

وابستگی رگرسیونی و توابع مفصل

محمد حسین علامت‌ساز^۱

فائزه شفیعی‌فر^۱

چکیده

مفاهیم وابستگی از ابزارهای مهم درک روابط بین رویدادهای تصادفی چندمتغیره هستند. در مقاله حاضر، به بررسی مفهوم وابستگی رگرسیونی خواهیم پرداخت. به ویژه، نشان خواهیم داد که چگونه می‌توان از توابع مفصل برای بررسی این نوع وابستگی سود جست. همچنین، به نوعی آزمون فرض در این مورد اشاره خواهیم کرد و با مثال‌های واقعی و عددی موضوع را تشریح خواهیم کرد. **واژه‌های کلیدی:** وابستگی رگرسیونی، توابع مفصل، استقلال.

۱. مقدمه

مفاهیم اولیه وابستگی مثبت و منفی، توسط لهن (۱۹۶۶) بیان شد. مفاهیمی که توسط وی ارایه شد، همگی در حالت دو متغیره طرح شده بودند. [۴] و در سال‌های اخیر بزرگ‌نیا و امینی [۱۰] مفاهیم وابستگی را در حالت چندمتغیره مورد بررسی قرار دادند.

توابع مفصل توسط اسکالر در سال ۱۹۵۹ معرفی شدند. مفصل تابعی است که تابع توزیع چندمتغیره را به تابع توزیع‌های حاشیه‌ای تک‌متغیره آن پیوند می‌دهد. امروزه، توابع مفصل به عنوان راهی برای مطالعه اندازه‌های آزاد از مقیاس و وابستگی پیشنهاد می‌شوند، زیرا ویژگی‌های ناپارامتری بودن، ناوردایی و وابستگی توابع توزیع متغیرهای تصادفی را به خوبی بیان می‌کنند.

در بخش ۲، وابستگی رگرسیونی مثبت و منفی را معرفی و تشریح می‌کنیم. در بخش ۳، توابع مفصل مورد بررسی قرار گرفته و برخی

خواص مهم آنها بیان خواهد شد. کاربرد توابع مفصل را در تعیین وابستگی رگرسیونی در بخش ۴ مورد توجه قرار خواهیم داد و بالاخره در بخش ۵ به آزمون فرضیه لی‌من در مورد وابستگی رگرسیونی در مقابل فرضیه استقلال اشاره خواهیم کرد.

۲. وابستگی رگرسیونی

۲-۱- وابستگی رگرسیونی دو متغیره

تعریف ۲-۱: زوج تصادفی (X, Y) را وابسته رگرسیونی منفی گوئیم، هرگاه $P[Y \leq y | X = x]$ تابعی غیر نزولی بر حسب x باشد. در این صورت، X و Y را وابسته رگرسیونی منفی^۲ (NRD) گوئیم و از نماد $NRD(Y | X)$ برای نمایش آن استفاده می‌کنیم. خانواده تمام توزیع‌های (X, Y) را که در آنها (X, Y) وابسته رگرسیونی منفی هستند، با \mathcal{g} نشان می‌دهیم.

^۱ گروه آمار، دانشگاه اصفهان

^۲ Negative Regression Dependence

اثبات:

چون $(X, U) \in g_+$ ، لذا $P[U \leq u | X = x]$ تابعی غیر نزولی بر حسب x است. برای اثبات کافی، نشان دهیم که

$$P[Y \leq y | X = x]$$

$\forall x_1 < x_2$:

$$\begin{aligned} P[Y \leq y | X = x_1] &= P[U \leq y - \alpha - \beta X \leq y | X = x_1] \\ &= P[U \leq y - \alpha - \beta x_1 | X = x_1] \\ &\leq P[U \leq y - \alpha - \beta x_2 | X = x_1] \\ &\text{با توجه به اینکه } (X, Y) \in g_+ \text{ و } \beta \leq 0 \\ &\leq P[U \leq y - \alpha - \beta x_2 | X = x_2] \\ &= P[U \leq y - \alpha - \beta X | X = x_2] \\ &= P[Y \leq y | X = x_2] \end{aligned}$$

در نتیجه $P[Y \leq y | X = x]$ تابعی غیر نزولی از x است [۱۰].

۲-۲- تعمیم چندمتغیره

سه مفهوم وابستگی که می‌توانند به عنوان تعمیم چندمتغیره وابستگی

رگرسیونی فرض شوند، عبارتند از

۱- وابستگی مثبت از طریق ترتیب تصادفی^۴ (PDS)

۲- صعودی شرطی در دنباله^۵ (CIS)

۳- نزولی شرطی در دنباله^۶ (CDS)

تعریف ۲-۳: بردار تصادفی (X_1, \dots, X_n) را (PDS) گوییم، اگر $\{X_i : i \neq j\}$ مشروط بر $X_j = x$ به ازای هر j با افزایش یافتن x به طور تصادفی افزایش یابد. به عبارت ساده‌تر، یعنی اگر بدانیم یکی از X ها در بردار (X_1, \dots, X_n) افزایش یافته است، آنگاه سایر X ها هم افزایش می‌یابند.

تعریف ۲-۴: بردار تصادفی (X_1, \dots, X_n) را (CIS) گوییم، هرگاه برای هر $x_i \in \mathcal{R}$:

$$P(X_i > x_i | X_j = x_j, j = 1, \dots, i-1) \quad i = 2, \dots, n \quad (2)$$

تابعی صعودی بر حسب X_1, \dots, X_{i-1} باشد. به عبارت ساده‌تر، یعنی اگر بدانیم X_1, \dots, X_{i-1} افزایش یافته‌اند، احتمال اینکه x_i افزایش

تعریف ۲-۲: X و Y را وابسته رگرسیونی مثبت^۷ (PRD) گوییم، هرگاه تابع $P[Y \leq y | X = x]$ به ازای هر $x, y \in \mathcal{R}$ بر حسب x غیر صعودی باشد؛ و از نماد $PRD(Y | X)$ برای نمایش آن استفاده می‌کنیم. توزیع‌های (X, Y) را که در آنها (X, Y) وابسته رگرسیونی مثبت هستند، با F_+ نشان می‌دهیم.

در اینجا به بیان چند خاصیت مقدماتی وابستگی رگرسیونی که توسط امینی [۱۰] اثبات شده می‌پردازیم.

قضیه ۲-۱:

الف) زوج (X, Y) ، NRD است، اگر و تنها اگر $(-X, Y)$ ، PRD باشد.

ب) زوج (X, Y) ، NRD است، اگر و تنها اگر $(X, -Y)$ ، PRD باشد.

مثال ۲-۱: فرض کنید $Y = \alpha + \beta X + U$ و X و U مستقل باشند، آنگاه،

الف: $(X, Y) \in g_+ \Leftrightarrow \beta \leq 0$

ب: $(X, Y) \in f_+ \Leftrightarrow \beta \geq 0$

اثبات:

$\forall x, y \in \mathcal{R}$:

$$\begin{aligned} P[Y \leq y | X = x] &= P[\alpha + \beta X + U \leq y | X = x] \\ &= P[U \leq y - \alpha - \beta X | X = x] \\ &\text{(بنابر استقلال } U \text{ و } X) \\ &= P[U \leq y - \alpha - \beta x] \\ &= F_u(y - \alpha - \beta x) \end{aligned}$$

مشاهده می‌کنیم که $F_u(y - \alpha - \beta x)$ تابعی غیر نزولی از x است، اگر و تنها اگر $\beta \leq 0$. در نتیجه (الف) ثابت می‌شود و به همین ترتیب تابعی غیر صعودی از x است، اگر و تنها اگر $\beta \geq 0$. در نتیجه (ب) هم ثابت می‌شود.

مثال ۲-۲: هرگاه $Y = \alpha + \beta X + U$ ، $\beta \leq 0$ و $(X, U) \in g_+$ ؛ آنگاه $(X, Y) \in g_+$.

^۴ Positive Dependence through the Stochastic ordering

^۵ Conditional Increasing in sequence

^۶ Conditional Decreasing in sequence

^۷ Positive Regression Dependence

برعکس، برای هر دو تابع توزیع تک‌متغیره F و G و هر مفصل C ، تابع H بالا یک تابع توزیع دو بعدی با توزیع حاشیه‌ای F و G تعریف می‌کند. همچنین، اگر F و G پیوسته باشد، C یکتاست.

فرض کنید X و Y متغیرهای تصادفی پیوسته با تابع توزیع توأم $H(x, y)$ و توزیع‌های حاشیه‌ای به ترتیب $F(x)$ و $G(y)$ باشند. متغیرهای تصادفی $U = F(X)$ و $V = G(Y)$ را در نظر می‌گیریم. بدیهی است که هر کدام از متغیرهای U و V دارای توزیع یکنواخت بر $[0, 1]$ تابع توزیع توأم $C(u, v)$ هستند. ایده اساسی و در عین حال ساده تبدیل هر کدام از متغیرهای تصادفی به این دلیل ارائه شده است که توزیع‌های حاشیه‌ای، توزیع یکنواخت بر فاصله‌ای به طول یک باشد و نیز گشتاورها و سایر پارامترهای توزیع دو متغیره «ناوردای مقیاس» باشند. بنابراین،

$$C(u, v) = H(F^{-1}(u), G^{-1}(v)) \quad , \quad u, v \in [0, 1] = I \quad (5)$$

بدیهی است که بنابر تعریف و بر اساس نتایج قضیه اسکالر داریم: X و Y مستقل هستند، اگر و تنها اگر $C(u, v) = uv$. در این صورت آن را مفصل حاصلضرب گوئیم و با نماد $\prod(u, v)$ نشان می‌دهیم.

کران‌های تابع مفصل توسط فرشه به صورت زیر مشخص گردید،

$$W(u, v) = \max\{u + v - 1, 0\} \quad (6)$$

$$M(u, v) = \min\{u, v\}$$

$$W(u, v) \leq C(u, v) \leq M(u, v)$$

این کران‌ها که به ترتیب، کران بالا و پایین فرشه گفته می‌شوند، خود توابع مفصل دو متغیره هستند.

طبق قضیه اسکالر مفصل C یک تابع توزیع می‌باشد. اگر آن را با C نشان دهیم، داریم؛

$$c(u, v) = \frac{\partial^2 C(u, v)}{\partial u \partial v} \quad (7)$$

در نتیجه

$$H(x_1, x_2) = C(F(x_1), G(x_2)) \quad (8)$$

$$H(x_1, x_2) = C(u, v)$$

$$\frac{\partial^2 H(x_1, x_2)}{\partial x_1 \partial x_2} = \frac{\partial^2 C(u, v)}{\partial u \partial v} \frac{\partial u}{\partial x_1} \frac{\partial v}{\partial x_2}$$

$$h(x_1, x_2) = c(u, v) f(x_1) g(x_2)$$

۴-۲- مفصل چندمتغیره

برای توزیع‌های چندمتغیره پیوسته، حاشیه‌های تک‌متغیره و چندمتغیره با ساختار وابستگی توسط مفصل‌ها قابل نمایش می‌باشند. در این حالت،

یابد بیشتر می‌شود. توجه داشته باشید در حالت دو متغیره مفهوم صعودی شرطی در دنباله معادل با وابستگی رگرسیونی مثبت می‌باشد.

تعریف ۲-۵: بردار تصادفی (X_1, \dots, X_n) را (CDS) گوئیم، هرگاه برای هر $x_i \in \mathcal{R}$:

$$P(X_i > x_i | X_j = x_j, \quad j = 1, \dots, i-1) \quad i = 2, \dots, n \quad (3)$$

تابعی نزولی بر حسب X_1, \dots, X_{i-1} باشد در حالت $n = 2$ مفهوم نزولی شرطی در دنباله معادل با وابستگی رگرسیونی منفی می‌باشد [۴].

۳. تابع مفصل^۷

مفصل‌ها، تابع‌هایی هستند که تابع‌های توزیع یک متغیره را به فرم تابع‌های توزیع چندمتغیره پیوند می‌دهند. کلمه مفصل برای اولین بار توسط اسکالر^۸ به کار رفت و از کلمه لاتین «به هم پیوستن» گرفته شده است. جو [۴] و نلسن [۹، ۸، ۷ و ۶] و امبرجتز [۳] گام‌های مهمی را برای اشاعه کاربرد توابع مفصل در زمینه مفاهیم وابستگی دو متغیره و چندمتغیره برداشتند.

۴-۱- مفصل‌های دو متغیره

یک مفصل دو بعدی تابعی مانند $I = [0, 1] \rightarrow I^2 : C$ با ویژگی‌های زیر می‌باشد؛

(الف)

$$C(0, t) = C(t, 0) = 0 \quad , \quad C(1, t) = C(t, 1) = t \quad \forall t \in [0, 1]$$

$$u_1 \leq u_2 \quad , \quad v_1 \leq v_2 \quad \forall u_1, u_2, v_1, v_2 \in I \quad (ب)$$

$$C(u_1, v_1) - C(u_1, v_2) - C(u_2, v_1) + C(u_2, v_2) \geq 0 \quad (۴)$$

به سادگی ثابت می‌شود که تابع مفصل C نسبت به هر یک از متغیرهای خود، غیر نزولی و پیوسته است [۷].

اهمیت مفصل‌ها در آمار ریاضی از قضیه‌ای موسوم به قضیه اسکالر مشخص می‌شود که بیان می‌دارد؛ هرگاه H یک تابع توزیع دو بعدی با تابع توزیع‌های حاشیه‌ای F و G باشد، در نتیجه یک مفصل C وجود دارد، به طوری که $H(x, y) = C(F(x), G(y))$.

^۷ Copula function

^۸ Sklar, A.

^۹ Copulare

$$\begin{aligned} G_i(x) &= P[T_i(X_i) \leq x] = P[X_i \leq T_i^{-1}(x)] = F_i(T_i^{-1}(x)) \\ C_T &= C_T[G_1(x_1), \dots, G_n(x_n)] \quad , \quad x \in \mathfrak{R} \\ &= P[T_1(X_1) \leq x_1, \dots, T_n(X_n) \leq x_n] \\ &= P[X_1 \leq T_1^{-1}(x_1), \dots, X_n \leq T_n^{-1}(x_n)] \\ &= C[F_1(T_1^{-1}(x_1)), \dots, F_n(T_n^{-1}(x_n))] \\ &= C[G_1(x_1), \dots, G_n(x_n)] \end{aligned}$$

چون X_1, \dots, X_n پیوسته هستند و دامنه G_i ها $[0, 1]$ است، می توان نتیجه گرفت که در $[0, 1]^n$ ، $C_T = C$ [۳].

فرض کنید X_1, \dots, X_n متغیرهای تصادفی و به ترتیب دارای توابع توزیع F_1, \dots, F_n و چگالی های f_1, \dots, f_n باشند. همچنین تابع چگالی توأم آنها برابر h و تابع مفصل چگالی برابر c باشد. آنگاه با توجه به (۹) داریم،

$$h(x_1, \dots, x_n) = c(u_1, \dots, u_n) f_1(x_1) \dots f_n(x_n) \quad (11)$$

۴. تابع مفصل و مفهوم وابستگی رگرسیونی

قضیه ۴-۱: فرض کنید X و Y دو متغیر تصادفی با تابع توزیع توأم $H(x, y)$ و به ترتیب حاشیه های $F(x)$ و $G(y)$ باشند و C تابع مفصل باشد. آنگاه

الف: $PRD(Y|X)$ اگر و تنها اگر $\frac{\partial}{\partial u} C(u, v)$ به ازای تمام v ها بر حسب u نزولی باشد.

ب: $NRD(Y|X)$ اگر و تنها اگر $\frac{\partial}{\partial u} C(u, v)$ به ازای تمام v ها بر حسب u صعودی باشد.

اثبات: فرض کنید تابع C یک مفصل متغیرهای تصادفی X و Y باشد و همچنین $F(x) = u$ و $G(y) = v$. آنگاه

$$\begin{aligned} P[Y \leq y | X = x] &= \frac{\frac{\partial}{\partial x} H(x, y)}{\frac{d}{dx} F(x)} \\ &= \frac{\frac{\partial}{\partial u} C(u, v) \frac{du}{dx}}{\frac{du}{dx}} = \frac{\partial}{\partial u} C(u, v) \end{aligned}$$

مفصل یک توزیع چندمتغیره است که تمام حاشیه های تک متغیره آن دارای توزیع یکنواخت $[0, 1]$ می باشند. بنابراین، اگر C یک مفصل باشد، تابع توزیع یک بردار تصادفی یکنواخت چندمتغیره است.

یک توزیع n متغیره H که از مین تک متغیره آن دارای تابع توزیع حاشیه ای F_j است را در نظر بگیرید. مفصل مربوط به H یک تابع توزیع به صورت $[0, 1]^n \rightarrow [0, 1]$ است که در رابطه زیر صدق می کند،

$$H(x) = C(F_1(x_1), \dots, F_n(x_n)) \quad (9)$$

$$x = (x_1, x_2, \dots, x_n) \in \mathfrak{R}^n$$

قضیه ۳-۱: اگر H یک تابع توزیع n متغیره پیوسته با حاشیه ای های تک متغیره F_1, \dots, F_n باشد و تابع چندک به صورت $F_1^{-1}, \dots, F_n^{-1}$ باشد، آنگاه یک انتخاب یکتا برای (۹) به صورت زیر است.

$$C(u) = H(F_1^{-1}(u_1), \dots, F_n^{-1}(u_n)) \quad (10)$$

$$u = (u_1, \dots, u_n) \in [0, 1]^n$$

$$F_i^{-1}(u_i) = \inf \{x | F(x_i) \geq u_i\} \quad , \quad u_i \in [0, 1]$$

اثبات:

برای $i = 1, \dots, n$ فرض می کنیم $X_i \sim F_i$ ، لذا $U_i = F_i(X_i) \sim U(0, 1)$

$$\begin{aligned} C(u_1, \dots, u_n) &= P[U_1 \leq u_1, \dots, U_n \leq u_n] \\ &= P[F_1(X_1) \leq u_1, \dots, F_n(X_n) \leq u_n] \end{aligned}$$

بنابر یکنوا بودن توابع توزیع

$$\begin{aligned} &= P[X_1 \leq F_1^{-1}(u_1), \dots, X_n \leq F_n^{-1}(u_n)] \\ &= F(F_1^{-1}(u_1), \dots, F_n^{-1}(u_n)) \end{aligned}$$

قضیه ۳-۲: فرض کنید (X_1, \dots, X_n) برداری از متغیرهای تصادفی پیوسته با مفصل C باشد و T_1, \dots, T_n توابعی پیوسته و صعودی روی دامنه X_1, \dots, X_n باشند. در این صورت، $(T_1(X_1), \dots, T_n(X_n))$ نیز دارای مفصل C است.

اثبات: فرض کنید F_1, \dots, F_n به ترتیب توابع توزیع متغیرهای تصادفی X_1, \dots, X_n و G_1, \dots, G_n به ترتیب تابع توزیع های تبدیل های $(T_1(X_1), \dots, T_n(X_n))$ فرض کنید که (X_1, \dots, X_n) دارای مفصل C بوده و $(T_1(X_1), \dots, T_n(X_n))$ مفصل C_T داشته باشد. داریم

۰/۱۰، ۰/۵۰ و ۰/۹۰ چگالی‌های تقریبی را با جدول (۲) مقایسه کرد و نتایج به دست آمده را در جدول (۳)، شکل (۳) خلاصه نمود.

مرحله بعدی، انتخاب یک تابع مفصل مناسب است که وابستگی متغیرهای تصادفی را دربرمی‌گیرد. اگرچه در ابتدا ارتباط متغیرها بسیار پیچیده به نظر می‌رسید، اما بررسی بر روی بسیاری از یافته‌های شهودی و داده‌های قبلی نشان داد که رابطه قوی وابستگی رگرسیونی بین متغیرها برقرار است. بررسی ریلی نشان داد که اگر X_i و X_j متغیرهای بحرانی باشند، $F(x_j | x_i)$ غیر صعودی (غیر نزولی) بر حسب x_j برای تمام x_i ها می‌باشد. در این مدل، وابستگی رگرسیونی مثبت و منفی وجود دارد و ریلی اصطلاح وابستگی رگرسیونی یکنوا را به کار برد. تعدادی از خانواده‌های مفصل را می‌شناسیم که می‌توانند وابستگی رگرسیونی را دربرگیرند. ولی انتخاب مفصل چندمتغیره نرمال بسیار معقول می‌باشد. زیرا اولاً برای مدلی که وابستگی از طریق ضریب همبستگی‌ای مثل اسپیرمن بیان شده است، مناسب می‌باشد. ثانیاً تغییرپذیری و راحتی تحلیل مفصل چندمتغیره نرمال، امتیاز بالایی را برای انتخاب شدن به این مفصل می‌دهد. در ادامه، چگونگی یافتن تابع مفصل برای داده‌های شرکت هوایی را شرح می‌دهیم.

خانواده گاوسی^{۱۰} (نرمال)

مفصل توزیع نرمال π متغیره با ماتریس همبستگی R به صورت زیر است.

$$C_R^{Ga}(u) = \Phi_R''(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_n)) \quad (12)$$

Φ_R'' نماد تابع توزیع توأم π متغیر از تابع توزیع نرمال استاندارد با ماتریس همبستگی R است. Φ^{-1} نماد معکوس تابع توزیع نرمال استاندارد تک‌متغیره می‌باشد. در این صورت داریم،

$$C_n(F_1, \dots, F_n) = \Phi_R''(\Phi^{-1}(F_1), \dots, \Phi^{-1}(F_n)) \quad (13)$$

به طوری که Φ_R'' تابع توزیع π متغیره نرمال است که متغیرهای تصادفی آن دارای همبستگی خطی R_n می‌باشند. R_n ماتریس مربعی است که عناصر قطر ۱ آن و عناصر غیر قطری برابر $(\frac{1}{2} \sin(\pi \rho_{ij} / 6))$ می‌کنیم که ρ_{ij} ضریب همبستگی رتبه‌ای اسپیرمن می‌باشد و F_i تابع توزیع متغیر π است [۳].

با توجه به روابط بالا و تعریف وابستگی رگرسیونی هر دو حکم ثابت می‌شود [۸ و ۶].

با استفاده از ویژگی‌های هندسی تابع مفصل، برای مفاهیم وابستگی می‌توان یک تفسیر هندسی ارائه داد. بدیهی است که اگر رویه $Z = C(u, v)$ منطبق بر رویه $Z = \prod(u, v)$ قرار گیرد، X و Y مستقل هستند.

قضیه ۴-۴: $PRD(Y | X)$ اگر و تنها اگر $Z = C(u, v)$ به ازای تمام u ها تابعی مقعر از u باشد.

اثبات: چون تابع $\mathcal{R} \rightarrow [0, 1]: f$ مقعر است، اگر و تنها اگر مشتق‌های یک طرفه آن f'_+ و f'_- نزولی باشند. در نتیجه، تابع $C(u, v)$ نسبت به u مقعر است، اگر و تنها اگر $\frac{\partial C(u, v)}{\partial u}$ بر حسب u ، به ازای تمام v ها نزولی باشد، لذا با توجه به قضیه ۴-۱ حکم ثابت می‌شود [۶].

مثال ۴-۱: در سال ۱۹۹۷ یک شرکت هواپیمایی تصمیم به خرید یک هواپیمای جدید گرفت. صاحبان این شرکت می‌خواستند بدانند خرید هواپیمای سود بیشتری به همراه می‌آورد یا یک سرمایه‌گذاری پولساز دیگر؟ آنها به آماردانی به نام ریلی مراجعه کردند تا سیستم ریسک و تصمیم را مدیریت کند و یک نمودار از عوامل مؤثر در تصمیم‌گیری و جدولی که سطح پایین، پایه و بالای متغیرهای ورودی را نشان می‌داد به او ارائه کردند. این نمودار و جدول در شکل (۲) و جدول (۱) نمایش داده شده است. در جدول ۱۰ امین، ۵۰ امین و ۹۰ امین درصد از توزیع احتمال هر کدام به عنوان سطوح متغیرها معرفی شده است [۲].

برای ریلی اولین مرحله ارائه یک تحلیل حساس، تعیین متغیرهای بحرانی بود. او پس از مطالعه و بررسی، چهار متغیر بحرانی و با نفوذ معرفی کرد. روش انتخاب این چهار متغیر، بر اساس فرآیند تصمیم‌گیری حساس دو مرحله‌ای بوده است که مورد بحث ما نمی‌باشد و از آن صرف‌نظر می‌کنیم. این چهار متغیر عبارتند از:

(H): ساعت پرواز ، (O): هزینه عملیات

(P): قیمت ، (C): گنجایش

اطلاعات مربوط به این متغیرهای در جدول (۲) نمایش داده شده است. مرحله دوم، مدل‌بندی توزیع توأم این متغیرها بود. او ابتدا چگالی‌های حاشیه‌ای متغیرها را مدل‌بندی کرد. برای اثبات صحت برازش، کسرهای

به طوری که $Y = (y_1, \dots, y_p)'$ و I ماتریس همانی 4×4 می باشد. دوباره عبارت (۱۲) را به کار می بریم و می نویسیم،

$$f(P, H, C, O|R) = f_1(P) \dots f_p(O) \quad (18)$$

$$\times \frac{1}{|R|^{1/2}} \exp\left\{-\frac{1}{2} Y^T (R^{-1} - I) Y\right\}$$

$$= \frac{f_1(P) \dots f_p(O)}{|R|^{1/2}}$$

$$\times \exp\left\{-\frac{1}{2} [\Phi^{-1}(F_1(P)), \dots, \Phi^{-1}(F_p(O))]^T \times (R^{-1} - I) [\Phi^{-1}(F_1(P)), \dots, \Phi^{-1}(F_p(O))]\right\}$$

این تابع چگالی توأم دارای حاشیه های مشخص $f_1(P)$ و ... و $f_p(O)$ می باشد. به دلیل آنکه ضریب همبستگی اسپیرمن تحت تبدیل یک به یک ناوردا است، بنابراین R همان ضریب همبستگی متغیرهای O, C, H, P می باشد. ریلی با محاسبه تابع چگالی توأم، به راحتی توانست توابع چگالی شرطی و همچنین امید ریاضی شرطی را به دست آورد و وارد مرحله آخر، یعنی مرحله تصمیم گیری شود [۲].

۵. آزمون استقلال در مقابل PRD

لی من [۵] آزمونی را برای فرض استقلال در مقابل فرض (PRD)

ارایه داد. آزمون پیشنهادی لی من به ترتیب زیر می باشد

$$\begin{cases} H_0: & Y \text{ و } X \text{ مستقل هستند} \\ H_1: & PRD(Y|X) \end{cases}$$

فرض کنید نمونه تصادفی $(x_1, y_1), \dots, (x_n, y_n)$ در دست باشد، آماره آزمون عبارت است از؛

$$T = \frac{1}{n(n-1)} \sum_{i \neq j} Sgn(x_i - x_j) Sgn(y_i - y_j) \quad (19)$$

مشخص است که آماره آزمون، برآورد τ ی کندال می باشد. فرض کنید تحت فرض استقلال ناحیه رد با اندازه α به صورت زیر باشد.

$$P(|T| \geq t_{\alpha/2} | H_0) = \alpha$$

که نماد $t_{\alpha/2}$ نمایانگر نقطه بحرانی توزیع دقیق T تحت فرض صفر می باشد. می توان ثابت کرد، تحت فرض استقلال آماره T دارای توزیع پارامتری مقارن است. این توزیع را می توان با محاسبه تعداد

فرض کنید تابع چگالی متغیر P با $f_1(P)$ و تابع توزیع این متغیر با $f_1(P)$ نمایش داده شوند. به همین ترتیب $f_2(H), f_2(C), f_2(O)$ و توابع چگالی و $F_2(H), F_2(C), F_2(O)$ توابع توزیع H, C, O هستند. $F_1(P) = u_1, u_2 = F_2(H), u_3 = F_2(C), u_4 = F_2(O)$ را در نظر بگیرید. مفصل نرمال عبارت است از؛

$$C^*(u) = \Phi_R^*(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_4)) \quad (14)$$

به طوری که Φ_R^* نمایانگر تابع توزیع توأم 4 متغیره نرمال استاندارد با ماتریس همبستگی خطی R می باشد و Φ^{-1} معکوس تابع توزیع نرمال استاندارد تک متغیره است. برای به دست آوردن R از رابطه $r_{ij} = 2 \sin(\pi \rho_{ij} / 6)$ استفاده می شود. بنابر رابطه (۱۲) می دانیم،

$$f(y_1, \dots, y_p) = f_1(y_1) \dots f_p(y_p) c(F_1(y_1), \dots, F_p(y_p)) \quad (15)$$

در اینجا، یک ترفند ریاضی - کار گرفته می شود. $y_i = \Phi^{-1}(u_i)$ دارای توزیع نرمال استاندارد است و می توان نوشت

$$f_N(y_1, \dots, y_p) = f_N(y_1) \dots f_N(y_p) c(F_N(y_1), \dots, F_N(y_p))$$

که F_N و f_N به ترتیب تابع توزیع و چگالی نرمال استاندارد تک متغیره هستند. در نتیجه،

$$c(F_N(y_1), \dots, F_N(y_p)) = \frac{f_N(y_1, \dots, y_p)}{f_N(y_1) \dots f_N(y_p)} \quad (17)$$

$$= \frac{1}{(\sqrt{2\pi})^p |R|^{1/2}} \exp\left\{-\frac{1}{2} Y^T R^{-1} Y\right\}$$

$$= \frac{1}{\prod_{i=1}^p \sqrt{2\pi} \exp\left\{-\frac{1}{2} y_i^2\right\}}$$

$$= \frac{1}{|R|^{1/2}} \exp\left\{-\frac{1}{2} Y^T R^{-1} Y\right\}$$

$$= \frac{1}{\exp\left\{-\frac{1}{2} Y^T Y\right\}}$$

$$= \frac{1}{|R|^{1/2}} \exp\left\{-\frac{1}{2} Y^T (R^{-1} - I) Y\right\}$$

و برای n های بزرگ، تحت فرض استقلال X و Y داریم [۱]

(۲۱)

$$Z = \frac{T - 0}{\sqrt{VarT}} = \frac{\sqrt[3]{n(n-1)}}{\sqrt{2(n+5)}} T \xrightarrow{d} N(0,1)$$

مثال ۲-۵: به عنوان یک مثال ساده برای آزمون وابستگی مثبت، داده‌های بیمه حوادث آتش سوزی [۱۱ و صفحات ۱۴۲-۱۳۸] را مورد بررسی قرار می‌دهیم. در این بیمه، دو نوع هزینه وجود دارد. هزینه‌های نوع A^{۱۱} که هزینه‌هایی هستند که مختص عملیات شرکت بیمه می‌باشند و شامل دستمزد مشاور حقوقی، هزینه‌های رسیدگی به ادعاها و ... می‌باشند.

هزینه‌های نوع L^{۱۲} که هزینه‌های شامل خساراتی که بیمه ادعا می‌کند، می‌باشد. در اینجا، از داده‌های واقعی بیمه آتش سوزی که دفتر خدمات بیمه‌ای در اختیار ما قرار داده، استفاده کرده‌ایم و از آنها نمونه‌های تصادفی به حجم ۱۵۰۰ مورد در نظر گرفته‌ایم. با استفاده از نرم‌افزار S-Plus آزمون لی‌من را بر این نمونه تصادفی انجام داده که برنامه و خروجی آن در ادامه می‌باشد. برآورد τ کندال عبارت است از $T = 0.3109967$. مقدار آماره آزمون برابر است با $Z = 18.0499$.

همچنین مقدار احتمال عبارت است از؛

$$P-V = P(\tau \text{ کندال جامعه حداقل به بزرگی } \tau \text{ مشاهده شده} | H_0) = 0$$

مشاهده می‌کنیم، مقدار احتمال از هر سطح معنی‌دار دلخواهی کوچک‌تر باشد. بنابراین، فرض استقلال به نفع وابستگی رگرسیونی مثبت رد می‌شود و نتیجه می‌گیریم متغیرهای A و L در داده‌های بیمه آتش سوزی وابستگی رگرسیونی مثبت می‌باشند.

جایگشت‌های رتبه‌های X و Y به دست آورد. دو مفهوم پایه‌ای هماهنگی و ناهماهنگی را یادآوری می‌کنیم.

تعریف ۱-۵: دو زوج مستقل (X_i, Y_i) و (X_j, Y_j) از متغیرهای تصادفی که از یک توزیع دو متغیره پیروی می‌کنند را در نظر بگیرید. اگر داشته باشیم؛ $X_i < X_j$ وقتی که $Y_i < Y_j$ یا $X_i > X_j$ وقتی که $Y_i > Y_j$ رابطه فوق یک هماهنگی (توافق) نامیده می‌شود.

تعریف ۲-۵: برای هر دو زوج مستقل، اگر داشته باشیم؛ $X_j < X_i$ وقتی که $Y_i > Y_j$ یا $X_i > X_j$ وقتی که $Y_i < Y_j$ ، این رابطه را یک ناهماهنگی (عدم توافق) نامیده می‌شود.

در جدول پیوست کتاب آمار ناپارامتری [۱] چندک‌های توزیع دقیق آماره $W = N_c - N_d$ آورده شده است که متغیر N_c معرف تعداد زوج‌های هماهنگی و N_d معرف تعداد زوج‌های ناهماهنگ می‌باشد.

مثال ۱-۵: توزیع دقیق آماره کندال را برای $n=3$ تحت فرض H_0 به دست می‌آوریم. تعداد جایگشت‌های متمایز رتبه X ها و Y ها برابر $3! = 6$ می‌باشد. این جایگشت‌ها را در جدول ۴ نشان می‌دهیم و برآورد آماره (t) را محاسبه می‌کنیم. که در آن R_{X_i} و R_{Y_i} به ترتیب رتبه x_i در بین X ها و رتبه y_i در بین Y ها و n_c و n_d به ترتیب تعداد زوج‌های هماهنگ و ناهماهنگ می‌باشند. توزیع دقیق آماره کندال برای $n=3$ عبارت است از

t	$P[\tau = t]$
۱	۱/۶
۱/۳	۲/۶
-۱/۳	۲/۶
-۱	۱/۶

توزیع احتمال T برای زوج‌های نمونه از هر جامعه دو متغیره به طور مجانبی نرمال است. بنابراین، اگر بتوان میانگین و واریانس T را تعیین کرد، T در نمونه‌های بزرگ برای استنباط در مورد فرض استقلال بسیار مناسب است. چون تحت فرض استقلال X و Y آماره T متقارن بوده و همواره $-1 \leq T \leq +1$ است، پس $E(T) = 0$ می‌توان ثابت کرد،

$$Var(T) = \frac{2(2n+5)}{9n(n-1)} \quad (20)$$

Allocated Loss Adjustment Expenses (ALAE's)^{۱۱}

Losses^{۱۲}

جدول ۱

متغیر (X)	پایین	پایه	بالا
درصد:	۰/۱۰	۰/۵۰	۰/۹۰
نسبت مؤسس	٪۴۵	٪۵۰	٪۷۰
گنجایش	۴۰۵	٪۵۰	٪۶۰
سطح قیمت	\$۹۵	\$۱۰۰	\$۱۰۸
ساعت پرواز	۵۰۰	۸۰۰	۱۰۰۰
هزینه عملیات در هر ساعت	\$۲۳۰	\$۲۴۵	\$۲۶۰
درصد سرمایه گذارها	٪۳۰	٪۴۰	٪۵۰
نرخ بهره	٪۱۰/۱۵۰	٪۱۱/۱۵۰	٪۱۳/۱۰۰
بیمه	\$۱۸۰۰۰	\$۲۰۰۰۰	\$۲۵۰۰۰
مالیات	\$۸۵۰۰۰	\$۸۷۵۰۰	\$۹۰۰۰۰

جدول ۲

متغیر	ضریب همبستگی						
	پایین	پایه	بالا				
متغیر	درصد:	۰/۱۰	۰/۵۰	۰/۹۰	P	H	C
P	\$۹۵	\$۱۰۰	\$۱۰۸				
H	۵۰۰	۸۰۰	۱۰۰۰	-۰/۱۵۰			
C	٪۴۰	٪۵۰	٪۶۰	-۰/۲۵	۰/۱۵۰		
O	\$۲۳۰	\$۲۴۵	\$۲۶۰	.	.	۰/۲۵	

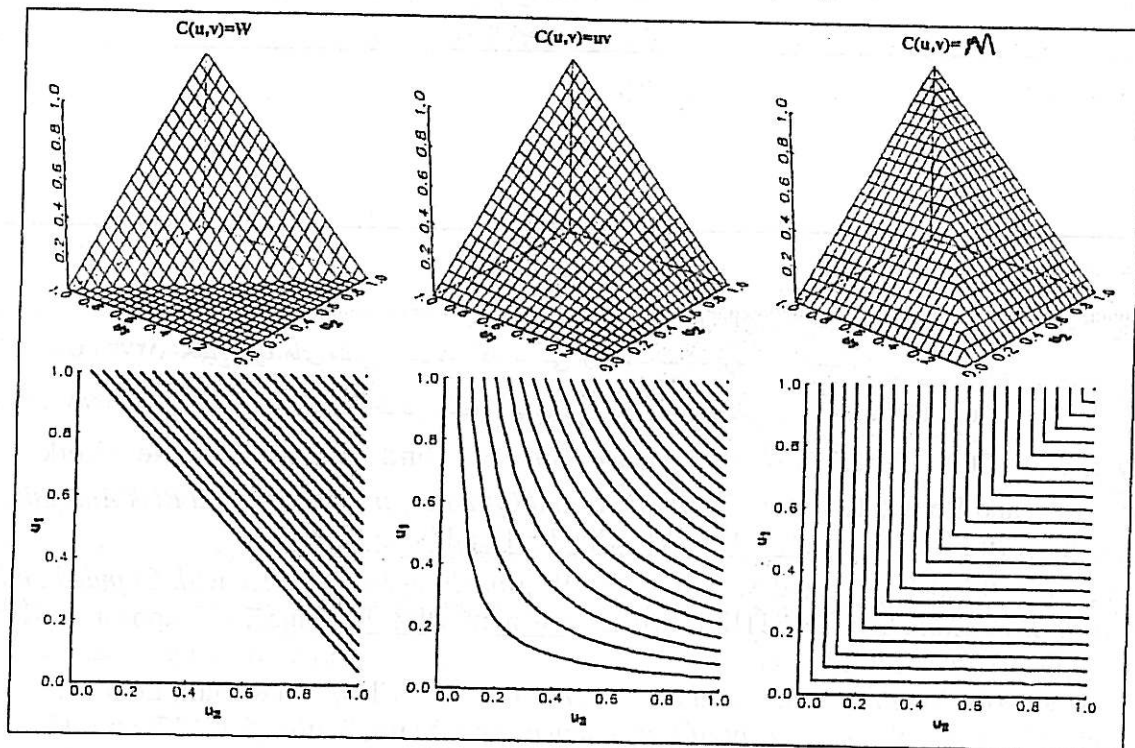
جدول ۳

متغیر	توزیع	پارامترها	دامنه
P	Scaled beta	=۹، =۱۵	[\$۸۱/۹۴، \$۱۳۳/۹۶]
H	Scaled beta	=۴، =۲	[۶۶/۹۱، ۱۱۵۳/۲۶]
C	Beta	=۲۰، =۲۰	[۰، ۱]
O	Normal	=۲۴۵، =۱۱/۷۲	(-∞، ∞)

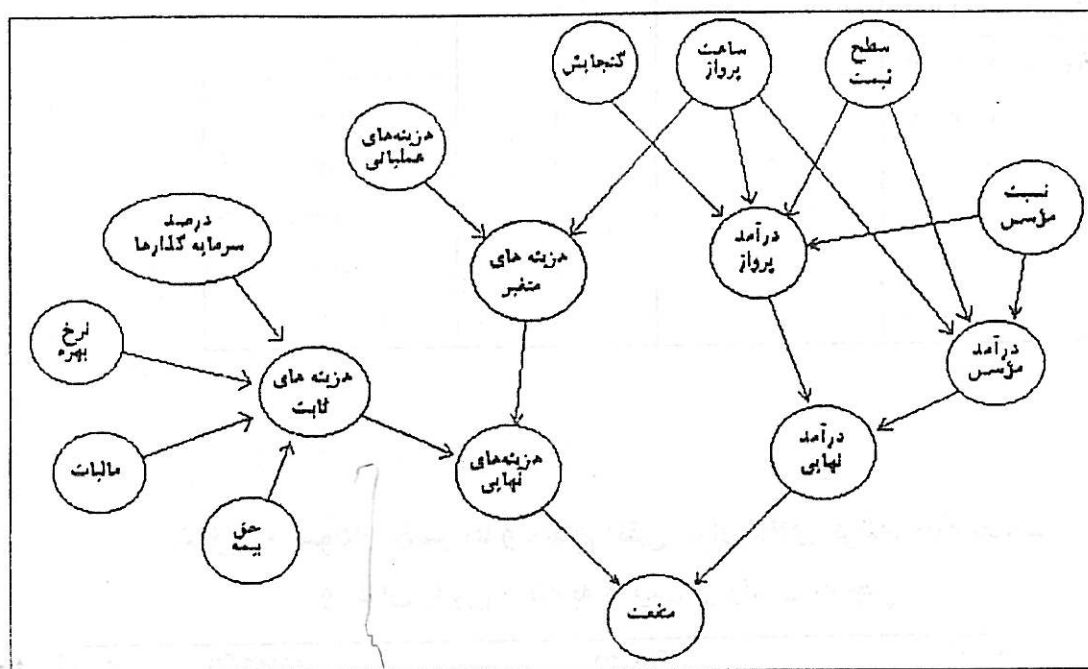
جدول ۴

$i = 1 \ 2 \ 3$	n_c	n_d	$t = (n_c - n_d) / 3$
$R_{Xi} = 1 \ 2 \ 3$			
$R_{Yi} = 1 \ 2 \ 3$	۳	۰	۱
۱ ۳ ۲	۲	۱	۱/۳
۲ ۱ ۳	۲	۱	۱/۳
۲ ۳ ۱	۱	۲	-۱/۳
۳ ۱ ۲	۱	۲	-۱/۳
۳ ۲ ۱	۰	۱	-۱

شکل ۱- نمودار مفصل‌ها و مقطع افقی کران بالای فرشه، حاصلضرب و کران پایین فرشه به ترتیب از راست به چپ



شکل ۲



مراجع

[۱] امینی دهک، محمد؛ ۱۳۷۳، تحلیل پیوندهای منفی، پایان‌نامه کارشناسی ارشد، دانشگاه فردوسی مشهد.

[۲] شفیعی فر، فائزه؛ ۱۳۸۱، اندازه‌های وابستگی و ارتباط در مدل‌های تصادفی، پایان‌نامه کارشناسی ارشد، دانشگاه اصفهان.

[3] Conover, W.T., 2000. *Practical Nonparametric Statistics*, Third Edition, Wiley, New York.

[4] Clemen, R.T. and Reilly, T., 1999. *Correlations and Copulas for Decision and Risk Analysis*, www.bus.utexas.edu/Faculty/jim.Dver/DA-WP/WP970012.pdf.

[5] Embrechts, P., Lindskog, F. and Mcneil, A., 2001. *Modelling Dependence with Copulas and Applications to Risk Management*, ETHZ, Zurich, www.math.ethz.ch/finance, will appear in: Handbook of Heavy Tailed Distributions in Finance.

[6] Joe, H., 1997. *Multivariate Models and Dependence Concepts*, Chapman & Hall, London.

[7] Lehmann, E.L., 1966. *Some Concept of Dependence*, Ann. Math. Statist., Vol. 37, pp. 1137-1153.

[8] Nelsen, R.B., 1992. *On Measures of Association as Measures of Positive Dependence*, Statistics and Probability Letters, Vol. 14, pp. 269-274.

[9] Nelsen, R.B., 1995. *Copulas, Characterization, Correlation and Counterexamples*, Mathematics Magazine, Vol. 68, pp. 193-198.

[10] Nelsen, R.B., 1999. *An Introduction to Copulas*, Springer, New York.

[11] Nelsen, R.B., 2001. *Concordance and Copulas: A Survey. Distributions with Given Marginals and Statistical Modelling*, (Cuadras, C., Fortiana, J., Rodriguez Lallena, J.A., Eds.), Kluwer Academic publishers, Dordrecht, in press.