

دومین کارگاه قابلیت اعتماد و کاربردهای آن، ۱۸ و ۱۹ خرداد ۱۳۹۰، قطب علمی داده‌های ترتیبی و فضایی

برآورد ضریب همبستگی کندال برای داده‌های بقای وابسته از طریق مدل‌های مفصل ارشمیدسی دو متغیره تحت سانسور تصادفی

هادی جباری نوقابی *

استادیار گروه آمار دانشگاه فردوسی مشهد Jabbarinh@um.ac.ir

اعظم مختاری

کارشناس ارشد آمار دانشگاه آزاد اسلامی واحد مشهد

حسنعلی آذرنوش

استاد گروه آمار دانشگاه آزاد اسلامی واحد مشهد، قطب علمی داده‌های ترتیبی و فضایی دانشگاه فردوسی مشهد

چکیده

بردار زمان بقای (T_1, T_2) را که به وسیله مدل مفصل ارشمیدسی مدل‌بندی می‌شود، در نظر می‌گیریم. در این مقاله با برآورد مفصل کلایتون به داده‌ها و تعمیم ضریب همبستگی کندال برای داده‌های سانسور شده از راست، یک برآورد گشتاوری ساده برای پارامتر وابستگی مدل‌های مفصل ارشمیدسی به دست می‌آوریم. مفصل‌های ارشمیدسی، مفصل کلایتون، ضریب همبستگی τ کندال، داده‌های بقای وابسته

۱ مقدمه

تاکنون مدل‌های زیادی برای بررسی داده‌های بقای چند متغیره‌ی وابسته پیشنهاد شده‌اند. در میان آن‌ها، مدل‌های ارشمیدسی بسیار متداول هستند. در این زمینه باندین روچ و لیانگ (۱۹۹۶)، ونگ و ول (۲۰۰۰) و هوگارد (۲۰۰۰) کار کرده‌اند. براین اساس تابع بقای دو متغیره $S(t_1, t_2)$ ، با توابع حاشیه‌ای $S(t_1) = S(t_1, 0)$ و $S(t_2) = S(0, t_2)$ را می‌توانیم به وسیله مفصل ارشمیدسی با در نظر گرفتن وابستگی بین دو زمان بقا مدل‌بندی کنیم. در این صورت داریم:

$$S(t_1, t_2) = p[q(S_1(t_1)) + q(S_2(t_2))]$$

که در آن، تابع $q(\cdot)$ یک تابع محدب و نزولی بر بازه $[0, 1]$ است و $q(1) = 0$. همچنین $q(\cdot)$ معکوس تابع $p(\cdot)$ است. جنست و ریوست (۱۹۹۳) و اکس (۱۹۸۹) در این زمینه تحقیق کرده‌اند. در این مقاله برای مدل‌بندی وابستگی بین داده‌های بقا از مفصل کلایتون استفاده می‌کنیم که در آن $p(S) = (1 + S)^{-\frac{1}{\alpha}}$. با توجه به مفصل کلایتون تابع بقای دو متغیره به صورت زیر خواهد بود:

$$S(t_1, t_2) = \left[\frac{1}{S_1(t_1)^{-\alpha} + S_2(t_2)^{-\alpha} - 1} \right]^{\frac{1}{\alpha}}$$

متغیرهای تصادفی و مستقل U و V به صورت زیر تعریف می‌کنیم:

$$U = \frac{q(S_1(T_1))}{q(S_1(T_1)) + q(S_2(T_2))} \quad (1)$$

$$V = S(T_1, T_2) = p[q(S_1(T_1)) + q(S_2(T_2))] \quad (2)$$

جنست و ریوست (۱۹۹۳) ثابت کردند که اگر (T_1, T_2) از مفصل ارشمیدسی با توابع بقای حاشیه‌ای $S_1(t_1)$ و $S_2(t_2)$ پیروی کند، آن‌گاه U دارای توزیع یکنواخت بر فاصله‌ی $[0, 1]$ و V دارای توزیع کندال با تابع چگالی $k(v) = \phi(v) \frac{\phi'(v)}{\phi(v)^2}$ است که بر فاصله‌ی $[0, 1]$ تعریف می‌شود و در آن $\phi(\cdot)$ و $\phi'(\cdot)$ به ترتیب مشتقات مرتبه اول و دوم تابع $\phi(\cdot)$ است.

فرض می‌کنیم بردار (C_1, C_2) نشان‌دهنده‌ی سانسور تصادفی راست با تابع توزیع دلخواه در ازای بردار (T_1, T_2) باشد. در بخش دوم، تابع توزیع شرطی $V = S(T_1, T_2)$ هنگامی که T_1 و T_2 هر دو یا یکی از آن‌ها سانسور شده‌اند، بررسی می‌شود. در بخش سوم برآوردگر مناسبی برای ضریب همبستگی τ کندال در معرض داده‌های سانسور شده معرفی می‌شود. در نهایت در بخش چهارم نتایج حاصل از برآوردگر ضریب همبستگی τ کندال که با استفاده از شبیه‌سازی بدست آمده است، ارائه می‌گردد.

۲ مواد و روش‌ها

در این بخش دو قضیه و یک نتیجه در مورد خواص متغیرهای تصادفی زمان بقا در معرض مشاهدات سانسور شده از راست بیان می‌شود.

فرض کنید (T_1, T_2) یک جفت متغیر تصادفی باشد طوری که تابع توزیع آن‌ها توسط یک مفصل ارشمیدسی تولید شده است. همچنین فرض کنید (T_1, T_2) بوسیله بردار (C_1, C_2) از راست سانسور شده باشد که دارای توزیع توأم پیوسته دلخواهی است. قضیه زیر توزیع شرطی $V = S(T_1, T_2)$ را با فرض زمان‌های بقای T_1 و T_2 ارائه می‌کند.

قضیه ۱: (ونگ و اکس، ۲۰۰۸)

۱. تابع توزیع شرطی $(V | T_1 > c_1, T_2 > c_2)$ برابر است با

$$F_1(v, c_1, c_2) = \frac{1}{S(c_1, c_2)} \left[v - \frac{\phi(v) - \phi(S(c_1, c_2))}{\phi'(v)} \right]; \quad 0 \leq v \leq S(c_1, c_2)$$

۲. تابع توزیع شرطی $(V | T_1 > c_1, T_2 = t_2)$ برابر است با

$$F_2(v, c_1, t_2) = \frac{\phi'(S(c_1, t_2))}{\phi'(v)}; \quad 0 \leq v \leq S(c_1, t_2)$$

۳. تابع توزیع شرطی $(V | T_1 = t_1, T_2 > c_2)$ برابر است با

$$F_3(v, t_1, c_2) = \frac{\phi'(S(t_1, c_2))}{\phi'(v)}; \quad 0 \leq v \leq S(t_1, c_2)$$

: تحت شرایط قضیه ۱ نتایج زیر بدست می آید:

$$i) E(V | T_1 > c_1, T_2 > c_2) = \frac{S(c_1, c_2)}{\gamma} - q[S(c_1, c_2)] \int_0^1 \frac{du}{\dot{q}[uS(c_1, c_2)]} + \int_0^1 \frac{q[uS(c_1, c_2)]}{\dot{q}[uS(c_1, c_2)]} du.$$

$$ii) E(V | T_1 > c_1, T_2 = t_2) = S(c_1, t_2) - S(c_1, t_2) \dot{q}[S(c_1, t_2)] \int_0^1 \frac{du}{\dot{q}[uS(c_1, t_2)]}.$$

$$iii) E(V | T_1 = t_1, T_2 > c_2) = S(t_1, c_2) - S(t_1, c_2) \dot{q}[S(t_1, c_2)] \int_0^1 \frac{du}{\dot{q}[uS(t_1, c_2)]}.$$

برهان :: تحت شرایط قضیه ۱، برای هر $0 < v < S(c_1, c_2)$ در مورد قسمت (i) داریم:

$$\begin{aligned} E(V | T_1 > c_1, T_2 > c_2) &= \int_v^{S(c_1, c_2)} [1 - F(v, c_1, c_2)] dv \\ &= \int_0^{S(c_1, c_2)} dv - \int_0^{S(c_1, c_2)} \frac{v}{S(c_1, c_2)} dv \\ &\quad + \frac{\int_0^{S(c_1, c_2)} \frac{\phi(v)}{\dot{\phi}(v)} dv}{S(c_1, c_2)} \\ &\quad - \frac{\int_0^{S(c_1, c_2)} \frac{\phi(S(c_1, c_2))}{\dot{\phi}(v)} dv}{S(c_1, c_2)} \end{aligned} \quad (3)$$

باتغییر متغیر $u = \frac{v}{S(c_1, c_2)}$ داریم:

$$\begin{aligned} E(V | T_1 > c_1, T_2 > c_2) &= v \Big|_0^{S(c_1, c_2)} - \frac{v^2}{2S(c_1, c_2)} \Big|_0^{S(c_1, c_2)} + \int_0^1 \frac{\phi(uS(c_1, c_2))}{\dot{\phi}(uS(c_1, c_2))} du \\ &\quad - \phi(S(c_1, c_2)) \int_0^1 \frac{du}{\dot{\phi}(uS(c_1, c_2))} \\ &= S(c_1, c_2) - \frac{S(c_1, c_2)}{\gamma} + \int_0^1 \frac{\phi(uS(c_1, c_2))}{\dot{\phi}(uS(c_1, c_2))} du \\ &\quad - \phi(S(c_1, c_2)) \int_0^1 \frac{du}{\dot{\phi}(uS(c_1, c_2))}. \end{aligned}$$

برای اثبات (ii) داریم:

$$\begin{aligned} E(V | T_1 > c_1, T_2 = t_2) &= \int_v^{S(c_1, t_2)} [1 - F(v, c_1, t_2)] dv \\ &= \int_0^{S(c_1, t_2)} dv - \left[\int_0^1 \frac{du}{\dot{\phi}(uS(c_1, t_2))} \right] \dot{\phi}(S(c_1, t_2)) S(c_1, t_2). \end{aligned} \quad (4)$$

برای اثبات (iii) داریم:

$$\begin{aligned}
 E(V | T_1 = t_1, T_2 > c_2) &= \int_v [1 - F(V, t_1, c_2)] dv \\
 &= \int_0^{S(c_1, c_2)} dv - \left[\int_0^1 \frac{du}{\phi(uS(t_1, c_2))} \right] \phi(S(t_1, c_2)) S(t_1, c_2)
 \end{aligned}
 \tag{5}$$

قضیه ۲: (ونگ و اکس، ۲۰۰۸) فرض کنید (T_1, T_2) یک بردار تصادفی باشد که توسط مفصل کلایتون ۱ مدل‌بندی می‌شود. آنگاه بنا به قضیه ۱ و فرع ۱ روابط زیر را حاصل می‌شوند:

$$\begin{aligned}
 i) \quad E(V | T_1 > c_1, T_2 > c_2) &= \frac{S(c_1, c_2)}{2} \left(\frac{\alpha + 1}{\alpha + 2} \right) \\
 ii) \quad E(V | T_1 > c_1, T_2 = t_2) &= S(c_1, t_2) \left(\frac{\alpha + 1}{\alpha + 2} \right) \\
 iii) \quad E(V | T_1 = t_1, T_2 > c_2) &= S(t_1, c_2) \left(\frac{\alpha + 1}{\alpha + 2} \right)
 \end{aligned}
 \tag{6}$$

پس در مدل کلایتون داریم:

$$E(V | T_1 > x_1, T_2 > x_2) = \frac{E(V | T_1 > x_1, T_2 = x_2)}{2}$$

و

$$E(V | T_1 > x_1, T_2 > x_2) = \frac{E(V | T_1 = x_1, T_2 > x_2)}{2}$$

واضح است که وقتی $\alpha = 0$ و T_1 و T_2 مستقل باشند، روابط (۶) به صورت زیر نوشته می‌شوند:

$$\begin{aligned}
 i) \quad E(V | T_1 > c_1, T_2 > c_2) &= \frac{S(c_1, c_2)}{4} = \frac{S_1(c_1)S_2(c_2)}{4} \\
 ii) \quad E(V | T_1 > c_1, T_2 = t_2) &= \frac{S_1(c_1)S_2(t_2)}{2} \\
 iii) \quad E(V | T_1 = t_1, T_2 > c_2) &= \frac{S_1(t_1)S_2(c_2)}{2}
 \end{aligned}$$

بنابراین هنگامی که α افزایش می‌یابد، امید ریاضی‌های شرطی نیز افزایش می‌یابند. دلیل این خاصیت این است که تابع بقا $S(t_1, t_2)$ در مدل مفصل کلایتون یک تابع صعودی از α است. همچنین هنگامی که $\alpha \rightarrow \infty$

$$E(V | T_1 > c_1, T_2 > c_2) \rightarrow \frac{S(c_1, c_2)}{2}$$

و

$$E(V | T_1 > c_1, T_2 = t_2) \rightarrow S(c_1, t_2)$$

و

$$E(V | T_1 = t_1, T_2 > c_2) \rightarrow S(t_1, c_2).$$

۳ کاربرد

یک راه ساده برای برآورد پارامترهای مدل‌های مفصل ارشمیدسی، استفاده از ضریب همبستگی τ کندهال است. این ضریب همبستگی به صورت زیر تعریف می‌شود:

$$\tau = E[\text{sign}\{T_{1i} - T_{1j}\}\text{sign}\{T_{2i} - T_{2j}\}], \quad 0 < \tau < 1$$

که در آن (T_{1i}, T_{2i}) و (T_{1j}, T_{2j}) مشاهدات مستقل از (T_1, T_2) بوده و sign تابع علامت بوده و به صورت زیر تعریف می‌شود:

$$\text{sign}(T_{1i} - T_{1j}) = \begin{cases} 1 & (T_{1i} - T_{1j}) > 0 \\ 0 & (T_{1i} - T_{1j}) = 0 \\ -1 & (T_{1i} - T_{1j}) < 0 \end{cases}$$

اگر δ_1 و δ_2 به ترتیب وضعیت سانسور T_1 و T_2 را نشان دهد، رابطه‌ی بین τ و $V = S(T_1, T_2)$ با استفاده از فرمول زیر بیان می‌شود:

$$\begin{aligned} \tau &= g(\theta) = 4E[V] - 1 \\ &= 4E[V(1 - \delta_1)(1 - \delta_2)] + E[V\delta_1(1 - \delta_2)] + E[V(1 - \delta_1)\delta_2] \\ &\quad + E[V\delta_1\delta_2] - 1 \end{aligned}$$

که در آن $\delta_1 = I(T_1 \leq C_1)$ و $\delta_2 = I(T_2 \leq C_2)$ و با استفاده از روابط (۳)، (۴) و (۵) داریم:

$$\begin{aligned} g(\theta) &= \frac{4}{n} \sum_i [\frac{S(C_{1i}, C_{2i})}{4} - q(S(C_{1i}, C_{2i})) \int_0^1 \frac{du}{\dot{q}(uS(C_{1i}, C_{2i}))} \\ &\quad + \int_0^1 \frac{q(uS(C_{1i}, C_{2i}))}{\dot{q}(uS(C_{1i}, C_{2i}))} du] (1 - \delta_{1i})(1 - \delta_{2i}) \\ &\quad + \frac{4}{n} \sum_i [S(C_{1i}, T_{2i}) - S(C_{1i}, T_{2i}) \dot{q}(S(C_{1i}, T_{2i})) \\ &\quad \times \int_0^1 \frac{du}{\dot{q}(uS(C_{1i}, T_{2i}))}] \delta_{1i}(1 - \delta_{2i}) \\ &\quad + \frac{4}{n} \sum_i [S(T_{1i}, C_{2i}) - S(T_{1i}, C_{2i}) \dot{q}(S(T_{1i}, C_{2i})) \\ &\quad \times \int_0^1 \frac{du}{\dot{q}(uS(T_{1i}, C_{2i}))}] (1 - \delta_{1i})\delta_{2i} \\ &\quad + \frac{4}{n} \sum_i S(T_{1i}, T_{2i}) \delta_{1i}\delta_{2i} - 1 \end{aligned} \tag{7}$$

بنا به نتایج جنست و ریوست (۱۹۹۳) در مدل مفصل کلایتون داریم: $\tau = \frac{\alpha}{\alpha+2}$. پس با تعریف $X_{1i} = \min(T_{1i}, C_{1i})$ و $X_{2i} = \min(T_{2i}, C_{2i})$ و جایگذاری آن‌ها در رابطه (۷) داریم:

$$\begin{aligned} \left(\frac{\alpha}{\alpha+2}\right) + 1 &= \frac{4}{n} \sum_i S(X_{1i}, X_{2i})(1-\delta_{1i})(1-\delta_{2i}) \left[\left(\frac{1}{2}\right) + \frac{\delta_{1i}}{(1-\delta_{1i})} + \frac{\delta_{2i}}{(1-\delta_{2i})} + \frac{\delta_{1i}\delta_{2i}}{(1-\delta_{1i})(1-\delta_{2i})} \left(\frac{\alpha+2}{\alpha+1}\right) \right] \\ &= \frac{4}{n} \left(\frac{\alpha+1}{\alpha+2}\right) \sum_i S(X_{1i}, X_{2i})(1-\delta_{1i})(1-\delta_{2i}) \left[\frac{(\alpha+1) + (\alpha+1)\delta_{1i} + (\alpha+1)\delta_{2i} + (1-\alpha)\delta_{1i}\delta_{2i}}{2(1-\delta_{1i})(1-\delta_{2i})(\alpha+1)} \right] \end{aligned}$$

در نتیجه با کمی محاسبات ساده داریم:

$$\sum_i S(X_{1i}, X_{2i}) \left[1 + \delta_{1i} + \delta_{2i} + \frac{1-\alpha}{1+\alpha} \delta_{1i}\delta_{2i} \right] - n = 0 \quad (7)$$

حال اگر در رابطه‌ی (7) به جای α از برآورد آن به صورت $\hat{\alpha}_n = \frac{2\hat{\tau}_n}{1-\hat{\tau}_n}$ استفاده کنیم، می‌توانیم برآوردی برای ضریب همبستگی τ ی‌کنندال به دست آوریم زمانی که داده‌ی سانسور شده داریم. با استفاده از برآورد تابع بقای دو متغیره نیز که توسط دابروسکا (۱۹۸۸) بیان شده است، داریم:

$$\hat{\tau}_n = \frac{\sum_i \hat{S}(X_{1i}, X_{2i})(1 + \delta_{1i} + \delta_{2i} + \delta_{1i}\delta_{2i}) - n}{\sum_i \hat{S}(X_{1i}, X_{2i})(-1 - \delta_{1i} - \delta_{2i} + 3\delta_{1i}\delta_{2i}) + n}$$

بدیهی است هنگامی که هیچ داده‌ی سانسور شده‌ای نداشته باشیم (یعنی به ازای هر i ، $\delta_{1i} = \delta_{2i} = 1$)، خواهیم داشت:

$$\hat{\tau}_n = \frac{4}{n} \sum_{i=1}^n \hat{S}(T_{1i}, T_{2i}) - 1 = \frac{4}{n} \sum_{i=1}^n \hat{V}_i - 1$$

۴ بحث و نتیجه‌گیری

در این مقاله با استفاده از قضیه ۱ و نتیجه‌ی حاصل از آن و مدل مفصل کلایتون یک برآورد گشتاوری برای τ ی‌کنندال زمانی که داده‌ی سانسور شده داریم، ارائه شد. همچنین با استفاده از قضیه ۱ می‌توان پارامتر مجهول مدل‌های مفصل ارشمیدسی را برآورد کرد. در مثال زیر کاربرد برآوردگر ضریب همبستگی ی‌کنندال را زمانی که با مشاهدات سانسور شده‌ی تصادفی از راست سرو کار داریم، توضیح می‌دهیم.

مثال ۱: در یک مطالعه شبیه‌سازی شده، ابتدا یک نمونه‌ی تصادفی ۱۰۰ تایی از مدل مفصل کلایتون با پارامترهای $\alpha = 0.5, 1.33, 3, 8$ ($\tau = 0.2, 0.4, 0.6, 0.8$)، بردارهای زمان را تولید می‌کنیم. بردارهای سانسور را نیز از توزیع یکنواخت در بازه (۰، ۱) تولید می‌کنیم. پس از تولید داده‌ها، مشاهده می‌کنیم که ۲۰ درصد آن‌ها سانسور شده‌اند. برای ارزیابی دقت برآوردگر ضریب همبستگی τ ی‌کنندال در حالت داده‌های سانسور شده، اریبی و انحراف معیار برآوردها را به صورت تجربی با تکرار مستقل فرآیند شبیه‌سازی به اندازه‌ی ۱۰۰۰ بار محاسبه می‌کنیم. جدول ۱ خلاصه‌ی نتایج شبیه‌سازی با نرم افزار R را نشان می‌دهد.

جدول ۱: برآورد اریبی و انحراف معیار برآوردهای ضریب همبستگی τ ی کندال

| α | ۰/۵ | ۱/۳۳ | ۳ | ۸ |
|-------------------------|---------|---------|---------|--------|
| اریبی برآورد شده | -۰/۷۵۹۸ | -۰/۵۸۴۸ | -۰/۳۹۲۴ | -۰/۱۹۷ |
| انحراف معیار برآورد شده | ۰/۰۲۷۴ | ۰/۰۰۹۲ | ۰/۰۰۴۹ | ۰/۰۰۱۷ |

با مقایسه داده‌های جدول ۱ مشاهده می‌کنیم که با افزایش α ، قدر مطلق اریبی و انحراف معیار برآوردهای ضریب همبستگی کندال کاهش می‌یابد.

- Bandeem-Roche, K. J., Liang, K.Y. (1996). Modelling failure-time associations in data with multiple levels of clustering. *Biometrika.*, **83**, 29–39.
- Dabrowska, D.M., (1988). Kaplan - Meier estimate on the plane. *Ann. Statist* , **16**, 1475–1489.
- Genest, C., Rivest, L.P., (1993). Statistical inference procedures for bivariate Archimedean copulas. *J.Amer. Statist. Assoc.*, **88**, 1034–1043.
- Hougaard, P., (2000). Analysis of Multivariate Survival Data. Springer-Verlag, New York, Inc.
- Wang, A., Oakes, D., (2008). Some properties of the Kendall distribution in bivariate Archimedean copula models under censoring. *Statistics. Probability.*, **78**, 2578–2583.
- Wang, J., Zafra, P., (2009). Estimating Bivariate Survival Function by Volterra Estimator Using Dynamic Programming Techniques. *Journal of Data Science.*, **7**, 365–380.
- Wang, W., Wells, M.T., (2000a). Model selection and semiparametric inference for bivariate failure-time data. *J. Amer. Statist. Assoc.*, **449**, 62–72.
- Wang, W., Wells, M.T., (2000b). Estimation of kendalls tau under censoring. *Statist. Sinica.*, **10**, 1199–1218.

Archimedean Copula, Clayton Copula, Kendall τ

Jabbari, H.¹, Mokhtari, A.², and Azarnoosh, H. A.³

Department of Statistics, Ferdowsi University of Mashhad, mashhad, Iran.

Department of Statistics, Azad Islamic University of Mashhad, mashhad, Iran.

Department of Statistics, Azad Islamic University of Mashhad, mashhad, Iran.

Abstract: Survival times vector (T_1, T_2) can be Modelled by an Archimedean Copula Model.

In this paper by fitting clayton copula for data and kendall correlation coefficients for censored data and We obtained, a simple Moment estimator of the dependence parameter in Archimedean Copula Models is proposed.

Keywords: Archimedean Copula, Clayton Copula, Kendall τ .

Mathematics Subject Classification (2000): 62N02