

مدل‌های رگرسیون براساس تابع مفصل

هادی جباری نوقایی، عاطفه جاوید نظام دوست

گروه آمار، دانشگاه فردوسی مشهد، مشهد، ایران

چکیده: در تجزیه و تحلیل مدل‌های وابستگی از جمله مدل‌های رگرسیونی، پیدا کردن تابعی که وابستگی بین متغیر پاسخ و متغیرهای توضیحی را نشان دهد، اهمیت دارد. وابستگی در این تابع به رفتار کناری متغیرها و رفتار توأم آن‌ها مربوط است. چون در تابع مفصل وابستگی به کناری‌ها بستگی ندارد، در این مقاله روشی برای جدا ساختن این دو نوع وابستگی از یکدیگر از طریق تابع مفصل ارائه می‌شود.

واژه‌های کلیدی: تابع مفصل، رگرسیون مفصل، مفصل FGM.

۱ مقدمه

تجزیه و تحلیل رگرسیون یکی از روش‌های آماری با کاربردهای وسیع است که از نظر تاریخی به سال ۱۸۷۵ و فعالیت‌های پژوهشی فرانسویس گالتون (۱۹۹۸) باز می‌گردد. در سال‌های اخیر برای تسهیل تحلیل رگرسیون از تابع مفصل که دارای ویژگی‌های بارزی است، استفاده شده است. برخی از مؤلفین در که در زمینه رگرسیون مفصل تحقیقاتی انجام داده‌اند عبارتند از: چروبینی (۲۰۰۴)، مک نیل و همکاران (۲۰۰۵)، کولیف و همکاران (۲۰۰۸)، و همچنین کتاب نلسن (۲۰۰۶) یک مرجع بسیار مفید برای مفصل و وابستگی می‌باشد.

ایده اصلی این روش، تشخیص رابطه بین چند متغیر است. ساختار مدل رگرسیونی به این صورت است که یک متغیر وابسته یا پاسخ (Y) با چند متغیر مستقل (X_1, \dots, X_n) مرتبط در نظر گرفته می‌شود. مساله اصلی این است که تعیین کنیم توزیع متغیر پاسخ چگونه با متغیرهای مستقل در ارتباط است. در این مقاله با تکیه بر مفهوم تابع توزیع شرطی، تابع رگرسیون مفصل را تعریف می‌کنیم. ساختار مقاله به شرح زیر است:

در بخش دوم پس از بیان تعاریف و مفاهیم اولیه، ویژگی‌های تابع رگرسیون مفصل، تابع رگرسیون خطی و همچنین یک مثال را بیان می‌کنیم. در بخش سوم روش برآورد با استفاده از مفصل را در تجزیه و تحلیل رگرسیون با ذکر یک مثال بیان می‌کنیم.

۱.۱ تعاریف و مفاهیم اولیه

فرض می‌کنیم X و Y متغیرهای پیوسته، به ترتیب با توابع توزیع کناری G و H و تابع توزیع توام F و مفصل C باشند. آن‌گاه $U = G(X)$ و $V = H(Y)$ دارای توزیع یکنواخت روی $(0, 1)$ است و توزیع توام آن‌ها مفصل C است.

تعریف ۱.۱. اگر بردار تصادفی (U, V) دارای توزیع‌های کناری یکنواخت روی $[0, 1]$ و مفصل C باشد، آن‌گاه $E_c[V | U = u]$ تابع رگرسیون مفصل V روی U نامیده می‌شود و آن را با $r_c(u)$ نشان می‌دهیم.

چند نکته **چروینی (۲۰۰۴)** را که در ادامه استفاده خواهد شد، در اینجا ذکر می‌کنیم:

۱- تابع توزیع شرطی V به شرط $U = u$ یعنی $C_u(v)$ به صورت زیر است:

$$\mathbb{P}(V \leq v | U = u) = \frac{\partial C(u, v)}{\partial u} = C_u(v) \quad (1)$$

۲- تابع توزیع شرطی Y به شرط $X = x$ برابر است با:

$$\begin{aligned} \mathbb{P}(Y \leq y | X = x) &= \mathbb{P}(V \leq H(y) | U = G(x)) \\ &= \frac{\partial C(u, v)}{\partial u} \Big|_{u=G(x), v=H(y)} \\ &= C_u(v) \Big|_{u=G(x), v=H(y)} \end{aligned} \quad (2)$$

۳- تابع رگرسیون مفصل V روی U به صورت زیر است:

$$E[V_u] = r_C(u) = 1 - \int_0^1 C_u(v) dv \quad (3)$$

۴- تابع رگرسیون Y روی X به صورت زیر است:

$$\begin{aligned} E[Y | X = x] = \hat{y} &= H^{[-1]} \left(1 - \int_0^1 C_u(v) dv \right) \\ &= H^{[-1]}(r_C(G(x))) \end{aligned} \quad (۴)$$

که در آن $H^{[-1]}$ معکوس تابع توزیع متغیر تصادفی Y است. شایان ذکر است که هر تبدیل اکیداً صعودی و یکنوا از X و Y توزیع حاشیه‌ای را تغییر خواهد داد، بدون اینکه بر رفتار توزیع توام X و Y تاثیری داشته باشد.

گزاره ۱.۱.۱. فرض کنید C یک مفصل و $r_C(u)$ مطابق رابطه (۳) باشد، بنابراین:

الف) اگر $C^*(u, v) = uv$ آنگاه:

$$r_*(u) = \frac{1}{4}$$

ب) اگر $C^+(u, v) = M(u, v) = \min\{u, v\}$ آنگاه:

$$r_+(u) = r_M(u) = u$$

ج) اگر $C^-(u, v) = W(u, v) = \max\{u + v - 1, 0\}$ آنگاه:

$$r_-(u) = r_W(u) = 1 - u$$

برهان. الف) از اینکه $C_*(v) = \frac{\partial C^*(u, v)}{\partial u} = v$ پس به ازای هر $u \in (0, 1)$ داریم:

$$r_*(u) = 1 - \int_0^1 C_*(v) dv = 1 - \int_0^1 v dv = 1 - \frac{1}{2} = \frac{1}{2}$$

$$C_U^+(v) = \frac{\partial C^+(u, v)}{\partial u} = \begin{cases} 0 & ; v \leq u \\ 1 & ; v > u \end{cases} \quad , \quad C^+(u, v) = \begin{cases} v & ; v \leq u \\ u & ; v > u \end{cases} \quad \text{ب) چون}$$

نتیجه:

$$r_+(u) = r_M(u) = 1 - \int_0^1 C_U^+(v) dv = 1 - \int_u^1 1 dv = 1 - 1 + u = u$$

$$\text{ج) با توجه به فرض، } C^-(u, v) = \begin{cases} 0 & ; v \leq 1-u \\ u+v-1 & ; v > 1-u \end{cases}$$

$$C_U^- = \frac{\partial C^-(u, v)}{\partial u} = \begin{cases} 0 & ; v \leq 1-u \\ 1 & ; v > 1-u \end{cases}$$

و در نتیجه:

$$r_-(u) = 1 - \int_0^1 C^-(v) dv = 1 - \int_{1-u}^1 1 dv = 1 - (1 - (1-u)) = 1-u$$

□

۲ تابع رگرسیون خطی

در این بخش تابع رگرسیون خطی را با ارائه یک مثال مورد بررسی قرار می‌دهیم. همچنین ضریب همبستگی پیرسون را محاسبه می‌کنیم. فرض کنید C_t کلاس مفصل‌هایی باشد که تابع رگرسیون خطی دارند، یعنی:

$$C_t = \{C : I^2 \rightarrow I, r_C(u) = \alpha + \beta u\}$$

قضیه ۱.۲. مفصل C دارای تابع رگرسیون خطی است، $(C \in C_t)$ اگر و فقط اگر $r_C(u)$ به صورت

$$r_C(u) = \alpha + (1-2\alpha)u \text{ یا } r_C(u) = \frac{1-\beta}{2} + \beta u \text{ باشد.}$$

برهان. از اینکه:

$$r_C(u) = 1 - \int_0^1 \frac{\partial C(u, v)}{\partial u} dv = \alpha + \beta u$$

نتیجه می‌شود:

$$\frac{\partial}{\partial u} \int_0^1 C(u, v) dv = 1 - \alpha - \beta u$$

و با انتگرال گیری از طرفین معادله به راحتی این قضیه را می توان اثبات کرد. شبیه قضیه ۱۰۲ به راحتی می توان قضیه بعدی را اثبات کرد. □

قضیه ۲۰۲. اگر مفصل C دارای تابع رگرسیون خطی باشد، آنگاه ضریب همبستگی پیرسون

$$\rho_C = 1 - 2\alpha \text{ است و تابع رگرسیون مفصل برابر با } r_C(u) = \frac{1-\rho_C}{4} + \rho_C u \text{ می باشد.}$$

مثال ۱۰۲. مفصل C از خانواده مفصل های FGM را به صورت زیر در نظر می گیریم:

$$C(u, v) = uv [1 + \theta(1 - u)(1 - v)]$$

که در آن $\theta \in [-1, 1]$. چون $\frac{\partial C(u, v)}{\partial u} = v\theta - v^2\theta - 2uv\theta + 2uv^2\theta + v$

$$\begin{aligned} r_C(u) &= 1 - \int_0^1 C_u(v) dv = 1 - \int_0^1 \frac{\partial C(u, v)}{\partial u} dv \\ &= \int_0^1 (1 + \theta - 2u\theta)v dv + \int_0^1 v^2(\theta + 2u\theta) dv \\ &= 1 - \frac{3 + \theta - 2u\theta}{6} = \frac{3 - \theta}{6} + \frac{\theta}{3}u \end{aligned}$$

با مقایسه با فرم کلی تابع رگرسیون خطی $\alpha = \frac{3-\theta}{6}$ ، تابع رگرسیون مفصل اصلاح شده به صورت زیر است:

$$r_C(u) = \alpha + (1 - 2\alpha)u$$

و از این که $\rho_C = 1 - 2\alpha$ در این مثال داریم:

$$\rho_C = 1 - 2\frac{3-\theta}{6} = \frac{\theta}{3}$$

۳ برآورد با استفاده از مفصل در تجزیه و تحلیل رگرسیون

تابع مفصل می تواند شامل یک بردار از پارامترها باشد، که درجه آن از وابستگی بین توابع توزیع حاشیه ای تک متغیره گرفته می شود. در زمینه رگرسیون هر توزیع حاشیه ای را می توان توسط یک بردار از متغیرهای کمکی مشخص کرد. برآورد حاصله توسط انتخاب (کولیف و پایوا، ۲۰۰۹) مفصل مناسب $C(\cdot, \cdot; \theta)$

که برداری از پارامترهای θ و توزیع‌های حاشیه‌ای $F_1(y_1 | x_1, \beta_1)$ و $F_2(y_2 | x_2, \beta_2)$ است. که در آن‌ها x_1 و x_2 متغیرهای کمکی و β_1 و β_2 پارامترهای مجهول هستند. سپس طبق روش‌های درست‌نمایی ماکسیمم و استاندارد کردن، تابع توزیع توام را به صورت زیر داریم:

$$F(y_1, y_2 | x_1, x_2, \beta_1, \beta_2) = C(F_1(y_1 | x_1, \beta_1), F_2(y_2 | x_2, \beta_2); \theta).$$

برای یک بردار از پارامترهای $\Theta = (\beta_1, \beta_2, \theta)$ ، تابع چگالی آن عبارت است از:

$$f(y_{1i}, y_{2i}, \Theta) = c(F_1(y_{1i}, \beta_1), F_2(y_{2i}, \beta_2), \theta) f_1(y_{1i}, \beta_1) f_2(y_{2i}, \beta_2). \quad (5)$$

مثال ۱.۳. فرض کنید توابع حاشیه‌ای y_{1i} و y_{2i} دارای توزیع نمایی به صورت زیر باشند:
 $Y_{1i} \sim \exp(\beta_1)$ و $Y_{2i} \sim \exp(\beta_2)$ و فرض کنید مفصل C از خانواده FGM به نحوه زیر باشد:

$$C(u, v; \theta) = uv + \theta u(1-u)(1-v) \quad \theta \in [-1, 1]$$

و تابع چگالی مفصل

$$c(u, v; \theta) = 1 + \theta(1-2u)(1-2v) \quad (6)$$

باشد. طبق فرمول (5) داریم:

$$f(y_{1i}, y_{2i}, \theta, \beta_1, \beta_2) = \beta_1 e^{-\beta_1 y_{1i}} \cdot \beta_2 e^{-\beta_2 y_{2i}} \cdot c(F_1(y_{1i}, \beta_1), F_2(y_{2i}, \beta_2), \theta)$$

با استفاده از رابطه (6) نتیجه می‌گیریم:

$$f(y_{1i}, y_{2i}, \theta, \beta_1, \beta_2) = \beta_1 \beta_2 e^{-\beta_1 y_{1i} - \beta_2 y_{2i}} [1 + \theta(1 - 2(1 - e^{-\beta_1 y_{1i}}))(1 - 2(1 - e^{-\beta_2 y_{2i}}))].$$

حال می‌خواهیم پارامترهای β_1 ، β_2 ، θ را برآورد کنیم. با توجه به این که $\tau_{X,Y} = \frac{\theta}{4}$ است، می‌توان θ را به صورت زیر برآورد کرد:

$$\hat{\theta} = \frac{9}{4} \hat{\tau}$$

بنابراین داریم:

$$f(y_{1i}, y_{2i}, \theta, \beta_1, \beta_2) = \beta_1 \beta_2 e^{-\beta_1 y_{1i} - \beta_2 y_{2i}} \left[1 + \hat{\theta} (1 - \gamma (1 - e^{-\beta_1 y_{1i}})) (1 - \gamma (1 - e^{-\beta_2 y_{2i}})) \right].$$

حال با استفاده از روش درستنمایی ماکسیمم می‌خواهیم پارامترهای β_1 و β_2 را برآورد کنیم. بنابراین داریم:

$$L(\beta_1, \beta_2) = \sum_{i=1}^n \ln(f \{y_{1i}, y_{2i}, \hat{\theta}, \beta_1, \beta_2\}) \\ = \sum_{i=1}^n \left\{ \ln \beta_1 + \ln \beta_2 - \beta_1 y_{1i} - \beta_2 y_{2i} + \ln \left[1 + \hat{\theta} (1 - \gamma (1 - e^{-\beta_1 y_{1i}})) (1 - \gamma (1 - e^{-\beta_2 y_{2i}})) \right] \right\}.$$

حال نسبت به β_1 مشتق می‌گیریم:

$$\frac{\partial L(\beta_1, \beta_2)}{\partial \beta_1} = \sum_{i=1}^n \frac{\partial \ln(f)}{\partial \beta_1}$$

که در آن:

$$\sum_{i=1}^n \frac{\partial \ln(f)}{\partial \beta_1} = \sum_{i=1}^n \left[\frac{\beta_2 e^{-\beta_1 y_{1i} - \beta_2 y_{2i}} - \beta_1 \beta_2 e^{-\beta_1 y_{1i} - \beta_2 y_{2i}} \gamma (1 + \hat{\theta} (-1 + \gamma e^{-\beta_1 y_{1i}})) (-1 + \gamma e^{-\beta_2 y_{2i}})}{\beta_1 \beta_2 e^{-\beta_1 y_{1i} - \beta_2 y_{2i}} (1 + \hat{\theta} (-1 + \gamma e^{-\beta_1 y_{1i}})) (-1 + \gamma e^{-\beta_2 y_{2i}})} \right. \\ \left. - \frac{\gamma \beta_1 \beta_2 e^{-\beta_1 y_{1i} - \beta_2 y_{2i}} \gamma y_{1i} e^{-\beta_1 y_{1i}} (-1 + \gamma e^{-\beta_2 y_{2i}})}{\beta_1 \beta_2 e^{-\beta_1 y_{1i} - \beta_2 y_{2i}} (1 + \hat{\theta} (-1 + \gamma e^{-\beta_1 y_{1i}})) (-1 + \gamma e^{-\beta_2 y_{2i}})} \right].$$

با توجه به این که توزیع متقارن است مشتق لگاریتم تابع درستنمایی نسبت به β_2 نیز مشابه β_1 است. بنابراین برای برآورد پارامترهای β_1 و β_2 باید دستگاه معادلات زیر حل شود:

$$\begin{cases} \sum_{i=1}^n \frac{\partial \ln(f)}{\partial \beta_1} = 0 \\ \sum_{i=1}^n \frac{\partial \ln(f)}{\partial \beta_2} = 0 \end{cases} \quad (7)$$

حال با ارایه یک مثال عددی، نتایج مثال ۱.۳ را بکار می‌گیریم. در این مثال داده‌ها را از صفحه اقتصادی روزنامه خراسان جمع‌آوری کرده‌ایم به این صورت که در هر روز از شش ماهه دوم سال ۱۳۹۰ نرخ دلار برحسب ریال استخراج شده است و بیشترین مقدار و کمترین مقدار نرخ دلار در هر هفته نیز مد نظر است. داده‌ها شامل ۲۵۰ زوج مشاهده از بیشترین نرخ دلار (y_1) با میانگین ۱/۰۲۳ و واریانس

۰/۹۸۷ و کمترین نرخ دلار (y_2) با میانگین ۱/۱۲۱ و واریانس ۱/۲۶۲ است. ابتدا با استفاده از آزمون نیکویی برازش کلموگروف-اسمیرنوف برای متغیر y_1 با $p - value = ۰/۹۷۵۳$ و y_2 با $p - value = ۰/۵۳۵۶$ ملاحظه می‌شود که هر دو متغیر دارای توزیع نمایی با پارامتر یک هستند. همچنین مقادیر ضرایب همبستگی پیرسون (۰/۰۳۹۵۵۴۴۴)، کندال (۰/۰۱۱۴۵۷۲۹) و اسپیرمن (۰/۰۱۷۸۳۵۴۵) تاییدی بر وابستگی ضعیف در بین این دو مجموعه داده است، بنابراین مفصل FGM برای این داده‌ها مناسب است. همچنین با انجام آزمون نیکویی برازش (gofCopula) در پکیج (copula) در R، با $p - value = ۰/۹۸۶۵$ مشخص می‌شود که مفصل FGM با پارامتر θ برای داده‌ها مناسب است.

$$\hat{\theta} = ۰/۰۵۵۸۹۹۰۲$$

بنابراین به راحتی با حل دستگاه (۷) می‌توان پارامترهای β_1 و β_2 را به صورت زیر برآورد کرد:

$$\hat{\beta}_1 = ۰/۹۷۶۴۴۴۹$$

$$\hat{\beta}_2 = ۰/۸۹۱۸۵۳۶.$$

بحث و نتیجه‌گیری

در این مقاله تلاش کردیم تا نگاه عمیق‌تری به معادله رگرسیونی نسبت به حالت تابعی ساده آن داشته باشیم و با مثالی نشان دادیم که فرم تابع خط رگرسیون، به رفتار توأم متغیرها (که مفصل تعیین کننده آن است) و رفتار حاشیه‌ای (که با تابع توزیع حاشیه‌ای شکل می‌گیرد) وابسته است. هدف اصلی در تجزیه و تحلیل رگرسیون، نمایش تابع رگرسیون به صورت خطی است، که در این مقاله نیز بررسی شد. همچنین نحوه برآورد با استفاده از مفصل در تجزیه و تحلیل رگرسیون را با ذکر مثالی نشان دادیم.

مراجع

Cherubini, U. Luciano, E. Vecchiato, W. (2004), *Copula Methods in Finance*. New York: John Wiley and Sons, p. 182

- Frees, E., and E. Valdez, "Understanding Relationships Using Copulas," *North American Actuarial Journal*, *Journal 2*, 1112, pp. 1–25.
- Kolev, N., Anjos, U. and Mendes, B. (2006). *Copulas: a review and recent developments*. *Statistical Models* 22, 617-660.
- Kolev, N., and Paiva, D., (2009), *Copula-based regression models*. *Journal of Statistical Planning and Inference* 139, 3847-3856
- McNeil, A., Frey, R. and Embrechts, P. (2005). *Quantitative Risk Management: Concepts, Techniques and Tools*. Princeton University Press.
- Nelsen, R. (2006). *An Introduction to Copulas*, 2nd Edition, Springer: New York.