

استفاده از مدل بی ثبات نرخ خطر متناسب در تحلیل داده‌های واقعی

محمد ملانوری^۱، حبیب نادری^۲، حامد احمدزاده^۳، سلمان ایزدخواه^۴

تاریخ دریافت: ۹۹/۱/۱۳

تاریخ پذیرش: ۹۹/۱/۱

چکیده:

جمعیت‌های زیادی در آنالیز بقا اغلب با مشکل ناهمگنی مواجه هستند. افراد در نحوه برخورد با علت مرگ، واکنش به معالجه و تحت تأثیر عوامل خطر قرار گرفتن، انعطاف‌پذیر هستند. نادیده گرفتن این ناهمگنی می‌تواند باعث به دست آمدن نتایج نادرست شود. برای برطرف کردن این مشکلات، مدل بی ثبات نرخ خطر متناسب را معرفی می‌کنیم. در این مقاله، ضمن معرفی مدل بی ثبات نرخ خطر متناسب، ویژگی‌های آن را مورد مطالعه قرار می‌دهیم. برازش مدل بی ثبات به داده‌های سانسور شده از راست را در حضور متغیرهای توضیحی (متغیرهای قابل مشاهده) بررسی می‌کنیم و در قالب یک مثال کاربردی برای برازش مدل بی ثبات به داده‌ها با در نظر گرفتن توزیع پایه و ایبول و نمایی در توابع درستنمایی، از آن‌ها برای برآورد پارامترهای مدل استفاده کرده و با معیارهای مختلف به مقایسه مناسب مدل‌ها می‌پردازیم.

واژه‌های کلیدی: برآورد ماکسیمم درستنمایی، توزیع پایه، تابع بقای غیرشرطی، مدل بی ثبات

۱ مقدمه

کرمانی [۶] مقایسه‌های تصادفی در مدل‌های بی ثبات را از طریق تأثیر مقایسه دو متغیر تصادفی پایه بر روی متغیرهای بی ثبات متناظر مطالعه کرده و نتایج مفیدی به دست آوردند. ژو و لی [۱۸] نیز متغیر پایه را مستقل از متغیر بی ثبات فرض کردند. در مدل آنان متغیر خروجی (یا متغیر جمعیت) در مدل بی ثبات است که با متغیر بی ثبات وابستگی دارد. کیاد و همکاران [۱۳] برخی از نتایج ترتیب تصادفی نسبی از دو مدل بی ثبات را مورد مطالعه قرار دادند. علاوه بر این محققین دیگری در زمینه مدل‌های بی ثبات تحقیقاتی را انجام داده‌اند. به‌عنوان مثال می‌توان به آلن [۱، ۲، ۳]، هودگارد [۷، ۸، ۹، ۱۰، ۱۱]، نلسون [۱۶] و مولر و استویان [۱۵] اشاره کرد.

انتخاب توزیع مناسب برای متغیر بی ثبات از دیدگاه بیزی همان انتخاب توزیع پیشین برای پارامتر بی ثبات است که نمی‌توان گفت هیچ دیدگاهی برای انتخاب توزیع پیشین وجود ندارد. در عمل بنا به تجربه معمولاً توزیع گاما، وایبول، لگک نرمال و توزیع گوسی وارون را برای متغیر بی ثبات انتخاب می‌کنند که دلیل اصلی آن انعطاف‌پذیر بودن این توزیع‌هاست. در واقع می‌توان گفت بیشتر نرم‌افزارها به انتخاب توزیع گاما برای متغیر بی ثبات محدود می‌شوند. بنابراین علاقه‌مندیم مدل نرخ خطر متناسب زمانی که عامل بی ثبات دارای توزیع گاما است را مورد بحث قرار دهیم. در این حالت با

یک مطالعه تحلیل بقا را در نظر بگیرید که زیرگروه‌های جامعه دارای یک متغیر تصادفی مشترک مشاهده نشده باشند. چنین متغیرهایی می‌توانند یک عامل ژنتیکی، اثر محیطی در سال‌های اولیه زندگی فرزندان یک خانواده یا اثر الگوی زندگی در زوجها باشد. به این متغیرهای مشترک، خطرهای مشترک و به مدل‌هایی که دربرگیرنده این ناهمگنی تصادفی هستند، مدل‌های بی ثبات می‌گویند. این مدل‌ها به‌منظور تبیین تغییرات ناشی از عوامل خطر مشاهده نشده به کار می‌روند. مدل‌های بی ثبات (مدل اثر تصادفی) قادرند ناهمگنی بین زمان‌های بقا در گروه‌های مختلف و همبستگی ایجادشده بین زمان‌های بقای هر گروه را به حساب آورند. مدل‌های بی ثبات مدل‌های نسبتاً جدیدی در تحلیل بقا بوده و به‌طور گسترده‌ای در سه دهه‌ی اخیر مورد مطالعه قرار گرفته‌اند. مدل‌های بی ثبات کلاسیک اولین بار توسط واپل و همکاران [۱۷]، برای مطالعه برخی ناهمگنی‌هایی که به‌وسیله متغیرهای کمکی قابل اندازه‌گیری نیستند، معرفی شدند. کاربرد این مدل‌ها در تحلیل بقا برای ناهمگنی مشاهده نشده (پنهان) بین افراد، با فرض اینکه نرخ خطر افراد متأثر از یک پارامتر مشخص به نام پارامتر شکنندگی (بی ثبات) و یک مخاطره پایه در نظر گرفته شود، مورد توجه بسیاری از محققان بوده است. گوپتا و

^۱ دانشجوی کارشناسی ارشد دانشگاه سیستان و بلوچستان

^۲ عضو هیئت علمی دانشگاه سیستان و بلوچستان h.h.naderi@gmail.com

^۳ عضو هیئت علمی دانشگاه سیستان و بلوچستان

^۴ عضو هیئت علمی دانشکده علوم پایه و فنی-مهندسی بیجار

به صورت زیر است:

$$h(t|v) = -\frac{d}{dt} \ln \bar{F}(t|v) = v \left(-\frac{d}{dt} \ln \bar{F}_0(t) \right) = v h_0(t) \quad (2)$$

که در آن $h_0(t)$ تابع نرخ خطر پایه و مستقل از v است. در واقع اگر T معرف طول عمر یک مؤلفه باشد آنگاه $h(t|v)$ معرف نرخ مخاطره در لحظه t برای مؤلفه‌ای با بی‌ثباتی v خواهد بود.

با استفاده از رابطه (۲) تابع بقای شرطی متغیر T به شرط V با توجه به اینکه $\bar{F}_0(t) = \exp(-\int_0^t h_0(x) dx)$ ، به صورت زیر به دست می‌آید:

$$\bar{F}(t|v) = [\bar{F}_0(t)]^v = \exp\left(-v \int_0^t h_0(x) dx\right) \quad (3)$$

که در آن $\bar{F}_0(t)$ تابع بقای پایه است. تابع چگالی احتمال شرطی متناظر با رابطه (۳) به صورت زیر است:

$$\begin{aligned} f(t|v) &= -\frac{\partial \bar{F}(t|v)}{\partial t} = -v(\bar{F}_0(t))^{v-1} \times (-f_0(t)) \quad (4) \\ &= v \frac{f_0(t)}{\bar{F}_0(t)} \cdot (\bar{F}_0(t))^v \\ &= v h_0(t) \exp\left(-v \int_0^t h_0(x) dx\right) \end{aligned}$$

با استفاده از رابطه (۳) تابع توزیع غیر شرطی (جمعیت) متغیر T پس از متوسط گیری نسبت به تأثیر متغیر پنهان V به صورت زیر به دست می‌آید:

$$\begin{aligned} \bar{F}(t) &= \int_0^\infty \bar{F}(t|v) g(v) dv \quad (5) \\ &= \int_0^\infty (\bar{F}_0(t))^v g(v) dv \\ &= E[(\bar{F}_0(t))^V]; t > 0 \end{aligned}$$

که در آن $g(\cdot)$ تابع چگالی متغیر تصادفی V است. با توجه به رابطه (۵) تابع چگالی غیر شرطی T به صورت زیر به دست می‌آید:

$$\begin{aligned} f(t) &= \int_0^\infty f(t|v) g(v) dv \quad (6) \\ &= \int_0^\infty v h_0(t) \exp\left(-v \int_0^t h_0(x) dx\right) g(v) dv \\ &= h_0(t) \int_0^\infty v \exp\left(-v \int_0^t h_0(x) dx\right) g(v) dv \end{aligned}$$

۳ مدل بی‌ثبات نرخ خطر متناسب در

برآورد پارامترهای مدل

قبل از شروع مباحث اصلی نیاز به بیان مفاهیمی مقدماتی است که در ابتدا به آن اشاره می‌کنیم.

انتگرال‌گیری از تابع درست‌نمایی شرطی نسبت به متغیر بی‌ثبات می‌توان به یک عبارت صریح و ساده برای درست‌نمایی حاشیه‌ای دست‌یافت. درست‌نمایی حاشیه‌ای که شامل پارامترهای مورد نظر است را ماکسیم می‌کنیم و برآورد پارامترها و خطای استاندارد آن‌ها را به دست می‌آوریم. در بخش دوم مدل بی‌ثبات نرخ خطر متناسب را معرفی و شاخص‌هایی برای آن را مورد مطالعه قرار می‌دهیم. در بخش سوم به برازش مدل بی‌ثبات به داده‌های سانسور شده از راست در حضور متغیرهای توضیحی (قابل مشاهده) می‌پردازیم. در بخش چهارم در قالب یک مثال کاربردی مدل بی‌ثبات نرخ خطر متناسب را به داده‌های سانسور شده از راست برازش می‌دهیم. در پایان نتیجه‌گیری را ارائه می‌دهیم.

۲ معرفی مدل بی‌ثبات نرخ خطر متناسب

در این بخش به معرفی مدل بی‌ثبات نرخ خطر متناسب و شرایط معادل با آن می‌پردازیم. قبل از معرفی این مدل لازم است تعریف زیر را بیان کنیم که از طریق آن می‌توانیم از یک توزیع پایه، توزیع دیگری به دست آوریم.

تعریف ۱.۲. برای هر $\alpha > 0$ و تابع توزیع پایه F_0 ، تابع $F(t|\alpha) = 1 - (1 - F_0(t))^\alpha$ یک تابع توزیع است که یک خانواده از توزیع‌های با نرخ خطر متناسب را تشکیل می‌دهد.

اکنون برای معرفی این مدل فرض کنید متغیر تصادفی نامنفی و پیوسته طول عمر T و متغیر تصادفی نامنفی V ، متغیر بی‌ثبات و $F_0(x)$ یک توزیع پایه باشد آنگاه بنا به تعریف (۱.۲) و با توجه به ارتباط بین تابع بقا و تابع نرخ خطر، تابع بقای شرطی شده روی متغیر بی‌ثبات (یعنی به شرط $V = v$) برای یک مؤلفه به صورت زیر بیان می‌شود:

$$\bar{F}(t|v) = (1 - F_0(t))^v = (\bar{F}_0(t))^v \quad (1)$$

بنابراین اعضای جمعیت با $v > 1$ در مقایسه اعضای با $v = 1$ ، دارای تابع بقای نزولی و نرخ خطر صعودی (با توجه به فرمول تابع نرخ خطر، تابع نرخ خطر و تابع بقا رابطه‌ی عکس باهم دارند) و اعضای جمعیت با $v < 1$ در مقایسه با اعضای $v = 1$ ، دارای تابع بقای صعودی و نرخ خطر نزولی هستند.

با به کار بردن رابطه $\bar{F}(x) = \exp\{-H(x)\}$ می‌توان نوشت:

$$\bar{F}(t|v) = \exp\{-vH_0(t)\}$$

که در آن $H_0(t) = \int_0^t h_0(x) dx = -\ln \bar{F}_0(t)$ ، تابع نرخ خطر تجمعی متناظر با $\bar{F}_0(x)$ می‌باشد. با استفاده از رابطه (۱) مدل بی‌ثبات نرخ خطر متناسب

اطلاعات حقیقی برای مؤلفه‌ی i ام، به ازای $i = 1, 2, \dots, n$ شامل زوج (y_i, δ_i) است به طوری که y_i مینیمم زمان پیشامد t_i و زمان سانسور c_i است. یعنی $y_i = \min(t_i, c_i)$ و δ_i نشانگر سانسور است و اگر پیشامد مورد نظر اتفاق افتد، مقدار یک را می‌گیرد و در غیر این صورت، مقدار صفر را خواهد گرفت.

$$\delta_i = \begin{cases} 1, & t_i \leq c_i \\ 0, & t_i > c_i \end{cases}$$

برای داده‌های سانسور شده، تابع درستنمایی بقا کاملاً از تابع درستنمایی کلاسیک برای داده‌های مستقل بدون سانسور متفاوت است. فرم عبارت درستنمایی با نوع داده‌هایی که در دسترس هستند تعیین می‌شود. به طور کلی در اینجا، داده‌هایی که در نظر می‌گیریم یا داده‌های کامل هستند (همه‌ی زمان‌های پیشامد مشاهده شده‌اند) یا داده‌های سانسور شده از راست هستند (زمان پیشامد مؤلفه‌ها برای سانسور راست) و یا داده‌های سانسور شده‌ی بازه‌ای هستند (زمان‌های پیشامد معلوم هستند و در یک بازه مشاهده می‌شوند). برای داده‌های سانسور شده از راست با سانسور تصادفی، سهم درستنمایی زمان پیشامد $(y_i = t_i, \delta_i = 1)$ به صورت زیر داده می‌شود:

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \frac{1}{\sqrt{\epsilon}} P(y_i - \epsilon < Y_i < y_i + \epsilon, \delta_i = 1) \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\sqrt{\epsilon}} P(y_i - \epsilon < Y_i < y_i + \epsilon, T_i \leq C_i) \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\sqrt{\epsilon}} \int_{y_i - \epsilon}^{y_i + \epsilon} \int_t^{\infty} dG(t) dF(t) \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\sqrt{\epsilon}} \int_{y_i - \epsilon}^{y_i + \epsilon} (1 - G(t)) dF(t) = (1 - G(y_i)) f(y_i) \end{aligned}$$

از طرف دیگر برای مشاهدات سانسور شده‌ی راست $(y_i = c_i, \delta_i = 0)$ سانسور تصادفی سهم درستنمایی به صورت زیر داده می‌شود:

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \frac{1}{\sqrt{\epsilon}} P(y_i - \epsilon < Y_i < y_i + \epsilon, \delta_i = 0) \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\sqrt{\epsilon}} P(y_i - \epsilon < Y_i < y_i + \epsilon, T_i > C_i) \\ &= (1 - F(y_i)) g(y_i) \end{aligned}$$

تحت سانسور تصادفی راست، داده بقا شامل ترکیب زمان‌های پیشامد و مشاهدات سانسور شده‌ی راست است. بنابراین طبق توضیحات بالا درستنمایی برای یک نمونه به اندازه‌ی n به صورت زیر به دست می‌آید:

$$L = \prod_{i=1}^n [(1 - G(y_i)) f(y_i)]^{\delta_i} [(1 - F(y_i)) g(y_i)]^{1 - \delta_i}$$

اگر فرض کنیم که توزیع زمان‌های سانسور به پارامترهای مرتبط با تابع بقا وابسته نیستند، آنگاه گوییم سانسور حاوی اطلاعات مفید نمی‌باشد. (لیانگ [۱۴]، فلمینگ و هرینگتون [۵]) و عوامل $(1 - G(y_i))^{\delta_i}$ و $(1 - F(y_i))^{1 - \delta_i}$ برای

⁵Right Censoring

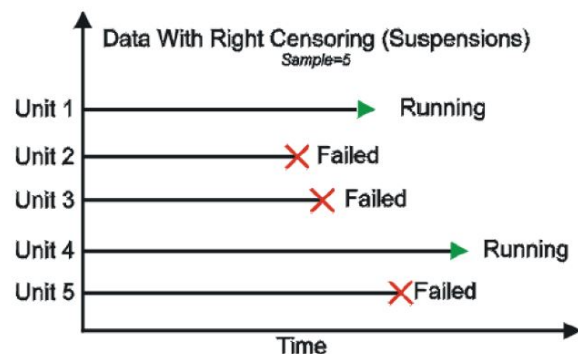
۱.۳ داده‌های سانسور شده از راست

رایج‌ترین نوع سانسور، متعلق به داده‌های سانسور شده از راست می‌باشد. در داده‌های مربوط به طول عمر، این مجموعه داده‌ها شامل واحدهایی هستند که دچار شکست نمی‌شوند. اگر C_1, C_2, \dots, C_n را متغیرهای تصادفی مستقل و هم‌توزیع با تابع توزیع G در نظر بگیریم. در اینجا C_i زمان سانسور با توجه به زمان شکست t_i (زمان شکست i امین واحد) می‌باشد. در این حالت سانسوری، داده‌ها به صورت زیر نمایش داده می‌شود

$$(Y_1, \delta_1), (Y_2, \delta_2), \dots, (Y_n, \delta_n)$$

که در آن $\delta_i = I(T_i \leq C_i)$ و $Y_i = \min(T_i, C_i)$ می‌باشد. $\delta_1, \dots, \delta_n$ حاوی اطلاعات مربوط به سانسور شدن زمان‌های شکست واحدهای مورد نظر می‌باشند. یعنی اگر $\delta_i = 0$ شد، مشاهده مورد نظر (Y_i) مربوط به زمان شکست سانسور شده می‌باشد و اگر $\delta_i = 1$ شد، (Y_i) برابر زمان دقیق شکست است.

مثال ۱.۳. اگر ۵ واحد را مورد بررسی قرار دهیم و فقط ۳ تا از آن‌ها تا پایان بررسی به زمان شکست برسند. با توجه به ۲ واحدی که به زمان شکستشان نرسیده‌اند، ما داده‌های سانسور شده از راست خواهیم داشت. واژه "از راست سانسور شده" بدین مفهوم است که زمان رخداد پیشامد مورد علاقه (زمان شکست) در سمت راست نقاط داده‌ها (برای ۳ واحد مشاهده) قرار دارد. به عبارت دیگر اگر واحدها به کار خود ادامه دهند، زمان شکست این دو واحد، زمانی بعد از سایر واحدهای مشاهده شده و یا سمت راست مقیاس زمانی آن‌ها قرار خواهد گرفت.



شکل ۱: نمودار سانسور از راست

۲.۳ تابع درستنمایی بقا برای داده‌های سانسور شده از راست

در حالت کلی، فرض می‌کنیم که زمان سانسور و زمان بقا از نظر آماری، متغیرهای تصادفی مستقل هستند. برای داده‌های سانسور شده از راست،

در اینجا می‌خواهیم پارامترهای θ و پارامترهای توزیع $g(\cdot)$ (پارامترهای توزیع بی‌ثبات) و نیز بردار ضرایب رگرسیونی θ را برآورد کنیم. با استفاده از رابطه (۹) و نیز با فرض

$$H_0(t; \theta) = \int_0^t h(u; \theta) du = -\log \bar{F}_0(t; \theta)$$

و همچنین طبق (۷) تابع درستنمایی شرطی بر اساس مشاهدات به صورت زیر است:

$$\begin{aligned} L(\beta, \theta | \mathbf{x}, v) &= \prod_{i=1}^n (f(t_i | x_i, v_i))^{\delta_i} (\bar{F}(t_i, x_i, v_i))^{1-\delta_i} \quad (10) \\ &= \prod_{i=1}^n (h(t_i | x_i, v_i))^{\delta_i} \bar{F}(t_i, x_i, v_i) \\ &= \prod_{i=1}^n \left(v_i \exp(\beta' x_i) g_0(t_i, \theta) \right)^{\delta_i} \\ &\quad \times \exp\left(-v_i \exp(\beta' x_i) R_0(t_i, \theta)\right) \end{aligned}$$

برای برازش مدل بی‌ثبات گاما فرض می‌کنیم متغیر بی‌ثبات V دارای توزیع گاما با میانگین یک و واریانس σ ($V \sim \Gamma(\frac{1}{\sigma}, \sigma)$) باشد، یعنی،

$$g(v) = \frac{v^{\frac{1}{\sigma}-1}}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \exp\left(-\frac{v}{\sigma}\right) \quad v > 0, \sigma > 0 \quad (11)$$

بنابراین تابع درستنمایی غیرشرطی بر اساس مشاهدات و روابط (۱۰) و (۱۱) به صورت زیر است:

$$\begin{aligned} L(\beta, \theta, \sigma | \mathbf{x}) &= \int_0^\infty L(\beta, \theta | \mathbf{x}, v) h(v) dv \quad (12) \\ &= \int_0^\infty \prod_{i=1}^n \left(v_i \exp(\beta' x_i) h_0(t_i, \theta) \right)^{\delta_i} \exp\left(-v_i \exp(\beta' x_i) R_0(t_i, \theta)\right) \frac{v_i^{\frac{1}{\sigma}-1}}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \exp\left(-\frac{v_i}{\sigma}\right) dv_i \\ &= \prod_{i=1}^n \int_0^\infty \left(v_i \exp(\beta' x_i) h_0(t_i, \theta) \right)^{\delta_i} \exp\left(-v_i \exp(\beta' x_i) R_0(t_i, \theta)\right) \cdot \frac{v_i^{\frac{1}{\sigma}-1}}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \exp\left(-\frac{v_i}{\sigma}\right) dv_i \\ &= \prod_{i=1}^n \frac{\left(\exp(\beta' x_i) h_0(t_i, \theta) \right)^{\delta_i}}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \int_0^\infty v_i^{\delta_i + \frac{1}{\sigma} - 1} \exp\left(-\frac{v_i}{\sigma} - v_i \exp(\beta' x_i) R_0(t_i, \theta)\right) dv_i \\ &= \prod_{i=1}^n \frac{\left(\exp(\beta' x_i) h_0(t_i, \theta) \right)^{\delta_i}}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \int_0^\infty v_i^{\delta_i + \frac{1}{\sigma} - 1} \exp\left(-\frac{v_i}{\sigma} - \frac{v_i}{\sigma} \sigma \exp(\beta' x_i) R_0(t_i, \theta)\right) dv_i \\ &= \prod_{i=1}^n \frac{\left(\exp(\beta' x_i) h_0(t_i, \theta) \right)^{\delta_i}}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \int_0^\infty v_i^{\delta_i + \frac{1}{\sigma} - 1} e^{-\frac{v_i}{\sigma} (1 + \sigma \exp(\beta' x_i) R_0(t_i, \theta))} dv_i \\ &= \prod_{i=1}^n \frac{\left(\exp(\beta' x_i) h_0(t_i, \theta) \right)^{\delta_i}}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \frac{\Gamma(\delta_i + \frac{1}{\sigma})}{\left(1 + \sigma \exp(\beta' x_i) R_0(t_i, \theta) \right)^{\delta_i + \frac{1}{\sigma}}} \\ &= \prod_{i=1}^n \frac{\Gamma(\delta_i + \frac{1}{\sigma})}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \left(\frac{\exp(\beta' x_i) h_0(t_i, \theta)}{1 + \sigma \exp(\beta' x_i) R_0(t_i, \theta)} \right)^{\delta_i} \times \left(1 + \sigma \exp(\beta' x_i) R_0(t_i, \theta) \right)^{-\frac{1}{\sigma}} \end{aligned}$$

استنباط درباره‌ی تابع بقا مفید نیستند و می‌توان آن‌ها را از درستنمایی حذف کرد، لذا داریم:

$$\begin{aligned} L &= \prod_{i=1}^n (f(y_i))^{\delta_i} (\bar{F}(y_i))^{1-\delta_i} \quad (7) \\ &= \prod_{i=1}^n (h(y_i))^{\delta_i} \bar{F}(y_i) \end{aligned}$$

در این بخش مدل بی‌ثبات $h(t | v) = v h_0(t)$ برای داده‌های سانسور تصادفی از راست در حضور متغیرهای توضیحی (متغیرهای قابل مشاهده) بررسی می‌شود. به این منظور، فرض کنید $d_i = 1, 2, \dots, n$ متغیرهای پیوسته نشان‌دهنده زمان بقای مؤلفه i مستقل و هم توزیع با تابع توزیع F و تابع چگالی f باشند و همچنین فرض کنید C_i ، $i = 1, 2, \dots, n$ ، زمان سانسور مؤلفه i ام باشد. در این صورت، متغیرهای تصادفی قابل مشاهده عبارت‌اند از: $Y_i = \min(T_i, C_i)$ و $\delta_i = I(T_i \leq C_i)$ که در آن $I(\cdot)$ نشانگر وقوع حادثه موردنظر است. هدف، برازش مدل بی‌ثبات پارامترهای نرخ خطر متناسب به نمونه تصادفی یادشده است. اکنون فرض کنید مدل بی‌ثبات $h(t | v) = v h_0(t)$ به صورت زیر در نظر گرفته شود:

$$h(t | \mathbf{x}, v) = v g(\mathbf{x}, \beta) h_0(t; \theta) \quad (8)$$

که در آن \mathbf{x} مقدار مشاهده شده بردار متغیرهای توضیحی $\mathbf{x}_{p \times 1}$ و $\beta_{p \times 1}$ بردار پارامترهای رگرسیونی، $g(\cdot)$ تابعی نامنفی از \mathbb{R}^+ و همچنین θ بردار پارامترهای توزیع پایه فرض شده است. با فرض $g(\mathbf{x}, \beta) = \exp(\beta' \mathbf{x})$ مدل (۸) به صورت زیر درمی‌آید:

$$h(t | \mathbf{x}, v) = v \exp(\beta' \mathbf{x}) h_0(t; \theta) \quad (9)$$

$$= \prod_{i=1}^n \frac{\Gamma(\delta_i + \frac{1}{\sigma})}{\sigma^{\frac{1}{\sigma}} \Gamma(\frac{1}{\sigma})} \left(\frac{\exp(\underline{\beta}' x_i) h_{\circ}(t_i, \theta)}{1 - \sigma \exp(\underline{\beta}' x_i) \log \bar{F}_{\circ}(t_i, \theta)} \right)^{\delta_i} \times \left(1 - \sigma \exp(\underline{\beta}' x_i) \log \bar{F}_{\circ}(t_i, \theta) \right)^{-\frac{1}{\sigma}}$$

که تحت فرض H_0 برای n های بزرگ دارای توزیع کای دو با یک درجه آزادی است. بنابراین در سطح α فرض H_0 رد می شود اگر

$$\chi_{(LR)}^2 > \chi_{1, \alpha}^2$$

بنابراین برای بررسی تأثیر متغیر توضیحی بر متغیر پاسخ یعنی آزمون

$$\begin{cases} H_0 : \beta = 0 \\ H_1 : \beta \neq 0 \end{cases}$$

آماره آزمون نسبت درستنمایی محاسبه و با مقدار توزیع کای دو با یک درجه آزادی مقایسه می شود. در واقع کافی است تفاوت دو مقدار $2 \log L$ در حالتی که متغیر توضیحی حذف شده است محاسبه شده و با مقدار توزیع کای دو با یک درجه آزادی مقایسه شود که آماره آزمون نسبت درستنمایی و p -مقدار در جدول (۱) آورده شده است. به این نتیجه رسیدیم که نوع درمان تأثیر معنی داری بر طول عمر بیماران دارد. در این مثال توزیع وایبول را با توزیع نمایی مقایسه می کنیم. نتایج به دست آمده از برازش مدل ها، برآورد ماکسیمم درستنمایی پارامترهای مدل مطرح شده، برآورد انحراف استاندارد برآوردگرها و معیار اطلاع آکائیکه و معیار اطلاع بیزی در جدول (۲) ارائه شده است.

نکته: در مورد معیار اطلاع آکائیکه (AIC) باید گفت معیاری برای سنجش نیکویی برازش است. این معیار بر اساس مفهوم آنتروپی بنا شده است و نشان می دهد که استفاده از یک مدل آماری به چه میزان باعث از دست رفتن اطلاعات می شود. به عبارت دیگر، این معیار تعادلی میان دقت مدل و پیچیدگی آن برقرار می کند. این معیار توسط هیروتسوگو آکائیک [۴] برای انتخاب بهترین مدل آماری پیشنهاد شد.

جدول ۱: مقادیر آماره آزمون و p -مقدار برای دو مدل نمایی و وایبول

توزیع پایه	آماره آزمون	p -مقدار
توزیع نمایی	۳۸/۷۹۴	۰/۰۰۰۴
توزیع وایبول	۴۹/۲۰۴	۰/۰۰۰۲

۴ مثال کاربردی

اکنون برازش مدل های (۸) به داده های سرطان ریه (کلیفلش و پرنیس [۱۲]) بررسی می شود. در آزمایش مربوط به بیماران سرطان ریه، ۱۳۷ مرد مبتلا به سرطان ریه پیشرفته به طور تصادفی تحت درمان معمولی و یا شیمی درمانی قرار می گیرند. نقطه ی پایان آزمایش، زمان فوت بیمار در نظر گرفته شده است. در اینجا برای برازش مدل (۹) به داده ها با در نظر گرفتن دو حالت زیر برای توزیع پایه در توابع درستنمایی (۱۲) از آن برای برآورد پارامترهای مدل استفاده می شود.

الف) توزیع پایه نمایی باشد یعنی تابع نرخ خطر آن به صورت زیر باشد:

$$g(t; c, \gamma) = \frac{\lambda e^{-\lambda t}}{e^{-\lambda t}} = \lambda$$

ب) توزیع پایه وایبول باشد یعنی تابع نرخ خطر آن به صورت زیر باشد:

$$g(t; \lambda) = \frac{c \gamma^c x^{c-1} e^{-(\gamma x)^c}}{e^{-(\gamma x)^c}} = c \gamma^c x^{c-1}$$

که در آن $\theta = (c, \gamma)$ است. در این مطالعه فرض شده است بیمارانی که از شروع آزمایش دارای طول عمر بیشتر از ۱۵ ماه بوده اند، به عنوان داده های سانسور شده از راست باشند. همچنین تنها نوع درمان بیمار به عنوان متغیر توضیحی در نظر گرفته شده است. فرض می کنیم یک عامل ژنتیکی که قابل مشاهده نیست در درمان بیماری افراد مؤثر باشد، که در اینجا از آن به عنوان متغیر بی ثبات یاد می شود.

تذکره: یکی از شیوه ها برای آزمون فرض

$$\begin{cases} H_0 : \beta = 0 \\ H_1 : \beta \neq 0 \end{cases}$$

آزمون نسبت درستنمایی است. آماره این آزمون عبارت است از:

$$\chi_{(LR)}^2 = 2 \left(\log L(\hat{\beta}) - \log L(\beta_0) \right)$$

پارامتر، افزایش پیدا کند و این امر ممکن است منجر به بیش برآورد شود که در این شرایط AIC و BIC ، این مسئله را با اعمال محدودیت روی تعداد پارامترهای مدل حل می‌کنند که محدودیت BIC در تعداد پارامترها نسبت به AIC بیشتر است. هنگامی که در جستجوی بهترین برازش برای داده‌ها در بین سایر توزیع‌ها هستیم توزیع با کوچک‌ترین مقدار AIC به‌عنوان بهترین برازش انتخاب می‌شود. باید توجه داشت که در مبحث تشخیص مدل معیار AIC با استفاده از مشاهدات به‌دست آمده به انتخاب مدل مناسب می‌پردازد. ممکن است با تغییرات اندکی در مشاهدات، این معیار مدل دیگری را به‌عنوان مدل درست مشخص کند. پس اگر داده‌ها در نقاط مناسبی جمع‌آوری نشده باشند این چنین معیارهایی به نتایج گمراه‌کننده منجر شده و مدل را به‌درستی انتخاب نمی‌کنند. بنابراین لازم است مشاهدات به‌گونه‌ای باشند که معیار انتخاب مدل در برابر آن‌ها پایدار بوده و تصمیم این معیارها در تعیین مدل درست با تغییر مشاهدات از ثبات کافی برخوردار باشد.

با توجه به داده‌ها، چند مدل آماری ممکن است با توجه به مقدار AIC رتبه‌بندی شوند و مدل دارای کمترین AIC ، بهتر است. در حالت کلی (AIC) برابر است با:

$$AIC = -2\ln(L) + 2k = -2\ell(\hat{\theta}) + 2k$$

که k تعداد پارامترهای مدل آماری است و L مقدار ماکسیمم تابع درستنمایی برای مدل برآورد شده است. معیار اطلاع بیزی (BIC) به‌صورت زیر است:

$$BIC = \ln(n)k - 2\ln(L) = \ln(n)k - 2\ell(\hat{\theta})$$

که k تعداد پارامترهای مدل آماری است و L مقدار ماکسیمم تابع درستنمایی برای مدل برآورد شده و n تعداد مشاهدات است. در توضیح معیار BIC باید گفت همانند AIC برای انتخاب مدل از بین چند مدل به کار می‌رود و بر اساس تابع درستنمایی است و لذا بسیار مرتبط با AIC می‌باشد. زمانی که مدل برازش می‌دهیم این امکان وجود دارد که تابع درستنمایی با اضافه کردن

جدول ۲: مقادیر MLE و لگاریتم درستنمایی (LL) و BIC و AIC برای دو مدل نمایی و وایبول

توزیع پایه	پارامترها	برآورد	انحراف استاندارد	AIC	BIC	$loglike$
توزیع نمایی	λ	۰/۰۳۷۴	۰/۰۱۰۹	۱۴۵۱/۱۸۸	۱۴۵۹/۹۲۶	-۷۲۲/۵۹۳۹
	β	۰/۰۳۶۶	۰/۰۴۸۶			
	σ	-۰/۰۲۴۵	۰/۰۰۹۱			
توزیع وایبول	c	۱/۱۳۵۸	۰/۱۳۹۴	۱۴۴۲/۷۰۰	۱۴۵۴/۳۵۱	-۷۱۷/۳۵۰۱
	γ	۰/۰۹۲۴	۰/۰۳۸۵			
	σ	۰/۲۵۷۹	۰/۲۰۵۱			
	β	-۰/۰۴۱۱	۰/۰۰۷۱			

۵ نتیجه‌گیری

در این مقاله مدل نرخ خطر متناسب را به‌عنوان مدل بی‌ثبات در نظر گرفتیم و ویژگی‌های توزیعی آن را مورد بررسی قرار دادیم. همچنین در قالب یک مثال کاربردی نشان دادیم که مدل بی‌ثبات نرخ خطر متناسب گاما با توزیع پایه وایبول با توجه به کمتر بودن معیارهای مقایسه مدل نسبت به توزیع نمایی، برای برازش به این داده‌ها مناسب است.

همان‌طور که مشاهده می‌شود توزیع وایبول دارای معیار اطلاع آکائیکه و بیزی کمتری نسبت به توزیع نمایی است بنابراین مدل بی‌ثبات نرخ خطر متناسب گاما با توزیع پایه وایبول نسبت به توزیع نمایی، برازش مناسب‌تری به این داده‌ها ارائه کرده است.

مراجع

- [1] Aalen, O.O. (1978). Nonparametric inference for a family of counting processes. *Ann. Stat.*, **6**, 701–726.
- [2] Aalen, O.O. (1998). Heterogeneity in survival analysis. *Statist. Med.*, **7**, 1121-1137.
- [3] Aalen, O.O. (1992). Modelling heterogeneity in survival analysis by the compound poisson distributhion. *Ann. Appl. Probab.*, **2**, 951-992.
- [4] Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, **19**, 716–723.
- [5] Fleming, T.R. and Harrington, D.P. (1991). *Counting Processes and survival analysis*. Wiley, New York.
- [6] Gupta, R.C. and Kirmani, S. (2005). Stochastic comparisons in frailty models. *J Stat Plan Inference.*, **136**, 3647–3658.
- [7] Hougaard, P. (1986). Survival models for heterogeneous populations derived from stable distributions. *Biometrika*, **73**, 387–396.
- [8] Hougaard, P. (1995). Frailty models for survival analysis. *Lifetime Data Anal.*, **1**, 255-273.
- [9] Hougaard, P. (2000). *Analysis of Multivariate Survival Data*. Springer-Verlag, New York.
- [10] Hougaard, P. (1984). Life table methods for heterogeneous population: distribution describing the heterogeneity. *Biometrika*, **71**, 75-83.
- [11] Hougaard, P. (1991). Modelling heterogeneity in survival analysis. *J. Appl. Prpbab.*, **28**, 695-701.
- [12] Kalbfleisch, D. and Prentice, R. L. (2011), *The Statistical Analysis of Failure Time Data*, John Wiley, New York.
- [13] Kayid, M., Izadkhah, S., & Zuo, M. J. (2017). Some results on the relative ordering of two frailty models. *Statistical Papers*, **58(2)**, 287-301.
- [14] Liang, K. Y., Self, S. G., Bandeen-Roche, K. J. and Zeger, S. L. (1995), Some Recent Developments for Regression Analysis of Multivariate Failure Time Data, *Lifetime Data Analysis*, **1**, 403-415.
- [15] Muller, A. and Stoyan, D. (2002). *Comparison Methods for Stochastic Models and Risks*. Wiley, NewYork.
- [16] Nelsen, R. B. (2006). *An Introduction to Copulas*. Lectures Notes in Statistics, 139, Springer-Verlag, New York.
- [17] Vaupel, J.W., Manton, K.G. and Stallard, E., (1979). The impact of heterogeneity in individual frailty on the dynamics of mortality. *Demography*, **16**, 439–454.
- [18] Xu, M. and Li, X.,(2008). Negative dependence in frailty models. *J Stat Plan Inference.*, **138**, 1433–1441.