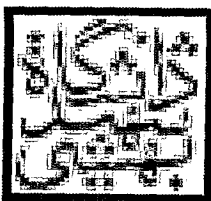


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

۸۷/۱۰۰۶۶۱

۸۷/۱۰/۱۴



دانشگاه شهید بهشتی
دانشکده علوم ریاضی
گروه آمار

رساله دوره دکتری در رشته آمار

استنباط آماری کلاسیک و بیزی در توزیعهای آمیخته متناهی

توسط:

محمد بهرامی

استاد راهنما:

دکتر محمد رضا مشکانی

استاد مشاور:

دکتر احمد خدادادی

بهمن ماه ۱۳۸۶

۱۰۷۷۲۱

کتابخانه
دانشگاه شهید بهشتی
گروه آمار

۱۳۸۷/۱۰/۰۶ - ۶

چکیده

یک توزیع آمیخته ترکیبی از دو یا چند توزیع آماری است. وقتی نمونه گیری از یک جامعه غیر همگن متشکل از دو یا چند زیرجامعه هر یک با توزیعی متفاوت انجام می گیرد، با توزیعی آمیخته سروکار پیدا می کنیم. برای مثال، توزیع زمان خرابی آمیخته ای از مؤلفه های سالم و معیوب، توزیع وزن حیوانات با گروههای سنی متفاوت یا مدت زمانی که بیماران قلبی پس از یک عمل جراحی در گروههای سنی متفاوت زندگی خواهند کرد. در این پایان نامه ابتدا توزیعهای آمیخته را در حالت کلی معرفی نموده و روشهای برآورد کردن پارامترها را در این گونه مدلها بررسی می کنیم. سپس درحالتهای خاص برای توزیعهای آمیخته بقا مانند دو توزیع نمایی و دو توزیع نرمال با به کارگیری روش نیوتون-رافسون و الگوریتم EM برآوردهای ماکسیمم درست نمایی را که در واقع برآوردهای کلاسیک پارامترها می باشند محاسبه می کنیم. در ادامه برآوردهای بیزی و بیزسلسله مراتبی پارامترهای موجود در توزیع های آمیخته را با به کارگیری الگوریتم نمونه گیری گیبز به دست آورده و در حالت کلیتر با مجهول در نظر گرفتن تعداد مؤلفه ها در یک توزیع آمیخته و با استفاده از فرایند زاد و مرگ و الگوریتم مربوط به جهش برگشت پذیر برآورد های بیزی و بیز سلسله مراتبی پارامترها و تعداد مؤلفه های نامعلوم در مدل آمیخته را مشخص می کنیم. در پایان عوامل موثر بر زمان بقای بیماران مبتلا به لوسمی حاد را که از بیمارستان امید اصفهان جمع آوری شده است بررسی کرده و از آنجا که داده های مذکور یک توزیع آمیخته را نشان می داد، با برآزش یک چگالی آمیخته شامل دو مؤلفه نرمال، میانگین زمان بقای این بیماران را به دست آورده ایم.

کلید واژه ها: الگوریتم EM ، برآورد بیزی، برآورد بیز سلسله مراتبی، برآورد ماکسیمم درست نمایی، بیماری لوسمی، تابع گاما، توزیع آمیخته، جهش برگشت پذیر، داده های کامل، فرایند زاد و مرگ، متغیرهای پنهان، مدل بقا، نمونه گیری گیبز.

تشکر و قدر دانی

سپاس پروردگار یکتا را که توفیق بندگی، تحصیل علم و انجام این پایان نامه را به من عطا فرمود. اکنون که در پرتو عنایات او موفق به انجام مراحل این رساله شده ام، بر خود لازم می دانم از تمامی عزیزانی که در این امر مرا یاری داده اند سپاسگزاری نمایم.

در ابتدا از استاد ارجمند جناب آقای دکتر مشکانی که مسئولیت راهنمایی اینجانب را برعهده داشته و در مراحل مختلف این پژوهش همواره از راهنماییهای ارزنده ایشان برخوردار بوده ام تشکر و قدر دانی می نمایم.

از استاد گرامی جناب آقای دکتر احمد خدادادی که در مقام استاد مشاور در این تحقیق مرا از راهنماییهای خود برخوردار ساختند، سپاسگزارم.

همچنین از آقایان دکتر محسن محمد زاده، دکتر حمید علوی مجد، دکتر عبدالرحیم شهلایی و دکتر محمد رضا فرید روحانی که زحمت داوری این رساله را تقبل کردند و با انتقادات و پیشنهادات خود موجب تقویت این رساله گردیدند، تشکر و قدردانی می نمایم.

از استاد گرامی جناب آقای دکتر محمد قاسم وحیدی اصل مدیر محترم وقت گروه آمار به خاطر بذل توجه و پیگیری های مداوم در انجام مراحل اداری کار سپاسگزارم.

تقدیم

به تنها دخترم، نازنین

فهرست مطالب

صفحه		عنوان
		پیشگفتار
		فصل اول
۱	تعاریف و مقدمات تحلیل بقا	
۲	مقدمه	۱-۱
۲	تعریف ها و مقدمات تحلیل بقا	۲-۱
۳	سانسور شدگی	۳-۱
۴	سانسور راست و چپ	۱-۳-۱
۵	سانسور بازه ای	۲-۳-۱
۵	تابع بقای تجربی	۴-۱
۶	مدل رگرسیون خطرات متناسب کاکس	۵-۱
۷	برآورد پارامترها در تحلیل بقا	۱-۵-۱
۹	مدلهای پارامتری	۶-۱
۱۳	برآوردهای بیزی برای پارامترهای مجهول	۷-۱
۱۳	تابع زیان	۱-۷-۱
۱۴	چگالی پیشینی و چگالی پسینی	۲-۷-۱
۱۵	برآوردگر بیزی	۳-۷-۱
۱۷	برآورد پارامترهای مجهول از دیدگاه بیز تجربی	۴-۷-۱
۲۰	توزیعهای آمیخته و برآورد های کلاسیک...	فصل دوم
۲۱	مقدمه	۱-۲
۲۲	معرفی مدل آمیخته	۲-۲
۲۶	برآورد نسبتهای آمیختگی دو توزیع معلوم	۳-۲
۲۶	برآورد جیمز	۱-۳-۲
۴۱	برآورد طبقه ای	۲-۳-۲
۴۴	برآورد بویز	۳-۳-۲
۵۱	برآورد گشتاوری	۴-۳-۲

۵۳	برآورد ماکسیمم درستنمایی	۵-۳-۲
۵۶	توزیعهای آمیخته پارامتری	۴-۲
۵۶	نمونه گیری از یک جامعه ناهمگن	۱-۴-۲
۵۸	توزیع بردار نشانگر Z_r	۲-۴-۲
۵۹	تابع بقاء در توزیعهای آمیخته	۳-۴-۲
۵۹	مدل آمیخته از دو توزیع نمایی	۵-۲
۶۰	برآوردهای ماکسیمم درستنمایی پارامترهای آمیخته دو توزیع نمایی	۱-۵-۲
۶۱	برآوردهای ماکسیمم درستنمایی با استفاده از الگوریتم EM	۲-۵-۲
۶۳	برآوردهای گشتاوری پارامترها در آمیخته دو توزیع نمایی	۳-۵-۲
۶۵	برآورد پارامترهای مدل آمیخته از دو توزیع نمایی با داده های راست سانسوریده	۶-۲
۶۶	برآورد پارامترها	۱-۶-۲
۷۵	آزمون نیکویی برازش	۷-۲
۷۷	استنباط بیزی در مورد پارامترهای مدل های آمیخته	فصل سوم
۷۸	مقدمه	۱-۳
۷۸	معرفی توزیع آمیخته با توزیعهای پیشینی مربوط به پارامترهای آن	۲-۳
۸۰	برآوردهای بیزی پارامترهای آمیخته دو توزیع نمایی	۳-۳
۸۳	برآوردهای بیز تجربی پارامترهای موجود در آمیخته دو توزیع نمایی	۴-۳
۸۴	برآوردهای بیز سلسله مراتبی پارامترهای مدل آمیخته از دو توزیع نمایی	۵-۳
۹۴	نتیجه گیری	۶-۳
۹۵	برآوردهای بیزی پارامترهای موجود در مدل های آمیخته با تعداد مولفه نامعلوم	۷-۳
۹۶	مدل آمیخته و مدل سلسله مراتبی	۱-۷-۳
۹۷	استنباط آماری به وسیله $MCMC$	۲-۷-۳
۹۸	ساختن یک زنجیره مارکوف به وسیله فرآیند نقطه ای	۳-۷-۳
۹۸	فرآیند زادو مرگ برای مولفه های یک مدل آمیخته	۴-۷-۳
۱۰۷	ساختار یک زنجیره مارکوف	۵-۷-۳

۱۲۳	فصل چهارم کاربرد توزیعهای آمیخته در پزشکی	
۱۲۴	مقدمه	۱-۴
۱۲۵	درمان لوسمی حاد	۲-۴
۱۲۵	عوامل موثر بر زمان بقای بیماران لوسمی حاد	۳-۴
۱۲۹	برازش بهترین مدل رگرسیون خطی	۱-۳-۴
۱۲۹	برآورد پارامترهای مدل	۲-۳-۴
۱۳۰	کیفیت برازش و کیفیت پیش بینی	۳-۳-۴
۱۳۰	برازش مدل مناسب	۴-۳-۴
۱۳۳	رگرسیون خطرات متناسب کاکس	۴-۴
۱۳۴	برآورد میانگین بقا در بیماران مبتلا به لوسمی حاد	۵-۴
۱۳۵	برآورد ماکسیمم درستنمایی پارامترها	۱-۵-۴
۱۳۵	تابع بقا در آمیخته دو توزیع نرمال	۲-۵-۴
۱۳۶	برآوردهای ماکسیمم درستنمایی توزیع آمیخته	۳-۵-۴
۱۳۸	آزمون نیکویی برازش توزیع آمیخته	۴-۵-۴
۱۴۰	برآوردهای گشتاوری پارامترهای موجود در توزیع آمیخته نرمال	۶-۴
۱۴۱	برآوردهای بیزی توزیع آمیخته نرمال	۷-۴
۱۴۲	برآوردهای بیزی پارامترهای توزیع آمیخته با توزیعهای پیشینی معلوم	۱-۷-۴
۱۴۵	بحث و نتیجه گیری	۲-۷-۴
۱۴۷	روش نیوتون-رافسن،.....	پیوست
۱۴۸	مقدمه	پ-۱
۱۴۸	روش نیوتون-رافسون	پ-۲
۱۵۱	الگوریتم-EM	پ-۳
۱۵۲	الگوریتم-EM در چگالیهای آمیخته	پ-۳-۱
۱۵۷	الگوریتم نمونه گیری گیبز	پ-۴
۱۶۱	مدلهای سلسله مراتبی	پ-۵
۱۶۲	همگرایی در نمونه الگوریتم گیری گیبز	پ-۶
۱۶۴	همگرایی چگالی پسین	پ-۶-۱
۱۶۹	برنامه های رایانه ای مربوط به برآورد پارامترهای توزیعهای آمیخته	پ-۷

۱۷۸	جدول داده ها	پ-۸
۱۸۰	لغت نامه فارسی به انگلیسی	
۱۸۲	مراجع فارسی	
۱۸۳	مراجع انگلیسی	
۱۸۸	چکیده انگلیسی	

پیشگفتار

در مراکز درمانی یکی از پرسش‌هایی که غالباً بیماران مطرح می‌کنند این است که "چه آینده‌ای در انتظار من است؟" به این سؤال نمی‌توان با اطمینان پاسخ داد، زیرا همیشه پیشگویی‌های لازم برای پاسخ به اینگونه سؤالات یک جزء نامعلوم دارد. معمولاً بهترین راهنمای پیشگویی‌ها تجربه سایر بیماران است. حتی وقتی پیامد نهایی را با قدری اطمینان بتوان پیش‌گویی کرد، ترتیب واقعی وقایع می‌تواند بین بیماران بسیار متفاوت باشد. مثلاً وضع بیماری را در نظر بگیرید که اخیراً برای او تشخیص ایدز داده شده است. در این مورد شانس بهبودی واقعاً صفر است، در نتیجه به دوران قابل پیش‌بینی بقا توجه می‌شود. برای پاسخ به این سؤال پزشک به پژوهش‌های انجام شده درباره پیشرفت بیماری ایدز مراجعه می‌کند. معمولاً اینگونه اطلاعات برای گروه‌های زیادی از بیماران جمع‌آوری می‌شود. با جمع‌آوری و بررسی زمان وقایع بحرانی هر بیمار (مثلاً، تاریخ تشخیص، بروز تظاهرات بیشتر بیماری و مرگ) روند پیشرفت بیماری را می‌توان به مراحل تقسیم کرد. یک روش تعیین سیر طبیعی یک بیماری، برآورد متوسط دوره زمانی موسوم به زمان بقا است، که از زمان تشخیص تا مرگ بیمار ادامه دارد. به طور کلی داده‌های بقا، داده‌هایی هستند که برای اندازه‌گیری زمان وقوع تصادفی از یک پیشامد آغازین (مانند تولد یک انسان) تا یک پیشامد نهایی (مانند مرگ همان انسان) در نظر گرفته می‌شود. برای مقایسه میانگین زمانهای بقا در داده‌های کامل در صورت نرمال بودن داده‌ها یا تبدیل آنها به نرمال از روشهای متداول آماری (مانند آزمون t - یا تحلیل واریانس) استفاده می‌شود. اما وجه تمایز داده‌های بقا با سایر داده‌ها وجود مشاهدات سانسوریده¹ است. مثلاً ممکن است پیشامد نهایی برای بعضی از موارد تحت بررسی مشاهده نشود، که این حالت را سانسور راست می‌نامند. زمان لازم تا وقوع یک پیشامد یک متغیر تصادفی پیوسته مثبت با یک تابع چگالی است. البته تعریف یک نقطه شروع نیز ضروری است که

معمولاً آن را زمان صفر در نظر می‌گیریم. مثلاً در مورد پیشامد مرگ، زمان صفر همان لحظه تولد است.

توزیع متغیر تصادفی زمان بقا را در حالت پارامتری، معمولاً به صورت یک توزیع نمایی، گاما یا وایبول در نظر می‌گیرند. اما حالاتی نیز وجود دارند که در مطالعه بقا با داده‌هایی سروکار داریم که به دلیل ناهمگن بودن جمعیت آماری نمی‌توان با توزیع‌های متداول برای آنها تابع بقای مناسبی را به دست آورد. از اینرو به توزیع‌هایی روی می‌آوریم که به صورت ترکیبی خطی از دو یا چند توزیع آماری باشند. اینگونه توزیع‌ها همان توزیع‌های آمیخته‌اند. در این رساله استنباط آماری در باره زمانهای بقایی را که به صورت آمیخته‌ای از دو یا چند توزیع نمایی پارامتری و دو یا چند توزیع نرمال پارامتری باشند بررسی می‌کنیم. ابتدا برآورد ماکسیمم درست نمایی پارامترهای چنین مدل‌های آمیخته را با نوشتن برنامه‌هایی در نرم افزار *MATLAB* با بکارگیری الگوریتم *EM*- به دست می‌آوریم.

در مرحله بعد مسئله را از دیدگاه استنباط بیزی بررسی کرده و شیوه‌های استنباطی لازم را بسط می‌دهیم. به عنوان کاربردی از روش پیشنهادی و با استفاده از نرم افزار پیشرفته *WINBUGS* برآوردهای بیزی و برآوردهای بیز سلسله مراتبی این پارامترها را در مورد داده‌های مربوط به زمان بقا بیماران لوسمی حاد که طی سال‌های ۱۳۷۳ تا ۱۳۸۲ به بیمارستان امید اصفهان مراجعه کرده‌اند، را محاسبه می‌کنیم.

فصل اول

تعاريف و مقدمات

تحليل بقا

۱-۱- مقدمه

به منظور روشن بودن تعریفها و مفهوم هایی که در فصلهای آینده به کار خواهیم برد، در این فصل به کوتاهی به بیان تعریفها و مفهوم هایی که در تحلیل بقا به کار می روند می پردازیم. در این راستا نمادگذاریها و فرمولهای را که برای توزیعها به کار می روند بیان می کنیم و خلاصه ای از روشهای برآورد را نیز توضیح می دهیم.

۱-۲- تعریف ها و مقدمات تحلیل بقا

بقا^۱: مدت زمان لازم از یک نقطه شروع (مانند تشخیص بیماری) تا زمان وقوع پیشامدی خاص (مانند وقوع مرگ) را بقا گویند.

تابع بقا: برای محاسبه احتمال بقا تا زمانی معین از تابعی استفاده می کنیم که به آن تابع بقا گفته می شود. فرض کنید T یک متغیر تصادفی با مقادیر مثبت باشد. اگر F_T تابع توزیع متغیر تصادفی T باشد، آنگاه تابع بقا که آن را با S_T نشان می دهیم به صورت زیر تعریف می شود.

$$S_T(t) = P(T > t) = 1 - F_T(t) \quad (1,2,1)$$

تابع خطر^۲: فرض کنید f_T تابع چگالی متغیر تصادفی مثبت T باشد. تابع خطر یا نرخ شکست لحظه ای به شرط زنده بودن فرد تا زمان t که آن را با $h(t)$ نشان می دهیم عبارت است از:

$$h(t) = \frac{f(t)}{S(t)} \quad (2,2,1)$$

صورت دقیق تر تابع خطر به شکل زیر است:

$$h(t) = \lim_{\Delta t \rightarrow \infty} \frac{P(t < T \leq t + \Delta t | t \leq T)}{\Delta t} \quad (3,2,1)$$

و عبارت است از احتمال اینکه نابودی در بازه $[t, t + \Delta t]$ رخ دهد به شرط آنکه فرد یا شیء مورد نظر تا زمان t زنده مانده باشد. به عبارت دیگر تابع خطر مشخص کننده احتمال شکست آنی در لحظه $T=t$ می باشد. برای زمان بقا پیوسته T می توان نوشت:

^۱-Survival

^۲-Hazard Function

$$h(t) = \frac{f(t)}{S(t)} = -\frac{S'(t)}{S(t)} \quad (4,2,1)$$

و با توجه به اینکه $S(0) = 1$ داریم:

$$S(t) = \exp\left(-\int_0^t h(u)du\right) \quad (5,2,1)$$

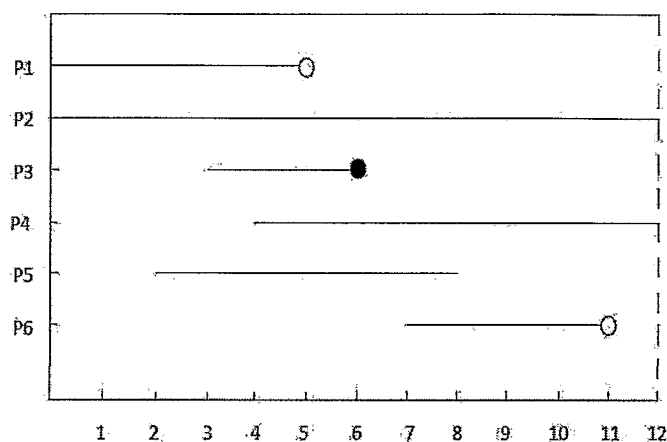
۱-۳ سانسور شدگی

یکی از مسائل مهم در تحلیل بقا وجود سانسورشدگی^۳ یا گمشدگی در داده هاست. بدین معنا که بعضی از افراد تحت مطالعه در طول دوره مورد مطالعه ممکن است به نوعی از مطالعه خارج شوند. مثلاً ممکن است دیگر تمایلی به شرکت در آزمایش پزشکی نداشته باشند یا اینکه به دلیلی قادر به ثبت وضعیت آنها بر حسب وقوع پیشامد نباشیم. برای مثال به دلیلی غیر از پیشامد تعریف شده بیمار یا فرد مورد مطالعه فوت کند و یا اینکه بیمار یا فرد مورد نظر پس از اتمام دوره پیگیری هنوز زنده باشد و پیشامد مورد نظر برایش رخ نداده باشد. بنابراین مشخص نیست که پیامد مورد نظر برای این دسته از افراد چه زمانی رخ می دهد و یا اینکه اصلاً رخ خواهد داد. تنها چیزی که می دانیم این است که تا پایان مطالعه پیشامد هنوز اتفاق نیفتاده است. این موارد را سانسور شدگی می نامیم. برای نشان دادن سانسورشدگی برای فرد i -ام در یک نمونه n تایی تحت مطالعه می توان از تابع نشانگر زیر استفاده کرد.

$$v_i = \begin{cases} 1 & T_i \leq c \\ 0 & T_i > c \end{cases}, \quad i = 1, 2, \dots, n$$

اگر پیشامد در طول دوره مورد مطالعه رخ دهد، مقدار v_i برابر یک و در غیر این صورت صفر است. متغیر T_i زمان بقای مربوط به فرد i -ام مورد مطالعه و c زمان بقای سانسوریده آن می باشد. نمودار ۱-۱ به عنوان مثال وضعیت شش بیمار مبتلا به ایدز را در یک دوره ۱۲ ماهه مورد مطالعه نشان می دهد. فرض کنید علامت 0 به منزله مرگ بیمار بر اثر بیماری ایدز و نشان دهنده مرگ بیمار بر اثر بیماری و علامت • به معنی مرگ بیمار بر اثر سایر عوامل باشد.

^۳-Censoring



نمودار ۱-۱ سانسورشدگی در داده های بقاء

بیمار P_1 : پس از ۵ ماه مطالعه بر اثر بیماری ایدز فوت شده است.

بیمار P_2 : سانسور شدگی پس از ۱۲ ماه رخ داده است.

بیمار P_3 : از ماه سوم وارد مطالعه شده و بر اثر حادثه ای غیر از بیماری ایدز در ماه ششم فوت شده

که باز هم نوعی سانسورشدگی است.

بیمار P_4 : از ابتدای ماه چهارم وارد مطالعه شده و تا پایان دوره، پیشامدی برای او اتفاق نیفتاده است.

بیمار P_5 : از ابتدای ماه دوم وارد مطالعه شده و در پایان ماه هشتم از مطالعه خارج شده است.

بیمار P_6 : از ماه هفتم وارد مطالعه شده و در پایان ماه یازدهم بر اثر بیماری ایدز فوت کرده است.

۱-۳-۱ سانسور راست و چپ

اگر از دست دادن فرد نمونه مدت زمانی پس از ورود او به مطالعه اتفاق افتد، این وضعیت

را سانسور راست می گوئیم، یعنی $T_i \geq t_0$ و اگر زمان بقا واقعی فرد در نمونه قبل از زمان شروع

مطالعه باشد، یعنی پیشامد قبل از ورود فرد به مطالعه اتفاق افتاده باشد و او را از دست داده باشیم،

به آن سانسور چپ می گوئیم یعنی: $T_i < t_0$. در نمودار ۱-۱ مشاهدات مربوط به بیماران P_4, P_5, P_6

و P_8 همگی سانسور راست هستند.

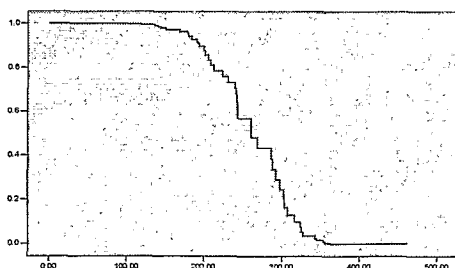
۱-۳-۲ سانسور بازه ای

این سانسورشده‌گی تلفیقی از سانسورشده‌گی راست و چپ است و در حقیقت حالتی است که پیشامد مورد نظر در یک بازه زمانی مثلاً از t_a تا t_b رخ می‌دهد. بنابراین $t_a < T_i < t_b$ که در آن T_i زمان رخ دادن پیشامد برای فرد i - ام است.

۴-۱ تابع بقای تجربی

تابع بقاء $S_T(t)$ نسبت افرادی در نمونه را نشان می‌دهد که از مبدأ زمانی خود که همان لحظه ورودشان به مطالعه است تا t واحد زمانی که پس از آن زنده مانده‌اند. مقدار آن در لحظه t برابر است با نسبت تعداد افرادی که در زمان t زنده مانده‌اند به تعداد کل افراد شرکت کننده در مطالعه. در شروع مطالعه هنوز هیچ اتفاقی رخ نداده است اما با گذشت زمان پیشامد ها یکی پس از دیگری رخ می‌دهند. اگر مدت زمان مطالعه نامتناهی باشد، بالاخره هیچ فردی زنده نخواهد ماند، یعنی خم بقا نهایتاً به صفر می‌گراید. اما چون در عمل مدت زمان مطالعه متناهی است، احتمالاً در پایان دوره هنوز تعدادی از افراد وجود دارند که پیشامد مورد نظر برای آنها اتفاق نیفتاده است و سانسور نشده‌اند.

از سوی دیگر چون تابع بقا بین دو شکست متوالی مقداری ثابت است، پس از مشاهده شکست در زمان t ، بلافاصله کاهش می‌یابد و تابع پله ای به دست می‌آید. این تابع پله ای در واقع همان تابع بقای تجربی^۴ است که آن را با $\hat{S}(t)$ نشان می‌دهیم و نمودار آن در شکل ۴-۱ رسم شده است.



شکل ۴-۱- نمودار تابع بقای تجربی در دوره مورد مطالعه

^۴-Empirical Survival Function

۵-۱ مدل رگرسیون خطرات متناسب کاکس

مانند تمام مباحث در استنباط آماری، در تحلیل بقا نیز فرض بر این است که زمان بقا برای تمام افراد یا اشیاء مورد آزمایش متغیرهای تصادفی مستقل و هم توزیع هستند. اما علاوه بر زمان بقا هر واحد، تعدادی متغیر دیگر نیز اندازه گیری می شوند که این متغیرها می توانند بر زمان بقا هر واحد تأثیر داشته باشند. این متغیرها را متغیرهای تبیینی^۵ می نامیم. در مسئله تحلیل بقا مهمترین موضوعی که به آن توجه می شود تابع خطراست که آن را با $h(t)$ نشان دادیم. واضح است که اگر زمان بقا از توزیعی خاص پیروی کند که به بردار پارامتری θ بستگی داشته باشد، آنگاه تابع خطر نیز تابعی از پارامتر θ به صورت $h(t, \theta)$ است.

چون زمان بقا یعنی T یک متغیر تصادفی پیوسته با مقادیر غیر منفی است، باید مدل‌هایی را در نظر گرفت که در آنها متغیر غیر منفی باشد، مانند توزیع نمایی، گاما و وایبول. البته در بسیاری از موارد نمی توان یک مدل پارامتری مناسب برای داده های بقا تعیین کرد، زیرا چولگی داده های بقا گریز ناپذیر است. مدل خطرات متناسب در واقع یک مدل نیمه پارامتری در تحلیل داده های بقا می باشد. یعنی علاوه بر پارامترهای توزیع، فرض می کنیم تابع خطر به متغیرهای تبیینی مانند $Z = (Z_1, \dots, Z_K)'$ نیز بستگی داشته باشد، یعنی فرم تابع خطر به صورت $h(t; Z, \theta)$ باشد. یکی از مهمترین مدل‌های نیمه پارامتری در تحلیل داده های بقا مدل کاکس (۱۹۷۲) است، که به آن مدل خطرات متناسب کاکس^۶ گفته می شود. در مدل کاکس متغیرهای تبیینی اثر ضربی بر روی تابع خطر دارد. یعنی:

$$h(t; Z) = h_0(t) \psi(Z, \beta) \quad (1, 5, 1)$$

که در آن $h_0(t)$ را تابع خطر پایه یا مبنا می نامیم و به متغیرهای تبیینی بستگی ندارد. چنانچه داشته باشیم: $Z = 0$ ، آن را شرط استاندارد می نامیم و به واحد آماری (فرد یا شیء) که برای آن $Z = 0$ است، مولفه استاندارد گفته می شود. بردار β ، بردار پارامترهای مجهول است که باید

^۵ - Explanatory Variables

^۶ - Cox Proportional Hazard Model

برآورد^۷ گردد و $\psi = (\underline{Z}, \underline{\beta})$ نیز تابعی از متغیرهای توضیحی است که به زمان بقا t بستگی ندارد. در ابتدا فرض بر این است که متغیرهای تبیینی \underline{Z} در طول زمان ثابت اند و تغییر نمی کنند. اما در مسئله کلی تر به این موضوع اشاره می شود که ممکن است برخی از متغیرهای تبیینی در طول زمان ثابت نمانند و بنابراین در تجزیه و تحلیل مؤثر باشند.

در بکارگیری مدل کاکس لازم است که نسبت خطر هر دو واحدی که در نمونه واقع اند در طول زمان ثابت باشد. فرض کنید $\underline{Z}_1, \underline{Z}_2$ بردار متغیرهای تبیینی برای دو واحد مورد نظر باشند. اگر فرض کنیم نسبت خطرات این دو واحد HR است آنگاه برآورد این نسبت در مدل خطرات متناسب عبارت است از:

$$\hat{HR} = \frac{\hat{h}(t; \underline{Z}_1)}{\hat{h}(t; \underline{Z}_2)} = \frac{\hat{h}_0(t) \cdot \psi(\underline{Z}_1, \underline{\beta})}{\hat{h}_0(t) \cdot \psi(\underline{Z}_2, \underline{\beta})} = \frac{\psi(\underline{Z}_1, \underline{\beta})}{\psi(\underline{Z}_2, \underline{\beta})} \quad (۲,۵,۱)$$

همان گونه که در (۲,۵,۱) ملاحظه می شود، این نسبت به t بستگی ندارد و نشان می دهد که تابع خطر مربوط به واحد اول متناسب با تابع خطر واحد دوم به صورت: $\hat{h}(t; \underline{z}_1) = \hat{k} \cdot \hat{h}(t; \underline{z}_2)$ می باشد. یک فرم بسیار ساده برای $\psi = (\underline{z}, \underline{\beta})$ عبارت است از:

$$\psi(\underline{Z}, \underline{\beta}) = \exp(Z_1\beta_1 + Z_2\beta_2 + \dots + Z_k\beta_k) \quad ; k \geq 1 \quad (۳,۵,۱)$$

اگر $\underline{Z} = \underline{0}$ باشد آنگاه $\psi(\underline{0}, \underline{\beta}) = 1$ و واحد در حالت استاندارد قرار دارد. اگر $\underline{\beta} = \underline{0}$ آنگاه $\psi(\underline{Z}, \underline{0}) = 1$ و بدین مفهوم است که متغیرهای تبیینی تأثیری بر مقدار خطر ندارند.

۱-۵-۱ برآورد پارامترها در تحلیل بقا

یکی از روشهای متداول و نسبتاً کارا در برآورد پارامترهای مدل‌های بقا روش ماکسیمم درست نمایی است. برآوردهای ماکسیمم درست نمایی دارای خاصیت ناوردایی^۸ بوده و تحت شرایط نظم عمومی بطور مجانبی کاراترین برآوردگر می باشند و توزیع مجانبی آنها نیز نرمال است (مود و

^۷-Estimate
^۸-Invariant Property

همکاران ۱۹۷۴). با داشتن مقادیر p متغیر تبیینی و انتخاب یک نمونه n تایی می توان n ردیف مشاهده به صورت بردار $p+1$ -بعدی $(t_i, Z_{i1}, Z_{i2}, \dots, Z_{ip})'$ را داشته باشیم که t_i زمان شکست i - امین مؤلفه می باشد. از اطلاعات فوق می توان برای برآورد پارامترهای β استفاده کرد. فرض کنید زمانهای شکست $a_1 < a_2 < \dots < a_n$ باشند. اگر برای واحدی در زمان a_j شکست رخ دهد آن را I_j می نامیم. بنابراین $I_j = i$ اگر و تنها اگر $t_i = a_j$. مجموع واحدهای در معرض خطر در زمان a_j را با $R(a_j)$ نشان می دهیم. در برخی موارد به آن مجموعه ریسک نیز گفته شده و داریم $R(a_j) = \{i | t_i \geq a_j\}$. تعداد واحد های موجود در $R(a_j)$ را با r_j نشان می دهیم. چون از $h_0(t)$ اطلاعاتی در دسترس نیست، بنابراین زمانهای a_j به تنهایی هیچ گونه اطلاعاتی از مقادیر β را نمی دهند. پس اطلاعات مربوط به مقادیر β را از I_j ها به دست می آوریم. مثلاً در زمان a_j زمانهای قبلی a_1, a_2, \dots, a_{j-1} معلوم بوده و می دانیم در هر یک از این زمانها برای کدام واحد پیشامد روی داده است. بنابراین مجموعه ریسک معلوم می شود. اکنون با شرط این که بدانیم برای یک واحد این مجموعه در زمان a_j پیشامدی رخ داده است، احتمال شرطی اینکه واحد مورد نظر، واحد i - ام باشد به صورت زیر به دست می آید.

$$\frac{h(a_j, z_{\sim i})}{\sum_{k \in R(a_j)} h(a_j, z_{\sim k})} = \frac{\psi(z_{\sim i}, \beta)}{\sum_{k \in R(a_j)} \psi(z_{\sim k}, \beta)} \quad (4,5,1)$$

همان گونه که در (۴,۵,۱) مشاهده می شود تابع خطر مبنا یعنی $h_0(a_j)$ حذف گردیده، بنابراین توزیع توأم تابع نشانگرهای I_j که در واقع همان تابع درست نمایی می باشد، به صورت زیر به دست خواهد آمد:

$$L(\beta) = \prod_{j=1}^n \frac{\psi(z_{\sim j}, \beta)}{\sum_{k \in R(a_j)} \psi(z_{\sim k}, \beta)} \quad (5,5,1)$$

اکنون می توانیم با استفاده از (۵,۵,۱) برآوردهای ماکسیمم درست نمایی برای بردار پارامتری β را به دست آوریم. ابتدا از (۵,۵,۱) لگاریتم می گیریم، داریم -

$$l(\beta) = \log L(\beta) = \sum_{j=1}^n \log \psi(z_j, \beta) - \sum_{j=1}^n \log \left(\sum_{k \in R(a_j)} \psi(z_k, \beta) \right) \quad (6,5,1)$$

حال نسبت به β_p از $l(\beta)$ مشتق می گیریم. بنابراین می توان نوشت:

$$\frac{\partial l(\beta)}{\partial \beta_p} = \sum_{j=1}^n \frac{\frac{\partial}{\partial \beta_p} \psi(z_j, \beta)}{\psi(z_j, \beta)} - \sum_{j=1}^n \frac{\frac{\partial}{\partial \beta_p} \sum_{k \in R(a_j)} \psi(z_k, \beta)}{\sum_{k \in R(a_j)} \psi(z_k, \beta)} \quad (7,5,1)$$

همچنین داریم:

$$\frac{\partial^2 l(\beta)}{\partial \beta_p \partial \beta_s} = \sum_{j=1}^n \frac{\left[\frac{\partial^2}{\partial \beta_p \partial \beta_s} \psi(z_j, \beta) \right] \psi(z_j, \beta) - \frac{\partial^2}{\partial \beta_p} \psi(z_j, \beta) \cdot \frac{\partial^2}{\partial \beta_s} \psi(z_j, \beta)}{\left(\psi(z_j, \beta) \right)^2} \quad (8,5,1)$$

$$- \sum_{j=1}^n \frac{\left[\frac{\partial^2}{\partial \beta_p \partial \beta_s} \sum_{k \in R(a_j)} \psi(z_k, \beta) \right] \left[\sum_{k \in R(a_j)} \psi(z_k, \beta) \right] - \frac{\partial^2}{\partial \beta_p} \sum_{k \in R(a_j)} \psi(z_k, \beta) \cdot \frac{\partial^2}{\partial \beta_s} \sum_{k \in R(a_j)} \psi(z_k, \beta)}{\left(\sum_{k \in R(a_j)} \psi(z_k, \beta) \right)^2}$$

اکنون دسته معادلات (7,5,1) را مساوی صفر قرار داده و برآوردهای مربوط به بردار پارامتری β را به دست آورده و با بکارگیری معادله (8,5,1) شرط ماکسیم شدن درست نمایی را بررسی می کنیم. البته واضح است که این کار به وسیله روشهای عددی امکان پذیر است.

۱-۶- مدل های پارامتری

فرض کنید داده های بقا از یک توزیع خاص تبعیت کنند و T زمان بقا یا طول عمر فردی با توزیعی مانند $F_T(t|\theta)$ و چگالی $f_T(t|\theta)$ باشد، که در آنها $\theta = (\theta_1, \theta_2, \dots, \theta_h)$ یک بردار پارامتری مجهول است و باید از روی داده های بقا برآورد گردد. در این بخش برآوردهای ماکسیم درست نمایی و برآوردهای بیزی بردار پارامتری مجهول را برای این توزیعها بدست می آوریم.

تابع بقا را با $S(t|\theta)$ نشان می دهیم و عبارت است از:

$$S(t|\theta) = P(T > t) = 1 - F(t|\theta) \quad (1,6,1)$$

در بسیاری از موارد T از یک توزیع خاص تبعیت می کند.