانواده برای نخستین بار توسط [۲] ارائه شد. در طول بیش از سه دهه گذشته توزیع‌های چوله متقارن و چوله نرمال متعددی توسط پژوهشگران مختلف معرفی و مورد مطالعه قرار گرفته است. در این مقاله، برای تحلیل کواریانس تحت فرض چوله نرمال بودن توزیع متغیر پاسخ مورد توجه قرار گرفته است. برای محاسبه برآوردهای ماکسیمم درستنمایی پارامترهای مدل پیشنهادی یک الگوریتم EM توسعه داده شده است. به منظور ارزیابی برآوردهای پیشنهادی از یک مطالعه شبیه‌سازی استفاده است.

* سخنران



۲ مدل پیشنهادی

معمولاً رابطه بین میانگین متغیر پاسخ و مجموعه‌ای از تیمارها و متغیرهای تبیینی، از طریق ترکیب روش‌های مختلف تحلیل رگرسیونی و تحلیل واریانس مورد بررسی قرار می‌گیرد. فرض کنید مدل تحلیل واریانس و رابطه رگرسیونی به ترتیب از طریق ساختار ماتریس طرح X و ماتریس متغیرهای تبیینی Z مشخص می‌شوند. در این صورت، مدل تحلیل کواریانس را می‌توان به شکل $\mu = X\beta + Z\gamma = W\theta$ نوشت، که در آن بردار میانگین مشاهدات پاسخ، $\beta = (\beta_0, \beta_1, \dots, \beta_{s-1})'$ و $\gamma = (\gamma_1, \dots, \gamma_q)$ به ترتیب نشان‌دهنده بردارهای تیمارها و ضرایب رگرسیونی و $W = [X|Z]$ ماتریس مشاهدات را نشان می‌دهد. $\theta = (\beta, \gamma)'$ بردار ضرایب مدل تحلیل کواریانس است. هدف اصلی، برآورد θ با استفاده از مشاهدات $y = (y_1, \dots, y_{sr})'$ و W است. در این بخش تحلیل کواریانس تحت فرض چوله نرمال بودن توزیع متغیر پاسخ مورد توجه قرار گرفته است. متغیر تصادفی Y دارای توزیع چوله نرمال ساهو است $(Y \sim SSN(\mu, \sigma^2, \lambda))$ ، هرگاه تابع چگالی آن به صورت

$$f_Y(y) = \frac{1}{\sigma} \phi(y|\mu, \sigma^2 + \lambda^2) \Phi\left(\frac{\lambda}{\sigma} \frac{(y - \mu)}{(\sigma^2 + \lambda^2)^{\frac{1}{2}}}\right)$$

باشد، که در آن $\phi(\cdot|\mu, \sigma^2)$ و $\Phi(\cdot|\mu, \sigma^2)$ به ترتیب تابع چگالی و تابع توزیع تجمعی توزیع نرمال با میانگین μ و واریانس σ^2 هستند [۷]. این توابع به ازای $(\mu, \sigma) = (0, 1)$ به ترتیب با نمادهای $\phi(\cdot)$ و $\Phi(\cdot)$ نشان داده می‌شوند. این توزیع به ازای $\lambda < 0$ دارای چولگی منفی، به ازای $\lambda > 0$ دارای چولگی مثبت و به ازای $\lambda = 0$ متقارن بوده و به توزیع نرمال تبدیل می‌شود. مدل تحلیل کواریانس تحت فرض چوله نرمال بودن متغیر پاسخ را به صورت

$$Y_{ij}|w_{ij} \sim SSN(w_{ij}\theta - \sqrt{\frac{\lambda}{\pi}}\lambda, \sigma^2, \lambda) \quad i = 1, \dots, s; j = 1, \dots, r, \quad (1)$$

در نظر بگیرید. تابع درستنمایی متناظر با مدل (۱) به صورت

$$L(\theta, \sigma^2, \lambda|y, W) = \prod_{i=1}^s \prod_{j=1}^r \frac{1}{\sigma} \phi(y_{ij}|w_{ij}\theta - \sqrt{\frac{\lambda}{\pi}}\lambda, \sigma^2 + \lambda^2) \Phi\left(\frac{\lambda}{\sigma} \frac{(y_{ij} - w_{ij}\theta)}{(\sigma^2 + \lambda^2)^{\frac{1}{2}}}\right),$$

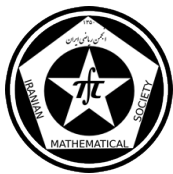
است. چون برآوردهای ماکسیمم درستنمایی پارامترهای تحلیل کواریانس تحت فرض چوله نرمال بودن توزیع متغیر پاسخ فرم بسته ندارند، برای استفاده از الگوریتم EM لازم است مسأله در قالب داده‌های کامل فرمول‌بندی شود. برای این منظور، توزیع چوله نرمال را می‌توان به صورت آمیخته‌ای از توزیع‌های نرمال و نیم‌نرمال به شکل $T_{ij} \sim HN(0, 1)$ و Y_{ij} ‌ها داده‌های ناقص و T_{ij} ‌ها داده‌های گم‌شده تلقی شوند، تابع چگالی توأم (Y_{ij}, T_{ij}) را می‌توان به صورت

$$f_{(Y_{ij}, T_{ij})}(y_{ij}, t_{ij}) = \frac{1}{\sqrt{\lambda\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_{ij} - w_{ij}\theta - \lambda(t_{ij} - \sqrt{\frac{\lambda}{\pi}}))^2\right\} \times \sqrt{\frac{\lambda}{\pi}} \exp\left\{-\frac{t_{ij}^2}{2}\right\},$$

نوشت. بنابراین تابع لگ درستنمایی داده‌های کامل به صورت

$$\begin{aligned} \ell_C(\theta, \sigma^2, \lambda|y, W, t) &= C - \frac{sr}{\lambda} \log \sigma^2 - \frac{1}{2\sigma^2} \left((y - W\theta)'(y - W\theta) - 2\lambda(y - W\theta)'t \right. \\ &\quad \left. - \lambda^2 \mathbf{1}'_{sr} (2\sqrt{\frac{\lambda}{\pi}}t - t^2) + 2\lambda\sqrt{\frac{\lambda}{\pi}} \mathbf{1}'_{sr} (y - W\theta) + sr\lambda^2 \frac{t^2}{\pi} \right), \end{aligned}$$

است، که در آن $\mathbf{1}_{sr}$ یک بردار ستونی sr بعدی از یک‌هاست. در مرحله E از الگوریتم EM ، محاسبه امید ریاضی تابع لگ درستنمایی داده‌های کامل به شرط داده‌های گم‌شده ضروری است. با استفاده



از روابط مربوط به گشتاورهای توزیع نرمال بریده شده (به عنوان نمونه [۳] را ببینید)، روابط

$$\begin{aligned} \hat{t}_{ij} &= E(t_{ij} | \boldsymbol{\theta} = \hat{\boldsymbol{\theta}}, y_{ij}) = \hat{\eta}_{ij} + \hat{\tau} \hat{\delta}_{ij}, \\ \hat{t}_{ij}^2 &= E(t_{ij}^2 | \boldsymbol{\theta} = \hat{\boldsymbol{\theta}}, y_{ij}) = \hat{\eta}_{ij}^2 + \hat{\tau}^2 + \hat{\tau} \hat{\delta}_{ij} \hat{\eta}_{ij}, \end{aligned}$$

به دست می‌آیند، که در آن‌ها $\hat{\delta}_{ij} = \frac{\phi(\frac{\hat{\eta}_{ij}}{\hat{\tau}})}{\Phi(\frac{\hat{\eta}_{ij}}{\hat{\tau}})}$ و $\hat{\tau}^2 = \frac{\sigma^2}{\sigma^2 + \lambda^2}$ ، $\hat{\eta}_{ij} = \frac{\lambda}{\sigma^2 + \lambda^2} (y_{ij} - \mathbf{w}_{ij} \hat{\boldsymbol{\theta}} + \sqrt{\frac{2}{\pi}} \lambda)$ از این رو، امید

ریاضی تابع لگ‌درست‌نمایی داده‌های کامل به شرط داده‌های گم شده به صورت

$$\begin{aligned} E(\ell_c(\boldsymbol{\theta}, \sigma^2, \lambda | \mathbf{y}, \mathbf{W})) &= C - \frac{sr}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \left((\mathbf{y} - \mathbf{W}\boldsymbol{\theta})' (\mathbf{y} - \mathbf{W}\boldsymbol{\theta}) - 2\lambda (\mathbf{y} - \mathbf{W}\boldsymbol{\theta})' \hat{\mathbf{t}} \right. \\ &\quad \left. - 2\lambda^2 \sqrt{\frac{2}{\pi}} \mathbf{1}'_{sr} \hat{\mathbf{t}} + \lambda^2 \mathbf{1}'_{sr} \hat{\mathbf{t}}^2 + 2\lambda \sqrt{\frac{2}{\pi}} \mathbf{1}'_{sr} (\mathbf{y} - \mathbf{W}\boldsymbol{\theta}) + sr \lambda^2 \frac{2}{\pi} \right), \end{aligned}$$

است. در مرحله M الگوریتم، مقادیری از پارامتر که امید ریاضی شرطی تابع لگ‌درست‌نمایی داده‌های کامل را به شرط داده‌های گم شده ماکسیم می‌سازند، به عنوان برآوردگرهای ماکسیم درست‌نمایی انتخاب می‌شوند. بر این اساس، در تکرار $(k+1)$ ام الگوریتم، پارامترها را براساس مقادیر آن‌ها در تکرار t ام به صورت زیر به دست می‌آیند:

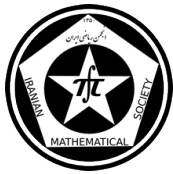
$$\begin{aligned} \hat{\boldsymbol{\theta}}^{(k+1)} &= (\mathbf{W}'\mathbf{W})^{-1} \mathbf{W}' \left(\mathbf{y} + \hat{\lambda}^{(k)} \left(\sqrt{\frac{2}{\pi}} \mathbf{1}_{sr} - \hat{\mathbf{t}}^{(k)} \right) \right), \\ \hat{\sigma}^{2(k+1)} &= \frac{1}{sr} \left((\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\theta}}^{(k)})' (\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\theta}}^{(k)}) - 2\hat{\lambda}^{(k)} (\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\theta}}^{(k)})' \hat{\mathbf{t}}^{(k)} \right. \\ &\quad \left. - \hat{\lambda}^{2(k)} \mathbf{1}'_{sr} \left(2\sqrt{\frac{2}{\pi}} \hat{\mathbf{t}}^{(k)} - \hat{\mathbf{t}}^{(k)} \right) + 2\hat{\lambda}^{(k)} \sqrt{\frac{2}{\pi}} \mathbf{1}'_{sr} (\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\theta}}^{(k)}) + sr \hat{\lambda}^{2(k)} \frac{2}{\pi} \right), \\ \hat{\lambda}^{(k+1)} &= \left(\sqrt{\frac{2}{\pi}} \mathbf{1}'_{sr} (\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\theta}}^{(k)}) - (\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\theta}}^{(k)})' \hat{\mathbf{t}}^{(k)} \right) \left(\mathbf{1}'_{sr} \left(2\sqrt{\frac{2}{\pi}} \hat{\mathbf{t}}^{(k)} - \hat{\mathbf{t}}^{(k)} \right) - sr \frac{2}{\pi} \right)^{-1} \end{aligned}$$

۳ مطالعه شبیه‌سازی

در این بخش، برای ارزیابی کارایی برآوردگرهای ماکسیم درست‌نمایی پیشنهادی، مطالعه‌ای شبیه‌سازی طراحی و اجرا شده است. برای این منظور، یک مدل تحلیل کواریانس با یک متغیر تبیینی و یک تیمار به صورت

$$\mu_{ij} = E(y_{ij} | \mathbf{X}, \mathbf{Z}) = \beta_0 + \beta_i + \gamma(z_{ij} - \bar{z}_i) \quad i = 1, 2; j = 1, \dots, r, \quad (2)$$

به ازای $\beta_0 = 2$ ، $\beta_1 = 5$ و $\beta_2 = 1$ در نظر گرفته شده است. متغیرهای پاسخ y_{ij} ، $i = 1, 2$ و $j = 1, \dots, r$ از توزیع $SSN(\mathbf{w}_{ij}\boldsymbol{\theta} - \sqrt{\frac{2}{\pi}}, \sigma^2, \lambda)$ شبیه‌سازی شده‌اند. به منظور ارزیابی قابلیت مدل پیشنهادی در مدل‌بندی مشاهدات دارای وضعیت‌های مختلف تقارن و چولگی، مقادیر مختلفی برای پارامتر چولگی در محدوده $\{-2, -1, 0, 1, 2\}$ لحاظ شده است. مقادیر ریشه میانگین توان دوم خطا ($RMSE$) برآوردگرهای ماکسیم درست‌نمایی پارامترهای مدل تحلیل کواریانس تحت شرایط بالا محاسبه و در جدول (۱) نشان داده شده است. در تمامی حالات برآوردگرهای حاصل از مدل نرمال نیز به عنوان برآوردگرهای سنتی و برای نشان دادن عدم استواری این برآوردگرها در مقابل انحراف از فرض نرمال بودن مشاهدات، به عنوان مدل رقیب مدنظر قرار گرفته است. به منظور لحاظ نمودن عدم قطعیت حاکم بر فرایند تولید نمونه‌های تصادفی و فراهم آوردن امکان تغییر احتمال به طریقه فراوانی نسبی، در تمام محاسبات تعداد تکرارها 5000 در نظر گرفته شده است. همانطور که انتظار می‌رود به ازای $\lambda = 0$ که مشاهدات متقارن بوده و توزیع‌های نرمال و چوله نرمال بر هم منطبق هستند،



جدول ۱: مقادیر ریشه میانگین توان دوم خطا براوردگرهای ماکسیمم درست‌نمایی ضرایب مدل تحلیل کواریانس چوله نرمال به همراه مقادیر متناظر برای توزیع نرمال به ازای اندازه نمونه‌های مختلف

| اندازه نمونه | λ | مدل تحلیل کواریانس | | | | | |
|--------------|-----------|--------------------|-----------|-----------|------------|-----------|-----------|
| | | نرمال | | | چوله نرمال | | |
| | | γ | β_1 | β_2 | γ | β_1 | β_2 |
| ۱۰۰ | -۲ | ۱۲/۶۷۷۱ | ۵/۰۰۸۴ | ۱/۰۱۸۶ | ۱۲/۲۶۰۳ | ۴/۷۳۳۴ | ۰/۹۱۲۶ |
| ۲۰۰ | -۲ | ۹/۶۶۵۶ | ۵/۰۰۲۹ | ۱/۰۰۷۷ | ۱۰/۱۲۴۹ | ۴/۸۴۸۵ | ۰/۹۲۰۶ |
| ۱۰۰ | -۱ | ۸/۶۴۵۳ | ۵/۰۰۶۸ | ۱/۰۰۶۸ | ۸/۲۷۰۹ | ۴/۷۳۷۰ | ۰/۶۹۹۵ |
| ۲۰۰ | -۱ | ۶/۱۴۴۲ | ۵/۰۰۴۰ | ۱/۰۰۱۷ | ۵/۷۷۴۹ | ۴/۷۹۸۲ | ۰/۷۱۷۰ |
| ۱۰۰ | ۰ | ۵/۴۰۹۳ | ۵/۰۰۱۳ | ۱/۰۰۳۶ | ۵/۴۰۹۳ | ۵/۰۰۱۳ | ۱/۰۰۳۶ |
| ۲۰۰ | ۰ | ۳/۶۸۲۶ | ۴/۹۹۸۱ | ۱/۰۰۱۹ | ۳/۶۸۲۶ | ۴/۹۹۸۱ | ۱/۰۰۱۹ |
| ۱۰۰ | ۱ | ۷/۸۱۴۳ | ۵/۰۰۳۲ | ۱/۰۰۶۷ | ۷/۴۶۲۵ | ۴/۹۵۹۵ | ۰/۶۷۱۱ |
| ۲۰۰ | ۱ | ۶/۲۸۷۷ | ۵/۰۰۰۲ | ۱/۰۰۳۷ | ۶/۱۸۷۵ | ۴/۹۹۸۴ | ۰/۹۹۸۹ |
| ۱۰۰ | ۲ | ۱۲/۹۵۱۱ | ۵/۰۰۹۰ | ۱/۰۱۳۱ | ۹/۶۹۹۴ | ۴/۱۷۵۰ | ۰/۵۳۰۸ |
| ۲۰۰ | ۲ | ۱۰/۲۶۷۶ | ۵/۰۰۱۲ | ۱/۰۱۰۰ | ۷/۶۱۰۹ | ۳/۸۸۴۵ | ۰/۴۴۰۱ |

به وضوح تفاوتی بین مدل‌های تحلیل کواریانس تحت فرض نرمال و چوله نرمال وجود ندارد. به ازای مقادیر مثبت و منفی پارامتر چولگی λ ، که به ترتیب متناظر با مشاهدات چوله به راست و چپ هستند، به وضوح مدل پیشنهادی از کارایی بیشتری برخوردار است.

۴ بحث و نتیجه‌گیری

به صورت سنتی یکی از فرض‌های اولیه در تحلیل کواریانس، نرمال بودن توزیع متغیر پاسخ است. این در حالی است که در بسیاری از کاربردها توزیع مشاهدات نامتقارن است. در چنین شرایطی استفاده از توزیع نرمال نتایج گمراه‌کننده‌ای به دنبال خواهد داشت. استفاده از توزیع چوله نرمال برای مدل‌بندی تغییرات متغیر پاسخ، کارایی مدل را به صورت قابل ملاحظه‌ای افزایش می‌دهد.

مراجع

- [1] D. Aigner, C. Lovell, P. Schmidt, *Formulation and Estimation of Stochastic Frontier Production Function Model*, Journal of Econometrics, 12 (1977), PP. 21-37.
- [2] A. Azzalini, *A Class of Distributions which Includes the Normal ones*, Scandinavian Journal of Statistics, 12 (1985), PP. 171-178.
- [3] R. Barr., E. Donald, and Sherril, *Mean and Variance of Truncated Normal Distribution*, The American Statistician, 53 (1999), PP. 357-361.
- [4] W. Cochran, *Analysis of Covariance: Its Nature and Uses*, Biometrics, 13 (1957), PP. 261-281.
- [5] R. A. Fisher, *Statistical Methods for Research Workers*, 4th ed., Oliver and Boyd, Edinburgh, 1932.
- [6] C. Roberts, *A Correlation Model Useful in The Study of Twins*, Journal of American Statistical Association, 61 (1966), PP. 1184-1190.
- [7] S. Sahu, D. Dey, and M. Branco, *A New Class of Multivariate Distributions with Applications to Bayesian Regression Models*, Canadian Journal of Statistics, 29 (2003), PP. 129-150.