

استنباط آماری بر روی داده های سانسور شده میانی با توزیع طول عمر نمایی

هنگامه حبیبی راد^۱، آرزو حبیبی راد و ابوالقاسم بزرگ نیا^۲

چکیده

در بعضی از مطالعات مربوط به داده های بقا این امکان وجود دارد که بعضی از واحدهای تحت آزمایش، به طور موقتی از مطالعه خارج شده باشند (برای مثال خروج موقتی یک فرد از شهر برای یک مدت نامشخص و برگشت مجدد او به شهر)، در این صورت با سانسور جدیدی به نام سانسور میانی مواجه می شویم که اولین بار توسط جامالاماداکا و مانگالام [۴] معرفی شده است. فرض کنید n واحد مشابه با زمان های بقای T_1, T_2, \dots, T_n تحت یک آزمایش بقا قرار گرفته باشند و برای i امین واحد یک بازه سانسور تصادفی (L_i, R_i) وجود داشته باشد، به طوری که زمان T_i برای i امین واحد در صورتی ثبت می شود که $T_i \notin [L_i, R_i]$ و در غیر این صورت سانسور می شود. یعنی اگر طول عمر یک واحد تحت آزمایش، در یک بازه تصادفی قرار گیرد، سانسور میانی رخ می دهد.

به دنبال ارائه این نوع سانسور آیر^۳ و همکاران [۳] تحلیل داده های سانسور شده میانی را با فرض نمایی بودن توزیع طول عمر مورد بررسی قرار داده اند. در این مقاله پس از معرفی کامل طرح سانسور میانی، ابتدا برآورد درستنمایی ماکسیمم آن را بررسی می کنیم و سپس به بیان نتایج تئوری و فضایی مربوطه می پردازیم.

واژه های کلیدی: برآورد درستنمایی ماکسیمم، توزیع طول عمر، توزیع نمایی، سانسور میانی.

۱- مقدمه

وقتیکه $\delta_i = 1$ ، مشاهدات سانسور نشده اند و مقدار واقعی T_i مشاهده می شود، در این شرایط ما (L_i, R_i) را مشاهده نمی کنیم. در شرایط دیگر، زمانی که $\delta_i = 0$ ، فقط فاصله سانسور شده $[L_i, R_i]$ را مشاهده می کنیم، در این صورت مشاهده سانسور شده است. برای i امین واحد داریم:

$$(Y_i, \delta_i) = \begin{cases} (T_i, 1), & T_i \notin [L_i, R_i], \\ ([L_i, R_i], 0), & o.w. \end{cases} \quad (1)$$

بنابراین، داده هایی که در حالت سانسور نشده بدست آمده اند با داده هایی که در شرایط فاصله سانسور شده به دست آمده اند یکسان نیستند. بر اساس مشاهدات، بدست آوردن تابع توزیع طول عمر T_i و گسترش روش های استنتاجی مورد نیاز مشکل به نظر می رسد.

در ابتدا به بیان چند مثال که در آن ها ممکن است سانسور میانی رخ دهد می پردازیم. فرض کنید در یک مطالعه طول عمر، موضوع مورد بررسی به طور موقت از مطالعه خارج

در این مقاله با داده هایی سروکار داریم که دارای سانسور میانی هستند، سانسور میانی زمانی رخ می دهد که داده در فاصله تصادفی قرار بگیرد و مقدار واقعی آن، مشاهده نشود.

طرح سانسور میانی به صورت زیر تعریف می شود:

فرض کنید n واحد یکسان مورد آزمایش قرار می گیرند و طول عمر این واحدها، T_1, T_2, \dots, T_n هستند. برای i امین واحد، T_i فقط زمانی مشاهده می شود که $T_i \notin [L_i, R_i]$ ، در غیر این صورت مشاهده نمی شود. فرض کنید $\delta_i = I(T_i \notin [L_i, R_i])$ ، که در آن $I(\cdot)$ تابع نشانگر در نظر گرفته شده است، یعنی

^۱ دانشجوی کارشناسی ارشد دانشگاه آزاد اسلامی مشهد

^۲ عضو هیئت علمی گروه آمار، دانشگاه فردوسی مشهد و دانشگاه آزاد

اسلامی مشهد

^۳ Jammalamadaka and Mangalam

^۴ Iyer

که فرض های متفاوتی برای طرح های سانسور ساخته شده، مراجعه کنید. به طور کلی دلایل مختلفی برای استقلال وجود دارد که به بعضی از آن ها اشاره می کنیم. در ابتدا، در بیشتر شرایطی که در واقعیت رخ می دهد احتمال بستگی داشتن فرایند سانسور شدن به طول عمر جامعه مورد نظر کم است و این مطلب اطلاعاتی در مورد تابع توزیع جامعه به ما می دهد. در مثال هایی که قبلاً بیان کردیم، سانسور میانی به دلایل خارجی رخ می دهد، که هیچ ارتباطی با متغیر طول عمر ندارد.

در ادامه مقاله، ابتدا در بخش دوم دو روش تکرار عددی و الگوریتم EM را برای برآورد پارامتر مجهول تابع درستنمایی ماکسیمم معرفی کرده و سپس در بخش سوم با کمک چند قضیه و چند لم خصوصیات مربوط به تابع درستنمایی معرفی شده در بخش دوم را، ثابت می کنیم و در ادامه در بخش چهارم نتایج گرفته شده از این مقاله بیان شده است.

۲- برآورد درستنمایی ماکسیمم برای توزیع طول عمر نمایی

ابتدا داده ها را با فرض نداشتن داده گم شده مرتب می کنیم، به این صورت که n_1 جمله اول را مشاهدات سانسور نشده و n_2 جمله بعدی را مشاهدات سانسور شده در نظر می گیریم. بنابراین مشاهدات به صورت زیر خواهد بود:

$$(3) \quad (T_1, 1), \dots, (T_{n_1}, 1), (L_{n_1+1}, R_{n_1+1}), \dots, (L_{n_1+n_2}, R_{n_1+n_2}),$$

که در آن $n_1 + n_2 = n$. بنابراین، برای n_1 مشاهده اول $T_i \in [L_i, R_i]$ و برای n_2 مشاهده بعدی $T_i \in [L_i, R_i]$.

لم ۱-۲: تابع درستنمایی داده های سانسور شده میانی به صورت زیر است

$$(4) \quad l(\theta) = c \theta^{n_1} e^{-\theta \sum_{i=1}^{n_1} t_i} \prod_{i=n_1+1}^{n_1+n_2} (e^{-\theta l_i} - e^{-\theta r_i}),$$

که در آن c یک مقدار ثابت بوده و تنها به α و β بستگی دارد. چون علاقه ای به برآورد α و β نداریم، در مورد آن بحث نمی کنیم.

اثبات: با توجه به تعریف سانسور میانی داریم

$$l(\theta) = \prod_{i=1}^{n_1} f_{T_i, \delta}(t_i, 1) \prod_{i=n_1+1}^{n_1+n_2} f_{L_i, Z_i, \delta}(l_i, z_i, 0),$$

از طرفی

شود (مثلاً، اگر فردی به طور موقت در یک دوره زمانی شهر را ترک کند و برگردد، در حالیکه هنوز زنده است) سانسور میانی رخ می دهد. همچنین سانسور میانی زمانی رخ می دهد که مثلاً وسیله اندازه گیری در یک دوره زمانی خراب شود یا اگر یک کلینیک که در آن مشاهدات جمع آوری می شود در یک دوره به دلیل شرایط ناگهانی مانند وقوع یک جنگ یا یک حادثه بسته شود. در چنین شرایطی، رخداد مورد نظر در مدت زمانی که یک مورد رخ نداده و یا مشاهده نمی شود، اتفاق می افتد.

در مقاله جامالاماداکا و مانگالام [۴] T_1, T_2, \dots, T_n متغیرهای تصادفی مستقل و هم توزیع با تابع توزیع نامعلوم $F(\cdot)$ در نظر گرفته شده اند، و $(L_1, R_1), \dots, (L_n, R_n)$ مستقل و هم توزیع هستند با تابع توزیع دو متغیره نامعلوم $G(\cdot, \cdot)$ و مستقل از T_i . البته در مقاله ذکر شده مسئله به صورت ناپارامتری مورد بررسی قرار گرفته است ولی در این مقاله فرمول های پارامتری مورد بررسی قرار می گیرند. برای این منظور فرض کنید که T_1, T_2, \dots, T_n متغیرهای تصادفی مستقل و هم توزیع نمایی با میانگین مثلاً $1/\theta$ با تابع چگالی (pdf) زیر باشد:

$$(2) \quad f(x; \theta) = \begin{cases} \theta e^{-\theta x} & , x > 0, \\ 0 & , o.w, \end{cases}$$

علاوه بر این، $(L_1, Z_1), \dots, (L_n, Z_n)$ نیز مستقل و هم توزیع هستند که L_i و $Z_i = R_i - L_i$ متغیرهای تصادفی مستقل نمایی و مستقل از T_i هستند. هم چنین فرض می شود که L_i و Z_i به ترتیب دارای میانگین های $1/\alpha$ و $1/\beta$ است و به θ بستگی ندارد. البته باید در نظر داشت فرایند سانسور شدن، مستقل از طول عمر جامعه مورد نظر است و هیچ اطلاعاتی در مورد طول عمر نمی دهد.

فرض هایی که در مورد استقلال، در تحلیل داده های طول عمر می شود کاملاً استاندارد هستند. برای مثال به کاپلان و میر [۶]، بابو و همکاران [۱]، جامالاماداکا و مانگالام [۴]، چو و تیچر [۲]، جامالاماداکا و آیر [۵] و مراجع ذکر شده در آن،

¹ Kaplan and Meier

² Babu, Rao and Rao

³ Chow and Teicher

⁴ Jammalamadaka and Iyer

زمانی متوقف می شویم که $|\theta^{(i)} - \theta^{(i+1)}| < \varepsilon$ و در آن ε مقدار کوچک مثبتی است که از قبل تعیین می شود. برای انتخاب اول θ می توان از $\theta^{(i)} = \frac{n_1}{\sum_{i=1}^n t_i}$ استفاده کرد. یکی از روش

های دیگر برای بدست آوردن MLE استفاده از الگوریتم EM است، که در این روش ابتدا مقدار $E(T|L < T < R)$ را بدست می آوریم، که در آن L و R مقادیر ثابت هستند و T دارای توزیع نمایی با میانگین $1/\theta$ است،

$$E(T|L < T < R) = \int_L^R f_{T|T \in (L,R)}(t|L < t < R) dt,$$

برای محاسبه رابطه فوق ابتدا θ را ثابت می کنیم.

لم ۲۰۲: تابع چگالی T ، به شرط آنکه $T \in (L, R)$ به صورت زیر است

$$f_{T|T \in (L,R)}(t|L < t < R) = \frac{\theta e^{-\theta t}}{e^{-\theta L} - e^{-\theta R}}.$$

اثبات: می دانیم

$$f_{T|T \in (L,R)}(t|L < t < R) = \frac{P(T = t, T \in (L, R))}{P(T \in (L, R))} = \frac{f_T(t)}{P(T \in (L, R))},$$

که در آن $f_T(t) = \theta e^{-\theta t}$ و

$$P(T \in (L, R)) = P(L < T < R) = e^{-\theta L} - e^{-\theta R}.$$

در نتیجه θ ثابت می شود. با کمک θ فوق داریم

$$E(T|L < T < R) = \int_L^R t f_{T|T \in (L,R)}(t|L < t < R) dt = \int_L^R \frac{t \theta e^{-\theta t}}{e^{-\theta L} - e^{-\theta R}} dt,$$

پس

$$E(T|L < T < R) = \frac{e^{-\theta L} \left(L + \frac{1}{\theta} \right)}{e^{-\theta L} - e^{-\theta R}} - \frac{e^{-\theta R} \left(R + \frac{1}{\theta} \right)}{e^{-\theta L} - e^{-\theta R}}. \quad (10)$$

توجه کنید که فرمول (۱۰) در محاسبه الگوریتم EM (برای آشنایی بیشتر با این الگوریتم می توانید به کتاب *Modeler's Approach* نوشته شده توسط جیمز تامپسون [۷] مراجعه کنید) استفاده می شود. در این صورت تابع شبه درستنمایی به شکل زیر در می آید:

$$f_{T_i, \delta}(t_i, 1) = P(\delta = 1 | T_i = t_i) f_{T_i}(t_i), \quad (5)$$

که در آن $f_{T_i}(t_i) = \theta e^{-\theta t_i}$

$$P(\delta = 1 | T_i = t_i) = P(T_i \notin (L_i, R_i) | T_i = t_i) = 1 - \frac{\alpha e^{-\beta t_i}}{\alpha - \beta} (1 - e^{-(\alpha - \beta)t_i}),$$

همچنین

$$f_{L_i, Z_i, \delta}(l_i, z_i, 0) = P(\delta = 0 | L_i = l_i, Z_i = z_i) = e^{-\theta l_i} - e^{-\theta r_i}. \quad (6)$$

در انتها با جایگذاری روابط (۵) و (۶) در فرمول (۴) لم فوق ثابت می شود.

در ادامه به دنبال محاسبه MLE برای پارامتر θ هستیم، که بر اساس فرمول (۴) لگاریتم درستنمایی به صورت زیر می باشد:

$$\ln l(\theta) = L(\theta) = \ln c + n_1 \ln \theta - \theta \sum_{i=1}^{n_1} t_i + \sum_{i=n_1+1}^{n_1+n_2} \ln(e^{-\theta l_i} - e^{-\theta r_i}). \quad (7)$$

با مشتق گیری از $L(\theta)$ و برابر صفر قرار دادن، داریم

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{n_1}{\theta} - \sum_{i=1}^{n_1} t_i + \sum_{i=n_1+1}^{n_1+n_2} \frac{(r_i - l_i)}{e^{\theta(r_i - l_i)} - 1} - \sum_{i=n_1+1}^{n_1+n_2} l_i = 0. \quad (8)$$

بنابراین $\hat{\theta}$ که از حل معادله (۸) حاصل می شود، یک MLE برای پارامتر θ است. ولی با استفاده از معادله (۸) نمی توان $\hat{\theta}$ را مستقیماً و به طور دقیق بدست آورد، لذا برای بدست آوردن MLE از روش های تکراری استفاده می کنیم. توجه کنید که برای این منظور رابطه (۸) را می توان به صورت زیر نوشت

$$h(\theta) = \theta, \quad (9)$$

که با استفاده از تعریف $z_i = r_i - l_i$ و با کمک رابطه (۸) واضح است که

$$h(\theta) = \frac{1}{\sum_{i=n_1+1}^{n_1+n_2} l_i + \sum_{i=1}^{n_1} t_i} \left[n_1 + \theta \sum_{i=n_1+1}^{n_1+n_2} \frac{z_i e^{-\theta z_i}}{1 - e^{-\theta z_i}} \right].$$

بنابراین یک روش ساده برای بدست آوردن MLE استفاده از روش تکرار برای معادله (۹) است. برای مثال می توانیم با مقدار اولیه $\theta^{(0)}$ شروع کنیم، سپس $\theta^{(1)} = h(\theta^{(0)})$ بدست می آید و به همین ترتیب ادامه می دهیم. در این روش تکراری

$$\frac{1}{n}L(\theta) = c' + \frac{n_1}{n} \ln \theta - \frac{\theta}{n} \sum_{i=1}^{n_1} t_i - \frac{\theta}{n} \sum_{i=n_1+1}^{n_1+n_2} l_i + \frac{1}{n} \sum_{i=n_1+1}^{n_1+n_2} \ln(1 - e^{-\theta t_i})$$

به سادگی داریم

$$\frac{1}{n}L(\theta) \rightarrow g(\theta) \text{ a.s.}$$

که در آن

$$g(\theta) = c' + p(\theta) \ln \theta + \theta \frac{(1-p(\theta))(\alpha + \beta + \gamma \theta)}{(\alpha + \theta)(\beta + \theta)} - \theta \frac{1}{\theta} - \theta \frac{(1-p(\theta))}{\alpha + \theta} - \frac{\alpha \beta}{\alpha + \theta} \left[\sum_{i=1}^{\infty} \frac{1}{i(\beta + i\theta)} - \sum_{i=1}^{\infty} \frac{1}{i(\beta + i\theta + \theta)} \right]$$

و

$$p(\theta) = \frac{\alpha \beta + \beta \theta + \theta^2}{(\alpha + \theta)(\beta + \theta)}, \quad c' = \frac{1}{n} \ln c, \quad (14)$$

اثبات: ابتدا تابع چگالی T را به شرط آن که $T \notin (L, R)$

$\alpha \neq \beta$ ، به صورت زیر محاسبه می کنیم. با فرض

$$\delta = \begin{cases} 1 & T \notin (L, R) \\ 0 & T \in (L, R) \end{cases}$$

داریم

$$f_{T|\delta}(t|\delta) = \frac{f_{T,\delta}(t)}{P(\delta=1)} = \frac{P(\delta=1|T=t)f_T(t)}{p(\theta)}$$

که در آن

$$P(\delta=1|T=t) = P(T \notin (L, R) | T=t) = 1 - P(L < t < R)$$

از طرفی می دانیم $R=L+Z$ ، پس

$$P(L < t < R) = \int_{t-L}^{\infty} \int_{t-L}^{\infty} f(l,z) dz dl = \frac{\alpha e^{-\beta t}}{\alpha - \beta} (1 - e^{-(\alpha - \beta)t})$$

و

$$P(\delta=1|T=t) = 1 - \frac{\alpha e^{-\beta t}}{\alpha - \beta} (1 - e^{-(\alpha - \beta)t})$$

در ادامه نشان می دهیم که $p(\theta) = \frac{\alpha \beta + \beta \theta + \theta^2}{(\alpha + \theta)(\beta + \theta)}$ ، برای

این منظور داریم

$$P_\theta(T \notin (L, R)) = 1 - P_\theta(T \in (L, R)) = 1 - [P_\theta(T < L + Z) - P_\theta(T < L)],$$

از طرفی

$$l(\theta) = \theta^{n_1+n_2} e^{-\theta \left(\sum_{i=1}^{n_1} T_i + \sum_{i=n_1+1}^{n_1+n_2} T_i^{(s)} \right)}, \quad (11)$$

که در آن

$$T_i^{(s)} = \frac{e^{-\theta t_i} \left(L_i + \frac{1}{\theta} \right) - e^{-\theta R_i} \left(R_i + \frac{1}{\theta} \right)}{e^{-\theta L_i} - e^{-\theta R_i}}. \quad (12)$$

مراحل الگوریتم EM عبارتند از:

مرحله ۱: فرض می کنیم θ_j ، z امین تکرار $\hat{\theta}$ باشد.

مرحله ۲: $T_{i(j)}^{(s)}$ را با استفاده از فرمول (۱۲) محاسبه می کنیم و

θ را توسط $\theta_{(j)}$ جایگزین می کنیم.

$$\text{مرحله ۳: } \theta_{(j+1)} = \frac{n_1 + n_2}{\sum_{i=1}^{n_1} T_i + \sum_{i=n_1+1}^{n_1+n_2} T_{i(j)}^{(s)}}$$

۳- نتایج نظری

قضیه ۱-۳: روش های تکراری مطرح شده همگرا هستند اگر

$$\sum_{i=n_1+1}^{n_1+n_2} r_i \leq 2 \sum_{i=1}^{n_1} t_i + 3 \sum_{i=n_1+1}^{n_1+n_2} l_i. \quad (13)$$

اثبات: با مشتق گیری از $h(\theta)$ داریم

$$\frac{\partial h(\theta)}{\partial \theta} = h'(\theta) = \frac{1}{\sum_{i=n_1+1}^{n_1+n_2} l_i + \sum_{i=1}^{n_1} t_i} \left[\sum_{i=n_1+1}^{n_1+n_2} \frac{z_i e^{-\theta z_i}}{1 - e^{-\theta z_i}} + \theta \sum_{i=n_1+1}^{n_1+n_2} \frac{-z_i^2 e^{-\theta z_i} (1 - e^{-\theta z_i}) - (z_i e^{-\theta z_i})^2}{(1 - e^{-\theta z_i})^2} \right]$$

پس

$$|h'(\theta)| = \frac{1}{\sum_{i=n_1+1}^{n_1+n_2} l_i + \sum_{i=1}^{n_1} t_i} \times \left| \sum_{i=n_1+1}^{n_1+n_2} \frac{z_i e^{-\theta z_i} (1 - e^{-\theta z_i} - \theta z_i)}{(1 - e^{-\theta z_i})^2} \right|$$

به سادگی می توان نشان داد که برای هر $x \geq 0$

$$\frac{|e^{-x} (1 - e^{-x} - x)|}{|1 - e^{-x}|} \leq \frac{1}{2}$$

بنابراین

$$|h'(\theta)| \leq \frac{\sum_{i=n_1+1}^{n_1+n_2} z_i}{\sum_{i=n_1+1}^{n_1+n_2} l_i + \sum_{i=1}^{n_1} t_i}$$

از طرفی می دانیم اگر روش های تکراری همگرا باشند

$|h'(\theta)| \leq 1$ (اسماعیلی [۹])، در این صورت رابطه (۱۳) برقرار

می باشد.

لم ۱۰۳: توجه کنید که بنابر رابطه

$$P_{\theta}(T \in (L, R)) = 1 - p(\theta),$$

و $f_L(l) = \alpha e^{-\alpha l}$ هم چنین

$$P(T \in (L, R) | L = l) = P(l < T < l + Z) \\ = P(T - Z < l) - P(T < l),$$

$$P(T - Z < l) = \int_0^{l+Z} \int_0^t (\beta e^{-\beta z} \theta e^{-\theta t}) dt dz \\ = 1 - \frac{\beta e^{-\theta l}}{(\beta + \theta)},$$

و

$$P(T < l) = \int_0^l \theta e^{-\theta t} dt = 1 - e^{-\theta l},$$

در نتیجه برای $l > 0$ داریم

$$f_{L|T \in (L, R)}(l) = \frac{1}{1 - p(\theta)} \times \frac{\alpha \theta}{(\beta + \theta)} e^{-(\alpha + \theta)l}.$$

در این صورت

$$E(L | T \in (L, R)) = \int_0^{\infty} l f_{L|T \in (L, R)}(l) dl \\ = \frac{1}{1 - p(\theta)} \times \frac{\alpha \theta}{(\beta + \theta)(\alpha + \theta)}.$$

همچنین تابع چگالی $Z = R - L$ را به شرط آن که $T \in (L, R)$ به صورت زیر بدست می آوریم

$$f_{Z|T \in (L, R)}(z) = \frac{P(Z = z, T \in (L, R))}{P(T \in (L, R))} \\ = \frac{P(T \in (L, R) | Z = z) f_Z(z)}{P(T \in (L, R))},$$

که در آن $f_Z(z) = \beta e^{-\beta z}$ و

$$P_{\theta}(T \in (L, R)) = 1 - p(\theta),$$

و

$$P(T \in (L, R) | Z = z) = P(L < T < L + z) \\ = P(T - L < z) - P(T < L),$$

که در رابطه فوق

$$P(T - L < z) = \int_0^{l+Z} \int_0^t (\alpha e^{-\alpha l} \theta e^{-\theta t}) dt dl \\ = 1 - \frac{\alpha e^{-\theta z}}{(\alpha + \theta)},$$

و

$$P(T < L) = \int_0^l \int_0^t \alpha e^{-\alpha l} \theta e^{-\theta t} dt dl \\ = \frac{\theta}{(\alpha + \theta)}.$$

در نتیجه برای $z > 0$ داریم

$$f_{Z|T \in (L, R)}(z) = \frac{1}{1 - p(\theta)} \times \frac{\alpha \beta e^{-\beta z}}{(\alpha + \theta)} (1 - e^{-\theta z}).$$

$$P_{\theta}(T < L + Z) = \int_0^{\infty} \int_0^{l+Z} \int_0^t (\alpha e^{-\alpha l} \beta e^{-\beta z} \\ \times \theta e^{-\theta t}) dt dl dz \\ = \frac{\alpha \theta + \beta \theta + \theta^2}{(\alpha + \theta)(\beta + \theta)},$$

و

$$P_{\theta}(T < L) = \int_0^l \int_0^t (\alpha e^{-\alpha l} \theta e^{-\theta t}) dt dl \\ = \frac{\theta}{(\alpha + \theta)},$$

در نتیجه

$$p(\theta) = P_{\theta}(T \notin (L, R)) = \frac{\alpha \beta + \beta \theta + \theta^2}{(\alpha + \theta)(\beta + \theta)}.$$

پس با فرض $\theta = \theta$ داریم

$$f_{T|T \notin (L, R)}(t) = \frac{1}{p(\theta)} \left\{ \theta e^{-\theta t} \left(1 - \frac{\alpha e^{-\beta t}}{\alpha - \beta} (1 - e^{-(\alpha - \beta)t}) \right) \right\}. \quad (15)$$

و به همین ترتیب در حالت $\alpha = \beta$ داریم

$$P(L < t < L + Z) = \int_0^t \int_0^{\infty} f(l, z) dz dl \\ = \alpha t e^{-\alpha t},$$

و

$$P(\delta = 1 | T = t) = 1 - \alpha t e^{-\alpha t},$$

در نتیجه

$$f_{T|T \in (L, R)}(t) = \frac{1}{p(\theta)} \left\{ \theta e^{-\theta t} (1 - \alpha t e^{-\alpha t}) \right\}. \quad (16)$$

با کمک فرمول های (15) و (16) داریم

$$E(T | T \notin (L, R)) = \int_0^{\infty} t f_{T|T \notin (L, R)}(t) dt \\ = \begin{cases} \frac{1}{p(\theta)} \left[\frac{1}{\theta} - \frac{2\alpha\theta}{(\alpha + \theta)^2} \right], & \alpha = \beta \\ \frac{1}{p(\theta)} \left[\frac{1}{\theta} - \frac{\alpha\theta}{(\alpha - \beta)(\beta + \theta)} - \frac{1}{(\alpha + \theta)} \right], & \alpha \neq \beta \end{cases}$$

در ادامه تابع چگالی L را به شرط آن که $T \in (L, R)$ باشد محاسبه می کنیم،

$$f_{L|T \in (L, R)}(l) = \frac{P(L = l, T \in (L, R))}{P(T \in (L, R))} \\ = \frac{P(T \in (L, R) | L = l) f_L(l)}{P(T \in (L, R))},$$

که با فرض $\theta = \theta$ داریم

$$\frac{\sum_{i=1}^{n_1} T_i}{n_1} \rightarrow E(T_i | T_i \notin (L_i, R_i)) \quad a.s.$$

و

$$\frac{\sum_{i=n_1+1}^{n_1+n_2} L_i}{n_1} \rightarrow E(L_i | T_i \in (L_i, R_i)) \quad a.s.$$

$$\frac{\sum_{i=n_1+1}^{n_1+n_2} \ln(1-e^{-\theta L_i})}{n_1} \rightarrow E(\ln(1-e^{-\theta L_i}) | T_i \in (L_i, R_i)) \quad a.s.$$

پس می توان به راحتی نشان داد که $\frac{1}{n}L(\theta) \rightarrow g(\theta) \quad a.s.$

لم ۳-۲: $g(\theta)$ یک تابع تک نمایی با ماکسیمم واحد است.

اثبات: با توجه به تعریف تابع $g(\theta)$ در رابطه (۱۴) داریم

$$g'(\theta) = \frac{p(\theta)}{\theta} + \frac{(1-p(\theta))(\alpha + \beta + \theta)}{(\alpha + \theta)(\beta + \theta)} - \frac{1}{\theta} - \frac{(1-p(\theta))}{\alpha + \theta} - \frac{\alpha\beta}{\alpha + \theta} \left[- \sum_{i=1}^{\infty} \frac{1}{(\beta + i\theta)} + \sum_{i=1}^{\infty} \frac{1}{(\beta + i\theta + \theta)} \right]$$

در نتیجه

$$g'(\cdot) = \infty$$

و

$$g'(\infty) = -\frac{\beta + \theta \cdot p(\theta)}{\theta \cdot (\beta + \theta)} < 0.$$

و

$$g''(\theta) = -\frac{p(\theta)}{\theta^2} - \frac{\alpha\beta}{\alpha + \theta} \times \left[\sum_{i=1}^{\infty} \frac{\gamma i}{(\beta + i\theta)^2} + \sum_{i=1}^{\infty} \frac{\gamma i}{(\beta + i\theta + \theta)^2} \right] < 0.$$

روابط فوق نشان می دهند که $g(\theta)$ فقط دارای یک ریشه است و چون مشتق دوم همواره منفی است پس آن ریشه یک نقطه ماکسیمم است. پس $g(\theta)$ یک تابع تک نمایی با یک نقطه ماکسیمم است.

لم ۳-۳: اگر $\hat{\theta}$ برآورد درستنمایی θ باشد، آنگاه $\hat{\theta} \rightarrow \theta^*$

که در آن θ^* جواب یکتای معادله غیر خطی زیر است

$$g'(\theta) = 0 \quad (17)$$

که $p(\theta)$ در (۱۴) تعریف شده است.

اثبات: بنابر لم ۲، رابطه (۱۷) تنها یک جواب مانند θ^* دارد زیرا $g'(\theta)$ محور x ها را تنها در یک نقطه قطع می کند و

از طرفی امید تابع $\ln(1-e^{-\theta})$ به شرط $T \in (L, R)$ به صورت زیر محاسبه می شود

$$E(\ln(1-e^{-\theta z}) | T \in (L, R)) = \int_{L_i}^{\infty} \ln(1-e^{-\theta z}) \times f_{z|T \in (L, R)}(z) dz = \frac{1}{1-p(\theta)} \times \frac{\alpha\beta}{(\alpha + \theta)} \times \left[\int_{L_i}^{\infty} e^{-\beta z} \ln(1-e^{-\theta z}) dz - \int_{L_i}^{\infty} e^{-(\beta+\theta)z} \ln(1-e^{-\theta z}) dz \right].$$

برای محاسبه مقدار انتگرال بالا از بسط $\ln(1-e^{-\theta})$ (عالم زاده [۱۰]) استفاده می کنیم

$$\ln(1-e^{-\theta z}) = \sum_{i=1}^{\infty} \frac{-e^{-i\theta z}}{i}, \quad |e^{-\theta z}| < 1$$

در نتیجه

$$\int_{L_i}^{\infty} e^{-\beta z} \ln(1-e^{-\theta z}) dz = \int_{L_i}^{\infty} e^{-\beta z} \sum_{i=1}^{\infty} \frac{-e^{-i\theta z}}{i} dz = - \sum_{i=1}^{\infty} \frac{1}{i(\beta + i\theta)},$$

و

$$\int_{L_i}^{\infty} e^{-(\beta+\theta)z} \ln(1-e^{-\theta z}) dz = \int_{L_i}^{\infty} e^{-(\beta+\theta)z} \times \sum_{i=1}^{\infty} \frac{-e^{-i\theta z}}{i} dz = - \sum_{i=1}^{\infty} \frac{1}{i(\beta + i\theta + \theta)}.$$

پس داریم

$$E(\ln(1-e^{-\theta z}) | T \in (L, R)) = -\frac{1}{1-p(\theta)} \times \frac{\alpha\beta}{(\alpha + \theta)} \times \left[\sum_{i=1}^{\infty} \frac{1}{i(\beta + i\theta)} - \sum_{i=1}^{\infty} \frac{1}{i(\beta + i\theta + \theta)} \right].$$

در ادامه بر اساس قانون قوی اعداد بزرگ با فرض

$$S_n = \sum_{i=1}^n X_i = n_1$$

$$\frac{n_1}{n} \rightarrow p(\theta) \quad a.s.$$

و با فرض $n_2 = n - n_1$

$$\frac{n_2}{n} \rightarrow 1 - p(\theta) \quad a.s.$$

همچنین

$$P(N(n)=i) = \frac{n!}{i!(n-i)!} p^i (1-p)^{n-i}; \quad i=0,1,\dots,n$$

پس برای $n \rightarrow \infty$

$$\frac{1}{\sqrt{N(n)}} \sum_{i=1}^{N(n)} U_i \xrightarrow{d} N(0,1).$$

اثبات: فرض کنید

$$Y_{N(n)} = \frac{1}{\sqrt{N(n)}} \sum_{i=1}^{N(n)} U_i,$$

تابع مشخصه^۲ $Y_{N(n)}$ را با $\phi_{N(n)}(t)$ نمایش داده که به صورت زیر است،

$$\begin{aligned} \phi_{N(n)}(t) &= E(e^{itY_{N(n)}}) \\ &= \sum_{k=0}^n E\left(e^{it \frac{1}{\sqrt{k}} \sum_{i=1}^k U_i} \mid N(n)=k\right) \\ &\quad \times \frac{n!}{i!(n-i)!} p^k (1-p)^{n-k}. \end{aligned}$$

می دانیم

$$\begin{aligned} E\left(e^{it \frac{1}{\sqrt{k}} \sum_{i=1}^k U_i}\right) &= E\left(e^{\sum_{i=1}^k it \frac{1}{\sqrt{k}} U_i}\right) \\ &= \left(\phi\left(\frac{t}{\sqrt{k}}\right)\right)^k, \end{aligned}$$

حال اگر $\phi_U(\cdot)$ تابع مشخصه U_i باشد، برای t ثابت داریم

$$\begin{aligned} \left| \phi_{N(n)}(t) - e^{-\frac{t^2}{2}} \right| &\leq \sum_{k=0}^n \left| E\left(e^{\sum_{i=1}^k it \frac{1}{\sqrt{k}} U_i} \mid N(n)=k\right) - e^{-\frac{t^2}{2}} \right| \\ &\quad \times \frac{n!}{i!(n-i)!} p^k (1-p)^{n-k} \\ &= \sum_{k=0}^n \left| \left(\phi\left(\frac{t}{\sqrt{k}}\right)\right)^k - e^{-\frac{t^2}{2}} \right| \\ &\quad \times \frac{n!}{i!(n-i)!} p^k (1-p)^{n-k}. \end{aligned}$$

بنابر قضیه حد مرکزی^۳

$$\frac{\bar{U} - \bar{\mu}}{\frac{\sigma}{\sqrt{k}}} \xrightarrow{d} N(0,1).$$

در نتیجه

$$\frac{\sum_{i=1}^k U_i}{\sqrt{k}} \xrightarrow{d} N(0,1),$$

پس $(\phi_{U_i}(t))^k \rightarrow e^{-\frac{t^2}{2}}$ یعنی

$g''(\theta) < 0$ به ازای هر θ . اثبات را در حالت اول: $\hat{\theta}_n$ برای همه n ها کراندار است، و در حالت دوم: $\hat{\theta}_n$ کراندار نیست ادامه داده و در آن $\hat{\theta}$ را با $\hat{\theta}_n$ نشان می دهیم.

حالت ۱: $\hat{\theta}_n$ برای همه n ها کراندار است.

فرض کنید $\hat{\theta}_n$ به θ^* همگرا نباشد. بنابراین، زیر مجموعه ای مانند $\{n_k\}$ از $\{n\}$ وجود دارد که برای $\tilde{\theta} \neq \theta^*$ داریم

$\tilde{\theta} \rightarrow \hat{\theta}_{n_k}$. از طرفی چون $\hat{\theta}_{n_k}$ MLE است، پس

$$\frac{1}{n_k} L(\hat{\theta}_{n_k}) \geq \frac{1}{n_k} L(\theta^*),$$

با حد گیری از طرفین این نامساوی داریم

$$g(\tilde{\theta}) \geq g(\theta^*).$$

و این با فرض θ^* تنها ماکسیمم $g(\theta)$ در تناقض بوده، پس $\hat{\theta}_n$ به θ^* همگراست.

حالت ۲: $\hat{\theta}_n$ کراندار نیست.

در این حالت نیز فرض کنید $\hat{\theta}_n$ به θ^* همگرا نباشد. بنابراین، زیر مجموعه ای مانند $\{n_k\}$ از $\{n\}$ وجود دارد $\hat{\theta}_{n_k} \rightarrow \infty$. پس داریم

$$\frac{1}{n_k} L(\hat{\theta}_{n_k}) \geq \frac{1}{n_k} L(\theta^*),$$

در نتیجه $\frac{1}{n_k} L(\hat{\theta}_{n_k}) \rightarrow -\infty$ (زیرا $\frac{1}{n_k} L(\hat{\theta}_{n_k}) \rightarrow g(\hat{\theta}_{n_k})$ و $g(\infty) = -\infty$) هم چنین $\frac{1}{n_k} L(\theta^*)$ همگرا به یک مقدار ثابت

است (زیرا $\frac{1}{n_k} L(\theta^*) \rightarrow g(\theta^*)$ که $g(\theta^*)$ یک مقدار ثابت است)، پس به این نتیجه می رسیم که $\hat{\theta}_{n_k}$ MLE نیست که با فرض مسئله در تناقض است. پس $\hat{\theta}_n$ به θ^* همگراست.

قضیه ۳-۲: برآورد درست‌نمایی ماکسیمم θ یک برآورد سازگار برای θ است.

اثبات: با کمک لم های لم ۲۰۳ و لم ۳۰۳ و نتیجه ۲۰۲ موجود در کتاب برآورد نقطه ای لهن^۱ [۸]، صفحه ۱۴، قضیه فوق به آسانی اثبات می شود.

لم ۳-۴: فرض کنید U_i هادنباله ای از متغیرهای تصادفی مستقل و هم توزیع با $E(U_i)=0$ ، $Var(U_i)=1$ و $N(n)$ دارای توزیع دوجمله ای (n,p) است که تابع احتمال تجمعی $N(n)$ به صورت زیر است

² Characteristic function

³ Central Limit Theorem

¹ Lehmann

$$L(\theta) = n_1 \ln \theta - \theta \sum_{i=1}^{n_1} T_i - \theta \sum_{i=n_1+1}^{n_1+n_2} L_i + \sum_{i=n_1+1}^{n_1+n_2} \ln(1 - e^{-\theta Z_i}).$$

پس

$$L'(\theta) = -\frac{n_1}{\theta} - \sum_{i=1}^{n_1} T_i - \sum_{i=n_1+1}^{n_1+n_2} L_i + \sum_{i=n_1+1}^{n_1+n_2} \frac{Z_i e^{-\theta Z_i}}{(1 - e^{-\theta Z_i})^2},$$

و

$$L''(\theta) = \frac{n_1}{\theta^2} + \sum_{i=n_1+1}^{n_1+n_2} \frac{Z_i^2 e^{-\theta Z_i}}{(1 - e^{-\theta Z_i})^3}.$$

با استفاده از قضیه مقدار میانگین،

$$L'(\bar{\theta}) - L'(\theta) = (\bar{\theta} - \theta) L''(\bar{\theta}),$$

که در آن $\bar{\theta}$ یک نقطه بین θ و $\bar{\theta}$ است. بنابراین،

$$\sqrt{n}(\bar{\theta} - \theta) = -\frac{\frac{1}{\sqrt{n}} L'(\theta)}{\frac{1}{n} L''(\bar{\theta})}.$$

اثبات کامل می شود اگر نشان دهیم

$$\frac{1}{\sqrt{n}} L'(\theta) \xrightarrow{d} N(0,1), \quad (18)$$

و

$$\frac{1}{n} L''(\bar{\theta}) \xrightarrow{a.s.} c. \quad (19)$$

بنابر لم ۴ رابطه (۱۸) برقرار است. از طرفی بنابر قانون قوی اعداد بزرگ $a.s. \bar{\theta} \rightarrow \theta$ ، در نتیجه فرمول (۱۹) نیز ثابت می شود.

۴- نتیجه

در این مقاله ابتدا سانسور میانی را به طور دقیق تعریف کردیم و سپس به بررسی استنباطی توزیع نمایی وقتی داده ها از سانسور میانی آمده اند پرداختیم. برای محاسبه MLE پارامتر θ از روش عددی و الگوریتم EM استفاده کردیم و سپس نشان دادیم روش های تکراری که در فرمول (۷) مطرح شده است همگراست اگر رابطه (۱۲) برقرار باشد. در ادامه ثابت کردیم $\frac{1}{n} L(\theta) \rightarrow g(\theta) \text{ a.s.}$ و $g(\theta)$ یک تابع تک نمایی با ماکسیمم واحد است. در بخش ۳، با بیان لم ۳-۳ ثابت کردیم که برآورد درستنمایی ماکسیمم θ یک برآورد سازگار

$$\lim_{k \rightarrow \infty} \left(\phi_U \left(\frac{t}{\sqrt{k}} \right) \right)^k = e^{-\frac{t^2}{\tau}}.$$

بنابراین برای $\varepsilon > 0$ داده شده، $N_1(t)$ را به اندازه کافی بزرگ اختیار می کنیم پس برای $k \geq N_1(t)$

$$\left| \left(\phi_U \left(\frac{t}{\sqrt{k}} \right) \right)^k - e^{-\frac{t^2}{\tau}} \right| \leq \varepsilon.$$

علاوه بر این، برای $N_1(t)$ ثابت، n را به اندازه کافی بزرگ اختیار می کنیم پس طبق قضیه حد مرکزی

$$\sum_{i=0}^{N_1(t)} \frac{n!}{i!(n-i)!} p^i (1-p)^{n-i} \leq \varepsilon.$$

بنابراین

$$\left| \phi_{N(n)}(t) - e^{-\frac{t^2}{\tau}} \right| \leq 3\varepsilon.$$

حال چون ε دلخواه است و $e^{-\frac{t^2}{\tau}}$ تابع مشخصه توزیع $N(0,1)$ است پس

$$\frac{1}{\sqrt{N(n)}} \sum_{i=1}^{N(n)} U_i \xrightarrow{d} N(0,1).$$

پس لم فوق ثابت می شود.

قضیه ۳-۳: برآورد درستنمایی ماکسیمم دارای توزیع مجانبی زیر است

$$\sqrt{n}(\bar{\theta} - \theta) \xrightarrow{d} N\left(0, \frac{\sigma^*}{c^*}\right),$$

که

$$\begin{aligned} \sigma^* &= \left[E \left\{ \left(T - \frac{1}{\theta} \right)^2 \mid T \notin (L, R) \right\} \right. \\ &\quad \left. - \left(E \left\{ T - \frac{1}{\theta} \mid T \notin (L, R) \right\} \right)^2 \right] \\ &\quad + [E\{L^* \mid T \in (L, R)\} - (E\{L \mid T \in (L, R)\})^2] \\ &\quad + [E\{V^* \mid T \in (L, R)\} - (E\{V \mid T \in (L, R)\})^2] \end{aligned}$$

و در آن

$$V = \frac{Ze^{-\theta Z}}{1 - e^{-\theta Z}},$$

و

$$\begin{aligned} c &= \frac{p(\theta)}{\theta} + (1 - p(\theta)) \\ &\quad \times \left\{ E \left\{ \frac{Z^2 e^{-\theta Z}}{(1 - e^{-\theta Z})^2} \mid T \in (L, R) \right\} \right\}. \end{aligned}$$

اثبات: می دانیم

برای θ است و نشان دادیم که $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N\left(0, \frac{\sigma^2}{c^2}\right)$ می باشد.

منابع

- [1]. Babu, G.J., Rao, C.R. and Rao, M.B. (1992). Nonparametric estimation of specific exposure rate in risk and survival analysis. *Journal of American Statistical Association*, 87, 84-89.
- [2]. Chow, Y.S. and Teicher, H. (1980). *Probability Theory*, Springer, New York.
- [3]. Iyer, S.K., Kundu, D. and Jammalamadaka, S.R. (2008). Analysis of Middle Censored Data with Exponential Lifetime Distribution. *Journal of Statistical Planning and Inference*, to be appear.
- [4]. Jammalamadaka, S.R. and Mangalam, V. (2003). Nonparametric estimation for middle censored data. *Journal of Nonparametric Statistics*, 15, 253-265.
- [5]. Jammalamadaka, S.R. and Iyer, S.K. (2004). Approximate self consistency for middle censored data, *Journal of Statistical Planning and Inference*, 124, 75-86.
- [6]. Kaplan, F.L. and Meier, P. (1958). Nonparametric estimation from incomplete observation, *Journal of American Statistical Association*, 63, 457-481.
- [7]. Thompson; J., (1999) *Simulation: A Modeler's Approach*. Wiley, New York, 306 pages.
- [8]. Lehmann, E.L. (1986). *Theory of point estimation*. Wiley.
- [۹]. اسماعیلی، ح. (۱۳۸۴)، آنالیز عددی مبتنی بر مطلب، انتشارات بو علی سینا، همدان.
- [۱۰]. عالم زاده، ع. (۱۳۷۴)، حساب دیفرانسیل و انتگرال با هندسه تحلیلی، موسسه نشر علوم نوین.