

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



انتشارات دانشگاه فردوسی مشهد، شماره ۳۱۲

تحلیل بقاء

نویسنده

ار. جی. میلر

مترجمان

ابوالقاسم بزرگ‌نیا - حجت رضایی پزند

۱۳۸۰

Miller, Rupert G.

میلر؛ روپرت

تحليل بقاء / نویسنده ار. جی. میلر؛ مترجمان ابوالقاسم بزرگ‌نیا، حجت رضایی بژند. - مشهد: دانشگاه فردوسی مشهد، ۱۳۸۰.

۱۸۶ ص: جدول، نمودار. - (انتشارات دانشگاه فردوسی مشهد؛ ۳۱۲).
۸۳۰۰ ریال

فهرست‌نویسی بر اساس اطلاعات فیبا.
عنوان به انگلیسی

Survival analysis.

واژه‌نامه.

کتابنامه.

۱. تحلیل داده‌های زمان افست. ۲. تحلیل بقاء (زیست‌سنجی). الف. بزرگ‌نیا، ابوالقاسم، مترجم. ب. رضایی بژند، حجت، مترجم. ج. دانشگاه فردوسی (مشهد). د. عنوان. ۱۳۱۲ -

۵۱۹/۲

۷ پ ۹ م / ۲۷۶ QA

۱۳۸۰

م ۸۰-۱۶۸۹۵

کتابخانه ملی ایران



تحليل بقاء

نوشته

ار. جی. میلر

مترجمان

ابوالقاسم بزرگ‌نیا، حجت رضایی بژند

وزیری، ۱۸۸ صفحه، ۱۰۰۰ نسخه، چاپ اول، پاییز ۱۳۸۰

امور فنی و چاپ: مؤسسه چاپ و انتشارات دانشگاه فردوسی

بها: ۸۳۰۰ ریال

(ISBN: 964-6335-83-7)

شابک ۹۶۴-۶۳۳۵-۸۳-۷

پیش‌گفتار نویسنده:

در بهار سال ۱۹۸۰، آموزش روشهای مفید و کاربردی در تحلیل بقاء به دانشجویان تحصیلات تکمیلی رشته آمار کاربردی در دانشگاه استانفورد، به این جانب پیشنهاد شد. با توجه به سفارشهای ارزشمند برادافرون، این کتاب تهیه و تنظیم شده است.

در سرفصلهای این نوشته، بیشتر از سمینارهای تحلیل بقاء در گروه غددشناسی کالیفرنیا شمالی - که به صورت زنجیره‌ای چاپ می‌شود - و همچنین از نوشته‌های آرت پترسن استفاده شده است.

بیل بران، ضمن تشویق ما، در تألیف این کتاب کمک کرده است. از نکات ارزشمند جری‌هالپرن و تری‌ترینو و نظریات ویرایشی الاین یونگ، در این کتاب استفاده شده است.

کارولاد کلو، با تایپ دقیق و سریع خود و همکاری در سرعت چاپ و ترسیمات بسیار عالی ماریاجد، به چاپ کتاب کمک نموده‌اند.

انتشار کتاب توسط مؤسسه تحقیقاتی علوم دارویی عمومی گرانت، پشتیبانی شده است.

رپرت - ج - میلر و جی‌ار

گیل گانگ

آلوارو مونوز

استانفورد، کالیفرنیا

جولای ۱۹۸۱

پیش‌گفتار مترجمان:

آمار حیاتی از شاخه‌های اصلی آمار به شمار می‌رود. بسیاری از روشهای آماری این شاخه به تحقیقات پزشکی وابسته است. رشد روزافزون تحقیقات پزشکی و نبود نوشتاری کاربردی-نظری نیاز به نگارش چنین کتابی را تأیید می‌کند. نوشتاری کوتاه و پرمحتوی به شیوه کتابهای درسی، می‌تواند این نیاز را تا اندازه‌ای برطرف کند. یکی از کتابهای بسیار مفید در این زمینه "تحلیل بقاء"، نوشته میسر است. چاپ دوباره آن در سال ۱۹۹۸، بیانگر استقبال از شیوه بیان و مطالب ارزشمند آن است. مثالهای خوب، مسائل حل شده در پایان کتاب، روشهای مناسب، ارائه قضایای لازم و اثبات روابط به صورت کوتاه، این کتاب را کاملاً کاربردی نموده است. در حالی که بیان نظری آن تا حد امکان حفظ شده است. کمی حجم و غنی بودن مطالب آن بیانگر توان بالای علمی مؤلف و انتقال مناسب آن به خوانندگان است. ذکر مراجع مختلف و مربوط به هر بحث، یکی از فواید این اثر و راه ارزنده‌ای برای دنبال کردن مطالب جانبی، اثبات قضایا و پژوهش در آمار کاربردی است. با توجه به موارد بالا به ترجمه این کتاب ارزشمند اقدام شده است. امید است بدین وسیله در پیشبرد آمار کاربردی خدمت کوچکی کرده باشیم.

مطالب این کتاب در حد دو تا سه واحد درسی برای دانشجویان کارشناسی رشته‌های آمار، داروسازی و رشته‌های مهندسی است. علاوه بر آن مرجعی مفید و کوتاه برای دانشجویان کارشناسی ارشد، دکترای پژوهشگران مختلف علوم پزشکی، مهندسی و علوم است. محققین و دانشجویان رشته‌های داروسازی، بیمه‌گری و آمارحیاتی بیشترین استفاده را از این کتاب می‌برند.

در پایان لازم می‌دانیم از زحمات آقایان مهندس تبرایی، کدکنی و فنائی (چاپخانه دانشگاه مشهد)، همچنین آقای مهندس امین تشکر و قدردانی نماییم. چاپ زیبا، سریع و دقیق آقای فردوسی مکان (آمارپژوه) کمک مهمی در ارائه این اثر داشته است. از ایشان قدردانی و تشکر می‌شود. طبیعی است که این ترجمه خالی از اشتباه و کاستی نیست، تذکر صاحب‌نظران ما را خوشحال و چاپ‌های بعدی را بهبود خواهد بخشید.

ابولقاسم بزرگ‌نیا

حجت رضایی پزند

فهرست مطالب

صفحه	عنوان
۱۱	فصل اول: مقدمه‌ای بر مفاهیم بقاء
۱۲	۱.۱ تابع بقاء و نرخ شکست
۱۳	۲.۱ انواع برش
۱۳	۱.۲.۱ برش نوع اول
۱۳	۲.۲.۱ برش نوع دوم
۱۴	۳.۲.۱ برش نوع سوم (برش تصادفی)
۱۵	۴.۲.۱ انواع دیگر برش
۱۶	مثال: کودکان آفریقایی
۱۶	نمادهای مورد استفاده
۱۹	فصل دوم: الگوهای پارامتری
۱۹	۱ توزیع‌ها
۱۹	۱.۱ توزیع نمایی
۱۹	۲.۱ توزیع گاما (دو پارامتری)
۲۰	۳.۱ توزیع وایبل

۲۰	۴.۱ توزیع رایلی
۲۱	۵.۱ توزیع لوگ نرمال (دو پارامتری)
۲۲	۶.۱ توزیع پارتو
۲۲	۷.۱ IFRA و IFR
۲۳	۲ برآورد کردن
۲۳	۱.۲ حداکثر درست‌نمایی
۲۴	روشهای نیوتن-رافسون و چوب خطی
۲۶	بازه‌های اطمینان و آزمونها
۲۷	مثال ۱: نمایی
۳۰	روش دلتا
۳۱	مثال ۲: وایبل
۳۲	برآورد $S(t)$
۳۳	۲.۲ ترکیب خطی آماره‌های مرتب
۳۵	توزیعهای فرین
۳۵	۳.۲ برآوردگرهای دیگر
۳۷	برآوردهای بیزی
۳۷	۳ الگوهای رگرسیونی
۳۸	۴ الگوهایی با کسرهای بقاء
۳۸	۱.۴ نمونه با حجم واحد
۳۸	۲.۴ رگرسیون
۴۲	فصل سوم: روشهای ناپارامتری (یک نمونه)
۴۲	۱ جدولهای طول عمر

۴۳	۱.۱ روش کاهش نمونه
۴۳	۲.۱ روش بیمه‌گری
۴۴	۳.۱ واریانس $\hat{S}(\tau_k)$
۴۴	۴.۱ انواع جدولهای طول عمر
۴۵	۲ برآوردگر حدی حاصل ضرب کاپلان-میر
۴۷	مثال: مطالعه درمان AML
۴۹	واریانس $\hat{S}(t)$
۴۹	۱.۲ الگوریتم تجدیدنظر در توزیع به راست
۵۰	۲.۲ خودسازگاری
۵۳	الگوریتم خودسازگاری
۵۴	۳.۲ برآوردگر حداکثر درست‌نمایی تعمیم‌یافته
۵۵	۴.۲ سازگاری
۵۸	۵.۲ نرمال مجانبی
۶۰	۳ برآوردگرهای تابع نرخ شکست
۶۱	نرمال مجانبی
۶۲	۴ برآوردگرهای تنومند
۶۳	۱.۴ میانگین
۶۴	۲.۴ برآوردگرهای L -
۶۵	۳.۴ برآوردگرهای M -
۶۶	۴.۴ میانه
۶۷	۵ برآوردگرهای بی‌زی
۶۹	برآوردگرهای تجربی بی‌زی

۷۱	فصل چهارم: روشهای ناپارامتری (دو نمونه)
۷۱	مثال: آزمایش بالینی فرضی
۷۱	۱ آزمون گهان
۷۴	۱.۱ میانگین و واریانس U
۷۵	۲.۱ روش محاسباتی مانند برای $\text{Var}_{0,p}^*(u)$
۷۶	۳.۱ مثال
۷۷	۴.۱ واریانس تحت H_0
۷۹	۲ آزمون مانند - هانزل
۷۹	۱.۲ جدول 2×2 ساده
۸۲	۲.۲ دنباله‌ای از جدولهای 2×2
۸۴	۳.۲ مثال
۸۴	۴.۲ نرمال مجانبی
۸۷	۳ رده آزمونهای تارون-وایر
۸۸	مثال
۸۸	۴ آزمون افرون
۹۱	فصل پنجم: روشهای ناپارامتری (k نمونه)
۹۱	۱ آزمون گهان تعمیم یافته (برسلو)
۹۲	انواع آزمونها
۹۴	۱.۱ ماتریس کوواریانس جایگشت
۹۵	۲.۱ توزیع تحت H_0
۹۵	۲ آزمون مانند - هانزل تعمیم یافته (تارون و وایر)

۹۶	انواع آزمونها
۹۹	فصل ششم: روشهای ناپارامتری: رگرسیون
۹۹	۱ الگوهای نرخ شکست متناسب کاکس
۱۰۰	۱.۱ تحلیل درست‌نمایی شرطی
۱۰۴	۲.۱ بررسی درست‌نمایی شرطی
۱۰۴	درست‌نمایی حاشیه‌ای برای رتبه‌ها
۱۰۶	درست‌نمایی جزئی
۱۰۷	۳.۱ بررسی نرمال مجانبی
۱۰۸	۴.۱ برآورد $S(t; \underline{x})$
۱۱۰	۵.۱ داده‌های گسسته یا طبقه‌ای
۱۱۲	۶.۱ متغیرهای وابسته به زمان
۱۱۳	۷.۱ مثال ۱: داده‌های پیوند قلب استانفورد
۱۱۳	۸.۱ مثال ۲: فرزند خواندگی و آبستنی
۱۱۳	۲ الگوهای خطی
۱۱۳	الگوهای زمانی شتاب داده شده
۱۱۴	۱.۲ آزمونهای رتبه خطی
۱۱۶	۲.۲ برآوردگرهای کمترین مربعات
۱۱۶	برآوردگرهای میلر
۱۱۹	برآوردگر باکلی - جیمز
۱۲۱	برآوردگر کول - سوسارلا - وان‌رایزن
۱۲۲	مثال: داده‌های پیوند قلب استانفورد

۱۲۹	فصل هفتم: نیکویی برازش
۱۲۹	۱ روشهای ترسیمی
۱۳۰	۱.۱ یک نمونه
۱۳۱	۲.۱ دو نمونه تا k نمونه
۱۳۱	مثال: مطالعه DNCB
۱۳۲	۳.۱ رگرسیون
۱۳۴	۲ آزمونها
۱۳۴	۱.۲ یک نمونه
۱۳۶	۲.۲ رگرسیون
۱۳۹	فصل هشتم: مباحث مختلف
۱۳۹	۱ برآوردگر دو منغیری کاپلان-مایر
۱۴۰	۲ نرخ شکست رقیب
۱۴۱	۳ برش وابسته
۱۴۲	۴ روش جک‌نایف و بوت‌استرپ
۱۴۵	فصل نهم: مسائل
۱۶۳	واژه‌نامه
۱۷۱	مراجع

فصل اول

مقدمه‌ای در مفاهیم بقاء

تحلیل بقاء از نظر علم آمار، عبارت است از: استفاده از فنون مختلف آماری در تحلیل متغیرهای تصادفی نامنفی. نوعاً، مقدار این متغیر تصادفی زمان شکست یک مؤلفه فیزیکی (مکانیکی یا الکتریکی) یا زمان مرگ یک واحد زنده (سلول، بیمار، حیوان و غیره) است. ممکن است این متغیر، زمان یادگیری یک مهارت باشد، یا حتی امکان دارد به زمان هیچ ارتباطی نداشته باشد. برای مثال، متغیر می‌تواند مبلغ پرداختی یک شرکت بیمه در وضعیت خاصی باشد.

در برخی موارد، یک بیمار بهبود یافته و مبلغ کل پرداختی بیمه او معلوم است. در موارد دیگر بیماری هنوز ادامه دارد و تنها مبلغ پرداختی تا آن زمان معلوم است. پیشینه و منشأ تحلیل بقاء کارهایی است که در گذشته در مورد جداول طول عمر انجام شده است. شکل جدید تحلیل بقاء از نیم قرن گذشته با کاربردهای مهندسی مورد توجه و مطالعه قرار گرفته است.

در جنگ جهانی دوم، علاقه‌مندی بسیاری به قابلیت اعتماد تجهیزات جنگی به وجود آمد. این علاقه‌مندی به قابلیت اعتماد، موضوع را به کارهای نظامی و فرآوردهای تجاری کشاند. قبلاً بیشتر تحقیقات آماری مورد استفاده در مهندسی بر الگوهای پارامتری متمرکز بود. در دو دهه اخیر تحقیقات آزمایشگاهی پزشکی افزایش یافته است. در این آزمایشها بیشتر به روشهای پارامتری توجه شده است. این کتاب هر دو روش پارامتری و ناپارامتری را مورد بحث قرار می‌دهد. ولی بیشتر به روشهای جدید ناپارامتری و کاربرد آنها در تحقیقات پزشکی تأکید شده است.

REFERENCE

Leavitt and Olshen, unpublished report (1974), give the insurance example.

۱.۱ تابعهای بقاء و نرخهای شکست

فرض کنید متغیر تصادفی $T > 0$ دارای تابعهای چگالی $f(t)$ و توزیع $F(t)$ باشد. تابع بقاء $S(t)$ ، به صورت زیر تعریف می‌شود:

$$S(t) = 1 - F(t) = P(T > t)$$

نرخ شکست یا تابع شکست $\lambda(t)$ ، به شکل زیر تعریف می‌شود:

$$\lambda(t) = \frac{f(t)}{1 - F(t)}$$

در علم امراض مسری، تابع $\lambda(t)$ را نرخ مرگ و میر نامند. نرخ شکست به صورت زیر قابل تفسیر است:

$$\lambda(t) dt \cong P(t < T < t + dt | T > t)$$

$$= P\left\{ \text{بقاء بیشتر از } t \mid \text{پیشامد در بازه } (t, t+dt) \text{ رخ دهد} \right\}$$

از $\lambda(t)$ انتگرال می‌گیریم، داریم:

$$\int_0^t \lambda(u) du = \int_0^t \frac{f(u)}{1 - F(u)} du = -\log \left[1 - F(u) \right]_0^t$$

$$= -\log[1 - F(t)] = -\log S(t)$$

در نتیجه می‌توان رابطه مهم زیر را نتیجه گرفت:

$$S(t) = e^{-\int_0^t \lambda(u) du}$$

توجه نمایید که $F(\infty) = 1$ (یعنی $S'(\infty) = 0$) خواهد بود، اگر و فقط اگر $\int_0^{\infty} \lambda(u) du = \infty$ باشد. مفاهیم بالا را می‌توان برای حالاتی که T دارای تابع چگالی نیست نیز تعمیم داد. به عبارت دیگر هنگامی که تابع توزیع F گسسته باشد. در بحث‌هایمان F را پیوسته در نظر می‌گیریم و هر جا لازم باشد، مفاهیم و رابطه‌ها را برای حالت گسسته اصلاح می‌کنیم.

۲.۱ انواع برشها (Censoring)

مطالب ارائه شده در این کتاب در بیشتر منابع آماری موجوداند. آنچه تحلیل بقاء را از سایر مباحث آماری جدا می‌کند، عمل برش است. به‌طور کلی، یک مشاهده بریده شده، فقط شامل قسمتی از اطلاعات مربوط به متغیرهای تصادفی مورد نظر است. در این کتاب سه نوع برش را مورد بحث و توجه قرار می‌دهیم.

فرض کنید T_1, T_2, \dots, T_n ، متغیرهای مستقل و هم‌توزیع (iid) با تابع توزیع F باشند. برشها را به شرح زیر بررسی می‌کنیم:

۱.۲.۱ برش نوع اول

فرض کنید t_c عددی ثابت باشد، که آن را زمان برش می‌نامیم. به جای مشاهده T_1 تا T_n (متغیرهای مورد توجه)، فقط متغیرهای Y_1 تا Y_n را به شرح زیر مشاهده می‌کنیم:

$$Y_i = \begin{cases} T_i & T_i \leq t_c \\ t_c & T_i > t_c \end{cases}$$

توجه کنید که تابع توزیع Y در $y = t_c$ ، دارای جرم مثبت $P(T > t_c) > 0$ است.

۲.۲.۱ برش نوع دوم

فرض کنید $r < n$ عددی ثابت و $T_{(1)} < T_{(2)} < \dots < T_{(n)}$ ، آماره‌های مرتب T_1 تا T_n باشند. اگر r امین شکست رخ دهد، مشاهده در این حالت پایان می‌پذیرد. بنابراین، فقط می‌توانیم $T_{(1)}$ تا $T_{(r)}$ را مشاهده کنیم. مشاهدات مرتب شده کامل به شرح زیراند:

$$\begin{aligned} Y_{(1)} &= T_{(1)} \\ Y_{(2)} &= T_{(2)} \\ &\vdots \\ Y_{(r)} &= T_{(r)} \\ Y_{(r+1)} &= T_{(r)} \\ &\vdots \\ Y_{(n)} &= T_{(r)} \end{aligned}$$

هر دو نوع، برش اول و دوم را در مهندسی به کار می‌برند. برای مثال، تعدادی

لامپ یا ترانزیستور موجود است. تمام آنها را در زمان $t=0$ مورد آزمایش قرار می‌دهیم و زمان شکست را یادداشت می‌کنیم. ممکن است طول عمر بعضی از ترانزیستورها طولانی باشد و نخواهیم به مدت طولانی منتظر بمانیم تا آزمایش به پایان برسد. بنابراین، امکان دارد، آزمایش را در زمان t_c ، که از قبل تعیین شده است، پایان دهیم. در این صورت برش نوع اول انجام شده است. در حالت دیگر از قبل ندانیم چه زمانی برای برش مناسب است. لذا، تصمیم بگیریم، که در انتظار بمانیم تا نسبت معینی از ترانزیستورها، مانند: $\frac{r}{n}$ بسوزند. در این جا، برش نوع دوم رخ داده است.

۳.۲.۱ برش نوع سوم (برش تصادفی)

این نوع برش که به برش تصادفی نیز معروف است، به صورت زیر تعریف می‌شود: فرض کنید، C_1 تا C_n ، متغیرهای iid با تابع توزیع G باشند. C_i زمان برش مربوط به T_i است. در این حالت فقط می‌توانیم $(Y_1, \delta_1), \dots, (Y_n, \delta_n)$ را مشاهده کنیم. که در آن Y_i و δ_i به شرح زیر تعریف می‌شوند:

$$Y_i = \min(T_i, C_i) = T_i \wedge C_i$$

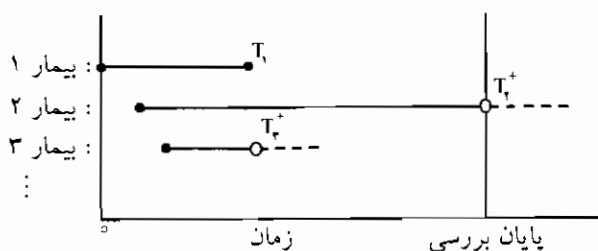
$$\delta_i = I(T_i \leq C_i) = \begin{cases} 1 & T_i \leq C_i & (\text{یعنی: } T_i \text{ بریده نمی‌شود}) \\ 0 & T_i > t_c & (\text{یعنی: } T_i \text{ بریده می‌شود}) \end{cases}$$

توجه شود که Y_1 تا Y_n ، متغیرهای iid با توزیع H هستند. علاوه بر این، δ_1 تا δ_n ، شامل اطلاعات مربوط به برش هستند. تذکر این که، در برش نوع اول و دوم می‌توانستیم ببینیم که چه مشاهداتی بریده می‌شوند. به این علت، تعریف δ_i ها به طور آشکار ضروری نبود. برش تصادفی در کارهای پزشکی و در مورد افراد یا آزمایشهای بالینی استفاده می‌شود. در یک آزمایش بالینی، ممکن است بیماران برای معالجه در زمانهای مختلف مراجعه کنند. سپس، هر بیمار به یکی از روشهای ممکن معالجه شود. حال اگر بخواهیم طول عمر آنها را مطالعه کنیم، عمل برش به یکی از صورتهای زیر انجام می‌شود:

۱. **عدم مراجعه:** ممکن است بیمار تصمیم بگیرد به پزشک دیگری مراجعه کند. در نتیجه، دیگر او را نخواهیم دید.

۲. **قطع معالجه:** ممکن است به علت عوارض بد جانبی، درمان آن را متوقف کنیم. یا این که، ممکن است بیمار هنوز در تماس باشد ولی از ادامه معالجه خودداری کند.

۳. اتمام بررسی: نمودار زیر یک بررسی ممکن را تشریح می‌کند:



در این‌جا، بیمار (۱) در زمان $t=0$ مورد بررسی قرار گرفته و در زمان T_1 فوت شده است. در نتیجه، یک مشاهده بریده نشده به دست آمده است. بیمار (۲) مورد بررسی قرار گرفته و در پایان بررسی هنوز زنده است. در نتیجه یک مشاهده بریده نشده T_1^+ به دست آمده است. بالاخره، بیمار (۳) مورد بررسی قرار گرفته و قبل از پایان بررسی، معالجه را قطع نموده است. در نتیجه، یک مشاهده بریده شده T_1^+ به دست آمده است. دربارهٔ برش تصادفی، فرض اساسی زیر را در نظر می‌گیریم:

فرض: متغیرهای تصادفی C_i و T_i مستقل‌اند.

بدون این فرض نتایج کمتری در دسترس است. معقول به نظر می‌رسد که فرض کنیم: مطالعه در زمانهای تصادفی شروع، به تصادف بررسی و قطع می‌شود. با این وجود، اگر دلیلی برای قطع بررسی در جریان معالجه وجود داشته باشد، در این صورت ممکن است بین T_i و C_i ، رابطه‌ای وجود داشته باشد.

۴.۲.۱ انواع دیگر برش

انواع دیگر برش در منابع مختلف مورد بررسی قرار گرفته‌اند. انواع قبلی برش به راست و چپ تقسیم می‌شوند. اگر متغیر مورد مطالعه بسیار بزرگ باشد و نخواهیم آن را به‌طور کامل مشاهده کنیم، آن را "راست برش" نامند. به‌طور مشابه "چپ برش" قابل تعریف است. برای مثال، در برش تصادفی چپ، تنها می‌توان $(Y_1, \varepsilon_1), \dots, (Y_n, \varepsilon_n)$ را مشاهده نمود، که Y_i و ε_i به شرح زیرند:

$$Y_i = \max(T_i, C_i) = T_i \vee C_i$$

$$\varepsilon_i = I(C_i \leq T_i)$$

مثال. کودکان آفریقایی: در این مثال هر دو برش راست و چپ انجام می‌شود. یک روان‌پزشک می‌خواهد سن گروه خاصی از کودکان آفریقایی را برای انجام یک عمل ویژه بداند. وقتی به روستای مورد نظر می‌رسد، تعدادی از بچه‌ها از قبل می‌دانند آن عمل را چگونه انجام دهند. در نتیجه، این کودکان مشاهدات بریده شده چپ را ارائه می‌کنند. بعضی از کودکان یادگیری را شروع کرده‌اند و سن یادگیری آنها را می‌توان ثبت کرد. هنگام بازگشت پژوهشگر، هنوز برخی از کودکان روستایی این عمل را یاد نگرفته‌اند. آنها مشاهدات بریده شده راست را نتیجه می‌دهند.

REFERENCE

Leiderman et al., Nature (1974).
Turnbull, JASA (1974).

برشهای راست و چپ؛ هر دو حالت‌های خاصی از برش فاصله‌ای هستند، که در آنها فقط مشاهداتی مورد توجه‌اند که در یک بازه قرار گیرند. اگر T_i برش تصادفی راست باشد، مشاهدات T_i در بازه $[C_i, \infty)$ قرار می‌گیرند. همچنین، اگر T_i برش تصادفی چپ باشد، مشاهدات T_i در بازه $(-\infty, C_i]$ قرار می‌گیرند.

در برابر برش (Censoring) فاصله‌ای، قطع مشاهدات (Truncation) قرار دارد. در این‌جا، اگر متغیر مورد علاقه در خارج بازه معینی قرار گیرد، از آن چشم‌پوشی می‌شود. برای مثال، فرض کنید بخواهیم توزیع و میانگین اندازه یک عنصر داخل سلول را تعیین کنیم. طبیعی است، به علت محدودیت وسایل اندازه‌گیری، اگر اندازه عنصر، کوچکتر از اندازه تعیین شده باشد، قابل تشخیص نیست.

نمادهای مورد استفاده:

متغیر تصادفی T_i را برای زمان بقاء و C_i را برای زمان برش در نظر می‌گیریم. متغیرهای تصادفی مشاهده شده، عبارت‌اند از: $Y_i = T_i \wedge C_i$ و $\delta_i = I(T_i \leq C_i)$. نمادهای دیگر مورد استفاده به شرح زیراند:

(۱) $X_i \sim F$: زمان بقاء و $Y_i \sim G$: زمان برش است، متغیرهای مشاهده‌ای به صورت $Z_i \simeq X_i \wedge Y_i \sim H$ و $\delta_i = I(X_i \leq Y_i)$ است.

این نمادها به نظر خوب می‌رسند. زیرا به راحتی می‌توان بین متغیر تصادفی و تابعهای توزیع تفاوت قائل شد. ولی در کاربردهای رگرسیونی که بعداً به آن می‌پردازیم،

X را برای متغیر مستقل به کار می‌بریم.

(۲) $X_i^0 \sim F^0$: زمان بقاء و $Y_i \sim G$: زمان برش است. متغیرهای مشاهده‌ای به صورت $Z_i \cong X_i^0 \wedge Y_i$ و $\delta_i = I(X_i^0 \leq Y_i)$ است.

در گزارش اعداد واقعی، مناسبتر است مشاهده بریده نشده را با (T_i) و مشاهده بریده شده را با (T_i^+) ، نشان دهیم. در این صورت ممکن است داده‌ها به شکل زیر باشند. ۵، ۱۱⁺، ۶٫۵، ۱۴⁺ و اعداد ۵ و ۶٫۵ بریده نشده و ۱۱ و ۱۴ بریده شده‌اند.

فصل دوم

الگوهای پارامتری

۱ توزیعها

۱.۱ توزیع نمایی:

الگوی نمایی دارای نرخ شکست ثابت است. یعنی: $\lambda(t) = \lambda > 0$ ، در نتیجه موارد

زیر را داریم:

$$\int_0^t \lambda(u) du = \lambda t$$

$$S(t) = e^{-\int_0^t \lambda(u) du} = e^{-\lambda t}$$

$$f(t) = -\frac{d}{dt} S(t) = \lambda e^{-\lambda t}$$

$$E(T) = \frac{1}{\lambda} \quad \text{و} \quad \text{Var}(T) = \frac{1}{\lambda^2}$$

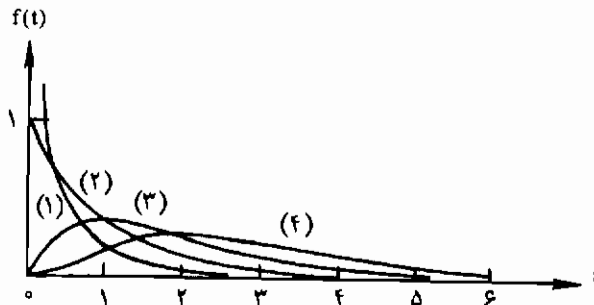
۲.۱ توزیع گاما (دو پارامتری)

الگوی گاما تعمیمی از الگوی نمایی است:

$$f(t) = \frac{\lambda^\alpha}{\Gamma(\alpha)} t^{\alpha-1} \cdot e^{-\lambda t} \quad \alpha, \lambda > 0$$

امید ریاضی و واریانس آن $E(T) = \frac{\alpha}{\lambda}$ و $\text{Var}(T) = \frac{\alpha}{\lambda^2}$ هستند. منحنی تابع چگالی گاما

با λ ثابت و به ازاء α های مختلف در نمودار (۱)، ارائه شده است:



نمودار ۱. تابع چگالی گاما برای $\lambda = 1$ و $\alpha = \frac{1}{4}$: (۱)، $\alpha = 1$: (۲)، $\alpha = 2$: (۳)، $\alpha = 3$: (۴).

متأسفانه الگوی گاما، عبارت بسته‌ای برای $S(t)$ و $\lambda(t)$ ندارد.

$$S(t) = 1 - \int_0^t f(u) du = 1 - \left(\frac{\text{تابع گاما ناقص}}{\text{تابع گاما کامل}} \right)$$

۳.۱ توزیع وایبل

الگوی وایبل تعمیم دیگری از توزیع نمایی است.

$$S(t) = e^{-(\lambda t)^\alpha} \quad \alpha, \lambda > 0$$

در نتیجه در موارد زیر به دست می‌آید:

$$\int_0^t \lambda(u) du = (\lambda t)^\alpha$$

$$\lambda(t) = \alpha \lambda (\lambda t)^{\alpha-1}$$

$$f(t) = \lambda(t) \cdot S(t) = \alpha \lambda (\lambda t)^{\alpha-1} \cdot e^{-(\lambda t)^\alpha}$$

در الگوی وایبل، $E(T)$ و $\text{Var}(T)$ دارای عبارت بسته‌ای نیستند. با این حال، شکل ساده‌ی $\lambda(t)$ و $S(t)$ ، الگوی وایبل را در تحلیل بقاء مفید ساخته است. نمودار (۲)، منحنی الگوی وایبل را نشان می‌دهد.

۴.۱ توزیع رایلی

این توزیع دارای مشخصات زیر است:

$$\lambda(t) = \lambda_0 + \lambda_1 t$$

$$\int_0^t \lambda(u) du = \lambda_0 t + \frac{1}{2} \lambda_1 t^2$$

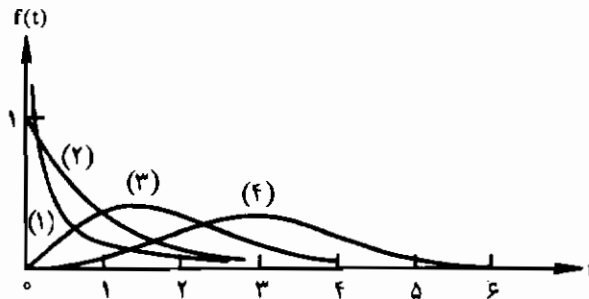
$$S(t) = \exp(-\lambda_0 t - \frac{1}{2} \lambda_1 t^2)$$

$$f(t) = (\lambda_0 + \lambda_1 t) \exp(-\lambda_0 t - \frac{1}{2} \lambda_1 t^2)$$

گشتاورهای رایلی عبارت بسته‌ای ندارند.

نرخ شکست خطی را می‌توان به صورت چندجمله‌ای عمومیت داد؛ داریم:

$$\lambda(t) = \sum_{i=0}^p \lambda_i t^i$$



نمودار ۲. چگالی و اویل در ازای (۱): $\lambda = 1$ و $\alpha = \frac{1}{2}$ ، (۲): $\lambda = 1$ و $\alpha = 1$ ، (۳): $\lambda = 0.5$ و $\alpha = 2$ ، (۴): $\lambda = 0.3$ و $\alpha = 3$.

۵.۱ توزیع لگ نرمال (دو پارامتری)

در این توزیع فرض می‌شود $\log T_i$ به صورت نرمال توزیع شده باشد.

$$\log T_i \approx N(\mu, \sigma^2)$$

$\lambda(t)$ و $S(t)$ دارای صورت بسته‌ای نیستند.

$$S(t) = 1 - P(T < t) = 1 - P\{\log T < \log t\}$$

$$= 1 - P\left\{\frac{\log T - \mu}{\sigma} < \frac{\log t - \mu}{\sigma}\right\} = 1 - \Phi\left(\frac{\log t - \mu}{\sigma}\right)$$

توزیع لگ نرمال می‌تواند برای داده‌های بریده نشده مفید باشد. یک تبدیل لگاریتمی

داده‌ها را به الگوی خطی استاندارد تبدیل می‌کند.

۶.۱ توزیع پارتو

این توزیع به شرح زیر است:

$$S(t) = \left(\frac{a}{t}\right)^\alpha I_{(a, \infty)}(t) \quad \alpha, a > 0$$

در نتیجه:

$$f(t) = \frac{\alpha a^\alpha}{t^{\alpha+1}} I_{[a, \infty)}(t)$$

$$\lambda(t) = \frac{\alpha}{t} I_{[a, \infty)}(t)$$

گشتاورهای این توزیع به آسانی محاسبه می‌شوند، ولی ممکن است نامتناهی باشند.

۷.۱ IFRA و IFR

اگر F یا f دارای تابع نرخ شکست صعودی باشند، آنها را IFR نامند. در این صورت $\lambda(t)$ تابعی صعودی است. اگر F یا f دارای میانگین نرخ شکست صعودی باشند، آنها را IFRA نامند. در این صورت تابع زیر صعودی است:

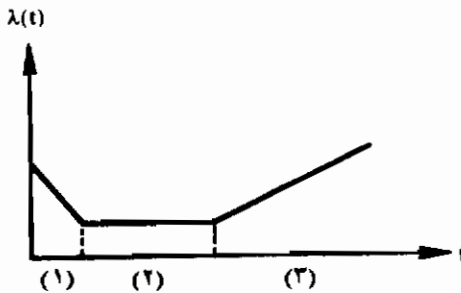
$$\frac{1}{t} \int_0^t \lambda(u) du$$

تعاریف مشابهی برای توابع DFR و DFRA وجود دارد.

DFR	IFR	FR ثابت
وایل ($\alpha < 1$)	وایل ($\alpha > 1$)	نمایی
گاما ($\alpha < 1$)	گاما ($\alpha > 1$)	
رایلی ($\lambda_1 < 0$)	رایلی ($\lambda_1 > 0$)	
پارتو ($t > a$)		

مفاهیم توزیعهای IFR و IFRA در کاربردهای مهندسی مفیدند. به ویژه در مطالعه دستگاههای متشکل از چندین مؤلفه. این دو معمولاً در آمار حیاتی کاربرد زیادی ندارند.

برای مثال در بررسی‌های علم بیماری‌های مسری، خطر بقاء دراز مدت، معمولاً به شکل وام حتم است، که در آن، زمان به سه دوره همانند شکل زیر تقسیم می‌شود:



- (۱): دوره خردسالی
 (۲): دوره بزرگسالی
 (۳): دوره کهنلت

REFERENCE

Barlow and Proschan, Statistical Theory of Reliability and Life Testing (1975).

۲ برآورد کردن

۱.۲ حداکثر درست‌نمایی

الگوی برش تصادفی را در نظر می‌گیریم. (توجه شود که در این جا برش نوع اول را با فرض $c_i \equiv t_c$ ، در نظر می‌گیریم. درست‌نمایی برای برش نوع دوم، مشابه نوع اول است. با این تفاوت که به علت منظور کردن ترتیب یک ضریب ثابت وجود دارد). تابع درست‌نمایی زوج مرتب (y_i, δ_i) به شرح زیر است:

$$L(y_i, \delta_i) = \begin{cases} f(y_i) & \delta_i = 1 \quad (\text{بدون برش}) \\ S(y_i) & \delta_i = 0 \quad (\text{بریده شده}) \end{cases}$$

$$= f(y_i)^{\delta_i} \cdot S(y_i)^{1-\delta_i}$$

تابع درست‌نمایی نمونه کامل نیز به شرح زیر است:

$$L = L(y_1, \dots, y_n, \delta_1, \dots, \delta_n) = \prod_{i=1}^n L(y_i, \delta_i)$$

$$= \left(\prod_u f(y_i) \right) \left(\prod_c S(y_i) \right)$$

در واقع تابعهای درست‌نمایی برای برش تصادفی، به شرح زیر است:

$$L(y_i, \delta_i) = \begin{cases} f(y_i)[1 - G(y_i)] & \delta_i = 1 \\ g(y_i)S(y_i) & \delta_i = 0 \end{cases}$$

$$L = \left(\prod_{\mathbf{u}} f(y_i) \right) \left(\prod_{\mathbf{c}} S(y_i) \right) \left(\prod_{\mathbf{c}} g(y_i) \right) \left(\prod_{\mathbf{u}} [1 - G(y_i)] \right)$$

اگر فرض شود که زمان برش با زمان بقاء ارتباط ندارد، دو حاصل ضرب آخری $\prod_{\mathbf{u}} [1 - G(y_i)]$ و $\prod_{\mathbf{c}} g(y_i)$ شامل پارامترهای مجهول نیستند. در نتیجه، این دو را می‌توان در بیشینه کردن L ، ثابت فرض کرد.

فرض کنید، $\underline{\theta} = (\theta_1, \dots, \theta_p)'$ بردار پارامتر باشد. محاسبه $\max L(\underline{\theta})$ با محاسبه جواب $\hat{\underline{\theta}}$ معادلات درست‌نمایی زیر یکسان است:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \log L(\underline{\theta}) &= \sum_{i=1}^n \frac{\partial}{\partial \theta_j} \log L_{\theta}(y_i, \delta_i) \\ &= \sum_{\mathbf{u}} \frac{\partial}{\partial \theta_j} \log f_{\theta}(y_i) + \sum_{\mathbf{c}} \frac{\partial}{\partial \theta_j} \log S_{\theta}(y_i) = 0 \quad j=1, 2, \dots, p \end{aligned}$$

معمولاً، محاسبه جواب به کمک رایانه و روشهای عددی امکان‌پذیر است.

روشهای نیوتن-رافسون و چوب‌خطی

نمادهای زیر را تعریف می‌کنیم:

$$L_i(\underline{\theta}) = L_{\theta}(y_i, \delta_i) \quad i=1, 2, \dots, n$$

$$\begin{aligned} \frac{\partial}{\partial \underline{\theta}} \log L(\underline{\theta}) &= \left(\frac{\partial}{\partial \theta_1} \log L(\underline{\theta}), \dots, \frac{\partial}{\partial \theta_p} \log L(\underline{\theta}) \right)' \\ \frac{\partial^2}{\partial \underline{\theta}^2} \log L(\underline{\theta}) &= \begin{pmatrix} \frac{\partial^2}{\partial \theta_1 \partial \theta_1} \log L(\underline{\theta}) & \dots & \frac{\partial^2}{\partial \theta_1 \partial \theta_p} \log L(\underline{\theta}) \\ \vdots & & \vdots \\ \frac{\partial^2}{\partial \theta_p \partial \theta_1} \log L(\underline{\theta}) & \dots & \frac{\partial^2}{\partial \theta_p \partial \theta_p} \log L(\underline{\theta}) \end{pmatrix} \end{aligned}$$

در نتیجه معادلات درست‌نمایی به صورت زیر خواهد بود:

$$\sum_i \frac{\partial}{\partial \theta_j} \log L_i(\underline{\theta}) = 0 \quad j=1, \dots, p$$

یا

$$\frac{\partial}{\partial \underline{\theta}} \log L(\underline{\theta}) = 0$$

فرض کنید بردار $\hat{\underline{\theta}} = (\hat{\theta}_1^\circ, \dots, \hat{\theta}_p^\circ)'$ یک مقدار اولیه برای جواب باشد. در صورت بسط حول $\hat{\underline{\theta}}^\circ$ داریم:

$$\begin{aligned} \sum_i \frac{\partial}{\partial \theta_j} \log L_i(\hat{\underline{\theta}}) &= \sum_i \frac{\partial}{\partial \theta_j} \log L_i(\hat{\underline{\theta}}^\circ) + \\ &+ \sum_k (\hat{\theta}_k - \hat{\theta}_k^\circ) \sum_i \frac{\partial^2}{\partial \theta_k \partial \theta_j} \log L_i(\hat{\underline{\theta}}^\circ) + \dots = 0 \quad j=1, \dots, p \end{aligned}$$

یا

$$\frac{\partial}{\partial \underline{\theta}} \log L(\hat{\underline{\theta}}) = \frac{\partial}{\partial \underline{\theta}} \log L(\hat{\underline{\theta}}^\circ) + \frac{\partial^2}{\partial \underline{\theta}^2} \log L(\hat{\underline{\theta}}^\circ) (\hat{\underline{\theta}} - \hat{\underline{\theta}}^\circ) + \dots = 0$$

فرض کنید $\hat{\underline{\theta}}^1$ جوابی باشد، که در آن از جملات درجه دوم و بیشتر چشم‌پوشی شده است.

$$\hat{\underline{\theta}}^1 = \hat{\underline{\theta}}^\circ + \left(-\frac{\partial^2}{\partial \underline{\theta}^2} \log L(\hat{\underline{\theta}}^\circ) \right)^{-1} \frac{\partial}{\partial \underline{\theta}} \log L(\hat{\underline{\theta}}^\circ) \quad (1)$$

بردار $\left(\frac{\partial}{\partial \underline{\theta}} \right) \log L(\hat{\underline{\theta}}^\circ)$ را بردار چوب‌خط (یا امتیاز) در $\hat{\underline{\theta}}^\circ$ نامند. ماتریس زیر را ماتریس اطلاع نمونه در $\hat{\underline{\theta}}^\circ$ نامند.

$$i(\hat{\underline{\theta}}^\circ) = -\frac{\partial^2}{\partial \underline{\theta}^2} \log L(\hat{\underline{\theta}}^\circ)$$

توجه شود که رابطه زیر صادق است. $I(\underline{\theta})$ ، اطلاع فیشر مشاهده‌ای نام است.

$$E(i(\hat{\theta})) = \left(-E \frac{\partial^2}{\partial \theta_k \partial \theta_j} \log L(\theta) \right) = I(\theta)$$

متذکر می‌شود که $I(\theta)$ اطلاع فیشر برای کل نمونه است.

$$I(\theta) = \sum_{i=1}^n I_i(\theta) = nI_1(\theta)$$

در این جا $I_i(\theta)$ اطلاع فیشر مربوط به مشاهده i ام است. روش تکراری که در (۱) به کار رفت، روش نیوتن-رافسون نامیده می‌شود. اگر اطلاعات نمونه در (۱) را با اطلاع فیشر جابه‌جا کنیم، داریم:

$$\hat{\theta}^1 = \hat{\theta}^0 + I^{-1}(\hat{\theta}^0) \frac{\partial}{\partial \theta} \log L(\hat{\theta}^0) \quad (2)$$

روش تکراری به کار رفته در (۲) را روش چوب‌خطی (یا امتیازی) گویند. روش (۲) باید در بعضی از موارد سریعتر همگرا شود. ولی در عمل، هنگامی که عمل برش وجود دارد، نمی‌توان از $I(\theta)$ در (۲) استفاده کرد.

REFERENCES

Rao, Linear Statistical Inference (1965), Section 5g.

Gross and Clark, Survival Distributions (1975), Chapter 6.

Kalbfleisch and The Statistical Analysis of Failure Time Data (1980), Section 3.7.

بازه‌های اطمینان و آزمونها

برای دو برش تصادفی نوع اول با شرایط همواری، رابطه زیر را داریم:

$$\hat{\theta} \underset{a}{\sim} N(\theta, I^{-1}(\theta))$$

معمولاً این نتیجه برای برش نوع دوم نیز برقرار است. ولی، اثبات آن مشکل است. نماد "a" به این معنی است، که رابطه به طور مجانبی برقرار است.

برای آزمون فرض: $H_0: \theta = \theta^0$ یا برای ساختن بازه‌های اطمینان، معمولاً از سه روش استفاده می‌شود:

۱ روش والد: تحت فرض H_0 ، داریم:

$$(\hat{\theta} - \theta^0)' I(\theta^0) (\hat{\theta} - \theta^0) \underset{a}{\sim} \chi_p^2$$

که در این رابطه از $I(\hat{\theta})$ به جای $I(\theta^0)$ می‌توان استفاده کرد.

۲ روش نیمن-پیرسن-نسبت درست‌نمایی و بیلکس: تحت فرض H_0 ، داریم:

$$-2 \log \frac{L(\hat{\theta}^0)}{L(\hat{\theta})} \underset{a}{\approx} \chi_p^2$$

۳ روش راثو: تحت فرض H_0 ، داریم:

$$\frac{\partial}{\partial \theta} \log L(\hat{\theta}^0)' I^{-1}(\hat{\theta}^0) \frac{\partial}{\partial \theta} \log L(\hat{\theta}^0) \underset{a}{\approx} \chi_p^2$$

توجه کنید که روش راثو از برآورد حداکثر درست‌نمایی (MLE) استفاده نمی‌کند. همچنین، محاسبات تکراری را لازم ندارد. علاوه بر آزمون، اغلب مایل به محاسبه $\hat{\theta}$ نیز هستیم. اگر $\hat{\theta}$ و $I(\hat{\theta}^0)$ را داشته باشیم، روش والد ساده‌تر خواهد بود. در حالت برش لازم است به جای $I(\hat{\theta})$ از $\dot{I}(\hat{\theta})$ استفاده شود. زیرا، معمولاً محاسبه $I(\hat{\theta})$ مشکل است. همچنین، افرن و هینکلی پیشنهاد می‌کنند، که استفاده از $\dot{I}(\hat{\theta})$ برای بازه اطمینان از $I(\hat{\theta})$ ، حتی زمانی که $I(\hat{\theta})$ قابل محاسبه باشد، بهتر است. با این وجود در این گونه موارد نظرها متفاوت است.

REFERENCES

Rao, Linear Statistical Inference (1965), Section 6e.

Efron and Hinkley, Biometrika (1978).

مثال ۱ نمایی: تحت برش تصادفی، فرض کنید n_u ، تعداد مشاهدات بریده نشده باشد، در این صورت داریم:

$$L = \lambda^{n_u} \exp \left(-\lambda \sum_u t_i - \lambda \sum_c c_i \right) = \lambda^{n_u} \exp \left(-\lambda \sum_{i=1}^n y_i \right)$$

$$\log L = n_u \log \lambda - \lambda \sum_{i=1}^n y_i$$

$$\frac{\partial}{\partial \lambda} \log L = \frac{n_u}{\lambda} - \sum_{i=1}^n y_i$$

$$\hat{\lambda} = \frac{n_u}{\sum_{i=1}^n y_i}$$

$$\frac{\partial^2}{\partial \lambda^2} \log L = \frac{-n_u}{\lambda^3}$$

$$i(\underline{\lambda}) = \frac{n_u}{\lambda^3}$$

توجه شود، که $\hat{\lambda} = n_u / \sum_{i=1}^n y_i$ برآورد MLE تحت برش نوع اول و دوم - به خوبی برآورد تحت برش تصادفی - نیز هست.

برای ساختن بازه‌های اطمینان و اجرای آزمونها، توزیع $\hat{\lambda}$ مورد نیاز است.

(الف) اگر مشاهدات بریده نشوند، داریم:

$$\hat{\lambda} = \frac{n}{\sum_{i=1}^n T_i} = \frac{1}{\bar{T}}$$

در رابطه بالا: T_1 تا T_n ، متغیرهای iid با توزیع نمایی و تابع چگالی زیراند:

$$f_{T_i}(t) = \lambda e^{-\lambda t}$$

در نتیجه، $S = \sum_{i=1}^n T_i$ دارای چگالی گاما به شرح زیر است.

$$f_S(t) = \frac{\lambda^n}{\Gamma(n)} t^{n-1} \cdot e^{-\lambda t}$$

بنابراین: $\sum_{i=1}^n T_i \sim \chi_{2n}^2$ یا به عبارت معادل:

$$\frac{2n\lambda}{\hat{\lambda}} \sim \chi_{2n}^2$$

یعنی: $\frac{2n\lambda}{\hat{\lambda}}$ یک آماره محوری است و آن را می‌توان برای آزمون و بازه اطمینان مسورد استفاده قرار داد (علامت " \sim " به معنی هم‌توزیع بودن است).

(ب) برای برش نوع دوم می‌توان نوشت:

$$\sum_{i=1}^n Y_i = T_{(1)} + T_{(2)} + \dots + (n-r)T_{(r)}$$

$$= nT_{(1)} + (n-1)[T_{(2)} - T_{(1)}] + \dots + (n-r+1)[T_{(r)} - T_{(r-1)}]$$

با استفاده از نتایج فرایند پواسن و زمانهای انتظار نمایی، داریم:

$$T_{(1)} = \{ \text{کمینه } (n-1) \text{ متغیر iid نمایی } T_i \text{ با پارامتر } \lambda \} \sim n\lambda e^{-n\lambda t}$$

$$nT_{(1)} \sim \lambda e^{-\lambda t}$$

$$T_{(2)} - T_{(1)} = \{ \text{کمینه } (n-1) \text{ متغیر iid نمایی } T_i \text{ با پارامتر } \lambda \} \sim (n-1)\lambda e^{-(n-1)\lambda t}$$

$$(n-1)[T_{(2)} - T_{(1)}] \sim \lambda e^{-\lambda t}$$

تا آخر. چون متغیرهای $nT_{(1)}$ ، $(n-1)[T_{(2)} - T_{(1)}]$ و $(n-r+1)[T_{(r)} - T_{(r-1)}]$ مستقل‌اند، پس داریم:

$$r\lambda \sum_{i=1}^n Y_i \sim \chi_{2r}^2$$

بنابراین، برای ساختن بازه‌های اطمینان و آزمونها، می‌توان از $r\lambda/\hat{\lambda}$ و ارتباط آن با توزیع χ^2 ، استفاده کرد. در این حالت، درجه آزادی دو برابر آماره‌های مرتب بریده نشده است.

(پ) اگر برش از نوع تصادفی یا نوع اول باشد، چاره‌ای جز به کارگیری نظریهٔ مجانبی نداریم. با توجه به محاسبات قبل، داریم:

$$\hat{\lambda} = \frac{n_u}{\sum_{i=1}^n y_i} \quad \text{و} \quad \frac{\partial^2}{\partial \lambda^2} \log L = \frac{-n_u}{\lambda^2}$$

در نتیجه: $\frac{\hat{\lambda} - \lambda}{\sqrt{\lambda^2/n_u}} \stackrel{a}{\approx} N(0, 1)$. در این حالت می‌توان به جای n_u از $E(n_u)$ استفاده نمود، به شرط آن که این میانگین معلوم باشد.

می‌توان تقریب نرمال را با تبدیل برآورد بهبود بخشید. به کمک روش دلتا (بعدها)

تعریف خواهد شد) و این که $\hat{\lambda} \stackrel{a}{\approx} N(\lambda, \frac{\lambda^2}{n_u})$ داریم:

$$\log \hat{\lambda} \underset{a}{\sim} N(\log \lambda, \frac{1}{n_u})$$

توجه شود که واریانس $\log \hat{\lambda}$ ، یعنی: $\frac{1}{n_u}$ ، به پارامتر مجهول بستگی ندارد. این یک حقیقت تجربی است که، تبدیل کردن یک برآورد - برای حذف وابستگی واریانس به پارامتر مجهول - معمولاً، تمایل دارد که با کاهش چولگی، همگرایی به نرمال را بهبود بخشد.

REFERENCE

Einstein and Sobel, JASA (1953), is a classic paper.

روش دلتا: فرض کنید، متغیر تصادفی Y دارای واریانس σ^2 و میانگین μ است (که به اختصار به صورت $Y \sim (\mu, \sigma^2)$ می‌نویسیم)، همچنین، فرض کنید توزیع تابع $g(Y)$ را می‌خواهیم. ابتدا، $g(Y)$ را پیرامون μ بسط می‌دهیم، داریم:

$$g(Y) = g(\mu) + (Y - \mu)g'(\mu) + \dots$$

اگر از جملات درجات بالاتر چشم‌پوشی کنیم، تقریب: $g(Y) \approx (g(\mu), \sigma^2(g'(\mu))^2)$ ، به دست می‌آید. در این عبارت " \approx " به معنی توزیع تقریبی است. علاوه بر این، اگر $Y \underset{a}{\sim} (\mu, \sigma^2)$ ، آن‌گاه داریم:

$$g(Y) \underset{a}{\sim} (g(\mu), \sigma^2(g'(\mu))^2)$$

روش چندمتغیره نیز به کار می‌رود. فرض کنید رابطه زیر را داریم:

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim \left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ & \sigma_y^2 \end{pmatrix} \right)$$

همچنین فرض می‌کنیم، مایل به یافتن توزیع $g(X, Y)$ هستیم. بسط آن به شرح زیر است:

$$g(X, Y) = g(\mu_x, \mu_y) + (X - \mu_x) \frac{\partial}{\partial x} g(\mu_x, \mu_y) + (Y - \mu_y) \frac{\partial}{\partial y} g(\mu_x, \mu_y) + \dots$$

در نتیجه، داریم:

$$g(X, Y) = \left(g(\mu_x, \mu_y), \sigma_x^2 \left(\frac{\partial}{\partial x} g \right)^2 + 2\sigma_{xy} \frac{\partial}{\partial x} g \frac{\partial}{\partial y} g + \sigma_y^2 \left(\frac{\partial}{\partial y} g \right)^2 \right)$$

علاوه بر این، اگر (X, Y) دارای توزیع مجانبی نرمال باشند، آنگاه $g(X, Y)$ نیز دارای توزیع مجانبی نرمال خواهد بود.

روش دلنا بسیار مفید است. مثلاً، می‌توان آن را برای محاسبه تقریبی $\text{Var}\left(\frac{X}{Y}\right)$ یا $\text{Var}(\bar{X}\bar{Y})$ به کار برد.

مثال ۲ وایبل: اگر توزیع وایبل را با پارامتر $\gamma = \lambda^\alpha$ بنویسیم، می‌توان مشتقات آن را ساده‌تر محاسبه کرد.

$$S(t) = e^{-(\lambda t)^\alpha} = e^{-\gamma t^\alpha}$$

$$f(t) = \gamma \alpha t^{\alpha-1} \cdot e^{-\gamma t^\alpha}$$

در این صورت، داریم:

$$L = (\gamma \alpha)^{n_u} \left(\prod_u t_i^{\alpha-1} \right) \exp\left(-\gamma \sum_u t_i^\alpha\right) \exp\left(-\gamma \sum_c c_i^\alpha\right)$$

$$= (\gamma \alpha)^{n_u} \left(\prod_u t_i^{\alpha-1} \right) \exp\left(-\gamma \sum_{i=1}^n y_i^\alpha\right)$$

$$\log L = n_u \log \gamma + n_u \log \alpha + (\alpha - 1) \sum_u \log t_i - \gamma \sum_{i=1}^n y_i^\alpha$$

$$\frac{\partial}{\partial \gamma} \log L = \frac{n_u}{\gamma} - \sum_{i=1}^n y_i^\alpha$$

$$\frac{\partial}{\partial \alpha} \log L = \frac{n_u}{\alpha} + \sum_u \log t_i - \gamma \sum_{i=1}^n y_i^\alpha \log y_i$$

بنابراین، برآورد حداکثر درست‌نمایی $(\hat{\alpha}, \hat{\gamma})$ به شرح زیر است:

$$\hat{\gamma} = \frac{n_u}{\sum_{i=1}^n y_i^{\hat{\alpha}}}$$

$$\frac{n_u}{\hat{\alpha}} + \sum_u \log t_i - \hat{\gamma} \sum_{i=1}^n y_i^{\hat{\alpha}} \log y_i = 0$$

این معادلات را باید به روش عددی حل کرد. روش نیوتن-رافسون به ماتریس اطلاع فیشر نمونه زیر - که در مسأله ۳ محاسبه شد - نیاز دارد.

$$-\frac{\partial^2}{\partial \theta^2} \log L = - \begin{pmatrix} \frac{\partial^2}{\partial \gamma^2} \log L & \frac{\partial^2}{\partial \gamma \partial \alpha} \log L \\ \frac{\partial^2}{\partial \alpha^2} \log L \end{pmatrix}$$

روش نیوتن-رافسون به مقادیر اولیه $\hat{\gamma}_0$ و $\hat{\alpha}_0$ نیز نیاز دارد. توجه نمایید، که برای محاسبه مقادیر اولیه معقول، روابط زیر قابل استفاده است:

$$S(t) = e^{-\gamma t^\alpha}$$

$$\log S(t) = -\gamma t^\alpha$$

$$\log(-\log S(t)) = \log \gamma + \alpha \log t$$

بنابراین، اگر برآورد $\hat{S}(t_i)$ را داشته باشیم، به کمک رگرسیون $\log(-\log \hat{S}(t_i))$ روی $\log t_i$ و محاسبه ضریب رگرسیون $\hat{\alpha}$ و مقدار ثابت $\log \hat{\gamma}$ ، نتیجه مطلوب حاصل می‌شود. مقدار انتخابی ممکن $\hat{S}(t_i)$ ، همان برآورد کاپلان-میر است، که بعداً بحث خواهد شد. همچنین می‌توان از یک تابع توزیع تجربی - که از برش چشم‌پوشی می‌کند - استفاده کرد.

REFERENCE

Cohen, Technometrics (1965), treats the MLE and gives additional references.

برآورد $S(t)$. برآورد تابع بقاء، یکی از اهداف اصلی تحلیل بقاء است.

$$S(t) = \exp\left(-\int_0^t \lambda(u) du\right)$$

برای مثال، یکی از مقیاسهای معالجه سرطان این است که احتمال زنده بودن بیمار را به مدت حداقل پنج سال حساب کنیم. در مطالعات مهندسی تابع بقاء را تابع اعتماد نامند و معمولاً آن را با $R(t)$ نشان می‌دهند. برای مثال، دانستن قابلیت اعتماد یک مؤلفه در یک دستگاه بعد از هزار ساعت کار آن دستگاه.

با داشتن MLE، برآورد تابع بقاء در دو حالت نمایی یا وایبل بسیار ساده است.

$$\hat{S}(t) = e^{-\hat{\lambda}t} \quad \text{حالت نمایی}$$

$$\hat{S}(t) = e^{-(\hat{\lambda}t)^{\hat{\alpha}}} = e^{-\hat{\gamma}t^{\hat{\alpha}}} \quad \text{حالت وایبل}$$

همچنین، برای هر مقدار ثابت t ، تابع بقاء $\hat{S}(t)$ خود تابعی از $\hat{\lambda}$ یا $(\hat{\gamma}, \hat{\alpha})$ است. بنابراین به دست آوردن تقریبی از توزیع $\hat{S}(t)$ به روش دلتا امکان پذیر است. علاوه بر این، می توان با تبدیل $\log \log$ ، همگرایی به نرمال را بهبود بخشید. در حالت نمایی، داریم:

$$S(t) = e^{-\lambda t}$$

$$\log[-\log S(t)] = \log \lambda + \log t$$

$$\log[-\log \hat{S}(t)] = \log \hat{\lambda} + \log t$$

$$\widehat{\text{Var}} \{ \log[-\log \hat{S}(t)] \} \cong \frac{1}{n_u}$$

در حالت وایبل، داریم:

$$S(t) = e^{-\gamma t^{\alpha}}$$

$$\log[-\log S(t)] = \log \gamma + \alpha \log t$$

$$\log[-\log \hat{S}(t)] = \log \hat{\gamma} + \hat{\alpha} \log t$$

$$\widehat{\text{Var}} \{ \log[-\log \hat{S}(t)] \} \cong \frac{\text{Var}(\hat{\gamma})}{\hat{\gamma}^2} + 2 \text{Cov}(\hat{\gamma}, \hat{\alpha}) \frac{\log t}{\hat{\gamma}} + \text{Var}(\hat{\alpha})(\log t)^2$$

۲.۲ ترکیب خطی آماره‌های مرتب

در این بخش فقط توزیع وایبل بررسی می شود. روش کار قابل تعمیم است. با تغییر پارامتر و تبدیل مناسب می توان مسأله برآورد λ و α را در توزیع وایبل به برآورد پارامتر مبدأ و مقیاس تبدیل کرد. با دوباره نویسی داریم:

$$P(Y > t) = e^{-(\lambda t)^{\alpha}} = \exp\{-\exp[\alpha(\log \lambda + \log t)]\} = \exp\left\{-\exp\left(\frac{\log t - \mu}{\sigma}\right)\right\}$$

در روابط بالا، $\mu = -\log \lambda$ و $\sigma = \frac{1}{\alpha}$ است. در این صورت داریم:

$$P(\log Y > t) = P(Y > e^t) = \exp\left\{-\exp\left(\frac{t-\mu}{\sigma}\right)\right\}$$

در معادله قبل دیده می‌شود که μ و σ همان پارامترهای موقعیت و مقیاس متغیر تصادفی $\log Y$ هستند. این نتیجه بسیار مهم است، زیرا در آمار قضایای زیادی برای برآورد دو پارامتر موقعیت و مقیاس وجود دارد.

فرض کنید، می‌خواهیم احتمال بقاء را برای زمان ثابت t_0 حساب کنیم، داریم:

$$S(t_0) = P(Y > t) = \exp\left\{-\exp\left(\frac{\log t_0 - \mu}{\sigma}\right)\right\}$$

با تعریف: $Y^0 = \frac{Y}{t_0}$ و $\mu_0 = \mu - \log t_0$ ، می‌توان عبارت را کوتاه‌تر نوشت. داریم:

$$S(t_0) = P(\log Y^0 > 0) = \exp\left\{-\exp\left(\frac{\mu_0}{\sigma}\right)\right\}$$

در این رابطه μ_0 و σ ، پارامترهای موقعیت و مقیاس متغیر تصادفی $\log Y^0$ هستند. اگر بتوانیم یک بازه اطمینان برای نسبت $\frac{\mu_0}{\sigma}$ پیدا کنیم، با دوبار استفاده از تبدیل نمایی، یک بازه اطمینان برای $S(t_0)$ به دست می‌آید. با استفاده از ترکیب خطی آماره‌های مرتب داریم:

$$\hat{\mu}_0 = \sum_{i=1}^n a_i \log Y_{(i)}^0 \quad \text{و} \quad \hat{\sigma}_0 = \sum_{i=1}^n b_i \log Y_{(i)}^0$$

در این دو رابطه $\sum_{i=1}^n a_i = 1$ و $\sum_{i=1}^n b_i = 0$ و همچنین b_1 تا b_n به گونه‌ای انتخاب می‌شوند که در یک شرط مجانبی بهینه صدق کنند. این روش به خصوص برای برش نوع دوم بسیار مناسب است، زیرا داریم:

$$a_{r+1} = \dots = a_n = 0$$

$$b_{r+1} = \dots = b_n = 0$$

بنابراین، برآوردها بر مبنای مشاهدات بریده نشده انجام می‌پذیرد.

REFERENCE

Johns and Lieberman, Technometrics (1966).

توزیع‌های فرین

تابع زیر یکی از سه توزیع حدی فرین است.

$$G_1 = \exp\{-\exp(-x)\} \quad -\infty < x < \infty$$

یک توزیع حدی فرین، توزیعی است مانند G ، به گونه‌ای که اگر X_1, \dots, X_n متغیرهای iid با توزیع F باشند، آنگاه $\max\{X_1, \dots, X_n\}$ - با فرض این که به طور مناسب استاندارد شده باشند - در توزیع به G همگراست. یکی دیگر از توزیع‌های حدی به شرح زیر است:

$$G_2(x) = \begin{cases} \exp\{-(-x)^\alpha\} & x < 0 \\ 1 & x > 0 \end{cases}$$

اگر دنباله فوقانی توزیع وایبل به طور مناسب مقیاس‌بندی شود، برابر با دنباله تحتانی G_2 است. همچنین، دنباله فوقانی توزیع لگاریتم متغیر وایبل (۳) برابر با دنباله تحتانی G_1 است. تابع توزیع G_2 از استاندارد کردن حد متغیر زیر به دست می‌آید:

$$\max\{X_1, \dots, X_n\} - x_0$$

در این رابطه x_0 نقطه قطع فوقانی است (یعنی: $F(x_0) = 1$ و $F(x_0^-) < 1$). در نتیجه داریم:

$$-\max\{X_1 - x_0, \dots, X_n - x_0\} = \min\{x_0 - X_1, \dots, x_0 - X_n\}$$

یک متغیر تصادفی وایبل را می‌توان به صورت کمینه (یعنی اولین شکست) تعداد زیادی زمانهای شکست بالقوه، تفسیر کرد. دستگاه با شکست اولین مؤلفه از کار می‌افتد.

۳.۲ برآوردگرهای دیگر

در این بخش فرض می‌شود که برآوردگرها دارای الگوی نمایی و بدون برش هستند.

برآوردگرهای ارباب اصلاح شده. در این مبحث، روش کار دارای اهمیت بیشتری از

نتایج به دست آمده است. فرض کنید احتمال بقاء را به صورت زیر برآورد می‌کنیم، که

$$\text{در آن } \bar{T} = \frac{1}{n} \sum_{i=1}^n T_i \text{ است.}$$

$$\hat{S}(t) = e^{-\hat{\lambda}t} = e^{-t/\bar{T}}$$

$$E[\hat{S}(t)] \neq e^{-\lambda t}$$

در نتیجه $\hat{S}(t)$ یک برآورد اریب است. می‌توان به روش دلتا اریبی را کاهش داد. با فرض: $\theta = E(T)$ ، داریم:

$$e^{-t/\bar{T}} = e^{-t/\theta} + (\bar{T} - \theta) \frac{t}{\theta^2} e^{-t/\theta} + \frac{1}{2} (\bar{T} - \theta)^2 \left[\left(\frac{t}{\theta^2} \right)^2 - \frac{2t}{\theta^3} \right] e^{-t/\theta} + \dots$$

$$E(e^{-t/\bar{T}}) = e^{-t/\theta} + \frac{1}{2} \frac{\theta^2}{n} \left[\left(\frac{t}{\theta^2} \right)^2 - \frac{2t}{\theta^3} \right] e^{-t/\theta} + \dots$$

$$= \left[1 + \frac{1}{2n} \left(\frac{t^2}{\theta^2} - \frac{2t}{\theta} \right) \right] e^{-t/\theta} + \dots$$

در نتیجه:

$$\tilde{S}(t) = \frac{e^{-\hat{\lambda}t}}{1 + \frac{1}{2n} (t^2 \hat{\lambda}^2 - 2t \hat{\lambda})}$$

$\tilde{S}(t)$ محاسبه شده در بالا باید اریبی کمتری از برآورد گسر $\hat{S}(t)$ داشته باشد. علاوه بر این، ثابت می‌شود که، میانگین مربع خطای $\tilde{S}(t)$ کمتر از $\hat{S}(t)$ است. برآورد جک‌نایف، که درباره آن بعداً بحث می‌شود، نیز تصحیح اریبی بالا را به همراه دارد.

برآوردگرهای نارایب با حداقل واریانس (UMVUE). در محاسبه UMVUE برای $S(t)$ از برآورد نارایب $U = I(T_1 > t)$ ، استفاده می‌شود. در این صورت با در نظر گرفتن آماره بسنده $S = \sum_{i=1}^n T_i$ و قضیه راثو-بلکول، UMVUE معادل $E(U|S)$ است. در

نتیجه داریم:

$$\tilde{S}(t) = E(U|S=s) = \left(1 - \frac{t}{s}\right)^{n-1} \cdot I(t < s)$$

برآوردهای بی‌زی. فقط یادآوری می‌کنیم که برآوردهای بیز را می‌توان با استفاده از توزیع پیشین گاما به دست آورد.

REFERENCES

Basu, *Technometrics* (1964), derives UMVUEs.

Zacks and Even, *JASA* (1966), compares mean square errors.

Gaver and Hoel, *Technometrics* (1970), look at estimators in the framework of sampling from a Poisson process.

۳ الگوهای رگرسیونی

در کاربردهای پزشکی امکان دارد که مدت بقاء به مقدار دارو یا تشعشع بستگی داشته باشد. همچنین، در کاربردهای مهندسی طول عمر یک لامپ ممکن است به دما یا عوامل دیگر وابسته باشد. فرض کنید Y متغیر وابسته و X متغیر مستقل باشد. دو الگوی زیر برای توزیع نمایی پیشنهاد می‌شود:

(الف) الگوی خطی

$$E(T) = \alpha + \beta X$$

در محاسبه برآوردها، از روش حداکثر درست‌نمایی استفاده می‌کنیم. عیب این الگو این است که اگر $\hat{\beta}$ منفی باشد، امکان دارد برآورد $E(T)$ منفی باشد.

(ب) الگوی خطی لگاریتمی

$$E(T) = \alpha e^{\beta X}$$

$$\log E(T) = \log \alpha + \beta X$$

در این جا نیز از برآورد حداکثر درست‌نمایی استفاده می‌شود.

REFERENCES

Feigl and Zelen, *Biometrics* (1965), discuss the uncensored case for both the linear and log – linear models.

Zippin and Armitage, *Biometrics* (1966), discusses the censored case

for the linear model.

Glasser, JASA (1967), discusses the censored case for the log – linear model.

Zippin and Lamborn, Stanford Univ. Tech. Report No. 20 (1969), discuss

the censored case for the log – linear model goodness of fit tests.

Mantel and Myers, JASA (1971), discuss the censored case for the multiple linear model.

۴ الگوهای با کسرهای بقاء

۱.۴ نمونه‌ای به حجم واحد.

فرض کنید:

$$p = P(\text{مرگ}) \quad \text{و} \quad 1-p = P(\text{بقاء})$$

نسبت $1-p$ را کسر بقاء نامند. با فرض: $P(T \leq t | \text{مرگ}) = 1 - e^{-\lambda t}$ ، تابع درست‌نمایی به صورت زیر است:

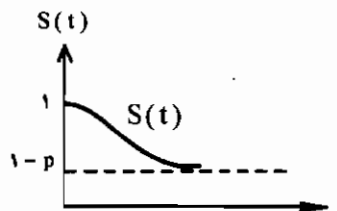
$$L(y, \delta) = \begin{cases} p\lambda e^{-\lambda y} & \delta = 1 \quad (\text{بدون برش}) \\ (1-p) + pe^{-\lambda y} & \delta = 0 \quad (\text{با برش}) \end{cases}$$

در محاسبه برآوردها از حداکثر درست‌نمایی استفاده می‌کنیم.

الگوهای با کسر بقاء را گاهی برای آزمایشهای کوتاه به کار می‌برند. در این

حالت فرض نمی‌شود که تابع بقاء $S(t)$ لزوماً به صفر نزدیک می‌شود. در عوض امکان

دارد $S(t)$ مطابق شکل زیر باشد:

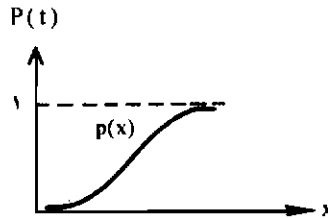


۲.۴ رگرسیون

$p(x)$ را به صورت زیر فرض می‌کنیم

$$p(x) = P(\text{مرگ} | x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

این تابع را لوجستیک گویند و به صورت نمودار زیر است:



همچنین فرض کنید: $P(T \leq t | \text{مرگ}) = 1 - e^{-\lambda t}$. در این صورت تابع درست‌نمایی به صورت زیر خواهد بود:

$$L(y, \delta, x) = \begin{cases} p(x)\lambda e^{-\lambda y} & \delta = 1 \quad (\text{بدون برش}) \\ 1 - p(x) + p(x)e^{-\lambda y} & \delta = 0 \quad (\text{با برش}) \end{cases}$$

برای به دست آوردن برآوردها از حداکثر درست‌نمایی استفاده می‌شود.

REFERENCE

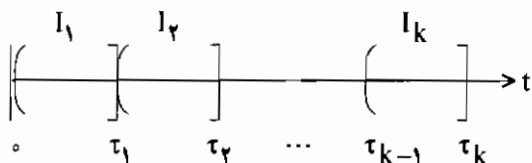
Farewell, Biometrika (1977).

فصل سوم

روشهای ناپارامتری (یک نمونه)

۱ جدولهای طول عمر

روش رسمی برآورد $S(t)$ در علوم بیماریهای مسری و بیمه‌گری، همان روش بیمه‌گری است که به جدول طول عمر بستگی دارد و در زیر آن را شرح می‌دهیم. فرض کنید زمان به دنباله‌های ثابتی از بازه‌های I_1, \dots, I_k افراز شده باشد. این بازه‌ها تقریباً همیشه (و نه لزوماً) مساوی در نظر گرفته می‌شوند. برای جمعیت انسانها این فاصله‌ها معمولاً یک سال است.



بنا به تعریف، در یک جدول طول عمر موارد زیر فرض می‌شوند:

$$n_i = \text{تعداد افراد زنده در شروع بازه } I_i.$$

$$d_i = \text{تعداد افراد مرده در طول بازه } I_i.$$

$$l_i = \text{تعداد گم شده‌ها در طول بازه } I_i.$$

$$w_i = \text{تعداد خارج شده‌ها در طول بازه } I_i.$$

$$p_i = \text{احتمال زنده بودن در طول بازه } I_i, \text{ در صورتی که در ابتدای بازه زنده بوده است.}$$

$$1 - p_i = q_i$$

جدول ۱، مثالی از یک جدول طول عمر است. I_1 تا I_5 ، هر کدام به مدت یکسال‌اند. ستون (۲) شامل n_i ، ستون (۳) شامل d_i ، ستون (۴) شامل l_i و ستون (۵) شامل w_i است. می‌خواهیم $S(5)$ را برآورد کنیم.

جدول ۱. محاسبه نرخ بقاء پنج ساله

سالهای بعد از تشخیص بیماری	تعداد زنده‌ها در ابتدای بازه	مردم در طول بازه	گم شده‌ها	خارج شده‌ها	تعداد مؤثر در معرض خطر مرگ	نسبت مرگ و میر	نسبت زنده‌ها	نسبت تجمعی زنده‌ها تا پایان بازه
(۱)	(۲)	(۳)	(۴)	(۵)	$(۲) - \frac{۱}{۲}[(۴)+(۵)]$	$(۳)/(۶)$	$۱ - (۷)$	$\prod_{i=1}^k (۸)_i$ (۹)
۰-۱	۱۲۶	۴۷	۴	۱۵	۱۱۶/۵	۰/۴	۰/۶۰	۰/۶۰
۱-۲	۶۰	۵	۶	۱۱	۵۱/۵	۰/۱	۰/۹۰	۰/۵۴
۲-۳	۳۸	۲	—	۱۵	۳۰/۵	۰/۰۷	۰/۹۳	۰/۵۰
۳-۴	۲۱	۲	۲	۷	۱۶/۵	۰/۱۲	۰/۸۸	۰/۴۴
۴-۵	۱۰	—	—	۶	۷	۰/۰۰	۱/۰۰	۰/۴۴

مرجع: Cutler and Ederer, J. Chronic Dis. (1958).

۱.۱ روش کاهش نمونه

برای برآورد $S(\tau_k)$ ، تنها افرادی را مورد توجه قرار می‌دهیم که در بازه $(0, \tau_k]$ در معرض خطر قرار دارند (یعنی: در بازه مورد توجه). موارد زیر را داریم:

$$n = n_1 - \sum_{i=1}^k \ell_i - \sum_{i=1}^k w_i$$

$$d = \sum_{i=1}^k d_i$$

$$\hat{S}(\tau_k) = 1 - \frac{d}{n}$$

در مثال جدول (۱)، داریم: $\hat{S}(5) = 1 - \frac{56}{60} = 0.067$ و $d = 56$ و $n = 126 - 12 - 56 = 60$ و عیب روش کاهش نمونه این است که از اطلاعات موجود در ℓ_i و w_i چشم‌پوشی می‌کند. این یک برآورد اریب (نقصانی) از $S(t)$ است.

۲.۱ روش بیمه‌گری

می‌توان احتمال بقاء $S(\tau_k)$ را به صورت حاصل ضرب چند احتمال تجزیه کرد. داریم:

$$\begin{aligned} S(\tau_k) &= P(T > \tau_k) \\ &= P(T > \tau_1) \cdot P(T > \tau_2 | T > \tau_1) \cdots P(T > \tau_k | \tau_{k-1}) \\ &= p_1 \cdot p_2 \cdots p_k \end{aligned}$$

در این رابطه: $p_i = P(T > \tau_i | T > \tau_{i-1})$.

روش بیمه‌گری، یک برآوردی جدا برای هر p_i ، ارائه می‌دهد. سپس، با ضرب این برآوردها، برآورد $S(\tau_k)$ نتیجه می‌شود.

در برآورد p_i ، می‌توان از $1 - \frac{d_i}{n_i}$ استفاده کرد، به شرطی که مشاهده‌ای در I_i گم یا خارج نشده باشد. با این وجود، اگر ℓ_i و w_i صفر نباشند، فرض می‌کنیم که به‌طور متوسط، افرادی که در بازه I_i حذف شده‌اند، در معرض خطر در نصف بازه بوده‌اند. بنابراین، حجم مؤثر نمونه را به صورت زیر تعریف می‌کنیم:

$$n'_i = n_i - \frac{1}{\nu}(\ell_i + w_i) \quad \text{و} \quad \hat{q}_i = \frac{d_i}{n'_i} \quad \text{و} \quad \hat{p}_i = 1 - \hat{q}_i$$

در نتیجه، برآورد روش بیمه‌گری معادل: $\hat{S}(\tau_k) = \prod_{i=1}^k p_i$ است. در جدول (۱)، ستون (۶) شامل n' ، ستون (۷) شامل \hat{q}_i ، ستون (۸) شامل \hat{p}_i و در ستون (۹) دیده می‌شود که $S(5) = 0.44$ است.

برای بهبودبخشی در پیدا کردن جانشینی برای حجم نمونه مؤثر، تلاشهای زیادی به انجام رسیده است. ولی اگر یک برآورد دقیقتری برای $S(t)$ لازم شود، حدّ حاصل ضرب برآوردگر کاپلان-میر، روش مناسبی خواهد بود.

۳.۱ واریانس $\hat{S}(\tau_k)$

برای برآورد واریانس $\hat{S}(\tau_k)$ ، از رابطه زیر استفاده می‌شود:

$$\log \hat{S}(\tau_k) = \sum_{i=1}^k \log \hat{p}_i$$

با فرض $n'_i \hat{p}_i \sim \text{Binomial}(n'_i, p_i)$ و استفاده از روش دلتا، داریم:

$$\begin{aligned} \text{Var}(\log \hat{p}_i) &\cong \text{Var}(\hat{p}_i) \left[\frac{d}{d p_i} (\log p_i) \right]^2 \\ &\cong \frac{p_i q_i}{n'_i} \cdot \frac{1}{p_i^2} = \frac{q_i}{n'_i p_i} \end{aligned}$$

فرض می‌شود $\log \hat{p}_1$ و \dots و $\log \hat{p}_k$ مستقل‌اند.

$$\text{Var}[\log \hat{S}(\tau_k)] \cong \sum_{i=1}^k \frac{q_i}{n'_i p_i}$$

$$\hat{\text{Var}}[\log \hat{S}(\tau_k)] = \sum_{i=1}^k \frac{\hat{q}_i}{n'_i \hat{p}_i} = \sum_{i=1}^k \frac{d_i}{n'_i (n'_i - d_i)}$$

حال با استفاده دوباره از روش دلتا، داریم:

$$\hat{\text{Var}}[\hat{S}(\tau_k)] \cong \hat{S}^2(\tau_k) \sum_{i=1}^k \frac{d_i}{n'_i (n'_i - d_i)}$$

تساوی فوق به "رابطه گرینوود" معروف است.

۴.۱ انواع جدولهای طول عمر

جدول (۱) و (۲)، مثالی از یک جدول طول عمر گروهی است. یک گروه یا دسته، مجموعه‌ای از افراد است که در جریان بررسی، مورد مطالعه قرار دارند. افراد در معرض خطر در ابتدای بازه I_i ، همان افرادی هستند که در بازه قبلی (I_{i-1}) زنده مانده‌اند (نمرده یا گم یا خارج نشده‌اند).

نوعی دیگر از جدول طول عمر، جدول طول عمر جاری است. در این جدول گروهی از افراد با سن τ_{i-1} در نظر گرفته می‌شوند، که در ابتدای بازه $I_i = (\tau_{i-1}, \tau_i]$ در معرض خطر قرار دارند. این گروه از افراد، کاملاً متفاوت از افرادی هستند که در بازه قبلی I_{i-1} در معرض خطر بوده‌اند. نوعاً، گروههای سنی مختلف در جمعیت در یک زمان بررسی می‌شوند.

REFERENCES

Berkson and Gage, Proc. Staff Meet. Mayo Clin. (1950).

Cutler and Ederer, J. Chronic Dis. (1958).

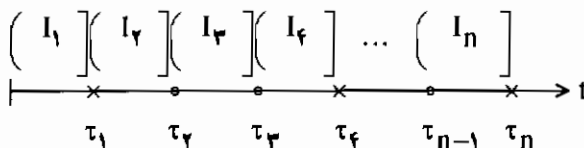
Elveback, JASA (1958).

Chiang, Stochastic Processes in Biostatistics (1968), Chapter 9.

Breslow and Crowley, Ann. Stat. (1974).

۲ برآوردگر حدی حاصل ضرب (کاپلان-میر)

برآوردگر حدی حاصل ضرب (PI)، همانند برآوردگر نیمه‌گری است با این تفاوت که طول بازه‌های I_i متغیراند. در واقع τ_i ، انتهای راست بازه I_i مشاهده بریده شده یا بریده نشده مرتبه i ام است.



نماد "x": نشانگر بریده نشده و نماد "o": نشانگر بریده شده هستند.

به خاطر داشته باشید، که مشاهدات تکراری وجود ندارند و زوجهای (Y_i, δ_i) تا (Y_n, δ_n) را مشاهده می‌کنیم. فرض کنید $Y_{(1)} < Y_{(2)} < \dots < Y_{(n)}$ ، آماره‌های مرتب

Y_1 تا Y_n باشند. همچنین، فرض کنید: $\delta_{(i)}$ مقدار δ متناظر با $Y_{(i)}$ باشد. یعنی: اگر $Y_{(i)} = Y_{j_i}$ ، آن گاه $\delta_{(i)} = \delta_{j_i}$ باشد. توجه شود که $\delta_{(1)}$ تا $\delta_{(n)}$ مرتب نیستند، فرض کنید $R(t)$ تابع خطر در زمان t باشد. یعنی: مجموعه افرادی که تا زمان t^- زنده هستند. همچنین n_i تعداد اعضای $R(Y_i)$ ، که در زمان $Y_{(i)}$ زنده هستند، d_i تعداد افرادی که در زمان $Y_{(i)}$ فوت شده و p_i ، احتمال زنده بودن در طول I_i باشد (به شرطی که در ابتدای بازه I_i زنده باشند). به عبارت معادل:

$$p_i = P(T > \tau_i | T > \tau_{i-1}) \quad \text{و} \quad q_i = 1 - p_i$$

برآوردهای \hat{p}_i و \hat{q}_i ، به شرح زیراند:

$$\hat{p}_i = 1 - \hat{q}_i = \begin{cases} 1 - \frac{1}{n_i} & \delta_{(i)} = 1 \quad (\text{بریده نشده}) \\ 1 & \delta_{(i)} = 0 \quad (\text{بریده شده}) \end{cases}$$

حال، برآورد PL، در صورت نبود تکرار، به شرح زیر است:

$$\begin{aligned} \hat{S}(t) &= \prod_{y_{(i)} \leq t} \hat{p}_i = \prod_{u: y_{(i)} \leq t} \left(1 - \frac{1}{n_i}\right) = \prod_{y_{(i)} \leq t} \left(1 - \frac{1}{n_i}\right)^{\delta_{(i)}} \\ &= \prod_{y_{(i)} \leq t} \left(1 - \frac{1}{n-i+1}\right)^{\delta_{(i)}} = \prod_{y_{(i)} \leq t} \left(\frac{n-i}{n-i+1}\right)^{\delta_{(i)}} \end{aligned}$$

REFERENCE

Kaplan and Meier, JASA (1958).

یادآوریهای لازم:

(الف) برای مشاهدات مکرر بریده نشده، فرض کنید قبل از زمان τ ، دقیقاً n فرد زنده هستند، و در زمان τ ، تعداد d تن فوت می‌شوند. فاصله d فوت را به فاصله‌های بی‌نهایت کوچک تقسیم می‌کنیم به گونه‌ای که عامل مربوط به d فوت در برآوردهای حدی حاصل ضرب به شکل زیر در آید:

$$\left(1 - \frac{1}{m}\right) \left(1 - \frac{1}{m-1}\right) \dots \left(1 - \frac{1}{m-d-1}\right) = \frac{m-d}{m} = 1 - \frac{d}{m}$$

(ب) اگر مشاهدات بریده شده و بریده نشده تکراری باشند، مشاهدات بریده نشده را قبل از مشاهدات بریده شده در نظر بگیرید.

(پ) اگر آخرین مشاهده (مرتب شده) $y(n)$ بریده شده باشد، آن گاه برای $\hat{S}(t)$ ، به طوری که در بالا تعریف شد، داریم:

$$\lim_{t \rightarrow \infty} \hat{S}(t) > 0$$

گاهی بهتر است که برای $t \geq y(n)$ تعریف کنیم $\hat{S}(t) = 0$ ، یا تصور کنیم که اگر $\delta(n) = 0$ باشد، $t > y(n)$ تعریف نشده است.

با توجه به یادآوریهای (الف) و (ب)، $Y'(1) < Y'(2) < \dots < Y'(r)$ ، زمانهای جداگانه حیات را نشان می‌دهند. همچنین:

$$\delta'_j = \begin{cases} 1 & \text{اگر مشاهده در زمان } y'(j) \text{ بریده نشده باشد} \\ 0 & \text{اگر مشاهده در زمان } y'(j) \text{ بریده شده باشد} \end{cases}$$

$$n_j = \text{تعداد } R(y'(j)) \text{ ها}$$

$$d_j = \text{تعداد فوت شده‌ها در زمان } y'(j)$$

با توجه به موارد بالا، برآورد PL برای حالت تکراری، به شرح زیر است:

$$\hat{S}(t) = \prod_{u: y'(j) \leq t} \left(1 - \frac{d_j}{n_j}\right) = \prod_{y'(j) \leq t} \left(1 - \frac{d_j}{n_j}\right)^{\delta'_j(j)}$$

مثال. مطالعه درمان AML: یک آزمایش بالینی به منظور ارزشیابی تأثیر یک روش شیمیایی روی بیماری (AML) صورت گرفته است. بعد از رسیدن به مرحله خاصی، بیماران به تصادف به دو گروه تقسیم شده‌اند. گروه اول تحت درمان شیمیایی قرار گرفته‌اند و گروه دوم یا شاهد، بدون درمان مورد مطالعه قرار گرفته‌اند. هدف این است که آیا درمان شیمیایی زمان شفایافتن را به تأخیر می‌اندازد. یعنی: آیا سبب افزایش زمان بهبودی می‌شود.

برای یک تحلیل اولیه در طول دوره آزمایش، داده‌ها (برحسب ۱۰/۷۴) به شرح صفحه بعد بوده است. طول بهبودی کامل برحسب هفته است.

۹, ۱۳, ۱۳⁺, ۱۸, ۲۳, ۲۸⁺, ۳۱, ۳۴, ۴۵⁺, ۴۸, ۱۶۱⁺

گروه تیمار

۵, ۵, ۸, ۸, ۱۲, ۱۶⁺, ۲۳, ۲۷, ۳۰, ۳۳, ۴۳, ۴۵

گروه شاهد

برآورد گر PL (کاپلان-میر) برای گروه تیمار، به صورت زیر محاسبه می‌شود:

$$\hat{S}(0) = 1$$

$$\hat{S}(9) = \hat{S}(0) \times \frac{10}{11} = 0,91$$

$$\hat{S}(13) = \hat{S}(9) \times \frac{9}{10} = 0,82$$

$$\hat{S}(18) = \hat{S}(13) \times \frac{7}{8} = 0,72$$

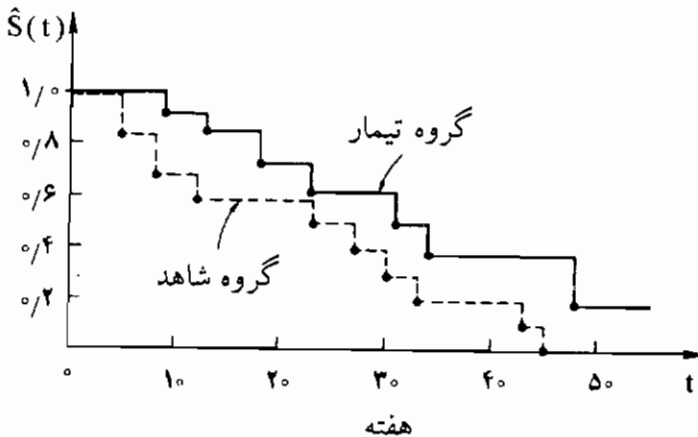
$$\hat{S}(23) = \hat{S}(18) \times \frac{6}{7} = 0,61$$

$$\hat{S}(31) = \hat{S}(23) \times \frac{4}{5} = 0,49$$

$$\hat{S}(34) = \hat{S}(31) \times \frac{3}{4} = 0,37$$

$$\hat{S}(48) = \hat{S}(34) \times \frac{1}{2} = 0,18$$

در نمودار (۳)، برآورد گرهای PL برای گروه تیمار و گروه شاهد نشان داده شده است.



نمودار ۳ برآورد تابع بقاء برای مطالعه درمان AML

REFERENCE

Embury et al., West. J. Med. (1977).

واریانس $\hat{S}(t)$: با استفاده از روش محاسبه واریانس بیمه‌گری و در حالت بدون تکرار، داریم:

$$\begin{aligned}\hat{\text{Var}}[\hat{S}(t)] &= \hat{S}^{\vee}(t) \sum_{y(i) \leq t} \frac{\hat{q}_i}{n_i \hat{p}_i} \\ &= \hat{S}^{\vee}(t) \sum_{y(i) \leq t} \frac{\delta(i)}{(n-i)(n-i+1)}\end{aligned}$$

در حالت وجود تکرار، داریم:

$$\hat{\text{Var}}[\hat{S}(t)] = \hat{S}^{\vee}(t) \sum_{y(i) \leq t} \frac{\delta'(j) d_j}{n_j (n_j - d_j)}$$

این تساویها به روابط گرین‌وود معروف است.

بررسی این روابط به روشنی جدول طول عمر نیست، زیرا تعداد جملات حاصل ضرب تصادفی است. همچنین، بیشتر جملات آن وابسته‌اند. با این وجود، بعداً به صورت تقریب مجانبی واریانس $\hat{S}(t)$ ، تحقیق خواهد شد.

توماس و میر، سه روش مختلف ساختن بازه اطمینان را مطالعه کرده‌اند. در یکی از

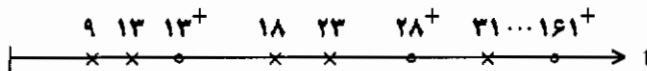
روشها از تقریب واریانس $\hat{\text{Var}}[\hat{S}(t)]$ ، استفاده شده است. همچنین به یادآوریهایی پایان فصل ششم، بخش (۴.۱)، مراجعه فرمایید.

REFERENCE

Thomas and Grunkemeier, JASA (1975).

۱.۲ الگوریتم تجدید نظر در توزیع به راست

افرون، روش دیگری را برای محاسبه برآوردگر PL معرفی کرد. آن را با مثال (AML) تشریح می‌کنیم. نمودار ($n=11$) زمان بقاء را رسم می‌کند.



برآورد معمولی $S(t)$ ، با فرض بریده نشدن به هر یک از زمانهای مشاهده شده، جرم $\frac{1}{11}$ را نسبت می‌دهد. اولین زمان بریده شده را در نظر می‌گیریم (13^+)، چون فوتی در 13^+

صورت نگرفته، لذا، در جایی و در سمت راست آن رخ می‌دهد. معقول به نظر می‌رسد که جرم $\frac{1}{11}$ نقطه ۱۳^+ را به‌طور مساوی بین تمام مشاهدات راست ۱۳^+ ، توزیع کنیم. بنابراین، جرم $(\frac{1}{11})$ را به ۱۸^+ ، ۲۳^+ ، ۲۸^+ ، ... اضافه می‌کنیم. حال دومین مشاهده بریده شده (۲۸^+) را در نظر می‌گیریم. جرم $(\frac{1}{11}) + \frac{1}{11}$ نقطه ۲۸^+ را بین تمام نقاط سمت راست ۲۸^+ توزیع می‌کنیم. در مورد بقیه مشاهدات بریده شده نیز به‌طور مشابه عمل می‌کنیم. نتایج مربوط به PL در جدول صفحه بعد آمده است:

REFERENCE

Efron, Proc. Fifth Berkeley Symp. IV (1967), pp. 831-853.

۲.۲ خودسازگاری

برای سادگی، فرض می‌کنیم تکرار وجود ندارد. برآوردگر $\hat{SC}(t)$ را خودسازگار می‌نامند، اگر به شرح زیر باشد:

$$\hat{SC}(t) = \frac{1}{n} \left[\sum_{i=1}^n (1) \cdot I(y(i) > t) + \sum_{i=1}^n (0) \cdot I(y(i) \leq t, \delta(i) = 1) + \sum_{i=1}^n \frac{\hat{SC}(t)}{\hat{SC}(y(i))} I(y(i) \leq t, \delta(i) = 0) \right] \quad (۴)$$

که در آن $\hat{SC}(t) / \hat{SC}(y(i))$ ، برآورد احتمال شرطی بقاء بعد از t است، با این شرط که در زمان $y(i)$ زنده بوده است. توجه کنید که (۴) معادل رابطه زیر است:

$$\hat{SC}(t) = \frac{1}{n} \left[N_y(t) + \sum_{y(i) \leq t} (1 - \delta(i)) \frac{\hat{SC}(t)}{\hat{SC}(y(i))} \right] \quad (۵)$$

که در آن $N_y(t) = \#(y_i > t)$.

(نماد " $\#$ ")، یعنی: تعدادی که در این شرط صدق می‌کنند.)

برآوردگر PL تنها برآوردگر خودسازگار برای $t < y(n)$ است. اثبات به قرار زیر

است:

$Y(i)$	جرم اولیه	جرم بعد از اولین تجدید	جرم بعد از دومین تجدید	جرم بعد از سومین تجدید	$\hat{S}(y(i))$
۹	$\frac{1}{11} = 0,09$	$0,09$	$0,09$	$0,09$	$0,91$
۱۳	$0,09$	$0,09$	$0,09$	$0,09$	$0,82$
13^+	$0,09$	0	0	0	
۱۸	:	$0,09 + (\frac{1}{8})(0,09) = 0,10$	$0,10$	$0,10$	$0,72$
۲۳		$0,10$	$0,10$	$0,10$	$0,61$
23^+		$0,10$	0	0	
۳۱	:	$0,10 + (\frac{1}{5})(0,10) = 0,12$	$0,12$	$0,12$	$0,49$
۳۴		$0,12$	$0,12$	$0,12$	$0,37$
34^+		$0,12$	0	0	
۴۸		$0,12 + (\frac{1}{4})(0,12) = 0,18$	$0,18$	$0,18$	$0,18$
48^+			$0,18$	$0,18$	

با توجه به رابطه (۵)، برآوردگر خودسازگار در شرایط زیر صدق می‌کند:

$$\hat{SC}(t) = \frac{N_y(t)}{n - \sum_{y(i) \leq t} \left(\frac{1 - \delta(i)}{\hat{SC}(y(i))} \right)}$$

$$= \begin{cases} 1 & t < y(1) \\ \frac{N_y(t)}{n - \sum_{i=1}^k \left(\frac{1 - \delta(i)}{\hat{SC}(y(i))} \right)} & y(k) \leq t < y(k+1) \\ & k=1, 2, \dots, n-1 \end{cases} \quad (6)$$

باید ثابت کرد که اگر $\hat{SC}(t)$ در رابطه (۶) صدق کند، آن‌گاه $\hat{SC}(t)$ با برآوردگر PL، $\hat{S}(t)$ برابر است. ابتدا، توجه کنید اگر $t < y(1)$ ، آن‌گاه

$$\hat{S}(t) = 1 - \hat{SC}(t)$$

همچنین، $\hat{S}(t)$ و $\hat{SC}(t)$ در بازه $(y(k), y(k+1))$ برای $k=1, 2, \dots, n-1$ ثابت است. بنابراین، تنها کافی است نشان دهیم که جهش تابع $\hat{SC}(t)$ در نقطه $y(k)$ برابر جهش تابع $\hat{S}(t)$ است.

(الف) اگر $\delta(k) = 0$ باشد، از (۶) نتیجه زیر به دست می‌آید:

$$N_y(y(k)) - 1 = N_y(y(k)) = \hat{SC}(y(k)) \left[n - \sum_{i=1}^k \left(\frac{1 - \delta(i)}{\hat{SC}(y(i))} \right) \right]$$

$$= \hat{SC}(y(k)) \left[n - \sum_{i=1}^{k-1} \left(\frac{1 - \delta(i)}{\hat{SC}(y(i))} \right) \right] - 1$$

$$= \hat{SC}(y(k)) \left[\frac{N_y(y(k) - 1)}{\hat{SC}(y(k) - 1)} \right] - 1$$

در نتیجه:

$$\hat{SC}(y(k)) = \hat{SC}(y(\bar{k}))$$

پس اگر $\delta(k) = 0$ باشد، $\hat{SC}(t)$ در نقطه $y(k)$ دارای جهش نیست. در نتیجه در $t = y(k)$ با $\hat{S}(t)$ برابر است.

(ب) اگر $\delta(k) = 1$ باشد، از (۶) نتیجه زیر به دست می‌آید:

$$\begin{aligned} \hat{SC}(y(k)) &= \frac{N_y(y(k))}{n - \sum_{i=1}^k \left(\frac{1 - \delta(i)}{\hat{SC}(y(i))} \right)} = \frac{N_y(y(k))}{N_y(y(\bar{k}))} \times \frac{N_y(y(\bar{k}))}{n - \sum_{i=1}^{k-1} \left(\frac{1 - \delta(i)}{\hat{SC}(y(i))} \right)} \\ &= \frac{n-k}{n-k+1} \hat{SC}(y(\bar{k})) \end{aligned}$$

در نتیجه، اگر $\delta(k) = 1$ باشد، $\hat{SC}(t)$ در نقطه $y(k)$ دارای جهش است لذا داریم:

$$\frac{\hat{SC}(y(k))}{\hat{SC}(y(\bar{k}))} = \frac{n-k}{n-k+1}$$

که در نقطه $t = y(k)$ با $\hat{S}(t)$ برابر است.

الگوریتم خودسازگاری

برآوردگر ساده زیر را در نظر می‌گیریم:

$$\hat{S}^0(t) = \frac{N_y(t)}{n}$$

این برآوردگر را با استفاده از رابطه بازگشتی زیر می‌توان بهبود بخشید:

$$\hat{S}^{(j+1)}(t) = \frac{1}{n} \left[N_y(t) + \sum_{y(i) \leq t} (1 - \delta(i)) \frac{\hat{S}^{(j)}(t)}{\hat{S}^{(j)}(y(i))} \right]$$

در واقع، $\hat{S}^{(j)}(t)$ به‌طور یکنواخت در چند مرحله متناهی به برآوردگر PL، همگرا

می‌شود. این الگوریتم محاسباتی می‌تواند در مسائل کلی برش مفید باشد.

REFERENCES

Efron, Proc. Fifth Berkeley Symp. IV (1967).

Turnbull, JASA (1974).

_____, JRSS B (1976).

۳.۲ برآوردگر حداکثر درست‌نمایی تعمیم یافته

در شرایط معمولی، فرض می‌کنیم که مشاهده \underline{X} دارای توزیع احتمال P_θ ، صادق در رابطه: $dP_\theta(\underline{x}) = f_\theta(\underline{x})d\mu(\underline{x})$ است. $\mu(\underline{x})$ ، اندازه کراندار برای رده $\{P_\theta\}$ است. به دست آوردن برآوردگر حداکثر درست‌نمایی معادل حداکثر کردن $L(\theta) = f_\theta(\underline{x})$ است. در حالت مورد بحث، فرض کنید اندازه احتمال مشاهده برابر P_F باشد، که به توزیع مجهول F بستگی دارد. رده $\{P_F\}$ دارای اندازه کراندار نیست، بنابراین به یک تعریف کلی‌تر برای حداکثر درست‌نمایی نیاز داریم.

کیفر و لفویتز، تعریف زیر را پیشنهاد نموده‌اند. فرض کنید $\mathcal{P} = \{P\}$ رده اندازه‌های احتمال باشد. برای دو عضو P_1 و P_2 در \mathcal{P} تابع زیر، که مشتق رادون-نیکودین P_1 نسبت به $P_1 + P_2$ است، را تعریف می‌کنیم.

$$f(\underline{x}; P_1, P_2) = \frac{dP_1(\underline{x})}{d(P_1 + P_2)}$$

اندازه احتمال \hat{P} را برآوردگر حداکثر درست‌نمایی تعمیم یافته (GMLE) گوئیم هرگاه رابطه زیر برای هر $P \in \mathcal{P}$ برقرار باشد:

$$f(\underline{x}; \hat{P}, P) \geq f(\underline{x}; P, \hat{P}) \quad (7)$$

این تعمیم، شامل تعریف معمولی MLE نیز هست.

برآوردگر کاپلان-میر PL، مقدار GMLE تابع F را ارائه می‌کند. اثبات به شرح

زیر است:

اگر یک اندازه احتمال \hat{P} به \underline{x} احتمال مثبت نسب دهد، آن‌گاه $f(\underline{x}; \hat{P}, P) = 0$ است، مگر این که P نیز به \underline{x} احتمال مثبتی نسبت دهد. پس، برای امتحان (7) به ازای $P \in \mathcal{P}$ ، کافی است، آن را برای P ‌هایی امتحان کنیم که $P\{\underline{x}\} > 0$ باشد. در نتیجه (7) به صورت زیر در می‌آید:

$$\hat{P}(\underline{x}) \geq P(\underline{x}) \quad (۸)$$

چون \hat{S} به نقطه $[\underline{x} = (y_1, \delta_1), \dots, (y_n, \delta_n)]$ جرم مثبتی را نسبت می‌دهد. تنها کافی است آن اندازه‌های احتمال P را در نظر بگیریم که جرم مثبت به این نقطه نسبت می‌دهند. سپس، نشان دهیم که \hat{S} مقدار احتمال $\{P\{(y_1, \delta_1), \dots, (y_n, \delta_n)\}$ را حداکثر می‌کند. برای چنین P ی، رابطه زیر را داریم:

$$L = P\{(y_1, \delta_1), \dots, (y_n, \delta_n)\} \\ = \prod_{i=1}^n P\{T=y_{(i)}\}^{\delta_{(i)}} P\{T>y_{(i)}\}^{1-\delta_{(i)}}$$

فرض کنید P احتمال p_i را به فاصله نیم‌باز $[y_{(i)}, y_{(i+1)})$ با شرط $y_{(n+1)} = +\infty$ ، نسبت می‌دهد. برای مقادیر ثابت p_1 تا p_n ، تابع درست‌نمایی در ازاء $P\{T=y_{(i)}\} = p_i$ با شرط $\delta_{(i)} = 1$ ، حداکثر می‌شود. همچنین اگر $\delta_{(i)} = 0$ باشد، L به ازاء $P\{y_{(i)} < T < y_{(i+1)}\} = p_i$ حداکثر می‌شود. پس برای مقادیر ثابت p_1 تا p_n ، مقدار حداکثر L به صورت زیر است:

$$\prod_{i=1}^n p_i^{\delta_{(i)}} \left(\sum_{j=1}^n p_j \right)^{1-\delta_{(i)}}$$

با توجه به مسأله (۸)، دیده می‌شود که رابطه (۹) به ازاء مقدار زیر حداکثر می‌شود:

$$\hat{p}_i = \prod_{j=1}^{i-1} \left(1 - \frac{\delta_{(i)}}{n-j+1} \right) \frac{\delta_{(i)}}{n-i+1}$$

که این متناظر با \hat{S} است. اثبات در حالت تکراری به‌طور مشابه انجام می‌شود.

REFERENCES

- Kiefer and Wolfowitz, Ann. Math. Stat. (1956).
Kaplan and Meier, JASA (1958).
Johansen, Scand. J. Stat. (1978).

۴.۲ سازگاری

می‌دانیم $S(t)$ به شرح زیر است:

$$S(t) = S_T(t) = P(T > t) = 1 - F(t)$$

تابع S^* را به صورت زیر تعریف می‌کنیم:

$$S^*(t) = S_Y(t) = P(Y > t) = 1 - H(t)$$

$$= [1 - F(t)][1 - G(t)]$$

توابع بقاء جزئی را به صورت زیر تعریف می‌کنیم:

$$S_u^*(t) = P\{Y > t, \delta = 1\} = \int_t^\infty [1 - G(u)] dF(u)$$

$$S_c^*(t) = P\{Y > t, \delta = 0\} = \int_t^\infty [1 - F(u)] dG(u)$$

در این صورت داریم:

$$S^*(t) = S_u^*(t) + S_c^*(t)$$

نشان خواهیم داد که $S(t)$ را می‌توان به صورت تابعی از $S_u^*(t)$ و $S_c^*(t)$ تعریف کرد.

(الف) فرض کنید $S_u^*(t)$ پیوسته است.

$$\begin{aligned} \int_0^t \frac{dS_u^*(u)}{S_u^*(u) + S_c^*(u)} &= \int_0^t \frac{-[1 - G(u)] dF(u)}{[1 - F(u)][1 - G(u)]} \\ &= \int_0^t \frac{-dF(u)}{1 - F(u)} = \log[1 - F(u)] \Big|_0^t = \log S(t) \end{aligned}$$

بنابراین داریم:

$$S(t) = \exp \left[\int_0^t \frac{dS_u^*(u)}{S_u^*(u) + S_c^*(u)} \right]$$

(ب) فرض کنید: S_u^* دارای جهش در t باشد. ولی S_c^* در t پیوسته باشد، داریم:

$$\begin{aligned} \log \frac{S_u^*(t^+) + S_c^*(t^+)}{S_u^*(t^-) + S_c^*(t^-)} &= \log \frac{[1 - F(t^)][1 - G(t^+)]}{[1 - F(t^-)][1 - G(t^-)]} \\ &= \log \frac{[1 - F(t^+)]}{[1 - F(t^-)]} = \log \frac{S(t^+)}{S(t^-)} \end{aligned}$$

(تساوی دوم از پیوستگی S_C^* در t به دست می‌آید، که در آن $G(t^+) = G(t^-)$)
بنابراین:

$$S(t^+) = S(t^-) = \exp \left\{ \log \left[\frac{S_U^*(t^+) + S_C^*(t^+)}{S_U^*(t^-) + S_C^*(t^-)} \right] \right\}$$

اگر توزیعهای F و G جهش مشترک نداشته باشند، آن‌گاه از (الف) و (ب) رابطه زیر نتیجه می‌شود:

$$S(t) = \exp \left\{ c \int_0^t \frac{dS_U^*(u)}{S_U^*(u) + S_C^*(u)} + d \sum_{u \leq t} \log \left[\frac{S_U^*(u^+) + S_C^*(u^+)}{S_U^*(u^-) + S_C^*(u^-)} \right] \right\} \quad (10)$$

در این رابطه c ، یعنی انتگرال روی بازه‌های پیوسته S_U^* و $d \sum$ ، یعنی مجموع روی نقاط جهش S_U^* است، عبارت (10) را رابطه پترسن نامند و نشان می‌دهد که $S(t)$ را می‌توان به صورت تابعی از S_U^* و S_C^* و t نشان داد؛ یعنی:

$$S(t) = \psi(S_U^*, S_C^*; t)$$

رابطه پترسن ثابت می‌کند، که برآوردگر PL مربوط به $\hat{S}(t)$ سازگار است. اثبات آن به شرح زیر است. دو تابع توزیع تجربی را به صورت زیر تعریف می‌کنیم:

$$\hat{S}_U^*(t) = \frac{1}{n} \sum_{i=1}^n I(Y_i > t, \delta_i = 1)$$

$$\hat{S}_C^*(t) = \frac{1}{n} \sum_{i=1}^n I(Y_i > t, \delta_i = 0)$$

می‌توان دید که برآوردگر PL برابر زیر است:

$$\hat{S}(t) = \psi(\hat{S}_U^*, \hat{S}_C^*; t)$$

به شرط آن‌که هر تکرار بین مشاهدات بریده شده و نشده به عنوان یک مشاهده بریده نشده بعد از برش در نظر گرفته شود. توجه شود که چون \hat{S}_U^* گسسته است، $\psi(\hat{S}_U^*, \hat{S}_C^*; t)$ تنها شامل مجموع روی جهشهای \hat{S}_U^* است.

بنابر قضیه گلیونکو-کانتلی، داریم:

$$\hat{S}_u^*(t) \xrightarrow{a.s.} S_u^*(t)$$

$$\hat{S}_c^*(t) \xrightarrow{a.s.} S_c^*(t) \quad \text{به طور یکنواخت در } t$$

(نماد " $\xrightarrow{a.s.}$ " همگرایی تقریباً همه جا را نشان می‌دهد). همچنین، ψ یک تابع پیوسته از S_u^* و S_c^* در اندازه کوچکترین کران بالاست. یعنی اگر داشته باشیم:

$$\|S_u^* - S_u^{**}\| = \sup_t |S_u^*(t) - S_u^{**}(t)| \rightarrow 0$$

و

$$\|S_c^* - S_c^{**}\| \rightarrow 0$$

آن گاه داریم:

$$\psi(S_u^*, S_c^*; t) \rightarrow \psi(S_u^{**}, S_c^{**}; t)$$

بنابراین، داریم:

$$\hat{S}(t) = \psi(\hat{S}_u^*, \hat{S}_c^*; t) \xrightarrow{a.s.} \psi(S_u^*, S_c^*; t) = S(t)$$

REFERENCE

Peterson, JASA (1977).

۵.۲ نرمال مجانبی

نشان می‌دهیم که اگر F و G بر $[0, T]$ پیوسته و $F(t) < 1$ باشد، با شرط $n \rightarrow \infty$ آن گاه:

$$Z_n(t) = \sqrt{n} [\hat{S}(t) - S(t)] \xrightarrow{w} Z(t)$$

که در آن $Z(t)$ یک فرایند گاوسی با گشتاورهای زیر است:

$$E[Z(t)] = 0$$

$$\begin{aligned} \text{Cov}[Z(t_1), Z(t_2)] &= S(t_1) \cdot S(t_2) \times \int_0^{t_1 \wedge t_2} \frac{dF_u(u)}{[1-H(u)]^2} \\ &= S(t_1) \cdot S(t_2) \times \int_0^{t_1 \wedge t_2} \frac{dF_u(u)}{[1-F(u)][1-H(u)]} \end{aligned}$$

که در آن:

$$F_u(t) = P(Y \leq t, \delta = 1) = \int_0^t [1 - G(u)] dF(u)$$

$$1 - H(u) = [1 - F(u)][1 - G(u)]$$

اثبات شامل تابع نرخ شکست است، که در بخش بعدی به آن می‌پردازیم.

توجه شود که $Z_n(t)$ ، به طو ضعیف (\xrightarrow{w}) به فرایند گاوسی $Z(t)$ ، میل می‌کند. یعنی: به ازاء هر t_1 تا t_k ، $Z_n(t_1)$ تا $Z_n(t_k)$ دارای توزیع مجانبی نرمال چندمتغیره هستند و دنباله اندازه‌های احتمال Z_n ، ملایم است، به گونه‌ای که $f(Z_n)$ در توزیع به $f(Z)$ - برای هر تابع پیوسته در اندازه کوچکترین کران بالا- همگراست. یک حالت خاص نتیجه بالا، به شرح زیر است:

$$\hat{S}(t) \underset{a}{\sim} N \left(S(t), \frac{S^{\vee}(t)}{n} \int_0^t \frac{dF_u(u)}{[1 - H(u)]^{\vee}} \right)$$

برای واریانس مجانبی $\hat{S}(t)$ می‌توان یک تقریب به دست آورد. زیرا $F_u(t) = P(Y \leq t, \delta = 1)$ و $H(t) = P(Y \leq t)$. موارد زیر را در نظر می‌گیریم (بدون تکرار)

$$d\hat{F}_u(y_{(i)}) = \frac{\delta_{(i)}}{n}$$

$$1 - \hat{H}(y_{(i)}) = 1 - \frac{i}{n} = \frac{n-i}{n}$$

$$1 - \hat{H}(y_{(i)}^-) = 1 - \frac{i-1}{n} = \frac{n-i+1}{n}$$

اگر در واریانس مجانبی به جای $[1 - H(u)]^{\vee}$ از $[1 - H(u)][1 - H(u^-)]$ استفاده کنیم و برآوردهای بالا را در آن قرار دهیم؛ داریم:

$$\begin{aligned} \widehat{AVar}[\hat{S}(t)] &= \frac{\hat{S}^{\vee}(t)}{n} \sum_{y_{(i)} \leq t} \frac{\delta_{(i)}/n}{[(n-i)/n][(n-i+1)/n]} \\ &= \hat{S}^{\vee}(t) \sum_{y_{(i)} \leq t} \frac{\delta_{(i)}}{(n-i)(n-i+1)} \end{aligned}$$

که دقیقاً رابطه گرینوود است. (\widehat{AVar} به معنی واریانس مجانبی است).

REFERENCES

- Billingsley, Convergence of Probability Measures (1968), for weak convergence.
Breslow and Crowley, Ann. Stat. (1974).

۳ برآوردگرهای تابع نرخ شکست

می‌دانیم تابع نرخ شکست به صورت: $\lambda(t) = \frac{f(t)}{1-F(t)}$ ، تعریف می‌شود. همچنین، مشکلات برآورد $\lambda(t)$ معادل برآورد تابع چگالی است. حالت ساده‌تر، برآورد تابع نرخ شکست تجمعی: $\Lambda(t) = \int_0^t \lambda(u) du$ ، است.

توابع Λ و S با رابطه $S(t) = e^{-\Lambda(t)}$ با هم ارتباط دارند. برای سادگی فرض می‌کنیم تکرار وجود ندارد. نلسون، $\Lambda(t)$ را به صورت زیر برآورد می‌کند:

$$\hat{\Lambda}(t) = \hat{\Lambda}_\nu(t) = \sum_{y(i) \leq t} \frac{\delta(i)}{n-i+1}$$

برآورد پترسن Λ به شرح زیر است:

$$\hat{\Lambda}_1(t) = \sum_{y(i) \leq t} -\log \left(1 - \frac{\delta(i)}{n-i+1} \right)$$

دو برآوردگر بالا بسیار نزدیک هم‌اند، زیرا، برای مقادیر کوچک x ، $\log(1-x) \cong -x$ است. برآوردگر پترسن متناظر برآوردگر PL تابع بقاء است. داریم:

$$\hat{S}_1(t) = e^{-\hat{\Lambda}_1(t)} = \prod_{y(i) \leq t} \left(1 - \frac{\delta(i)}{n-i+1} \right) = \hat{S}(t)$$

در حالی که، برآورد نلسون متناظر یک برآوردگر متفاوت دیگر تابع بقاست.

$$\hat{S}_\nu(t) = e^{-\hat{\Lambda}_\nu(t)}$$

فلمینگ و هارینگتن، $\hat{S}_\nu(t)$ را به عنوان یک برآوردگر دیگر تابع بقاء پیشنهاد می‌کنند و نشان می‌دهند که در بعضی از مواقع دارای میانگین مربع خطای کوچکتری است.

REFERENCES

Nelson, J. Qual. Tech. (1969).

_____, Technometrics (1972).

Peterson, JASA (1977).

Fleming and Harrington, unpublished manuscript (1979).

نرمال مجانبی

از نتایج استاندارد توابع توزیع، داریم:

$$\sqrt{n} [\hat{F}_u(t) - F_u(t)] \xrightarrow{w} Z_{F_u}(t)$$

$$\sqrt{n} [\hat{H}(t) - H(t)] \xrightarrow{w} Z_H(t)$$

که: Z_{F_u} و Z_H ، فرایندهای گاوسی‌اند. با بسط $\hat{A}(t)$ داریم:

$$\begin{aligned} \hat{A}(t) &= \int_0^t \frac{d\hat{F}_u(u)}{1 - \hat{H}(u^-)} \\ &= \int_0^t \left[\frac{1}{1-H} + \frac{\hat{H}-H}{(1-H)^2} + \dots \right] \left[dF_u + d(\hat{F}_u - F_u) \right] \\ &= \int_0^t \frac{dF_u}{1-H} + \int_0^t \frac{\hat{H}-H}{(1-H)^2} dF_u + \int_0^t \frac{d(\hat{F}_u - F_u)}{1-H} + \dots \\ &= \Lambda(t) + \int_0^t \frac{\hat{H}-H}{(1-H)^2} dF_u + \frac{(\hat{F}_u - F_u)(t)}{1-H(t)} - \int_0^t \frac{\hat{F}_u - F_u}{(1-H)^2} dH + \dots \end{aligned}$$

تساوی آخر از انتگرال جزء به جزء به دست می‌آید. با تبدیل و ضرب در \sqrt{n} داریم:

$$\begin{aligned} \sqrt{n} [\hat{A}(t) - \Lambda(t)] &= \int_0^t \frac{\sqrt{n} (\hat{H}-H)}{(1-H)^2} dF_u \\ &\quad + \frac{\sqrt{n} (\hat{F}_u - F_u)(t)}{1-H(t)} - \int_0^t \frac{\sqrt{n} (\hat{F}_u - F_u)}{(1-H)^2} dH + \dots \\ &\xrightarrow{w} \int_0^t \frac{Z_H}{(1-H)^2} dF_u + \frac{Z_{F_u}(t)}{1-H(t)} - \int_0^t \frac{Z_{F_u}}{(1-H)^2} dH = Z_\Lambda(t) \end{aligned}$$

حد $Z_{\Lambda}(t)$ برابر میانگین موزون فرایندهای گاوسی است و خود نیز یک فرایند گاوسی است. داریم:

$$E[Z_{\Lambda}(t)] = 0$$

$$\text{Cov}[Z_{\Lambda}(t_1), Z_{\Lambda}(t_2)] = \int_0^{t_1 \wedge t_2} \frac{dF_u}{(1-H)^2}$$

با استفاده از رابطه و تقریب $\hat{S}(t) = e^{-\hat{\Lambda}(t)}$ ، توزیع مجانبی $\hat{S}(t)$ به دست می‌آید.

$$e^{-\hat{\Lambda}(t)} = e^{-\Lambda(t)} - [\hat{\Lambda}(t) - \Lambda(t)]e^{-\Lambda(t)} + \dots$$

$$\hat{S}(t) \cong S(t) - [\hat{\Lambda}(t) - \Lambda(t)]S(t) + \dots$$

$$\sqrt{n} [\hat{S}(t) - S(t)] \cong -\sqrt{n} [\hat{\Lambda}(t) - \Lambda(t)]S(t) + \dots$$

$$\xrightarrow{w} Z(t)$$

به گونه‌ای که $Z(t)$ یک فرایند گاوسی با $E[Z(t)] = 0$ و کوواریانس زیر است:

$$\text{Cov}[Z(t_1), Z(t_2)] = S(t_1) \cdot S(t_2) \times \int_0^{t_1 \wedge t_2} \frac{dF_u}{(1-H)^2}$$

REFERENCES

Breslow and Crowley, Ann. Stat. (1974).

Aalen, Scand. J. Stat. (1976).

_____, Ann. Stat. (1978).

۴ برآوردهای تنومند

در مسائل برآورد، اغلب می‌توان پارامتر مورد علاقه را به صورت تابعی مانند:

$$\theta = T(F)$$

نشان داد، که در آن F ، تابع توزیع مربوط است.

اگر برش نباشد برآوردگر معمول به صورت: $\hat{\theta} = T(F_n)$ است، که F_n ، تابع توزیع

تجربی است، ولی اگر برش وجود داشته باشد، یک برآوردگر معقول به صورت:

$$\hat{\theta} = T(\hat{F})$$

است، که در آن $\hat{F} = 1 - \hat{S}$ و \hat{S} برآوردگر PL است.

۱.۴ میانگین

$$\theta = T(F) = \int_0^{\infty} x dF(x) = \int_0^{\infty} [1 - F(x)] dx = \int_0^{\infty} S(t) dt$$

در صورت نبود برش، داریم:

$$\hat{\theta} = T(F_n) = \int_0^{\infty} x dF_n(x) = \bar{x} = \int_0^{\infty} [1 - F_n(x)] dx$$

در صورت وجود برش، داریم:

$$\hat{\theta} = T(\hat{F}) = \int_0^{\infty} x d\hat{F}(x) = \int_0^{\infty} \hat{S}(t) dt$$

$$AVar(\hat{\theta}) = \frac{1}{n} \int_0^{\infty} \frac{1}{[1 - H(s)]^2} \left(\int_s^{\infty} S(u) du \right)^2 dF_u(s)$$

در حالت بدون تکرار، داریم:

$$A\hat{V}ar(\hat{\theta}) = \sum_{i=1}^n \left(\int_{y(i)}^{\infty} \hat{S}(u) du \right)^2 \frac{\delta(i)}{(n-i)(n-i+1)}$$

اگر $y(n)$ بریده شود، آن گاه با شرط $t \rightarrow \infty$ ، برآورد $\hat{S}(t)$ به صفر میل نمی کند. در نتیجه حاصل انتگرال بی نهایت خواهد شد. سه راه حل را بررسی می کنیم:

۱ تعریف دوباره آخرین مشاهده. با تعویض $\delta(n) = 0$ به $\delta(n) = 1$ ؛ از داده‌ها AML برای تشریح استفاده می کنیم:

$$\begin{aligned} \hat{\theta} &= 9 \times 0,091 + 13 \times 0,091 + 18 \times 0,102 + 23 \times 0,102 + 31 \\ &\quad \times 0,123 + 34 \times 0,123 + 48 \times 0,184 + (161 \times 0,184) \\ &= 23,011 + (29,624) = 52,635 \end{aligned}$$

دنباله و بخصوص آخرین مشاهده، وزن زیادی دارند. علت این است که برآورد گر PL وزنهای صعودی را به آخرین مشاهدات و به چولگی توزیع می دهد.

۲ میانگین محدود شده (مایر و سنדר). به ازای مقدار ثابت S_0 ، میانگینی را در بازه $[0, S_0]$ تعریف نموده و آنرا به شرح زیر برآورد می کنیم:

$$\hat{\theta} = \int_0^{S_0} \hat{S}(t) dt$$

۳ حد متغیر (سوسارلا و وان ریزین). $\theta = \int_0^{\infty} S(t) dt$ را با $\hat{\theta} = \int_0^{S_n} \hat{S}(t) dt$ برآورد می‌کنیم. در این برآورد، $\{S_n\}$ دنباله‌ای از اعدادی است که به طور یکنواخت به ∞ میل می‌کند.

متأسفانه، انتخاب مناسب S_n به F و G بستگی دارد و در این مورد روش کاربردی دقیق وجود ندارد.

REFERENCES

- Kaplan and Meier, JASA (1958).
 Meier, Perspectives in Prob. and Stat. (1975).
 Sander, Stanford Univ. Tech. Report No. 8. (1975).
 Susarla and Van Ryzin, Ann. Stat. (1980).

۲.۴ برآوردگرهای L

یک فرض اساسی برای استفاده از برآوردگرهای L ، این است که توزیع اولیه F نسبت به θ متقارن باشد. نوعاً، زمانهای بقاء توزیع متقارن ندارند. زیرا مثبت‌اند. ولی می‌توان قبل از برآورد داده‌ها را با استفاده از یک تبدیل متقارن کرد. برای مثال با لگاریتم‌گیری. یک برآوردگر L به صورت زیر است.

$$\hat{\theta} = \int_{-\infty}^{+\infty} xJ(\hat{F}(x))d\hat{F}(x)$$

در این رابطه J بر $[0, 1]$ تعریف شده، نسبت به $\frac{1}{4}$ متقارن و در رابطه زیر صدق می‌کند:

$$\int_0^1 J(u) du$$

یک برآوردگر مهم L ، با پیرایش کردن میانگین به شرح زیر به دست می‌آید:

$$J(u) = \frac{1}{1-2\alpha} I[\alpha, 1-\alpha](u)$$

با داده‌های بریده شده، واریانس مجانبی یک برآوردگر L به صورت زیر است:

$$A\text{Var}(\hat{\theta}) = \frac{1}{n} \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} S(t) J(S(t)) S(u) \cdot J(S(u)) \left\{ \int_{-\infty}^{t \wedge u} \frac{dF_u(s)}{[1-H(s)]^2} \right\} dt du$$

REFERENCES

Sander, Stanford Univ. Tech. Report No. 8. (1975).

Reid Ann. Stat. (1981).

۳.۴ برآوردگر - M

در این جا نیز یک فرض اساسی متقارن بودن F است. در نتیجه ابتدا باید داده‌ها را تبدیل کرد. برآوردگر - M مربوط به $\hat{\theta}$ جواب معادله زیر است:

$$\int_{-\infty}^{+\infty} \psi(x - \hat{\theta}) d\hat{F}(x) = 0$$

تابع $\psi(x - \theta)$ ، تعمیم $f'(x - \theta)/f(x - \theta)$ است. در نتیجه، برآوردگرهای -L، تعمیم برآوردگرهای حداکثر درست‌نمایی‌اند. برآوردگر دو وزنی توکی متناظر با تابع زیر است:

$$\psi(x) = \begin{cases} x(1-x^2)^2 & |x| \leq 1 \\ 0 & |x| > 1 \end{cases}$$

در کاربردهای واقعی باید داده‌ها با یک مقیاس برآورد شده، مقیاس‌بندی شوند. واریانس مجانبی یک برآوردگر - M در حالت برش، به شرح زیر است:

$$A\text{Var}(\hat{\theta}) = \frac{1}{n} \int_{-\infty}^{+\infty} \frac{1}{[1-H(s)]^2} \cdot \left(\int_s^{+\infty} \frac{1}{E\psi'} S(t) \psi'(t - \theta) dt \right)^2 dF_u(s)$$

$$E\psi' = \int_{-\infty}^{+\infty} \psi'(t - \theta) dF(t)$$

در این جا:

REFERENCE

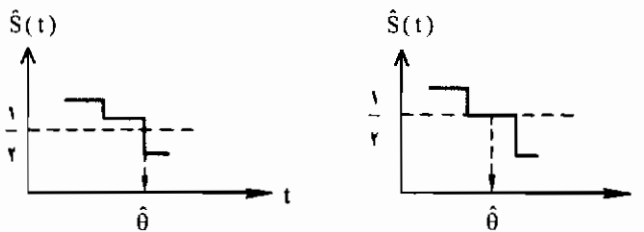
Reid Ann. Stat. (1981).

تا زمان حال برآوردگرهای -L و M با داده‌های بریده شده به‌طور آزمایشگاهی بررسی

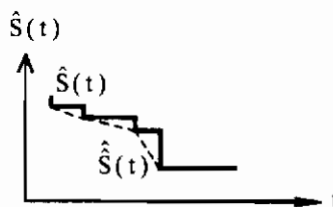
شده‌اند. عیوب و محاسن آنها بررسی نشده است. در حالی که برآوردگر میانه کاربرد عملی زیادی دارد.

۴.۴ میانه

با توجه به $\theta = S^{-1}(\frac{1}{p})$ ، یک برآوردگر معقول برای θ ، عبارت است از: اگر $\hat{\theta} = \hat{S}^{-1}(\frac{1}{p})$ دارای جواب منحصر به فرد نباشد، آن‌گاه $\hat{\theta}$ را نقطه میانی بازه شامل جوابها در نظر می‌گیریم:



شواهد تجربی نشان دهنده این است، که این برآوردگر سراسر است، خیلی بزرگ به نظر می‌رسد. برآوردگر PL با افزایش t ، جهش بزرگتری می‌دهد و به خاطر مشاهدات بریده کنار گذاشته شده، شکاف بین مشاهدات بریده نشده با t افزایش پیدا نمی‌کند. بنابراین $\hat{\theta}$ بزرگتر می‌شود. برای رفع این مشکل $\hat{S}(t)$ را به صورت خطی هموار از $\hat{S}(t)$ تعریف می‌کنیم و $\hat{\theta} = \hat{S}^{-1}(\frac{1}{p})$



برای مثال، در داده‌ها AML، داریم:

$$\hat{S}(23) = 0,614$$

$$\hat{S}(31) = 0,491$$

$$\hat{\theta} = 31 - \frac{\Lambda(0,009)}{(0,123)} = 30,415$$

باید واریانس $\hat{\theta}$ را محاسبه کرد. واریانس مجانبی به شرح: $AVar(\hat{\theta}) = \frac{AVar(\hat{S}(\theta))}{f^2(\theta)}$ است. $AVar(\hat{S}(\theta))$ را می‌توان به کمک رابطه گرین وود برآورد کرد. ولی، f یک تابع چگالی مجهول بوده و برآورد آن مشکل است.

REFERENCES

- Sander, Stanford Univ. Tech. Report No. 5 (1975), discusses the asymptotic variance.
- Földes, Rejto, and Winter, unpublished manuscript (1978), discuss density estimation using censored data.
- Reid, Ann. Stat. (1981), discusses the asymptotic variance..
- _____ and Iyengar, unpublished notes (1978), consider estimates of the variance.
- Efron, Stanford Univ. Tech. Report No. 53 (1980), uses the bootstrap to measure the variability of $\hat{\theta}$.

۵ برآوردگرهای بیزی

فرض می‌کنیم تکرار وجود ندارد. با در نظر گرفتن $N_y(t) = \#(y_i > t)$ ، داریم:

$$\begin{aligned}\hat{S}(t) &= \prod_{y(i) \leq t} \left[\frac{n-i}{n-i+1} \right]^{\delta(i)} \\ &= \prod_{y(i) \leq t} \left[\frac{n-i+1}{n-i} \right]^{-\delta(i)} \cdot \frac{1}{n} \left\{ \frac{n}{n-1} \cdot \frac{n-1}{n-2} \cdots \frac{N_y(t)+1}{N_y(t)} \right\} \frac{N_y(t)}{1} \\ &= \frac{N_y(t)}{n} \prod_{y(i) \leq t} \left[\frac{n-i+1}{n-i} \right]^{1-\delta(i)}\end{aligned}$$

سوسالار و وانرایزین، نشان داده‌اند که برآوردگر بیز $S(t)$ دارای صورت زیر نیز هست:

$$\hat{S}_\alpha(t) = \frac{\alpha(t, \infty) + N_y(t)}{\alpha(o, \infty) + n} \times \prod_{y(i) \leq t} \left[\frac{\alpha[y(i), \infty] + (n-i+1)}{\alpha[y(i), \infty] + (n-i)} \right]^{1-\delta(i)}$$

بر آوردگر $\hat{S}_\alpha(t)$ ، بر آوردگر بیز با تابع زیان زیر است:

$$L(\hat{\delta}, S) = \int_0^\infty [\hat{\delta}(t) - S(t)]^2 dw(t)$$

که در این رابطه، w هر تابع نامنفی صعودی و یا فرایند پیشین دیرکله \mathcal{P}_α با پارامتر α بر خانواده $\{P\}$ در تمام توزیعهای ممکن است. پارامتر α یک اندازه متناهی نامنفی بر $(0, \infty)$ است.

اندازه احتمال تصادفی P را دارای فرایند پیشین دیرکله با پارامتر α گوئیم، هرگاه برای هر افزایش اندازه پذیر β_1 تا β_k از $(0, \infty)$ داشته باشیم:

$$(P(\beta_1), \dots, P(\beta_k)) \sim \text{Dirichlet}(\alpha(\beta_1), \dots, \alpha(\beta_k))$$

توجه شود که توزیع دیرکله $(\alpha_1, \dots, \alpha_k)$ دارای چگالی زیر است:

$$f(x_1, \dots, x_k) \propto x_1^{\alpha_1-1} x_2^{\alpha_2-1} \dots x_k^{\alpha_k-1} \cdot I(x_i \geq 0, x_1 + \dots + x_k = 1)$$

توجه نمایید که در ازا $k=2$ ، توزیع دریکله دقیقاً همان توزیع بتاست.

فرض می کنیم مشاهده X دارای توزیع \mathcal{P}_θ است، که در آن θ طبق یک توزیع پیشین انتخاب می شود. در حالت ناپارامتری، متغیر T با توزیع P - که طبق توزیع \mathcal{P}_α انتخاب می شود - مشاهده می شود. به عبارت دیگر، زمان بقاء T توسط \mathcal{P}_α به دست می آید و سپس، P را تولید نموده و P نیز متغیر T را تولید می کند. می توان تساوی زیر را اثبات کرد:

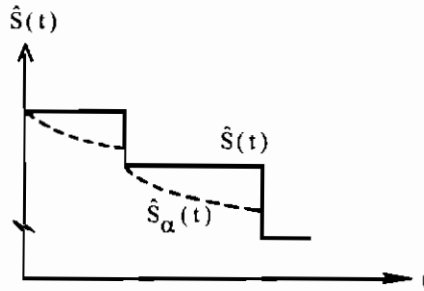
$$P\{T \in A\} = \frac{\alpha(A)}{\alpha(0, \infty)} \quad (11-1)$$

معادله (11-1)، تفسیری از پارامتر α را ارائه می کند. نسبت $\frac{\alpha(A)}{\alpha(0, \infty)}$ حدس اولیه بر احتمال مجموعه A است. برای مثال اگر فرض کنیم T دارای توزیع نمایی با میانگین $\frac{1}{\lambda_0}$ باشد، آن گاه:

$$\frac{\alpha(t, \infty)}{\alpha(0, \infty)} = e^{-\lambda_0 t}$$

همچنین، جرم کل $\alpha(0, \infty)$ ، شدت باور پیشین را نشان می دهد. برای مثال، $\alpha(0, \infty) = 10$ ، بیان می کند که باور پیشین ما ارزش ده مشاهده را دارد.

با توجه به رابطه $\frac{\alpha(t, \infty)}{\alpha(\cdot, \infty)} = e^{-\lambda_0 t}$ ، مقادیر $\hat{S}(t)$ و $\hat{S}_\alpha(t)$ مطابق شکل زیر قابل مقایسه‌اند:



سوسارلا و وان‌رایزین نشان دادند که در بسیاری از حالات، \hat{S}_α دارای میانگین مربع خطای کوچکتری از \hat{S} - حتی اگر حدس پیشین درست نباشد - است. برآورد بیزی در حالت تکرار مطابق زیر است:

$$\hat{S}_\alpha(t) = \frac{\alpha(t, \infty) + N_y(t)}{\alpha(\cdot, \infty) + n} \times \prod_{y'(j) \leq t} \left[\frac{\alpha[y'(j), \infty] + N_y(y'(j) -)}{\alpha[y'(j), \infty] + N_y(y'(j))} \right]^{1 - \delta'(j)}$$

REFERENCES

Ferguson, Ann. Stat. (1973), discusses the Dirichlet process prior.

Susarla and Van Ryzin, JASA (1976), derive the Bayes estimate in the censored case.

_____ and _____, Ann. Stat. (1978b), study the asymptotic behavior of Bayes estimates.

Ferguson and Phadia, Ann. Stat. (1979), examine more general prior distributions.

Rai, Susarla, and Van Ryzin, Comm. Stat. B (1980), look at mean square errors.

برآوردگرهای تجربی بیزی. به جای کاربرد یک حدس پیشین برای α ، می‌توان نمونه را برای برآورد α به کاربرد.

REFERENCES

Susarla and Van Ryzin, Ann. Stat. (1978a).

Phadia, Ann. Stat. (1980).

فصل چهارم

روشهای ناپارامتری (دو نمونه)

برای نمونه اول: فرض کنید T_1 تا T_m ، متغیرهای iid با توزیع F_1 و C_1 تا C_m ، متغیرهای iid با توزیع G_1 باشند. C_j زمان برش مربوط به T_j است. می‌توان داده‌های (X_1, δ_1) تا (X_m, δ_m) را به شرح زیر مشاهده نمود:

$$X_j = T_j \wedge C_j \quad \text{و} \quad \delta_j = I(T_j \leq C_j)$$

برای نمونه دوم: فرض کنید U_1 تا U_n ، متغیرهای iid با توزیع F_2 و D_1 تا D_n ، متغیرهای iid با توزیع G_2 باشند. D_j زمان برش مربوط به U_j است. می‌توان مشاهدات (Y_1, ε_1) تا (Y_n, ε_n) را به شرح زیر مشاهده نمود:

$$Y_j = U_j \wedge D_j \quad \text{و} \quad \varepsilon_j = I(U_j \leq D_j)$$

معمولاً، در مسأله دو نمونه‌ای آزمون فرض $H_0: F_1 = F_2$ ، مورد نظر است.

مثال. آزمایش بالینی فرضی: این آزمایش توسط بایرون ویام و بران جی آر، مطابق نمودارهای (۴. الف) و (ب)، بنا شده است. فرض نمایید مشاهدات X مربوط به تیمار A و مشاهدات Y مربوط به تیمار B باشند.

$$R_x A: 3, 5, 7, 9^+, 18$$

$$R_x B: 12, 19, 20, 20^+, 33^+$$

۱ آزمون گهان

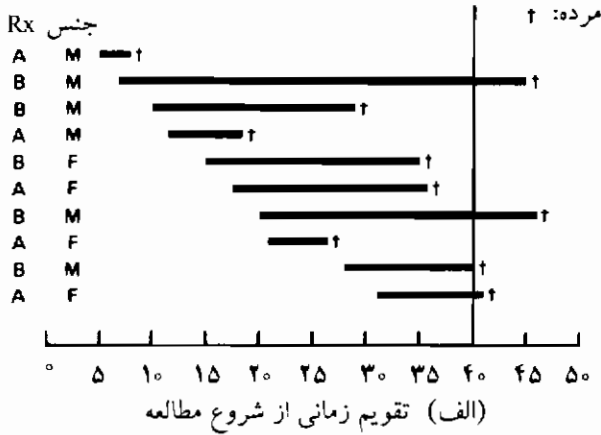
این آزمون تعمیمی از آزمون ویلکاکسن است. فرض کنید مشاهدات دو نمونه به

صورت زیر باشد:

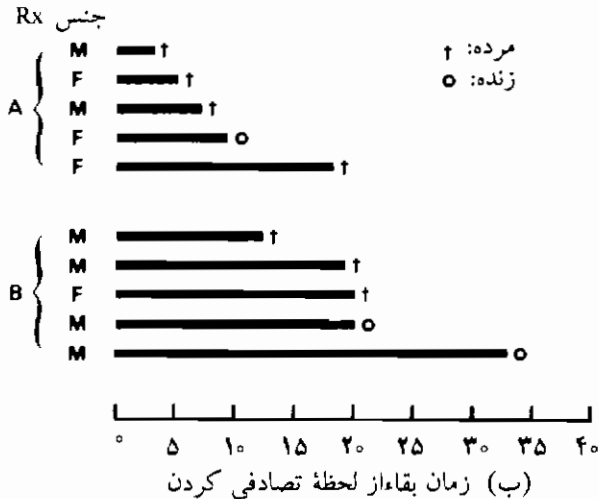
$$X_1, \dots, X_m; Y_1, \dots, Y_n$$

این ترکیب را مرتب می‌کنیم، داریم:

$$Z_{(1)}, Z_{(2)}, \dots, Z_{(m+n)}$$



نمودار ۴. (الف) بررسی زمان بقاء ده بیمار سرطانی که به تصادف تحت تیمار (A) و (B) قرار گرفته‌اند.



نمودار ۴. (ب) زمان بقاء از لحظه تصادفی کردن ده بیمار سرطانی جدا از زمان ختم مطالعه در $t = 40$.

R_1 را رتبه X_i و $R_1 = \sum_{i=1}^m R_{1i}$ در نظر می‌گیریم. فرض H_0 رد می‌شود، هرگاه R_1 خیلی کوچک یا خیلی بزرگ باشد. به کمک جدولهای مربوط، برای نمونه‌های کوچک و تقریب نرمال برای نمونه‌های بزرگ و به شرح زیر، آزمون را انجام می‌دهیم. در رابطه زیر $E_0(R_1)$ و $\text{Var}_0(R_1)$ ، گشتاورهای تحت فرض صفراند.

$$\frac{R_1 - E_0(R_1)}{\sqrt{\text{Var}_0(R_1)}} = \frac{R_1 - \frac{m(m+n+1)}{2}}{\sqrt{\frac{mn(m+n+1)}{12}}} \underset{a}{\sim} N(0, 1)$$

صورت من-ویتی آزمون ویلکاکسن مفید خواهد بود. با توجه به تعریف زیر می‌توان R_1 را مطابق بسط زیر نوشت:

$$U(X_i, Y_j) = U_{ij} = \begin{cases} +1 & X_i > Y_j \\ 0 & X_i = Y_j \\ -1 & X_i < Y_j \end{cases}$$

می‌توان نشان داد:

$$U = \sum_{i=1}^m \sum_{j=1}^n U_{ij}$$

در نتیجه:

$$R_1 = \frac{m(m+n+1)}{2} + \frac{1}{2}U$$

برای اثبات توجه شود، که اگر مشاهدات به صورت:

$$X_{(1)} < \dots < X_{(m)} < Y_{(1)} < \dots < Y_{(n)}$$

درآیند، آن‌گاه $R_1 = \frac{m(m+1)}{2}$ می‌شود. برای هر تعویض زوج X و Y ، R_1 به اندازه یک

واحد صعود می‌کند و تعداد این تعویض‌ها برابر $\sum_i \sum_j \frac{1}{2}(U_{ij}+1)$ است؛ بنابراین داریم:

$$\begin{aligned} R_1 &= \frac{m(m+1)}{2} + \sum_i \sum_j \frac{1}{2}(U_{ij}+1) = \frac{m(m+1)}{2} + \frac{mn}{2} + \frac{1}{2}U \\ &= \frac{m(m+n+1)}{2} + \frac{1}{2}U \end{aligned}$$

آزمون من-ویتنی فرض H_0 را در صورتی رد می‌کند، که U یا $|U|$ خیلی بزرگ باشد. در نمونه‌های کوچک از جدولها و در نمونه‌های بزرگ از تقریب نرمال استفاده می‌کنیم.

$$\frac{U - E_0(U)}{\sqrt{\text{Var}_0(U)}} = \frac{U}{\sqrt{\frac{mn(m+n+1)}{3}}} \stackrel{a}{\sim} N(0, 1)$$

برای داده‌های بریده شده، U_{ij} را به صورت زیر تعریف می‌کنیم:

$$U_{ij} = \begin{cases} +1 & t_i > u_j \text{ یا } (X_i > Y_j, \varepsilon_j = 1) \text{ یا } (X_i = Y_j, \delta_i = 0, \varepsilon_j = 1) \\ 0 & \text{در سایر جاها} \\ -1 & t_i < u_j \text{ یا } (X_i < Y_j, \delta_i = 1) \text{ یا } (X_i = Y_j, \delta_i = 1, \varepsilon_j = 0) \end{cases}$$

$$U = \sum_{i=1}^m \sum_{j=1}^n U_{ij}$$

فرض H_0 در صورت بزرگ بودن U یا $|U|$ رد می‌شود. آماره U به طور مجانبی دارای توزیع نرمال است. ولی، برای محاسبه معیارهای آزمون تنها گشتاورهای U کفایت می‌کند.

۱.۱ میانگین و واریانس U

در حالت بریده نشده، میانگین و واریانس را می‌توان با استفاده از قضیه جایگشت محاسبه کرد. تحت H_0 ، نمونه‌گیری از m مهره بدون جایگذاری را از جعبه‌ای شامل $m+n$ مهره به شماره‌های Z_1, \dots, Z_{m+n} در نظر بگیرید، که زیرنویسهای m مهره نمونه‌گیری شده به عنوان مقادیر X_1 تا X_m و زیرنویسهای n مهره باقی‌مانده به عنوان Y_1 تا Y_n در نظر گرفته می‌شوند. فرض کنید $E_{0,p}(U)$ و $\text{Var}_{0,p}(U)$ ، گشتاورهای تحت الگوی جایگشت باشند. در نتیجه:

$$E_{0,p}(U) = 0 = E_0(U)$$

$$\text{Var}_{0,p}(U) = \frac{mn(m+n+1)}{3} = \text{Var}_0(U)$$

در حالت برش، گهان قضیه جایگشت را تحت فرض قویتر زیر به کار می‌برد:

$$H_0^*: F_1 = F_2, G_1 = G_2$$

فرض کنید نمونه مرکب به صورت: (Z_1, ξ_1) تا (Z_{n+m}, ξ_{n+m}) باشد. نمونه‌ای شامل m مهره بدون جایگذاری را از جعبه‌ای شامل $n+m$ مهره به صورت (Z_1, ξ_1) تا (Z_{n+m}, ξ_{n+m}) در نظر می‌گیریم. زیرنویسهای نمونه m تایی را به صورت (X_1, δ_1) تا (X_m, δ_m) و زیرنویسهای n مهره باقی‌مانده را با (Y_1, ϵ_1) تا (Y_n, ϵ_n) در نظر می‌گیریم. در این صورت:

$$E_{0,p}^*(U) = 0$$

$$\text{Var}_{0,p}^*(U) = (\text{به مقاله گهان بخش ۳.۴ صفحه ۲۵۶ مراجعه شود})$$

مقدار واریانس به صورت عبارتی پیچیده در می‌آید، که از آن صرف‌نظر می‌شود. در عوض از روش مانتل مقدار $\text{Var}_{0,p}^*(U)$ به صورت ساده‌تری به دست می‌آید.

۲.۱ روش محاسباتی مانتل برای $\text{Var}_{0,p}^*(U)$

بنابه تعریف:

$$U_{k\ell} = U((Z_k, \xi_k), (Z_\ell, \xi_\ell))$$

$$U_{ij} = \begin{cases} +1 & (Z_k > Z_\ell, \xi_\ell = 1) \text{ یا } (Z_k = Z_\ell, \xi_k = 0, \xi_\ell = 1) \\ 0 & \text{در سایر جاها} \\ -1 & (Z_k < Z_\ell, \xi_k = 1) \text{ یا } (Z_k = Z_\ell, \xi_k = 1, \xi_\ell = 0) \end{cases}$$

$$U_k^* = \sum_{\substack{\ell=1 \\ \ell \neq k}}^{m+n} U_{k\ell}, \quad U = \sum_{k=1}^{m+n} U_k^* I(k \in I_1)$$

که در آن I_1 ، مجموعه اعداد صحیح در نمونه یک است. توجه شود که U معادل آماره گهان است، زیرا $U_{k_1 k_2} = -U_{k_2 k_1}$. در نتیجه اگر k_1 و k_2 عضو I_1 باشند، یکدیگر را حذف می‌کنند.

برای محاسبه توزیع جایگشت U ، فرض کنید U_1^* تا U_{n+m}^* ، تحت H_0 معلوم باشند. یک نمونه m تایی بدون جایگذاری از این U_k^* ها انتخاب می‌کنیم. سپس U ، مجموع این m مقدار را تشکیل می‌دهیم. با استفاده از نتایج نمونه‌گیری از جمعیت‌های متناهی، داریم:

$$\begin{aligned} \text{Var}_{\circ, P}^*(U) &= m \left(\frac{1}{m+n-1} \sum_{i=1}^{m+n} (U_i^*)^2 \right) \left(1 - \frac{m}{m+n} \right) \\ &= \frac{mn}{(m+n)(m+n-1)} \sum_{i=1}^{m+n} (U_i^*)^2 \end{aligned}$$

۳.۱ مثال. به کمک داده‌های مثال قبل (۴. الف و ب)، داریم:

Z	Rx	# < Z	# > Z	U*
۳	A	۰	۹	-۹
۵	A	۱	۸	-۷
۷	A	۲	۷	-۵
۹ ⁺	A	۳	۰	۳
۱۲	B	۳	۵	-۲
۱۸	A	۴	۴	۰
۱۹	B	۵	۳	+۲
۲۰	B	۶	۲	+۴
۲۰ ⁺	B	۷	۰	+۷
۲۳ ⁺	B	۷	۰	+۷

$$U = -۹ - ۷ - ۵ + ۳ + ۰ = -۱۸$$

$$E_{\circ, P}^*(U) = ۰ \quad \text{و} \quad \text{Var}_{\circ, P}^*(U) = \frac{(۵)(۵)(۲۸۶)}{(۱۰)(۹)} = ۷۹,۴۴$$

تحت فرض H_0 ، داریم:

$$\frac{U}{\sqrt{\text{Var}_{\circ, P}^*(U)}} = \frac{-۱۸}{\sqrt{۸,۹۱}} = -۲,۰۲ \approx N(۰, ۱)$$

بنابراین $P = 0.022$ ، مقدار P برای آزمون یک طرفه است.

REFERENCES

Gehan, *Biometrika* (1965).

Mantal, *Biometrics* (1967).

۴.۱ واریانس تحت H_0

نتایج مباحث قبل درباره واریانس با فرض: $G_1 = G_2$ ، $F_1 = F_2$ ، H_0^* ، به دست آمده‌اند. حال واریانس جایگشت تحت فرض $H_0: F_1 = F_2$ با طرح برش ثابت چیست؟ فرض کنید V_1, \dots, V_{m+n} ، نمونه مرکب از T_1, \dots, T_m و U_1, \dots, U_n باشد. تحت فرض H_0 ، از آن V_k را بدون جایگذاری نمونه‌گیری کرده و در عبارت زیر قرار می‌دهیم:

$$(-, D_n) \text{ و } \dots \text{ و } (-, D_1); (-, C_m) \text{ و } \dots \text{ و } (-, C_1)$$

به کمک این مشاهدات، صورت زیر مورد توجه است:

$$(X_1, \delta_1) \text{ و } \dots \text{ و } (X_m, \delta_m); (Y_1, \varepsilon_1) \text{ و } \dots \text{ و } (Y_n, \varepsilon_n)$$

$$(X_1, \delta_1), \dots, (X_m, \delta_m); (Y_1, \varepsilon_1), \dots, (Y_n, \varepsilon_n)$$

با توجه به این که تمام T_i, C_i, U_j و D_j ها قابل مشاهده نیستند، متأسفانه، نمی‌توان تمام (X_j, δ_j) و (Y_j, ε_j) ها را ساخت.

هاید $(\text{Var}_{0,P}^*(U))$ را با $\text{Var}_0(U)$ ، مقایسه می‌کند. داریم:

$$\begin{aligned} \text{Var}_0(U) &= E_0(U^2) = E \left\{ \left(\sum_{i=1}^m \sum_{j=1}^n U_{ij} \right)^2 \right\} \\ &= mn E_0(U_{ij}^2) + mn(n-1) E_0(U_{ij} U_{ij'}) + m(m-1)n E_0(U_{ij} U_{i'j}) \\ &\quad + m(m-1)n(n-1) E_0(U_{ij} U_{i'j'}) \end{aligned}$$

$E_0(\text{Var}_{0,P}^*(U)) =$ مجموع جملات مشابه

اگر با شرط $\lambda \rightarrow \frac{m}{m+n}$ و $0 < \lambda < 1$ ، n و m به بی‌نهایت میل کنند، آنگاه داریم:

$$\begin{aligned}
 R^{\lambda} &= p - \lim_{m,n \rightarrow \infty} \frac{\text{Var}_{o,p}^*(U)}{\text{Var}_o(U)} = \lim_{m,n \rightarrow \infty} \frac{E_o(\text{Var}_{o,p}^*(U))}{\text{Var}_o(U)} \\
 &= 3\lambda(1-\lambda) + \left\{ \lambda^3 P\{C_1 \wedge C_2 \wedge C_3 > T_1 \wedge T_2 \wedge T_3\} + \right. \\
 &\quad \left. + (1-\lambda)^3 P\{D_1 \wedge D_2 \wedge D_3 > T_1 \wedge T_2 \wedge T_3\} \right\} \\
 &\quad \times \left\{ \lambda P\{C_1 \wedge C_2 \wedge D_1 > T_1 \wedge T_2 \wedge T_3\} + \right. \\
 &\quad \left. + (1-\lambda) P\{C_1 \wedge D_1 \wedge D_2 > T_1 \wedge T_2 \wedge T_3\} \right\}^{-1} \quad (11)
 \end{aligned}$$

با توجه به رابطه (۱۱) داریم: $R^{\lambda} > 3\lambda(1-\lambda)$ اگر $\lambda = \frac{1}{4}$ ، آن گاه $R^{\lambda} > \frac{3}{4}$ یا $R > 0,875$. پس، اگر حجم نمونه‌ها برابر باشند، $SD_{o,p}^*(U)$ نمی‌تواند خیلی کوچکتر از $SD_o(U)$ باشد. در این جا نوع برش مهم نیست.

فرض کنید توزیعهای بریده شده به صورتهای مختلف لهن باشند. یعنی:

$$(1-G_1)^{T_1} = 1-F \quad \text{و} \quad (1-G_2)^{T_2} = 1-F$$

که در آن r_1 و r_2 به مقادیر زیر و به صورت $p_1 = \frac{1}{r_1+1}$ و $p_2 = \frac{1}{r_2+1}$ وابسته‌اند.

$$p_1 = P(C_1 < T_1) = P(\text{مشاهده در جمعیت ۱ بریده شده است})$$

$$p_2 = P(D_1 < U_1) = P(\text{مشاهده در جمعیت ۲ بریده شده است})$$

هاید در جدول ۲، مقادیر R را برای $\lambda = 0,5$ ، تحت توزیعهای لهن با تغییر سطوح بریده p_1 و p_2 محاسبه کرده است. جدول به خاطر تعیین حالت‌های $0,5 < |R-1|$ ، افزاز شده است. جدول ۳، با فرض $\lambda = 0,2$ مساوی جدول ۲ است. از جدول ۲ دیده می‌شود که آزمون گهان (طرحهای برش را مساوی فرض می‌کند) از یک انحراف معیار تقریباً صحیح - حتی وقتی که احتمالهای برش تفاوت زیادی دارند - استفاده می‌کند. علاوه بر این، هنگامی که یک نمونه چهار برابر دیگری باشد (جدول ۳)، آزمون گهان از یک انحراف معیار تقریباً صحیح برای دامنه وسیعی از احتمالهای بریده شده، استفاده می‌کند.

REFERENCES

Gilbert, Univ. Chicago thesis (1962), was the first to calculate $\text{Var}_0(U)$.

Hyde Stanford Univ. Tech. Report No. 30 (1977).

۲ آزمون مانتل-هانزل

۱.۲ جدول 2×2 . فرض کنید دو جمعیت داریم، به گونه‌ای که یک فرد می‌تواند در هر دو جمعیت یکی از دو مشخصه‌ای را دارا باشد. برای مثال، جمعیت ۱، شامل بیماران سرطانی تحت تیمار معین و جمعیت ۲، شامل بیماران تحت تیمار دیگر است. ممکن است، بیماران هر یک از دو جمعیت در عرض یک سال بمیرند یا زنده بمانند. می‌توان داده‌ها را در یک جدول 2×2 و به شرح زیر خلاصه نمود. همچنین p_1 و p_2 را به شرح زیر تعریف می‌کنیم:

	مردم	زنده	
جمعیت ۱	a	b	n_1
جمعیت ۲	c	d	n_2
	m_1	m_2	n

$$p_1 = P(\text{مردن} | \text{جمعیت ۱})$$

$$p_2 = P(\text{مردن} | \text{جمعیت ۲})$$

برای آزمون $H_0: p_1 = p_2$ ، از آماره زیر استفاده می‌شود، که در آن $\hat{p}_1 = \frac{a}{n_1}$ ، $\hat{p}_2 = \frac{c}{n_2}$ و $\hat{p} = \frac{m_1}{n}$ است:

$$\chi^2 = \left[\frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})(1/n_1 + 1/n_2)}} \right]^2 = \frac{n(ad-bc)^2}{n_1 n_2 m_1 m_2}$$

یا با توجه به تصحیح پیوستگی داریم:

$$\chi_c^2 = \frac{n(|ad-bc| - n/2)^2}{n_1 n_2 m_1 m_2}$$

جدول ۲. مقادیر R برای $\lambda = 0.5$ و توزیعهای بریده شده به صورتهای مختلف لپس هستند.

	P ₁									
	0.100	0.100	0.200	0.300	0.400	0.500	0.600	0.700	0.800	0.900
0.100	1.000	1.000	1.000	1.000	1.001	1.001	1.002	1.004	1.009	1.020
0.200	1.000	1.000	1.000	1.000	1.000	1.001	1.002	1.004	1.007	1.018
0.300	1.000	1.000	1.000	1.000	1.000	1.001	1.001	1.003	1.006	1.016
0.400	1.000	1.000	1.000	1.000	1.000	1.000	1.001	1.002	1.005	1.013
0.500	1.001	1.000	1.000	1.000	1.000	1.000	1.000	1.001	1.004	1.011
0.600	1.001	1.001	1.001	1.000	1.000	1.000	1.000	1.001	1.002	1.009
0.700	1.004	1.004	1.003	1.002	1.001	1.001	1.000	1.000	1.001	1.006
0.800	1.009	1.007	1.006	1.005	1.004	1.002	1.001	1.000	1.000	1.001
0.900	1.020	1.018	1.016	1.013	1.011	1.009	1.006	1.004	1.001	1.000

جدول ۳. مقادیر R برای $\lambda = 0.4$ و توزیعهای بریده شده به صورتهای مختلف لگس هستند.

P_2	۰/۰۰	۰/۱۰	۰/۲۰	۰/۳۰	۰/۴۰	۰/۵۰	۰/۶۰	۰/۷۰	۰/۸۰	۰/۹۰
۰/۰۰	۱/۰۰	۱/۰۱	۱/۰۲	۱/۰۴	۱/۰۶	۱/۰۹	۱/۱۴	۱/۲۱	۱/۳۳	۱/۶۵
۰/۱۰	۰/۹۹	۱/۰۰	۱/۰۱	۱/۰۳	۱/۰۵	۱/۰۸	۱/۱۲	۱/۱۸	۱/۲۹	۱/۵۹
۰/۲۰	۰/۹۸	۰/۹۹	۱/۰۰	۱/۰۱	۱/۰۳	۱/۰۶	۱/۰۹	۱/۱۵	۱/۲۶	۱/۵۳
۰/۳۰	۰/۹۷	۰/۹۸	۰/۹۹	۱/۰۰	۱/۰۲	۱/۰۴	۱/۰۷	۱/۱۲	۱/۲۲	۱/۴۷
۰/۴۰	۰/۹۶	۰/۹۷	۰/۹۷	۰/۹۹	۱/۰۰	۱/۰۲	۱/۰۵	۱/۰۹	۱/۱۸	۱/۴۰
۰/۵۰	۰/۹۵	۰/۹۵	۰/۹۶	۰/۹۷	۰/۹۸	۱/۰۰	۱/۰۲	۱/۰۶	۱/۱۴	۱/۳۳
۰/۶۰	۰/۹۴	۰/۹۴	۰/۹۵	۰/۹۶	۰/۹۷	۰/۹۸	۱/۰۰	۱/۰۳	۱/۰۹	۱/۲۶
۰/۷۰	۰/۹۳	۰/۹۳	۰/۹۳	۰/۹۴	۰/۹۵	۰/۹۶	۰/۹۷	۱/۰۰	۱/۰۵	۱/۱۸
۰/۸۰	۰/۹۲	۰/۹۲	۰/۹۲	۰/۹۳	۰/۹۳	۰/۹۴	۰/۹۵	۰/۹۷	۱/۰۰	۱/۰۹
۰/۹۰	۰/۹۲	۰/۹۲	۰/۹۲	۰/۹۲	۰/۹۲	۰/۹۲	۰/۹۲	۰/۹۳	۰/۹۵	۱/۰۰

متغیر χ^2 ، تقریباً دارای توزیع χ^2_1 است. این یک تقریب تحت فرض H_0 برای توزیع شرطی گسسته دقیق است. با معلوم بودن n_1, n_2, m_1 و m_2 ، متغیر تصادفی A ، عدد مربوط به خانه $(1, 1)$ جدول 2×2 ، دارای توزیع فوق هندسی به شرح زیر است:

$$P(A = a) = \frac{\binom{n_1}{a} \binom{n_2}{m_1 - a}}{\binom{n}{m_1}}$$

دو گشتاور اولیه توزیع فوق هندسی، به شرح زیر است:

$$E_0(A) = \frac{n_1 m_1}{n} \quad \text{و} \quad \text{Var}_0(A) = \frac{n_1 n_2 m_1 m_2}{n^2 (n-1)}$$

در نتیجه داریم:

$$ad - bc = n(a - E_0(A))$$

$$n_1 n_2 m_1 m_2 = n^2 (n-1) \text{Var}_0(A)$$

$$\chi^2 = \frac{(ad - bc)^2}{n_1 n_2 m_1 m_2} = \frac{n}{n-1} \left[\frac{a - E_0(A)}{\sqrt{\text{Var}_0(A)}} \right]^2$$

۲.۲ دنباله‌ای از جدولهای 2×2

فرض کنید، دنباله‌ای از جدولهای 2×2 در اختیار داریم. برای مثال، در k بیمارستان، بیماران تیمار ۱ یا تیمار ۲ را دریافت می‌کنند و رفتار آنها ثبت می‌شود. امکان دارد بیمارستانها دارای تفاوت باشند، لذا، نمی‌توانیم k جدول را در یک جدول 2×2 خلاصه کنیم. در نتیجه فرض زیر را آزمون می‌کنیم:

$$H_0: p_{11} = p_{12}, \dots, p_{k1} = p_{k2}$$

که p_{i1} و p_{i2} به شرح زیراند:

$$p_{i1} = P(\text{تیمار ۱} | \text{بیمارستان } i | \text{ مردن})$$

$$p_{i2} = P(\text{تیمار ۲} | \text{بیمارستان } i | \text{ مردن})$$

جداول به شرح زیراند:

	مرده	زنده	
تیمار ۱	a_1		n_{11}
تیمار ۲			n_{12}
	m_{11}	m_{12}	n_1

بیمارستان یکم

	مرده	زنده	
تیمار ۱	a_k		n_{k1}
تیمار ۲			n_{k2}
	m_{k1}	m_{k2}	n_k

بیمارستان kام

با استفاده از آماره مانتل-هانزل، داریم:

$$MH = \frac{\sum_{i=1}^k (a_i - E_o(A_i))}{\sqrt{\sum_{i=1}^k (\text{Var}_o(A_i))}}$$

در صورت استفاده از تصحیح پیوستگی، داریم:

$$MH_c = \frac{\left| \sum_{i=1}^k (a_i - E_o(A_i)) \right| - \frac{1}{2}}{\sqrt{\sum_{i=1}^k (\text{Var}_o(A_i))}}$$

اگر جدولها مستقل باشند یا این که k مقدار ثابتی بوده و $n_j \rightarrow \infty$ ، یا هنگامی که $k \rightarrow \infty$ میل کند، آن گاه: $MH \underset{d}{\sim} N(0, 1)$ است. همچنین، جدولها هم توزیع اند.

در تحلیل بقاء آماره MH به صورت زیر به کار برده می شود. فرض کنید: $(Z_{(1)}, \xi_{(1)})$ تا $(Z_{(m+n)}, \xi_{(m+n)})$ ، نمونه مرتب شده ترکیبی باشد. برای هر زمان بریده شده یک جدول 2×2 بسازید. آماره MH را برای این دنباله از جدولها برای آزمون $H_0: F_1 = F_2$ ، بسازید.

این جدولها مستقل نیستند. زیرا برای مثال، $\mathcal{R}(Z_{(1)})$ و $\mathcal{R}(Z_{(3)})$ تقریباً برابرند. ولی هنوز نرمال بودن مجانبی برقرار است. تغییر طرحهای برشی بر آماره MH بی تأثیر است. این آزمون را آزمون لگاریتم رتبه نیز گویند (به فصل ششم، بخش (۱.۲)، مراجعه کنید).

۳.۲ مثال

محاسبه آماره MH (مربوط به شکل‌های ۴. الف و ۴. ب)، در جدول شماره ۴ داده شده است. ستون z شامل مشاهدات مرتب بریده نشده است. چهار ستون بعدی n_1, m_1, n و a، جدول‌های (۲×۲) را می‌سازند. ستون بعدی مقدار $E_0(A) = \frac{n_1 m_1}{n}$ را محاسبه می‌کند. حاصل ضرب دو ستون آخر $\frac{m_1(n-m_1)}{(n-1)}$ و $\frac{n_1(1-\frac{n_1}{n})}{n}$ مقدار $Var_0(A)$ را می‌دهند. بهتر است، مقدار $Var_0(A)$ را به این طریق محاسبه کنیم. زیرا، مقدار $\frac{m_1(n-m_1)}{(n-1)}$ معمولاً برابر یک و $\frac{n_1(1-\frac{n_1}{n})}{n}$ برابر حاصل ضرب نسبت‌های دو نمونه است.

$$MH = \frac{\text{مجموع ستون } E_0(A) - \text{مجموع ستون } a}{\sqrt{\left(\frac{n_1(n-m)}{n-1} \text{ ستون} \times \frac{n_1(1-\frac{n_1}{n})}{n} \text{ ستون} \right) \text{مجموع}}} = \frac{2,31}{1,02} = 2,26$$

برای آزمون یک طرفه $P = 0,012$ و $MH_c = \frac{2,31-0,5}{1,02} = 1,77$ و $P = 0,038$ است. دنباله جدول‌های ۲×۲ مانند هانزل به شرح زیر است:

	x	o	y	o	o	x	→ t
	$Z(1)$	$Z(2)$	$Z(3)$	$Z(4)$	$Z(5)$	$Z(6), Z(7)$	
	$\xi(1) = 1$	$\xi(2) = 0$	$\xi(3) = 1$	$\xi(4) = 0$	$\xi(5) = 0$	$\xi(6) = 1, \xi(7) = 1$	

	D	A
X	۱	
Y	o	

 n_1

	D	A
X	o	
Y	۱	

 n_2

	D	A
X	۲	
Y	o	

 n_3

۴.۲ نرمال مجانبی

برای اثبات نرمال مجانبی، فرض می‌کنیم تکرار وجود ندارد و موارد زیر را در نظر

می‌گیریم:

جدول ۴. محاسبات مربوط به آماره مانتل-هانزک در آزمایش فرضی بالینی بران.

z	n	m_1	n_1	a	$E_0(A)$	$a - E_0(A)$	$\frac{m_1(n - m_1)}{n - 1}$	$\frac{n_1}{n} \left(1 - \frac{n_1}{n}\right)$
۳	۱۰	۱	۵	۱	۰٫۵۰	۰٫۵۰	۱	۰٫۲۵۰۰
۵	۹	۱	۴	۱	۰٫۴۴	۰٫۵۶	۱	۰٫۲۴۶۹
۷	۸	۱	۳	۱	۰٫۳۸	۰٫۶۲	۱	۰٫۲۳۴۴
۱۲	۶	۱	۱	۰	۰٫۱۷	-۰٫۱۷	۱	۰٫۱۳۸۹
۱۸	۵	۱	۱	۱	۰٫۲۰	۰٫۸۰	۱	۰٫۱۶۰۰
۱۹	۴	۱	۰	۰	۰	۰	۱	۰
۲۰	۳	۱	۰	۰	۰	۰	۱	۰
مجموع				۴	۱٫۶۹	۲٫۳۱		

$$N = n + m$$

$$\hat{H}(t) = \frac{1}{N} \sum_{i=1}^N I(Z_i \leq t)$$

$$\hat{H}_1(t) = \frac{1}{m} \sum_{i=1}^m I(X_i \leq t)$$

$$\hat{H}_u(t) = \frac{1}{N} \sum_{i=1}^N I(Z_i \leq t, \xi_i = 1)$$

$$\hat{H}_{1u}(t) = \frac{1}{m} \sum_{i=1}^m I(X_i \leq t, \delta_i = 1)$$

حال می‌توان صورت MH را به شرح زیر نوشت:

$$\sum_{i=1}^k (a_i - E_o(A_i)) = m \left\{ \int_0^{\infty} d\hat{H}_{1u}(s) - \int_0^{\infty} \frac{1 - \hat{H}_1(s^-)}{1 - \hat{H}(s^-)} d\hat{H}_u(s) \right\}$$

زیرا، $E(A_i) = \frac{m_{i1} n_{i1}}{n_i}$ که در آن a_i ، m_{i1} ، n_{i1} و n_i از جدول 2×2 ، متناظر با مشاهده i ام بریده نشده به دست می‌آید:

	D	A	
X	a_i		n_{i1}
Y			n_i
	m_{i1}		

چون فرض کرده‌ایم تکرار وجود ندارد، پس $m_{i1} = 1$. اگر s_i زمان i امین مشاهده بریده نشده باشد، داریم:

$$n_{i1} = \#(X \text{ های باقی مانده در زمان } s_i) = m(1 - \hat{H}_1(s_i^-))$$

$$n_i = \#(Z \text{ های باقی مانده در زمان } s_i) = N(1 - \hat{H}(s_i^-))$$

حال می‌توان صورت MH را برحسب توابع توزیع تجربی نوشت و می‌توان از روش

مشابه نرمال مجانبی برآوردگر PL، استفاده کرد.

REFERENCES

Mantel and Haenszel, J. Natl. Cancer Inst. (1962).

Crowley, JASA (1977).

Lininger et al., Biometrika (1979).

۳ رده آزمونهای تارون-وایر

بعد از ساختن جدول 2×2 برای هر مشاهده بریده شده، تارون و وایر، وزنی را برای هر جدول پیشنهاد می کنند، به گونه ای که:

$$\sum_{i=1}^k w_i [a_i - E_o(A_i)] = \sum_{i=1}^k w_i \left[a_i - \frac{m_{i1} n_{i1}}{n_i} \right] \quad (12)$$

و برای واریانس

$$\sum_{i=1}^k w_i^2 \text{Var}_o(A_i) = \sum_{i=1}^k w_i^2 \left[\frac{m_{i1}(n_i - m_{i1})}{n_i - 1} \right] \times \left[\left(\frac{n_{i1}}{n_i} \right) \left(1 - \frac{n_{i1}}{n_i} \right) \right] \quad (13)$$

سه حالت ویژه مهم وجود دارد:

(الف) $w_i = 1$ ، در این صورت آماره MH به دست می آید.

(ب) $w_i = n_i$ ، در این صورت آماره گهان به دست می آید.

(پ) $w_i = \sqrt{n_i}$ ، در این صورت آماره تارون-وایر به دست می آید.

یادآوریها:

(الف) از کدام آزمون استفاده کنیم؟ آماره گهان به مشاهدات مقدماتی وزن بیشتری می دهد، در حالی که آماره MH به تمام مشاهدات وزن یک می دهد. پیشنهاد تارون-وایر بین دو روش قبلی قرار دارد. بنابه اظهار آنان وزنه های $w_i = \sqrt{n_i}$ ، تأثیر فراوانی در دامنه تغییرات دارد.

(ب) هر چند (۱۲) با آماره گهان U برابر است، ولی $\hat{\text{Var}}_{\text{TW}}(U)$ که از

رابطه (۱۳) به دست می آید، برابر $\text{Var}_{o,p}^*(U)$ نیست. به طور مجانبی $\hat{\text{Var}}_{\text{TW}}(U)$ معادل واریانس U، تحت فرض H_o است، در حالی که $\text{Var}_{o,p}^*(U)$ ، واریانس تحت H_o^* است.

مثال. با مراجعه به جدول ۴، که در آن آمارهٔ MH را محاسبه کرده‌ایم، داریم:

$$\sum_{i=1}^k n_i (a_i - E_o(A_i)) = (10)(0,50) + (9)(0,56) + (8)(0,62) \\ + (6)(-0,17) + (5)(0,80) = 17,98$$

که برابر آمارهٔ U در آمارهٔ گهان است، البته بدون علامت و مقدار گرد شده، همچنین داریم:

$$\hat{\text{Var}}_{\text{TW}}(U) = \sum n_i^2 \left[\frac{m_{i1}(n_i - m_{i1})}{n_i - 1} \right] \left[\left(\frac{n_{i1}}{n_i} \right) \left(1 - \frac{n_{i1}}{n_i} \right) \right] \\ = (10^2)(0,25) + (9^2)(0,2469) + (8^2)(0,2344) \\ + (6^2)(-0,1389) + (5^2)(0,16) = 69$$

$$\text{Var}_{o,p}^*(U) = 79,44$$

در نتیجه داریم:

$$\sqrt{\hat{\text{Var}}_{\text{TW}}(U)} = 8,31 \quad \text{و} \quad \sqrt{\text{Var}_{o,p}^*(U)} = 8,91$$

REFERENCE

Tarone and Ware, Biometrika (1977).

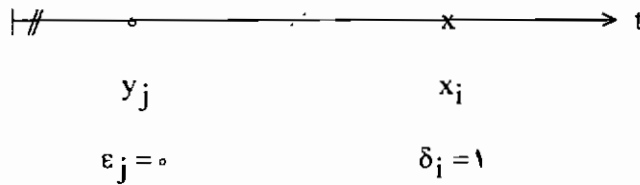
۴ آزمون افرون

به خاطر بیاورید که در ساختن آمارهٔ گهان، تابع امتیاز را به صورت زیر تعریف

کردیم:

$$U_{ij} = \begin{cases} +1 & t_i > u_j \\ 0 & \text{در سایر جاها} \\ -1 & t_i < u_j \end{cases}$$

فرض کنید، حالت زیر را داریم:



آزمون گهان، امتیاز $U_{ij} = 0$ را به این زوج - با صرفنظر از بزرگی X_i نسبت به Y_i - می‌دهد. افرون پیشنهاد می‌کند که امتیاز زیر را به آنها نسبت دهیم:

$$U_{ij} = \hat{P}\{T_i > U_j | (x_i, \delta_i), (y_j, \varepsilon_j)\}$$

برای حالت مورد نظر داریم:

$$U_{ij} = \hat{P}\{U_j < x_i | U_j > y_j\} = \frac{\hat{F}_T(x_i) - \hat{F}_T(y_j)}{1 - \hat{F}_T(y_j)}$$

در این رابطه \hat{F}_T برآوردگر PL کاپلان-مایر برای جمعیت T است. استفاده از این امتیازها و کاربرد 0 و 1 به جای 1 و -1 ، آماره زیر را نتیجه می‌دهد:

$$\int_0^{\infty} [1 - \hat{F}_T(u)] d\hat{F}_T(u) = \hat{P}\{T_i > U_j\} \quad (14)$$

برآوردگر $\hat{P}(T_i > U_j)$ همان GMLE احتمال $P(T_i > U_j)$ است، که پارامتر آماره ویلکاکسن در حالت بریده نشده است، یعنی:

$$\frac{1}{mn} U \xrightarrow{\text{a.s.}} P(X > Y)$$

در حالت بریده نشده: $U_{ij} = 0$ یا 1 .

برآوردگر (۱۴)، در دنباله پایدار نیست و مانع کاربرد وسیع آن می‌شود.

REFERENCE

Efron, Proc. Fifth Berkeley Symp. IV (1967).

فصل پنجم

روشهای ناپارامتری K نمونه

برای نمونه i ام ($i=1,2,\dots,K$)، فرض کنید T_{i1} تا T_{in_i} ، متغیرهای iid بوده، که توزیع هر یک F_i باشد. همچنین C_{i1} تا C_{in_i} ، متغیرهای iid با تابع توزیع G_i باشند. C_{ij} زمان برش مربوط به T_{ij} است. مشاهدات عبارت‌اند از: (X_{i1}, δ_{i1}) تا $(X_{in_i}, \delta_{in_i})$ ، که در آن:

$$X_{ij} = T_{ij} \wedge C_{ij} \quad \text{و} \quad \delta_{ij} = I(T_{ij} \leq C_{ij})$$

فرض زیر مورد نظر ماست:

$$H_0: F_1 = \dots = F_K$$

۱ آزمون گهان تعمیم یافته (برسلو)

با استفاده از تابع امتیاز مسأله دو نمونه‌ای، فرضهای زیر را در نظر می‌گیریم:

$$W_i = \sum_{j=1}^{n_i} \sum_{\substack{i'=1 \\ \neq i}}^K \sum_{j'=1}^{n_{i'}} U((X_{ij}, \delta_{ij}), (X_{i'j'}, \delta_{i'j'}))$$

$$\underline{W} = (W_1, \dots, W_K)'$$

برسلو ماتریس کوواریانس مجانبی \underline{W} ، تحت فرض محدود کننده زیر را به دست آورد:

$$H_0^*: F_1 = \dots = F_K \quad ; \quad G_1 = \dots = G_K$$

با تعریف $N = \sum_{i=1}^K n_i$ ، فرض می‌کنیم با شرط $N \rightarrow \infty$ و $\frac{n_i}{N} \rightarrow \lambda_i$ ، $i=1,2,\dots,K$

$$\underline{W} \underset{a}{\sim} N(\underline{\mu}_0^*, N^{\tau} \underline{\Sigma}_0^*)$$

که در این رابطه $\underline{\mu}_0^* = 0$ و

$$\underline{\Sigma}_0^* = \left(\int_0^{\infty} [1 - H(u)]^{\tau} dH_{\mathbf{u}}(u) \right) \times \begin{pmatrix} \lambda_1(1-\lambda_1) & & & \\ & \ddots & & \\ & & -\lambda_i \lambda_j & \\ & & -\lambda_i \lambda_j & \ddots \\ & & & & \lambda_K(1-\lambda_K) \end{pmatrix}$$

و

$$H_i(t) = P(X_{i1} \leq t)$$

$$H_{iu}(t) = P(X_{i1} \leq t, \delta_{i1} = 1)$$

$$H(t) = \lambda_1 H_1(t) + \dots + \lambda_K H_K(t)$$

$$H_{\mathbf{u}}(t) = \lambda_1 H_{1u}(t) + \dots + \lambda_K H_{Ku}(t)$$

چون ماتریس کوواریانس مجانبی به پارامترهای مجهول بستگی دارد، جانشین می‌کنیم.

$$\hat{\lambda}_i = \frac{n_i}{N}$$

$$\hat{H}(t) = \frac{1}{N} \sum_{i=1}^K \sum_{j=1}^{n_i} I(X_{ij} \leq t)$$

$$\hat{H}_{\mathbf{u}}(t) = \frac{1}{N} \sum_{i=1}^K \sum_{j=1}^{n_i} I(X_{ij} \leq t, \delta_{ij} = 1)$$

REFERENCE

Breslow, Biometrika (1970).

انواع آزمونها

۱ آزمون χ^2 امنیبوس: از آماره زیر استفاده می‌شود:

$$\frac{1}{N^{\tau} \int_0^{\infty} (1 - \hat{H})^{\tau} d\hat{H}_{\mathbf{u}}} \sum_{i=1}^K \frac{W_i^{\tau}}{\hat{\lambda}_i} \underset{a}{\sim} \chi_{K-1}^{\tau}$$

این آماره معادل $W'(\Sigma_0^*)^{-1}W$ است، که $(\Sigma_0^*)^{-1}$ معکوس تعمیم یافته Σ_0^* است. اگر برش نباشد، این آماره به طور مجانبی معادل آماره کروسکال-والیس است. یادآوری می‌شود که اگر R_{ij} مرتبه X_{ij} در میان N مشاهده باشد، داریم:

$$R_i = \sum_{j=1}^{n_j} R_{ij} \quad \text{و} \quad \bar{R}_i = \frac{1}{n_j} R_i \quad \text{و} \quad \bar{R}_\cdot = \frac{1}{N} \sum_{i=1}^K R_i$$

حال، آماره کروسکال-والیس، به شرح زیر است:

$$\frac{12}{N(N+1)} \sum_{i=1}^K n_i (\bar{R}_i - \bar{R}_\cdot)^2 = \left(\frac{12}{N(N+1)} \sum_{i=1}^K \frac{n_i^2}{n_i} \right) - 3N(N+1)$$

این آماره دارای توزیع χ_{K-1}^2 ، تحت فرض H_0 است.

۲ آزمون روند: فرض کنید، اگر جمعیت‌ها همگی برابر نباشد، آن‌گاه قابل مرتب شدن هستند. برای مثال می‌توان جمعیت‌ها را برحسب اندازه مقدار دارو به صورت: $d_1 < \dots < d_K$ ، مرتب کرد. d_i مقدار دارویی است که جمعیت i ام دریافت می‌کند. در حالت‌های دیگر، ممکن است از قبل بدانیم هرگاه جمعیت‌ها با هم متفاوت باشند، باید به طور یکنواخت تغییر کنند و این تغییر همبستگی عددی نیست.

$$\underline{\ell} = (d_1, \dots, d_K)'$$

در هنگامی که متغیرهای کمی در دسترس نباشد، از تعریف زیر استفاده می‌کنیم:

$$\underline{\ell} = (-(K-1), \dots, -3, -1, +1, +3, \dots, +(K-1))' \quad \text{زوج } K$$

$$\underline{\ell} = \left(-\frac{(K-1)}{2}, \dots, -1, 0, +1, \dots, +\frac{(K-1)}{2} \right)' \quad \text{فرد } K$$

آپسلون و توکی پیشنهاد می‌کنند که مقابله‌های خطی مرتبه ۲ یا مرتبه ۲ تا ۴، که برای مقادیر زوج K به شرح زیر تشریح می‌شود، مورد استفاده قرار گیرد.

$$\underline{\ell} = (-2(K-1), -(K-3), -(K-5), \dots, +(K-5), +(K-3), +2(K-1))'$$

$$\underline{\ell} = (-4(K-1), -2(K-3), -(K-5), \dots, +(K-5), +2(K-3), +4(K-1))'$$

برای نرمال کردن W ، می‌نویسیم: $\bar{W}_i = \frac{W_i}{n_i(N-n_i)}$ و $\underline{\bar{W}} = (\bar{W}_1, \dots, \bar{W}_K)'$

حال c را چنان در نظر بگیرید که، $c'W = \ell' \bar{W}$ باشد، در نتیجه داریم:

$$\frac{c'W}{\sqrt{N^* c' \Sigma_0^* c}} \underset{a}{\sim} N(0, 1)$$

آماره بالا را می‌توان برای آزمون فرض $H_0: F_1 = \dots = F_K$ در مقابل $H_1: F_1 < \dots < F_K$ به کار برد.

اگر کمیت‌های اندازه‌پذیر در دست باشند، روش‌های رگرسیونی قابل استفاده‌اند. آنها را در فصل ششم بررسی می‌کنیم.

۱.۱ ماتریس کوواریانس جایگشت

متغیر W^* را به صورت زیر تعریف می‌کنیم:

$$W_{ij}^* = \sum_{i'=1}^K \sum_{j'=1}^{n_{j'}} U((X_{ij}, \delta_{ij}), (X_{i'j'}, \delta_{i'j'})) \\ (i', j') \neq (i, j)$$

متغیرهای $W_{K1}^*, \dots, W_{Kn_K}^*$ و $W_{1n_1}^*, \dots, W_{1n_1}^*$ را به صورت W_1^*, \dots, W_N^* در نظر می‌گیریم. برای محاسبه توزیع جایگشت W ، فرض کنید W_1^* تا W_N^* معلوم‌اند. تحت فرض H_0 و بدون جایگذاری W_i^* ها را نمونه‌گیری می‌کنیم. فرض کنید W_1 مجموع اولین n_1 مشاهده نمونه، W_2 مجموع دومین n_2 نمونه و تا آخر باشد.

ماتریس کوواریانس $W = (W_1, \dots, W_K)'$ ، تحت این طرح نمونه‌گیری به شکل

زیر است:

$$\Sigma_{0,p}^* = \frac{1}{N} \left(\frac{\sum_{i=1}^K \sum_{j=1}^{n_j} (W_{ij}^*)^2}{N-1} \right) \times \begin{pmatrix} n_1(N-n_1) & & & \\ & \ddots & & \\ & & -n_i n_j & \\ & & & \ddots \\ -n_i n_j & & & & n_K(N-n_K) \end{pmatrix}$$

ماتریس $\Sigma_{0,p}^*$ را می‌توان به جای $N^* \hat{\Sigma}_0^*$ به کار برد. زیرا، به طور مجانبی معادل‌اند.

REFERENCE

Marcuson and Nordbrock, *Biom. Zeit. / Biom. J.* (1981).

۲.۱ توزیع تحت H_0

تحت فرض: $F_1 = \dots = F_K$ ، ثابت می‌شود که

$$\underline{W} \stackrel{d}{\sim} N(\underline{\mu}_0, N^2 \underline{\Sigma}_0)$$

که درایه‌های $\underline{\Sigma}_0$ ، به شرح زیراند:

$$\sigma_{ij}^0 = -\lambda_i \lambda_j \int_0^{\infty} (1-H_i)(1-H_j) dH_u \quad i \neq j$$

$$\sigma_{ii}^0 = \lambda_i \int_0^{\infty} [(1-H)(1-H_i) - \lambda_i(1-H)^2] dH_u$$

برای برآورد $\underline{\Sigma}_0$ ، به جای جایگذاری مقدار در H_i ، H_u و H ، ساده‌تر است که از آماره ML ، استفاده کنیم.

REFERENCE

Breslow, *Biometrika* (1970).

۲ آزمون مانتل-هانزل تعمیم یافته (تارون و وایر)

فرض کنید، نمونه ترکیبی مرتب شده به صورت: $(Z_{(1)}, \xi_{(1)}), \dots, (Z_{(N)}, \xi_{(N)})$

و همچنین به فرض $\mathfrak{R}(i) = \mathfrak{R}(Z_{(i)})$ باشد.

برای هر نقطه زمانی بریده نشده، جدول 2×2 را به شرح زیر می‌سازیم. این جدول به ازای $k=2$ ، ترانهاده جدولهای 2×2 فصل چهارم است.

	۱	۲	...	K	
مرده	$a_{i1} = 0$	$a_{i2} = 1$...	$a_{iK} = 0$	m_{i2}
زنده			...		m_{i1}
	n_{i1}	n_{i2}	...	n_{iK}	$N_i = \# \mathfrak{R}_u(i)$

تحت فرض $H_0: F_1 = \dots = F_K$ ، داریم:

$$E_o(\underline{A}_i) = (E_o(A_{i1}), \dots, E_o(A_{iK}))' = \left(\frac{m_{i1} n_{i1}}{N_i}, \dots, \frac{m_{i1} n_{iK}}{N_i} \right)'$$

$$\underline{\Sigma}_o(\underline{A}_i) = \left(\frac{m_{i1} n_{i1}}{N_i - 1} \right) \times \begin{pmatrix} \frac{n_{i1}}{N_i} \left(1 - \frac{n_{i1}}{N_i} \right) & & & \\ & & -\frac{n_{iK}}{N_i} \frac{n_{i1}}{N_i} & \\ & \dots & & \\ & & -\frac{n_{iK}}{N_i} \frac{n_{i1}}{N_i} & \\ & & & \frac{n_{iK}}{N_i} \left(1 - \frac{n_{iK}}{N_i} \right) \end{pmatrix}$$

بنا به تعریف داریم:

$$\underline{a} - E_o(\underline{A}) = \sum_i w_i (\underline{a}_i - E_o(\underline{A}_i))$$

$$\underline{\Sigma}_o = \sum_i w_i^2 \underline{\Sigma}_o(\underline{A}_i)$$

در این روابط، w_i وزنه‌های نسبت داده شده‌اند. سه حالت خاص به شرح زیر وجود دارد:

(الف) $w_i = 1$ ، در این صورت آزمون تعمیم یافته MH به دست می‌آید.

(ب) $w_i = N_i$ ، در این صورت آزمون تعمیم یافته گهان به دست می‌آید.

(پ) $w_i = \sqrt{N_i}$ ، کارایی بالایی روی دامنه تغییرات می‌دهد.

در بخش اول، یک ماتریس کوواریانس مجانبی از آماره تعمیم یافته گهان را تحت فرض محدود کننده H_0^* به دست آوردیم. این ماتریس تحت H_0 برابر $\underline{\Sigma}_o$ بوده و از نظر محاسباتی از قبلی ساده‌تر است.

REFERENCE

Tarone and Ware, *Biometrika* (1977).

انواع آزمونها

۱ آزمون χ^2 . چون $\underline{\Sigma}_o$ منفرد است، یکی از جمعیتها، مثلاً اولی را حذف می‌کنیم. بنابه تعریف $(\underline{A}_{-1}) - E_o(\underline{A}_{-1})$ و $\underline{\Sigma}_{o,-1}$ به ترتیب برابر $(\underline{A}_{-1}) - E_o(\underline{A}_{-1})$ و $\underline{\Sigma}_{o,-1}$ باشند، که در

آنها جمعیت اول حذف شده است. حال تحت فرض H_0 ، داریم:

$$W = (\underline{a}_{-1} - E_0(\underline{A}_{-1}))' \Sigma_{0,-1}^{-1} (\underline{a}_{-1} - E_0(\underline{A}_{-1})) \stackrel{a}{\sim} \chi_{K-1}^2$$

اگر هر کدام از جمعیتها را حذف کنیم، W تغییر نمی کند.

برای آزمون تعمیم یافته ماننل-هانزل ($w_i = 1$) آزمون تقریبی زیر را داریم:

$$\sum \frac{(O-E)^2}{E} = \sum_{k=1}^K \frac{(a_k - E_0(A_k))^2}{E_0(A_k)} \approx \chi_{K-1}^2$$

$$E_0(A_k) = \sum_i E_0(A_{ik}) = \sum_i \frac{m_{i1} n_{ik}}{N_i} \text{ و } a_k = \sum_i a_{ik}$$

هرچند این آزمون به علت نامساوی $\sum \frac{(O-E)^2}{E} \leq W$ تا حدودی محتاطانه است.

ولی، در کاربرد ساده تر است. زیرا در آن به محاسبه معکوس ماتریس نیازی نیست.

REFERENCES

Peto and Pike, Biometrics (1973).

Peto et al., British J. Cancer (1976, 1977).

۲ آزمون روند. فرض H_1 را به صورت $F_1 < \dots < F_K$ ، در نظر می گیریم. برای انتخاب ℓ مطابق بخش اول، از آماره $\ell'(\underline{a} - E_0(\underline{A}))$ ، که به طور مجانبی نرمال است، استفاده می کنیم.

REFERENCE

Tarone, Biometrika (1975).

فصل ششم

روشهای ناپارامتری: رگرسیون

۱ الگوهای نرخ شکست متناسب کاکس

فرض کنید T_1 تا T_n و C_1 تا C_n ، متغیرهای تصادفی مستقل باشند. متغیر C_i زمان برش مربوط به زمان بقاء T_i است. مشاهدات به صورت زیراند، که در آن $Y_i = T_i \wedge C_i$ و $\delta_i = I(T_i \leq C_i)$ است.

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

همچنین، مقادیر معلوم به صورت: x_1, \dots, x_n بوده که در آن $\underline{x}_i = (x_{i1}, \dots, x_{ip})'$ بردار متغیرهای وابسته به متغیر تابع T_i است. به خاطر بیاورید که تابع نرخ شکست را به صورت زیر تعریف کردیم، که در آن وابستگی T از طریق \underline{x} منظور شده است.

$$\lambda(t; \underline{x}) = \frac{f(t; \underline{x})}{1 - F(t; \underline{x})}$$

در این الگو فرض می شود که: $\lambda(t; \underline{x}) = e^{\beta' \underline{x}} \lambda_0(t)$ باشد، که بردار: $\underline{\beta} = (\beta_1, \dots, \beta_p)'$ ضرایب رگرسیون است. نرخ شکست برابر حاصل ضرب یک عدد در تابع $\lambda_0(t)$ است، که این ضریب عددی به ضرایب رگرسیون و کوواریانسها وابسته است. قضیه می تواند در این حالت نیز برقرار باشد که، به جای $e^{\beta' \underline{x}}$ ، هر مقدار محسوس $h(\beta' \underline{x})$ را، که در آن h مثبت باشد، جانشین کنیم. ضرایب رگرسیون $\underline{\beta}$ و تابع نرخ شکست $\lambda_0(t)$ ، هیچ کدام معلوم نیستند.

خانواده‌ای از توزیعها را "خانواده توزیمهای لهن" گوئیم، هرگاه تابع توزیمی مانند F_0 وجود داشته باشد، به گونه‌ای که در ازای هر F در ایسن خانواده

رابطه $1 - F = (1 - F_0)^\gamma$ برقرار باشد. در این رابطه γ یک عدد مثبت است. می توان رابطه را بر حسب تابع بقاء و به صورت $S = S_0^\gamma$ نیز نوشت. از الگوی نرخهای متناسب نتیجه می شود که توابع توزیع آنها تشکیل یک خانواده توزیعیهای لهن می دهند. اثبات آن به شرح زیر است:

$$\begin{aligned} S(t; \underline{x}) &= \exp\left\{-\int_0^t \lambda(u; \underline{x}) du\right\} = \exp\left\{-e^{-\beta' \underline{x}} \int_0^t \lambda_0(u) du\right\} \\ &= \exp\left\{-\int_0^t \lambda_0(u) du\right\} e^{-\beta' \underline{x}} = S_0(t) e^{-\beta' \underline{x}} \end{aligned}$$

که در رابطه بالا $S_0(t) = \exp\left\{-\int_0^t \lambda_0(u) du\right\}$ است. حالت خاص $p=1$ را در نظر می گیریم:

$$x_i = \begin{cases} 1 & \text{اگر مشاهده } i\text{ام از جمعیت } 1 \text{ باشد} \\ 0 & \text{اگر مشاهده } i\text{ام از جمعیت } 2 \text{ باشد} \end{cases}$$

$$e^{\beta x_i} = \begin{cases} e^\beta = \gamma & \text{اگر مشاهده } i\text{ام از جمعیت } 1 \text{ باشد} \\ 1 & \text{اگر مشاهده } i\text{ام از جمعیت } 2 \text{ باشد} \end{cases}$$

در نتیجه توابع بقاء برای جمعیت ۱ و ۲، به صورت زیر ارتباط دارند:

$$S_1(t) = S_2^\gamma(t)$$

۱.۱ تحلیل درست نمایی شرطی

در نوشته های کاکس داریم: فرض کنید $\lambda_0(t)$ اختیاری باشد. هیچ اطلاعاتی درباره β از فاصله های زمانی، که در آنها هیچ شکستی رخ نداده، به دست نمی آید. زیرا، امکان دارد، $\lambda_0(t)$ به طور قابل درکی با صفر متحد باشد. بنابراین، روی آن لحظه هایی که در آنها شکست رخ می دهد، شرطی می کنیم. در زمان گسسته نیز روی مشاهدات تکراری شرطی می کنیم. به هر حال، نیاز به روشی داریم که تمام $\lambda_0(t)$ را تحلیل ننماید.

در نظر گرفتن این توزیع شرطی اجباری به نظر می‌رسد. فرض کنید تکرار وجود نداشته باشد. حالتی که تکرار وجود دارد را بعداً بررسی می‌کنیم. زمانهای مشاهده‌ای را به صورت $y(1) < y(2) < \dots < y(n)$ مرتب می‌کنیم. فرض کنید، $\delta(i)$ تابع نشانگر و $x(i)$ ، متغیر وابسته به $y(i)$ باشد، می‌نویسیم:

$$\mathfrak{R}(i) = \mathfrak{R}(y(i)-)$$

برای هر زمان بریده نشده $y(i)$ ، داریم:

$$P\{[y(i), y(i) + \Delta y] \text{ یک مرگ در بازه } | \mathfrak{R}(i)\} \cong \sum_{j \in \mathfrak{R}(i)} e^{-\beta' x} \lambda_o(y(i)) \Delta y$$

$$P\{y(i) \text{ در زمان } | \mathfrak{R}(i) \text{ در یک مرگ در زمان } y(i)\} = \frac{e^{-\beta' x(i)}}{\sum_{j \in \mathfrak{R}(i)} e^{-\beta' x_j}}$$

اگر حاصل ضرب سه احتمال شرطی را حساب کنیم، درست‌نمایی شرطی به دست می‌آید.

$$L_c(\underline{\beta}) = \prod_u \frac{e^{-\beta' x(i)}}{\sum_{j \in \mathfrak{R}(i)} e^{-\beta' x_j}}$$

کاکس پیشنهاد می‌کند، که از درست‌نمایی شرطی به جای درست‌نمایی معمولی استفاده کنیم. به ویژه پیدا کردن برآورد حداکثر درست‌نمایی از بردار امتیاز و ماتریس اطلاعات نمونه، به شرح زیر استفاده شود:

$$\frac{\partial}{\partial \underline{\beta}} \log L_c(\underline{\beta}) = \left(\frac{\partial}{\partial \beta_1} \log L_c(\underline{\beta}), \dots, \frac{\partial}{\partial \beta_p} \log L_c(\underline{\beta}) \right)'$$

$$i(\underline{\beta}) = - \frac{\partial^2}{\partial \underline{\beta}^2} \log L_c(\underline{\beta}) = - \begin{pmatrix} \frac{\partial^2}{\partial \beta_1 \partial \beta_1} \log L_c(\underline{\beta}) & \dots & \frac{\partial^2}{\partial \beta_1 \partial \beta_p} \log L_c(\underline{\beta}) \\ \vdots & & \vdots \\ \frac{\partial^2}{\partial \beta_p \partial \beta_1} \log L_c(\underline{\beta}) & \dots & \frac{\partial^2}{\partial \beta_p \partial \beta_p} \log L_c(\underline{\beta}) \end{pmatrix}$$

می‌خواهیم معادلات زیر را، که معمولاً به روش تکرار منجر می‌شود، حل کنیم:

$$\frac{\partial}{\partial \underline{\beta}} \log L_c(\underline{\beta}) = 0$$

فرض کنید، $\hat{\underline{\beta}}^\circ$ یک حدس اولیه باشد، داریم:

$$\hat{\underline{\beta}}^1 = \hat{\underline{\beta}}^\circ + i^{-1}(\hat{\underline{\beta}}^\circ) \frac{\partial}{\partial \underline{\beta}} \log L_c(\hat{\underline{\beta}}^\circ)$$

اگر $\hat{\underline{\beta}}$ جواب معادله باشد، داریم:

$$\hat{\underline{\beta}} \underset{a}{\sim} N(\underline{\beta}, i^{-1}(\underline{\beta}))$$

با مشتق‌گیری از عبارت زیر، رابطه بردار امتیاز و ماتریس اطلاعات نمونه به دست می‌آید، داریم:

$$\log L_c(\underline{\beta}) = \sum_u \left[\underline{\beta}' \underline{x}_{(i)} - \log \left(\sum_{j \in \mathcal{R}(i)} e^{\underline{\beta}' \underline{x}_j} \right) \right]$$

$$\frac{\partial}{\partial \beta_k} \log L_c(\underline{\beta}) = \sum_u \left(x_{(i)k} - \frac{\sum_{j \in \mathcal{R}(i)} x_{jk} e^{\underline{\beta}' \underline{x}_j}}{\sum_{j \in \mathcal{R}(i)} e^{\underline{\beta}' \underline{x}_j}} \right)$$

$$i_{k\ell}(\underline{\beta}) = - \frac{\partial^2}{\partial \beta_k \partial \beta_\ell} \log L_c(\underline{\beta})$$

$$= \sum_u \left(\frac{\sum_{j \in \mathcal{R}(i)} x_{jk} x_{j\ell} e^{\underline{\beta}' \underline{x}_j}}{\sum_{j \in \mathcal{R}(i)} e^{\underline{\beta}' \underline{x}_j}} - \frac{\sum_{j \in \mathcal{R}(i)} x_{jk} e^{\underline{\beta}' \underline{x}_j}}{\sum_{j \in \mathcal{R}(i)} e^{\underline{\beta}' \underline{x}_j}} \times \frac{\sum_{j \in \mathcal{R}(i)} x_{j\ell} e^{\underline{\beta}' \underline{x}_j}}{\sum_{j \in \mathcal{R}(i)} e^{\underline{\beta}' \underline{x}_j}} \right)$$

برای آزمون $H_0: \underline{\beta} = 0$ ، کاکس از آماره نوع رانو به شرح زیر، که به طور مجانبی دارای توزیع χ_p^2 تحت فرض H_0 است، استفاده می‌کند.

$$\left(\frac{\partial}{\partial \underline{\beta}} \log L_c(\circ) \right)' i^{-1}(\circ) \left(\frac{\partial}{\partial \underline{\beta}} \log L_c(\circ) \right)$$

بردار امتیاز و ماتریس اطلاعات نمونه در $\underline{\beta} = \circ$ ، دارای صورت ساده زیر است:

$$\frac{\partial}{\partial \beta_k} \log L_c(\circ) = \sum_u (x_{(i)k} - \bar{x}_{(i)k})$$

$$i_{k\ell}(\circ) = L_c(\underline{\beta}) = \sum_u \left(\frac{1}{n_i} \sum_{j \in \mathcal{R}(i)} x_{jk} x_{j\ell} - \bar{x}_{jk} \bar{x}_{j\ell} \right)$$

$$= \sum_u \left(\frac{1}{n_i} \sum_{j \in \mathcal{R}(i)} (x_{jk} - \bar{x}_{(i)k}) \times (x_{j\ell} - \bar{x}_{(i)\ell}) \right)$$

که در این روابط، $n_i = \#\mathcal{R}(i)$ و $\bar{x}_{(i)} = \frac{\sum_{j \in \mathcal{R}(i)} x_j}{n_i}$ است.

ماتریس اطلاع نمونه عبارت است از: مجموع ماتریسهای کوواریانس برای مجموعه‌های نرخ شکست مشاهدات بریده نشده.

حالت خاص $p=1$ و x_i را به شرح زیر در نظر می‌گیریم:

$$x_i = \begin{cases} 1 & \text{اگر } i \text{ در نمونه } 1 \text{ باشد} \\ 0 & \text{اگر } i \text{ در نمونه } 2 \text{ باشد} \end{cases}$$

در این صورت با توجه به نمادهای آزمون MH، داریم:

$$\frac{\partial}{\partial \beta} \log L_c(\circ) = \sum_u (x_{(i)} - \bar{x}_{(i)}) = \sum_u \left(a_i - \frac{n_{i1}}{n_i} \right)$$

$$i(\circ) = \sum_u \left(\frac{1}{n_i} \sum_{j \in \mathcal{R}(i)} x_j^2 - \bar{x}_{(i)}^2 \right)$$

$$= \sum_u \frac{n_{i1}}{n_i} \left(1 - \frac{n_{i1}}{n_i} \right)$$

در نتیجه اگر تکرار نباشد، آزمون کاکس و MH یکسان هستند.

REFERENCES

Cox, JRSS B (1972).

Prentice and Kalbfleisch, Biometrics (1979), has a nice survey of the Cox procedure.

Kalbfleisch and Prentice, The Statistical Analysis of Failure Time Data (1980), is an excellent new text on the Cox approach.

۲.۱ بررسی درست‌نمایی شرطی

درست‌نمایی حاشیه‌ای برای رتبه‌ها. مجدداً فرض کنید تکرار وجود ندارد. با وجود تکرار استدلال زیر صحیح نیست. فرض نمایید داده‌ها بریده نشده‌اند و Y_1 تا Y_n مستقل و Y_i دارای توزیع F_i و چگالی f_i باشد. دو نماد $\underline{Y} = (Y_1, \dots, Y_n)$ و $\underline{R} = (R_1, \dots, R_n)$ ، که در آن R_i رتبه مشاهده Y_i است را در نظر می‌گیریم. در این صورت احتمال بردار رتبه \underline{r} به صورت زیر است:

$$p(\underline{r}) = \int \cdots \int_{u_1 < \cdots < u_n} \prod_{i=1}^n f_i(u_i) du_1 \cdots du_n$$

در این رابطه $f_i(i)$ چگالی متناظر $y(i)$ است. برای مثال با $n=3$ و $\underline{r} = (3, 1, 2)$ داریم:

$$p(\underline{r}) = P\{R_1 = 3, R_2 = 1, R_3 = 2\} = \int \int \int_{u_1 < u_2 < u_3} f_2(u_1) f_3(u_2) f_1(u_3) du_1 du_2 du_3$$

کلب‌فلیچ و پرنیس رابطه زیر را اثبات نموده‌اند:

$$F_i(t) = 1 - \exp\left(-e^{-\beta' \underline{x}_i} \int_0^t \lambda_o(u) du\right)$$

در نتیجه داریم:

$$p(\underline{r}) = \prod_{i=1}^n \frac{e^{-\beta' \underline{x}_i}}{\sum_{j \in \mathcal{R}(i)} e^{-\beta' \underline{x}_j}}$$

در حالت برش با استفاده از نماد $(Y_1, \delta_1), \dots, (Y_n, \delta_n)$ و فرض این که Y_i دارای توزیع F_i و چگالی f_i است، بردار رتبه را به صورت زیر تعریف می کنیم:

$$\underline{R}^{u/c} = (R_1^{u/c}, \dots, R_n^{u/c})$$

که در آن:

$$R_i^{u/c} = \begin{cases} \delta_i = 1 & \text{رتبه } Y_i \text{ در میان مشاهدات بریده نشده با شرط} \\ \delta_i = 0 & \text{رتبه مشاهده بریده نشده قبلی با شرط} \end{cases}$$

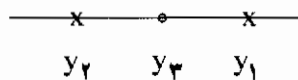
و بردار نشانگر $\underline{\delta}$ به صورت $\underline{\delta} = (\delta_1, \dots, \delta_n)$ است.

حال احتمال $(\underline{r}^{u/c}, \underline{\delta})$ به صورت زیر است:

$$p(\underline{r}^{u/c}, \underline{\delta}) = \int_{u_1 < \dots < u_{n_u}} \prod_{i=1}^{n_u} \left\{ f_{u(i)}(u_i) \times \prod_{j \in C_{i, i+1}} [1 - F_j(u_i)] \right\} du_1 \dots du_{n_u}$$

در این رابطه، $f_{u(i)}$ ، چگالی متناظر i امین مشاهده بریده نشده مرتب، $C_{i, i+1}$ ، مجموع اندیسهای متناظر مشاهدات بریده شده واقع در بین مشاهدات بریده نشده i ام و $(i+1)$ ام و n_u تعداد کل مشاهدات بریده نشده است.

برای مثال، $\underline{r} = (2, 1, 1)$ و $\underline{\delta} = (1, 1, 0)$ ، متناظر با شکل زیر است:



احتمال رتبه به شرح زیر است:

$$p((2, 1, 1), (1, 1, 0)) = \iint_{u_1 < u_2} f_2(u_1) [1 - F_3(u_1)] \times f_1(u_2) du_1 du_2$$

کلب فلیش و پرنتمس نشان دادند که اگر $F_i(t)$ به صورت زیر باشد، امتحان قابل محاسبه است.

$$F_i(t) = 1 - \exp \left(-e^{-\beta' x_i} \int_0^t \lambda_0(u) du \right)$$

$$p(\underline{r}^{u/c}, \underline{\delta}) = \prod_u \frac{e^{-\beta' \underline{x}(i)}}{\sum_{j \in \mathcal{R}(i)} e^{-\beta' \underline{x}_j}} = L_c$$

REFERENCE

Kalbfleisch and Prentice, *Biometrika* (1973).

درست‌نمایی جزئی. دنباله کمیتهای تصادفی زیر را در نظر بگیرید:

$$(X_1, S_1; X_2, S_2; \dots; X_m, S_m)$$

در رگرسیون با برش، فرض کنید: $y_{u(i)}$ ، شاهد i ام بریده نشده باشد. همچنین، فرض کنید: متغیر X_i دارای تمام اطلاعات بریده شده در $[y_{u(i-1)}, y_{u(i)}]$ - با این اطلاع که یک شکست در زمان $y_{u(i)}$ رخ داده - باشد. S_i را متغیری در نظر می‌گیریم که شامل مشاهده خاصی باشد، که با همبستگی $\underline{x}_{u(i)}$ در زمان $y_{u(i)}$ شکست می‌خورد. درست‌نمایی حاشیه‌ای S_1 تا S_m به شرح زیر است:

$$p(S_1, \dots, S_m | \underline{\beta}) = \prod_{i=1}^m p(S_i | S_1, \dots, S_{i-1}; \underline{\beta})$$

و درست‌نمایی شرطی S_1 تا S_m با فرض X_1 تا X_m نیز مطابق زیر است:

$$p(S_1, \dots, S_m | X_1, \dots, X_m; \underline{\beta}) = \prod_{i=1}^m p(S_i | S_1, \dots, S_{i-1}; X_1, \dots, X_m; \underline{\beta})$$

و بالاخره درست‌نمایی کامل به شرح زیر است:

$$p(X_1, \dots, X_m; S_1, \dots, S_m | \underline{\beta})$$

$$= \prod_{i=1}^m p(X_i, S_i | X_1, \dots, X_{i-1}; S_1, \dots, S_{i-1}; \underline{\beta})$$

$$= \prod_{i=1}^m p(X_i | X_1, \dots, X_{i-1}; S_1, \dots, S_{i-1}; \underline{\beta}) \times$$

$$\times \prod_{i=1}^m p(S_i | X_1, \dots, X_{i-1}; S_1, \dots, S_{i-1}; \underline{\beta})$$

کاکس، عبارت دوم حاصل ضرب (عبارت زیر) را درست‌نمایی جزئی نامیده است.

$$\prod_{i=1}^m p(S_i | X_1, \dots, X_{i-1}, X_i; S_1, \dots, S_{i-1}; \beta)$$

در رگرسیون با داده‌های بریده شده، درست‌نمایی جزئی با L_C برابر است، که ما آن را درست‌نمایی شرطی گوئیم. از مقایسه تعاریف درست‌نمایی‌های شرطی و جزئی، معلوم می‌شود که، درست‌نمایی جزئی، یک درست‌نمایی شرطی یسا حاشیه‌ای واقعی نیست. کاکس مدعی است که: درست‌نمایی جزئی شامل بیشترین اطلاعات درباره β برای رگرسیون با داده‌های بریده شده است. لذا، می‌توان از جمله اول حاصل ضرب، یعنی عبارت زیر چشم پوشی کرد.

$$\prod_{i=1}^m p(X_i | X_1, \dots, X_{i-1}; S_1, \dots, S_{i-1}; \beta)$$

بدون آن که چیز زیادی از دست برود، افرون، اطلاع فشر را در درست‌نمایی جزئی با اطلاع فشر در درست‌نمایی کامل برای تعدادی از الگوها مقایسه کرده است. معمولاً اطلاع در L_C بسیار بالاست. با کارایی بیش از ۹۰٪. در حالات نادری L_C ، حداقل به اندازه درست‌نمایی کامل شامل اطلاعات است.

REFERENCES

- Cox, *Biometrika* (1975).
 Efron, *JASA* (1977).
 Oakes, *Biometrika* (1977).

۳.۱ بررسی نرمال مجانبی

در مقاله ۱۹۷۲، کاکس بیان شده، که $\hat{\beta}$ ، جواب معادله $\frac{\partial}{\partial \beta} \log L_C(\beta) = 0$ ، به طور مجانبی دارای توزیع نرمال است. باکس در مقاله ۱۹۷۵ خود یک شناسه کاشف که مشابه شناسه درست‌نمایی حداکثر استاندارد است را ارائه می‌دهد.

تستیاتیس، اثباتی از نرمال مجانبی $\hat{\beta}$ با کاربرد انتگرال و فرایندهای تصادفی ارائه داده است. این اثبات مشابه اثباتی است که توسط برسلو و کرولی برای نرمال مجانبی برآورد گر PL و اثبات ارائه شده توسط کرولی برای آماره MH است.

ییلی با کاربرد تصویرهای هاجک یک شناسه ارائه داده است.

REFERENCES

Cox, Biometrika (1975).

Bailey, Univ. of Chicago thesis (1979).

Tsiatis, Ann. Stat. (1981).

۴.۱ برآورد $S(t; \underline{x})$

تحت الگوی نرخ شکست کاکس، داریم:

$$S(t; \underline{x}) = \exp\left(-e^{\beta' \underline{x}} \int_0^t \lambda_0(u) du\right) = \exp\left(-e^{\beta' \underline{x}} \Lambda_0(t)\right) = S_0(t) e^{\beta' \underline{x}}$$

که در آن $S_0(t) = e^{-\Lambda_0(t)}$ است.

برای برآورد $S(t; \underline{x})$ ، به جای β از $\hat{\beta}$ استفاده می‌کنیم. ولی $\Lambda_0(t)$ یا $S_0(t)$ را چگونه برآورد کنیم.

برسלו فرض می‌کند که $\lambda_0(t)$ بین مشاهدات بریده نشده ثابت است، داریم:

$$\hat{\lambda}_{0,B}(t) = \frac{1}{(y_{u(i)} - y_{u(i-1)}) \sum_{j \in \mathcal{R}_{u(i)}} e^{\beta' \underline{x}_j}}, \quad y_{u(i-1)} < t < y_{u(i)}$$

و $S_0(t)$ را به صورت زیر برآورد می‌کنیم:

$$\hat{S}_{0,B}(t) = \prod_{y(i) \leq t} \left(1 - \frac{\delta(i)}{\sum_{j \in \mathcal{R}(i)} e^{\beta' \underline{x}_j}} \right)$$

توجه شود که $\hat{\Lambda}_{0,B}(t) = \int_0^t \hat{\lambda}_{0,B}(u) du$ و $\hat{S}_{0,B}(t)$ ، حتی به ازای $t = y(i)$ سازگار

نیستند، یعنی:

$$\hat{S}_{0,B}(t) \neq e^{-\hat{\Lambda}_{0,B}(t)}$$

همچنین $\hat{S}_{0,B}(t)$ می‌تواند مقادیر منفی اختیار کند. تسیاتیس از مقدار زیر استفاده

می کند:

$$\hat{S}_{\circ,T}(t) = e^{-\hat{\Lambda}_{\circ,T}(t)}$$

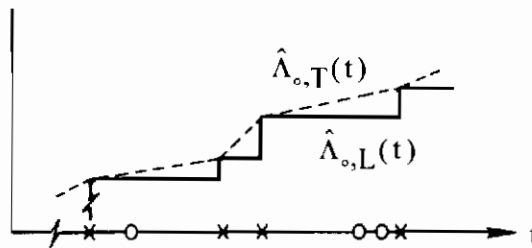
که در آن $\hat{\Lambda}_{\circ,B}(t) = \sum_{y(i) \leq t} \frac{\delta(i)}{\sum_{j \in \mathcal{R}(i)} e^{\beta' x_j}}$ است. ولی هنگامی که $\hat{\beta} = \circ$ باشد، $\hat{S}_{\circ,T}$

به برآوردگر PL، ساده نمی شود. توجه نمایید که $\hat{\Lambda}_{\circ}(t)$ یک تابع پله ای است.

اینک یک صورت خطی $\hat{\Lambda}_{\circ,B}(t) = \sum_{y(i) \leq t} \frac{\delta(i)}{\sum_{j \in \mathcal{R}(i)} e^{\beta' x_j}}$ را با تعریف زیر به

کار می برد، که معادل انتگرال برسلوی $\lambda_{\circ}(t)$ است.

$$\hat{S}_{\circ,L}(t) = e^{-\hat{\Lambda}_{\circ,L}(t)}$$



تسیاتیس و لینک هر دو $\text{Var}(\hat{S}_{\circ}(t))$ را محاسبه کرده اند. اوکی با استفاده از الگوی درست‌نمایی و دومی به روش دلتا. لینک با استفاده از روشهای مونت کارلو، فاصله‌های

اطمینان مربوط به $\hat{S}_{\circ,L}(t)$ و $\log \hat{S}_{\circ,L}(t)$ و $\log \frac{S_{\circ,L}(t)}{1 - \log \hat{S}_{\circ,L}(t)}$ را بررسی و اثبات می کند که احتمال پوشاندن $\hat{S}_{\circ,L}(t)$ خیلی کم، از $\log \hat{S}_{\circ,L}(t)$ خیلی زیاد و از $\log \hat{S}_{\circ,L}(t)$ تقریباً درست است. این نتایج مربوط به پوشش فاصله‌های اطمینان، برای برآوردگر PL نیز برقرار است.

برآوردگرهای دیگر $S(t; \underline{x})$ ، که از نظر محاسباتی پیچیده تراند توسط کاکس و افرون پیشنهاد شده است.

REFERENCES

Breslow, JRSS B (1972), in Discussion on Cox's paper.

_____, Biometrics (1974).

Tsiatis, Univ. Wisconsin Tech. Report No. 524 (1978).

_____, Ann. Stat. (1981).

Link, Stanford Univ. Tech. Report No. 45 (1979).

۵.۱ داده‌های گسسته یا طبقه‌ای

زمانهای متمایز و مرتب شده بقاء را به صورت: $y'(1) < \dots < y'(r)$ می‌نویسیم و فرض زیر را در نظر می‌گیریم:

$\mathfrak{R}(i) = y'(i)$ - نرخ شکست در زمان

$\wp(i) = y'(i)$ (مجموعه مرگهای جدا در زمان $y'(i)$ مرگ در زمان)

$d_i = \#(\wp(i))$

کاکس پیشنهاد می‌کند که از کمیته زیر استفاده شود:

$$L_c = \prod_{i=1}^r P\{\wp(i) \mid \mathfrak{R}(i), d_i\}$$

با:

$$P\{\wp(i) \mid \mathfrak{R}(i), d_i\} = \frac{\exp\left(\sum_{j \in \wp(i)} \beta' \underline{x}_j\right)}{\sum_{\wp^*(i)} \exp\left(\sum_{j \in \wp^*(i)} \beta' \underline{x}_j\right)}$$

در رابطه قبل، مجموع منخرج روی تمام زیر مجموعه‌های $\mathfrak{R}(i) \subset \wp^*(i)$ در نظر گرفته می‌شود، به گونه‌ای که $\# \wp^*(i) = d_i$ باشد. برای $i=1, 2, \dots, r$ ، تعداد $\binom{n_i}{d_i}$ زیرمجموعه وجود دارد، که حتی برای مجموعه داده‌های متوسط هم از نظر محاسباتی مناسب به نظر نمی‌رسد.

یک تابع درست‌نمایی دیگر به شرح زیر توسط برسلو و دیگران پیشنهاد شده است. اگر تعداد تکرارها زیاد نباشد، در عمل معقول به نظر می‌رسد:

$$L_C = \prod_{i=1}^r \frac{\exp\left(\sum_{j \in \mathcal{P}(i)} \beta'_j x_j\right)}{\left(\sum_{j \in \mathcal{R}(i)} e^{\beta'_j x_j}\right)^{d_i}}$$

پرنتمیس محور زمان را به صورت: $a_0 = 0 < a_1 < \dots < a_{r-1} < a_r = \infty$ ، $A_j = [a_{j-1}, a_j)$ ، افزایش می‌کند. اگر یک زمان بقاء در بازه A_j قرار گیرد، آن‌گاه زمان j را ثبت می‌کنیم.

با توجه به عبارت: $\alpha_j = \exp\left(-\int_{a_{j-1}}^{a_j} \lambda_0(t) dt\right)$ ، که احتمال شرطی یک فرد با

متغیر وابسته $x = x_j$ واقع در بازه A_j - با فرض این که در بازه A_{j-1} زنده بوده - است، می‌توان احتمال مشاهده i ام را که در ابتدای A_j زنده است، به صورت زیر محاسبه کرد:

$$P\{Y_i = j, \delta_i\} = \left(\prod_{k=1}^{j-1} \alpha_k e^{\beta'_k x_i}\right) \left(1 - \alpha_j e^{\beta'_j x_i}\right)^{\delta_i}$$

درست‌نمایی کامل به صورت: $L = \prod_{i=1}^n P\{Y_i = j, \delta_i\}$ بوده که تابعی از پارامترهای

مجهول β و α_1 تا α_r است.

برای برآورد این پارامترها از حداکثر درست‌نمایی استفاده می‌کنیم. توجه شود که

α_i ها به صورت زیر محدود می‌شوند:

$$0 < \alpha_j < 1 \quad \text{و} \quad j=1, \dots, r \quad \text{و} \quad \sum_{j=1}^r \alpha_j = 1$$

اگر α_r را حذف نموده و فرض کنیم: $\gamma_j = \log(-\log \alpha_j)$ ، با $j=1, \dots, r-1$ ، به طوری که:

$$-\infty < \gamma_j < +\infty \quad \text{و} \quad j=1, \dots, r-1$$

آن‌گاه حداکثر کردن، نسبت به γ_1 تا γ_{r-1} ، ساده‌تر است از این که نسبت به α_1 تا α_r حداکثر کنیم. زیرا، در این صورت نگران مقادیر مرزی نیستیم. همچنین روش نیوتن-رافسون سریعتر همگرا خواهد بود.

REFERENCES

Cox, JRSS B (1972).

Kalbfleisch and Prentice, JRSS B (1972), and

Peto, JRSS B (1972), in Discussion on Cox's paper.

Breslow, Biometrics (1974).

Prentice and Gloeckler, Biometrics (1978).

۶.۱ متغیرهای وابسته به زمان

به حالتی که متغیرها با زمان تغییر می‌کنند، تعمیم می‌دهیم. در نتیجه همراه با $Y_i = T_i \wedge C_i$ و $\delta_i = I(T_i \leq C_i)$ ، مقادیر $x_i(t)$ را مشاهده می‌کنیم. در الگوی نرخ شکست متناسب، فرض می‌شود، که تابع نرخ شکست λ_i آمین شاهد به شرح زیر باشد:

$$\lambda_i(t) = e^{\beta' x_i(t)} \lambda_0(t)$$

در نتیجه داریم:

$$P\{\text{یک مرگ در } \mathcal{R}(i) \text{ در زمان } y(i) \mid \text{مرگ } (i) \text{ در زمان } y(i)\} = \frac{e^{\beta' x_i(i)(y(i))}}{\sum_{j \in \mathcal{R}(i)} e^{\beta' x_j(i)(y(i))}}$$

درست‌نمایی شرطی به صورت زیر در می‌آید:

$$L_c = \prod_u \frac{e^{\beta' x_i(i)(y(i))}}{\sum_{j \in \mathcal{R}(i)} e^{\beta' x_j(i)(y(i))}}$$

در این حالت که به زمان بستگی دارد، هیچ‌گونه اثباتی برای نرمال مجانبی وجود ندارد. همچنین، برای حجم نمونه متوسط یا بزرگ محاسبات کند و پرهزینه خواهد بود.

۷.۱ مثال ۱: داده‌های پیوند قلب استانفورد. آیا بیماران قلبی که قلب خود را عوض می‌کنند از بیمارانی که چنین نمی‌کنند، عمر بیشتری دارند؟ نوعاً، یک بیمار، قلب خود را هنگامی عوض می‌کند که قلبی وجود داشته باشد. در این موقعیت می‌گوییم که بیمار از جمعیت بدون تعویض به جمعیت تعویض شده‌ها منتقل شده است. متغیر وابسته که انتقال را نشان می‌دهد از ۰ به ۱ تغییر می‌کند. متغیرهای وابسته دیگر شامل سن، زمان انتظار برای انتقال، زمان شروع معالجه و ... است.

REFERENCES

Turnbull, Brown, and Hu, JASA (1974).

Crowley and Hu, JASA (1977).

۸.۱ مثال ۲: فرزند خواندگی و آبستنی

آیا زوجهایی با دارا بودن مشکل نازایی و عقیمی، که بچه‌ای را به فرزند قبول کرده‌اند، بیشتر از زوجی که این عمل را انجام نداده‌اند، مایل به آبستنی هستند یا خیر؟ در این‌جا، امکان دارد زوجها از جامعه فرزند خوانده‌دار منتقل شوند. برش هنگامی رخ می‌دهد که زوج کوشش خود را در آبستن شدن متوقف کنند.

REFERENCES

Lamb and Leurgans, Amer. J. Obstet. Gyn. (1979).

Leurgans, Stanford Univ. Tech. Report No. 57 (1980).

۲ الگوی خطی

الگوی خطی استاندارد به صورت $T_i = \alpha + \beta x_i + e_i$ و یا به شکل زیر است:

$$T_i = \alpha + \beta' \underline{x}_i + e_i, \quad i = 1, \dots, n$$

که در آن e_1 تا e_n دارای خاصیت iid با توزیع مشترک F هستند. فرض کنید C_1 تا C_n مستقل باشند. C_i زمان برش مربوط به T_i است. مشاهدات به شرح زیرند:

$$Y_i = T_i \wedge C_i \quad \text{و} \quad \delta_i = I(T_i \leq C_i)$$

الگوهای زمانی شتاب داده شده: الگوهای خطی به وسیله الگوهای زمانی شتاب داده شده، به الگوهای خطی شکست مربوط می‌شوند. فرض کنید Z_0 زمان بقاء با تابع نرخ شکست زیر باشد:

$$\lambda_0 = \frac{f_0(z)}{1 - F_0(z)}$$

همچنین، فرض می‌کنیم که زمان بقاء یک فرد با متغیر وابسته x با متغیر زیر هم توزیع باشد.

$$Z_x = e^{\beta' x} Z_0$$

توجه شود که اگر $\beta' x < 0$ باشد، آن‌گاه Z_x از Z_0 کوتاهتر بوده و می‌گوییم که متغیر وابسته، زمان شکست را شتاب داده است. نرخ شکست Z_x ، به شرح زیر است:

$$\lambda_x(z) = \frac{f_x(z)}{1 - F_x(z)} = \frac{f_0(e^{-\beta' x} z) e^{\beta' x}}{1 - F_0(e^{-\beta' x} z)} = \lambda_0(e^{-\beta' x} z) e^{\beta' x}$$

با تعریف $T_x = \log Z_x$ ، داریم:

$$E(T_x) = \beta' x + E(\log Z_0) = \beta' x + \alpha$$

در نتیجه الگوی زمانی شتاب داده شده با یک الگوی لگاریتمی خطی، به شرح زیر برابر است:

$$T_x = \alpha + \beta' x + e$$

به گونه‌ای که $e = \log Z_0 - E(\log Z_0)$ است. در کاربرد روشهای الگوی خطی برای داده‌های بقاء، اغلب لازم است که داده‌ها را توسط لگاریتم گیری تبدیل کنیم تا متقارن شوند. در نتیجه از این دیدگاه، الگوی زمان شتاب داده شده مناسب است.

REFERENCES

Prentice and Kalbfleisch, *Biometrics* (1979), and
Kalbfleisch and Prentice, *The Statistical Analysis of Failure Time Data*
(1980), both discuss the accelerated time model.

۱.۲ آزمونهای رتبه خطی

اگر برش نباشد، تواناترین آماره رتبه برای آزمون $H_0: \beta = 0$ در مقابل $H_0: \beta \neq 0$ ،
در حالت $p=1$ به صورت: $\left. \frac{d}{d\beta} \log p(r) \right|_{\beta=0}$ است، که احتمال بردار رتبه r است.

همچنین، در صورت وجود برش می‌توان از آماره زیر استفاده کرد:

$$\left. \frac{d}{d\beta} \log p(\underline{r}^{u/c}, \underline{\delta}) \right|_{\beta=0}$$

که با توجه به بخش اول، داریم:

$$p(\underline{r}^{u/c}, \underline{\delta}) = \int \cdots \int_{u_1 < \cdots < u_{n_u}} \prod_{i=1}^{n_u} \left\{ f_{u(i)}(u_i) \times \prod_{j \in C_{i,i+1}} [\lambda - F_j(u_j)] \right\} du_1 \cdots du_{n_u}$$

$$f_i(u) = f(u - \beta x_i)$$

می‌توان نشان داد که:

$$\left. \frac{d}{d\beta} \log p(\underline{r}^{u/c}, \underline{\delta}) \right|_{\beta=0} = \sum_{i=1}^{n_u} \left\{ x_{u(i)} c_i + \left(\sum_{j \in C_{i,i+1}} x_j \right) C_i \right\}$$

که در این رابطه، داریم:

$$c_i = \left(\prod_{j=1}^{n_u} n_{u(j)} \right) \int \cdots \int_{u_1 < \cdots < u_{n_u}} \left\{ -\frac{d}{du_i} \log f(u_j) \right\} \\ \times \prod_{j=1}^{n_u} \left\{ f(u_j) [\lambda - F(u_j)]^{m_{u(j)}} \right\} du_1 \cdots du_{n_u}$$

$$C_i = \left(\prod_{j=1}^{n_u} n_{u(j)} \right) \int \cdots \int_{u_1 < \cdots < u_{n_u}} \left\{ -\frac{d}{du_i} \log [\lambda - F(u_i)] \right\} \\ \times \prod_{j=1}^{n_u} \left\{ f(u_j) [\lambda - F(u_j)]^{m_{u(j)}} \right\} du_1 \cdots du_{n_u}$$

$$m_{u(j)} = \# \text{ in } C_{j,j+1}$$

فرض کنید توزیع خطا، توزیع مقادیر فرین باشد. که تابع چگالی و بقاء آن به شرح زیر است:

$$f(t) = e^{t-e^t} \quad \text{و} \quad 1 - F(t) = e^{-e^t}$$

حال مقدار C_i ، به صورت زیر خواهد بود:

$$c_i = \sum_{j=1}^i \frac{1}{n_{u(j)}} - 1$$

$$C_i = \sum_{j=1}^i \frac{1}{n_{u(j)}}$$

در نتیجه تواناترین آماره رتبه در این حالت به صورت: $-\sum (x_{(i)} - \bar{x}_{(i)})$ ، است. این عبارت برابر منخرج آماره کاکس در آزمون $H_0: \beta = 0$ است (بخش اول) و آن را آزمون رتبه لگاریتمی می‌نامند.

REFERENCES

Peto and Peto, JRSS A (1972), introduce linear rank tests and coin the term "log rank test".

Latta, Biometrika (1977), establishes a connection between linear rank tests and Efron's test.

Morton, Biometrika (1978), discusses permutation theory for linear rank tests.

Prentice, Biometrika (1978), give the preceding derivation of the linear rank test statistic and calculates its variance.

Kalbfleisch and Prentice, The Statistical Analysis of Failure Time Data (1980), Chapter 6.

۲.۲ برآوردگرهای کمترین مربعات

برای سادگی فرض کنید $p = 1$ باشد، یعنی: $E(T_i) = \alpha + \beta x_i$. این برآوردگر را می‌توان به بیش از یک متغیر وابسته تعمیم داد.

برآوردگرهای میلر. برای مشاهدات بریده نشده، برآوردهای $\hat{\alpha}$ و $\hat{\beta}$ ، عبارت زیر را کمینه می‌سازند:

$$\sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 = n \int_{-\infty}^{+\infty} z^2 dF_n(z)$$

به طوری که F_n توزیع تجربی Z_1 تا Z_n است و $z_i = y_i - \alpha + \beta x_i$.

در صورت وجود برش، میلر پیشنهاد می کند که عبارت زیر را کمینه سازیم:

$$n \int_{-\infty}^{+\infty} z^{\dagger} d\hat{F}(z) = \sum_{i=1}^n \hat{w}_i(\beta)(y_i - \alpha - \beta x_i)^{\dagger}$$

که در آن \hat{F} برآوردگر PL برآوردگر مبنای $(z_1, \delta_1), \dots, (z_n, \delta_n)$ است و وزنهای $\hat{w}_1(\beta)$ تا $\hat{w}_n(\beta)$ جهشهای برآوردگر PL هستند. در نگاه اول می توان نتیجه گرفت که مجموع موزون مربعها به مشاهدات بریده شده بستگی ندارند. با این وجود که در اصل، برآوردگر PL و در نتیجه هر کوواریانس از وزنهای به مشاهدات بریده شده وابسته اند.

اگر آخرین مشاهده بریده شده باشد و $\delta(n) = 0$ ، آن را بریده نشده تغییر می دهیم.

در این صورت $\sum_{i=1}^n \hat{w}_i(\beta) = 1$ است.

وزنها را فقط به صورت تابعی از β نوشته ایم. زیرا، اضافه کردن ثابتی مانند α به هر T_i ، نتیجه اش فقط یک انتقال برآوردگر PL و جهشهای آن است و در نتیجه وزنهای α بستگی ندارند.

برای محاسبه $\hat{\alpha}$ و $\hat{\beta}$ ، نسبت به α مشتق می گیریم، داریم:

$$\hat{\alpha} = \sum_{i=1}^n \hat{w}_i(\beta) y_i - \beta \sum_{i=1}^n \hat{w}_i(\beta) x_i$$

اگر این مقدار را در مجموع موزون مربعات قرار دهیم، فقط تابعی از β به دست می آید.

$$f(\beta) = \sum \hat{w}_i(\beta)(y_i - \hat{\alpha} - \beta x_i)^{\dagger}$$

که با تجسس می توان آن را کمینه کرد.

چون تابع $f(\beta)$ پیوسته نیست، ممکن است روش تجسس خسته کننده باشد، به ویژه در ابعاد بالا. میلر روش بهتر زیر را به عنوان راه حل دیگری پیشنهاد نموده است. در این روش، برآورد اولیه را به صورت زیر تعریف می کنیم، که $\hat{\beta}^{\circ}$ همان شیب خط کمترین مربعات مشاهدات بریده نشده است.

$$\hat{\beta}^{\circ} = \frac{\sum y_i(x_i - \bar{x}_u)}{\sum (x_i - \bar{x}_u)^{\dagger}}$$

با این حدس اولیه $\hat{\beta}^{\circ}$ ، داریم: $\hat{z}_i^{\circ} = y_i - \hat{\beta}^{\circ} x_i$ ، $i = 1, \dots, n$. فرض کنید \hat{F}° ، برآوردگر PL

بر مبنای (\hat{z}_1, δ_1) تا (\hat{z}_n, δ_n) و $\hat{w}_1(\hat{\beta}^\circ)$ تا $\hat{w}_n(\hat{\beta}^\circ)$ ، جهشهای \hat{F}° باشند. برآورد جدید به صورت زیر تعریف می‌شود:

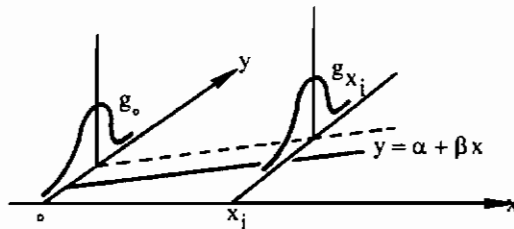
$$\hat{\beta}^1 = \frac{\sum_u \hat{w}_i^*(\hat{\beta}^\circ) y_i (x_i - \bar{x}_u^*)}{\sum_u \hat{w}_i^*(\hat{\beta}^\circ) (x_i - \bar{x}_u^*)^2}$$

$$\bar{x}_u^* = \sum_u \hat{w}_i^*(\hat{\beta}^\circ) x_i \quad \text{و} \quad \hat{w}_i^*(\hat{\beta}^\circ) = \frac{\hat{w}_i(\hat{\beta}^\circ)}{\sum_u \hat{w}_i(\hat{\beta}^\circ)}$$

با دوباره نرمال کردن وزنه‌های $\hat{w}_i^*(\hat{\beta}^\circ)$ ، به ما اجازه می‌دهد که از آخرین جمله مرتب شده \hat{z}_i° - در هر دو حالت بریده شده و نشده - چشم‌پوشی کنیم. تنها مشاهدات برید نشده در مجموع ظاهر می‌شوند. روش معمول تعریف دوباره آخرین جمله مرتب شده \hat{z}_i° - به گونه‌ای که بریده شود، اگر بریده شده است - نتایج کمتر پایداری را در برآورد تکراری β ارائه می‌کند. ولی هنوز می‌توان آن را برای برآورد α به کار برد. روش بالا را تکرار می‌کنیم با این امید که همگرا شود. با این وجود، امکان دارد دنباله برآورد گره‌های β در یک دور که بین دو مقدار نوسان می‌کند، قرار بگیرد. در این حالت متوسط دو مقدار را انتخاب می‌کنیم.

فرض نمایید که تغییرپذیری مربوط به وزنه‌های $\hat{w}_i^*(\hat{\beta})$ قابل حذف باشد. داریم:

$$\hat{\text{Var}}(\hat{\beta}) = \frac{\sum_u \hat{w}_i^*(\hat{\beta}^\circ) (y_i - \hat{\alpha} - \hat{\beta} x_i)}{\sum_u \hat{w}_i^*(\hat{\beta}^\circ) (x_i - \bar{x}_u^*)^2}$$



برای محاسبه برآورد واریانس بالا، مانند اثبات سازگاری برآوردهای $\hat{\alpha}$ و $\hat{\beta}$ ، به این فرض بستگی دارد که توزیع بریده شده مشاهده‌ی am به صورت: $G_{x_i}(c) = G_0(c - \beta x_i)$

باشد. برای تابع توزیعی مانند G_0 ، اگر چگالی آن (g_0) مانند شکل زیر باشد، آن گاه G_{X_i} دارای چگالی g_{X_i} خواهد بود. در واقع g_{X_i} ، انتقال یافته g_0 به اندازه βx_i است.

REFERENCE

Miller, Biometrika (1976).

بر آوردگر باکلی - جیمز. در این الگو فرض می شود: $E(T_i) = \alpha + \beta x_i$. متأسفانه، نمی توان T_i را مشاهده کنیم. ولی، فقط Y_i قابل مشاهده است، همچنین $E(Y_i) \neq \alpha + \beta x_i$. در نتیجه باکلی و جیمز متغیرهای کاذب زیر را تعریف می کنند:

$$Y_i^* = Y_i \delta_i + E(T_i | T_i > Y_i)(1 - \delta_i)$$

محاسبه $E(Y_i^*)$ به شرح زیر است:

$$\begin{aligned} E(Y_i^*) &= \int_0^{\infty} u(1 - G_i(u)) dF_i(u) + \int_0^{\infty} \left[\int_u^{\infty} \frac{s dF_i(s)}{1 - F_i(u)} \right] (1 - F_i(u)) dG_i(u) \\ &= \int_0^{\infty} u(1 - G_i(u)) dF_i(u) + \int_0^{\infty} \left[\int_0^s dG_i(u) \right] s dF_i(s) \\ &= \int_0^{\infty} u(1 - G_i(u)) dF_i(u) + \int_0^{\infty} G_i(s) s dF_i(s) \\ &= \int_0^{\infty} u dF_i(u) \\ &= \alpha + \beta x_i \end{aligned}$$

بنابراین، اگر بتوانیم y_1^* تا y_n^* را مشاهده کنیم، معقول به نظر می رسد که از برآوردگرهای زیر استفاده کنیم:

$$\hat{\alpha} = \bar{y}^* - \hat{\beta} \bar{x} \quad \text{و} \quad \hat{\beta} = \frac{\sum_{i=1}^n y_i^*(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (15)$$

چون نمی توانیم تمام y_1^* تا y_n^* را مشاهده کنیم، لذا، از برآوردگرهای مشاهدات گم شده

استفاده می‌کنیم. اگر $\delta_i = 0$ باشد، تعریف زیر را در نظر می‌گیریم:

$$\hat{E}(T_i | T_i > y_i) = \hat{\beta}x_i + \frac{\sum_{\hat{z}_k > \hat{z}_i} \hat{w}_k(\hat{\beta}) \hat{z}_k}{1 - \hat{F}(\hat{z}_i)} \quad (16)$$

که در این رابطه $\hat{z}_i = y_i - \hat{\beta}x_i$ و \hat{F} برآوردگر PL بر مبنای (\hat{z}_1, δ_1) تا (\hat{z}_n, δ_n) است. $\hat{w}_n(\hat{\beta})$ تا $\hat{w}_1(\hat{\beta})$ ، جهشهای \hat{F} است. در این صورت تعریف زیر را داریم:

$$\hat{y}_i^* = y_i \delta_i + \left[\hat{\beta}x_i + \frac{\sum_{\hat{z}_k > \hat{z}_i} \hat{w}_k(\hat{\beta}) \hat{z}_k}{1 - \hat{F}(\hat{z}_i)} \right] (1 - \delta_i) \quad (17)$$

با توجه به این که معادلات (۱۵)، $\hat{\beta}$ را به صورت تابعی از y_i^* و معادله (۱۷)، y_i^* را به صورت تابعی از $\hat{\beta}$ می‌دهد، لذا، باید روش تکراری را برای محاسبه آنها به کار ببریم. مانند برآورد میلر، امکان دارد دنباله برآوردهای β بین دو مقدار نوسان کند، در این حالت نیز مقدار متوسط این دو را انتخاب می‌کنیم.

باکلی و جیمز مدعی هستند که اگر برآوردهای β نوسان کنند، آن‌گاه تفاضل بین این دو مقدار، کوچکتر از برآورد میلر است. علاوه بر این اعتبار روش آنها به فرضهای توزیع بریده شده G_1 بستگی ندارد. باکلی و جیمز، برآورد واریانس را به صورت زیر ارائه می‌کنند:

$$\widehat{\text{Var}}(\hat{\beta}) = \frac{\hat{\sigma}_u^2}{\sum_u (x_i - \bar{x}_u)^2}$$

که در آن $\hat{\sigma}_u^2$ به صورت زیر است:

$$\hat{\sigma}_u^2 = \frac{1}{n_u - 2} \sum_u (y_i - \bar{y}_u - \hat{\beta}(x_i - \bar{x}_u))^2$$

در این جا اثبات آن ارائه نمی‌شود.

REFERENCE

Buckley and James, *Biometrika* (1979).

تذکرها:

(الف) روش باکلی و جیمز یک روش ناپارامتری و مشابه روش نرمال اشمی و

هان است. با توجه به تعریف $W_i = \frac{T_i - \alpha - \beta x_i}{\sigma}$ ، اگر F نرمال باشد. آن گاه داریم:

$$E(T_i | T_i > y_i) = E\left(\sigma W_i + \alpha + \beta x_i \mid W_i > \frac{y_i - \alpha - \beta x_i}{\sigma}\right)$$

$$= \alpha + \beta x_i + \frac{\sigma \int_{(y_i - \alpha - \beta x_i)/\sigma}^{\infty} w \phi(w) dw}{1 - \Phi\left(\frac{y_i - \alpha - \beta x_i}{\sigma}\right)} = \alpha + \beta x_i + \frac{\sigma \phi\left(\frac{y_i - \alpha - \beta x_i}{\sigma}\right)}{1 - \Phi\left(\frac{y_i - \alpha - \beta x_i}{\sigma}\right)}$$

در این رابطه ϕ و Φ به ترتیب تابعهای چگالی و توزیع نرمال استاندارد است. اشمی و هان به جای (۱۶)، برآورد زیر را به کار برده‌اند.

$$\hat{E}(T_i | T_i > y_i) = \hat{\alpha} + \hat{\beta} x_i + \frac{\hat{\sigma} \phi\left(\frac{y_i - \hat{\alpha} - \hat{\beta} x_i}{\hat{\sigma}}\right)}{1 - \Phi\left(\frac{y_i - \hat{\alpha} - \hat{\beta} x_i}{\hat{\sigma}}\right)}$$

REFERENCE

Schmee and Hahn, Technometrics (1979).

(ب) هر دو روش پارامتری و ناپارامتری، شبیه الگوریتم EM در نظریه حداکثر

درست‌نمایی‌اند.

REFERENCE

Dempster, Laird, and Rubin, JRSS B (1977).

برآوردگر کول-سوزارلا-وان رایزین

تعریف $Y_i^* = \frac{\delta_i Y_i}{1 - G(Y_i)}$ را در نظر می‌گیریم، داریم:

$$E(Y_i^*) = \int_0^{\infty} \frac{u}{1 - G(u)} (1 - G(u)) dF_i(u) = \int_0^{\infty} u dF_i(u) = \alpha + \beta x_i$$

بنابراین، اگر بتوانیم y_1^* تا y_n^* را مشاهده کنیم، می‌توانیم به طور معمول α و β را به کمک (۱۵) برآورد کنیم. متأسفانه، نمی‌توان همه y_1^* تا y_n^* را مشاهده کنیم، ولی، می‌توانیم برآوردها را جانشین نماییم. برای این کار، به جای G از برآوردگر PL استفاده شود، که در آن نقش متغیرهای بقاء و بریده شده عوض می‌شود. این پژوهشگران پیشنهاد می‌کنند که از برآوردگر بیز G استفاده شود.

امتیاز بزرگ این روش بی‌نیازی از تکرار است. همچنین بیشتر بر مبنای یک توزیع بریده شده عمل می‌کند تا با توزیعهای برشی انتقال یافته‌ای، مانند برآوردگر میلر. با این وجود، y_i^* ها مقادیر خاصی هستند. اینها یا صفراند یا مقادیر y_i ها را افزایش می‌دهند. رفتار این برآوردها در این جا ارزشیابی نمی‌شود.

REFERENCE

Koul, Susarla, and Van Ryzin, unpublished manuscript (1979).

مثال. داده‌های پیوند قلب استانفورد: روشهای کاکس-میلر و باکلی-جیمز را به کمک داده‌های جدول (۵)، مقایسه می‌کنیم. در رگرسیون اولی (شکل ۵) متغیر وابسته: (زمان بقاء) \log_{10} است، که در آن زمان بقاء تا زمان مرگ به علت پذیرفتن قلب اهدایی و متغیر کمکی امتیاز مربوط به طرز عمل است. در رگرسیون دومی (شکل ۶)، متغیر وابسته: (زمان بقاء) \log_{10} است، که در آن زمان بقاء تا لحظه مرگ -صرفنظر از این که به علت پذیرفتن قلب اهدایی یا علل دیگر باشد- است. متغیر کمکی سن است. اگر زمان بقاء صفر باشد آن را به ۱ تغییر می‌دهیم تا لگاریتم آن قابل محاسبه باشد. مقایسه‌های هر سه روش در جدولهای شماره ۶ و ۷، ارائه شده‌اند.

هر سه روش، نتایج پیچیده‌ای را در رگرسیون روی سن نشان می‌دهند: روش کاکس، نشانگر این است که اثر سن بسیار معنی‌دار است. روش میلر مدعی است که سن اثری ندارد و روش باکلی و جیمز، ادعای وسطی را بیان می‌کند. در این حالت می‌توان به کمک الگوی برش، برآوردهای میلر را کنار گذاشت.

همچنین درباره‌ی درجه اعتبار اثر ناهمخوانی امتیاز عدم توافق وجود دارد. باید تلاشهای بیشتری را در مورد این که کدام الگو (زمانی شتاب داده شده یا نرخ شکست متناسب) برای این داده‌ها مناسبتراند، انجام داد.

REFERENCES

Millre, *Biometrika* (1976).

Buckley and James, *Biometrika* (1979).

جدول ۵. داده‌های پیوند قلب استانفورد

روزهای تحمل	زمان بقاء	۱=مرده ۰=زنده	۱=نپذیرفتن ۰=پذیرفتن	مقادیر ناپارامتر T ₅	سن در Tx زمان
۳	۱۵	۱	۰	۱,۱۱	۵۴,۳
۴	۳	۱	۰	۱,۶۶	۴۰,۴
۷	۶۲۴	۱	۱	۱,۳۲	۵۱,۰
۱۰	۴۶	۱	۱	۰,۶۱	۴۲,۵
۱۱	۱۲۷	۱	۰	۰,۳۶	۴۸,۰
۱۳	۶۴	۱	۱	۱,۸۹	۵۴,۶
۱۴	۱۳۵۰	۱	۱	۰,۸۷	۵۴,۱
۱۶	۲۸۰	۱	۱	۱,۱۲	۴۹,۵
۱۸	۲۳	۱	۰	۲,۰۵	۵۶,۹
۲۰	۱۰	۱	۱	۲,۷۶	۵۵,۳
۲۱	۱۰۲۴	۱	۱	۱,۱۳	۴۳,۴
۲۲	۳۹	۱	۱	۱,۳۸	۴۲,۸
۲۳	۷۳۰	۱	۱	۰,۹۶	۵۸,۴
۲۴	۱۳۶	۱	۱	۱,۶۲	۵۲,۰
۲۵	۱۷۷۵	۰	۰	۱,۰۶	۳۳,۳
۲۸	۱	۱	۰	۰,۴۷	۵۴,۲
۳۰	۸۳۶	۱	۱	۱,۵۸	۴۵,۰
۳۲	۶۰	۱	۱	۰,۶۹	۶۴,۵
۳۳	۱۵۳۶	۰	۰	۰,۹۱	۴۹,۰
۳۴	۱۵۴۹	۰	۰	۰,۳۸	۴۰,۶
۳۶	۵۴	۱	۱	۲,۰۹	۴۹,۰
۳۷	۴۷	۱	۱	۰,۸۷	۶۱,۵
۳۸	۰	۱	۰	۰,۸۷	۴۱,۵
۳۹	۵۱	۱			۵۰,۵
۴۰	۱۳۶۷	۰	۰	۰,۷۵	۴۸,۶

ادامه جدول ۵.

روزهای تحمل	زمان بقاء	۱=مرده ۰=زنده	۱=نپذیرفتن ۰=پذیرفتن	مقادیر نا برابر T۵	سن در زمان Tx
۴۱	۱۲۶۴	۰	۰	۰٫۹۸	۴۵٫۵
۴۵	۴۴	۱	۰	۰٫۰	۳۶٫۲
۴۶	۹۹۴	۱	۱	۰٫۸۱	۴۸٫۶
۴۷	۵۱	۱	۱	۱٫۳۸	۴۷٫۲
۴۹	۱۱۰۶	۰	۰	۱٫۳۵	۳۶٫۸
۵۰	۸۹۷	۱			۴۶٫۱
۵۱	۲۵۳	۱	۱	۱٫۰۸	۴۸٫۸
۵۳	۱۴۷	۱			۴۷٫۵
۵۵	۵۱	۱	۱	۱٫۵۱	۵۲٫۵
۵۶	۸۷۵	۰	۰	۰٫۹۸	۳۸٫۹
۵۸	۳۲۲	۱	۱	۱٫۸۲	۴۸٫۱
۵۹	۸۳۸	۰	۰	۰٫۱۹	۴۱٫۶
۶۰	۶۵	۱	۱	۰٫۶۶	۴۹٫۱
۶۳	۸۱۵	۰	۰	۱٫۹۳	۳۲٫۷
۶۴	۵۵۱	۱	۰	۰٫۱۲	۴۸٫۹
۶۵	۶۶	۱	۱	۱٫۱۲	۵۱٫۳
۶۷	۲۲۸	۱	۰	۱٫۰۲	۱۹٫۷
۶۸	۶۵	۱	۱	۱٫۶۸	۴۵٫۲
۶۹	۶۶۰	۰	۰	۱٫۲۰	۴۸٫۰
۷۰	۲۵	۱	۱	۱٫۶۸	۵۳٫۰
۷۱	۵۸۹	۰	۰	۰٫۹۷	۴۷٫۵
۷۲	۵۹۲	۰	۰	۱٫۴۶	۲۶٫۷
۷۳	۶۳	۱	۱	۲٫۱۶	۵۶٫۴
۷۴	۱۲	۱	۰	۰٫۶۱	۲۹٫۲
۷۶	۴۹۹	۰	۰	۱٫۷۰	۵۲٫۲

ادامه جدول ۵.

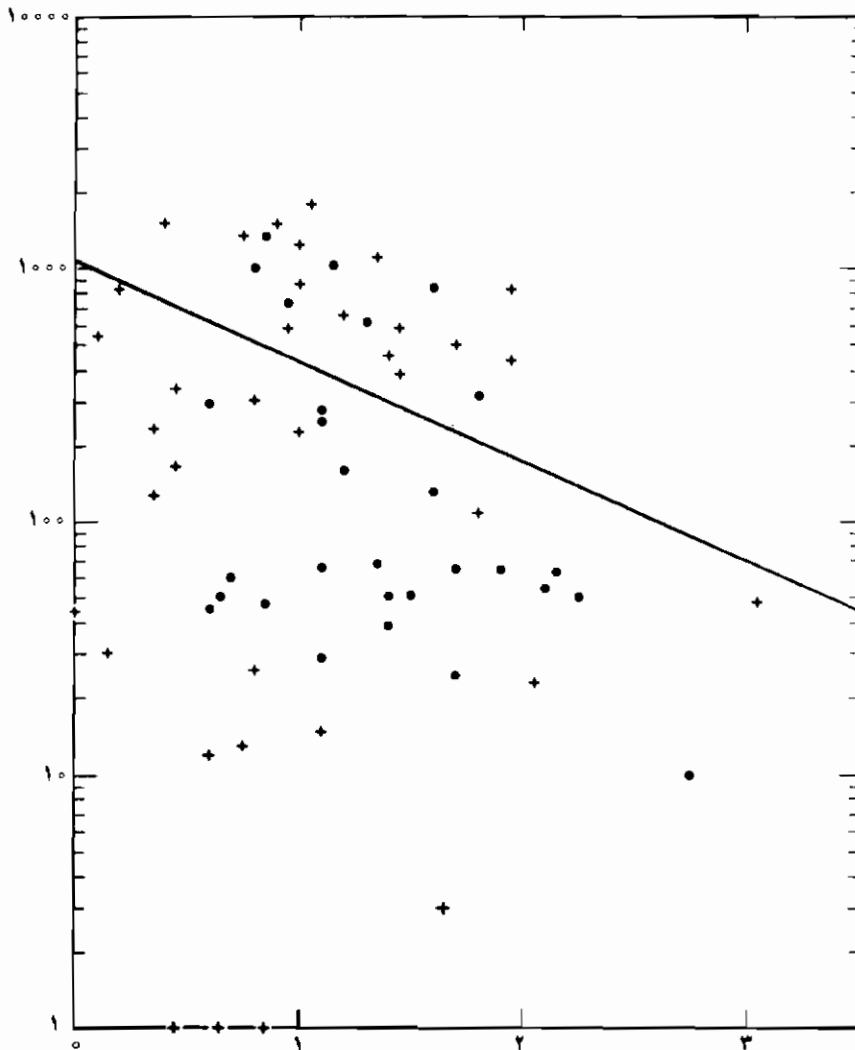
روزهای تحمل	زمان بقاء	۱=مرده ۰=زننده	۱=نپذیرفتن ۰=پذیرفتن	مقادیر نا برابر T5	سن در زمان Tx
۷۸	۳۰۵	۰	۰	۰٫۸۱	۴۹٫۳
۷۹	۲۹	۱	۱	۱٫۰۸	۵۴٫۰
۸۰	۴۵۶	۰	۰	۱٫۴۱	۴۶٫۵
۸۱	۴۳۹	۰	۰	۱٫۹۴	۵۲٫۹
۸۳	۴۸	۱	۰	۳٫۰۵	۵۳٫۴
۸۴	۲۹۷	۱	۱	۰٫۶۰	۴۲٫۸
۸۶	۳۸۹	۰	۰	۱٫۴۴	۴۸٫۹
۸۷	۵۰	۱	۱	۲٫۲۵	۴۶٫۴
۸۸	۳۳۹	۰	۰	۰٫۶۸	۵۴٫۴
۸۹	۶۸	۱	۱	۱٫۳۳	۵۱٫۴
۹۰	۲۶	۱	۰	۰٫۸۲	۵۲٫۵
۹۲	۳۰	۰	۰	۰٫۱۶	۴۵٫۸
۹۳	۲۳۷	۰	۰	۰٫۳۳	۴۷٫۸
۹۴	۱۶۱	۱	۱	۱٫۲۰	۴۳٫۸
۹۵	۱۴	۱			۴۰٫۳
۹۶	۱۶۷	۰	۰	۰٫۴۶	۲۶٫۷
۹۷	۱۱۰	۰	۰	۱٫۷۸	۲۳٫۷
۹۸	۱۳	۰	۰	۰٫۷۷	۲۸٫۹
۱۰۰	۱	۰	۰	۰٫۶۷	۳۵٫۲

نمودار ۵. بقاء در مقابل امتیازهای ناهم‌خوان T_1 .

"+" = زنده یا مرگ رد شده

"•" = مرگ رد شده

"-" = خط کمترین مربعات کاپلان-مایر

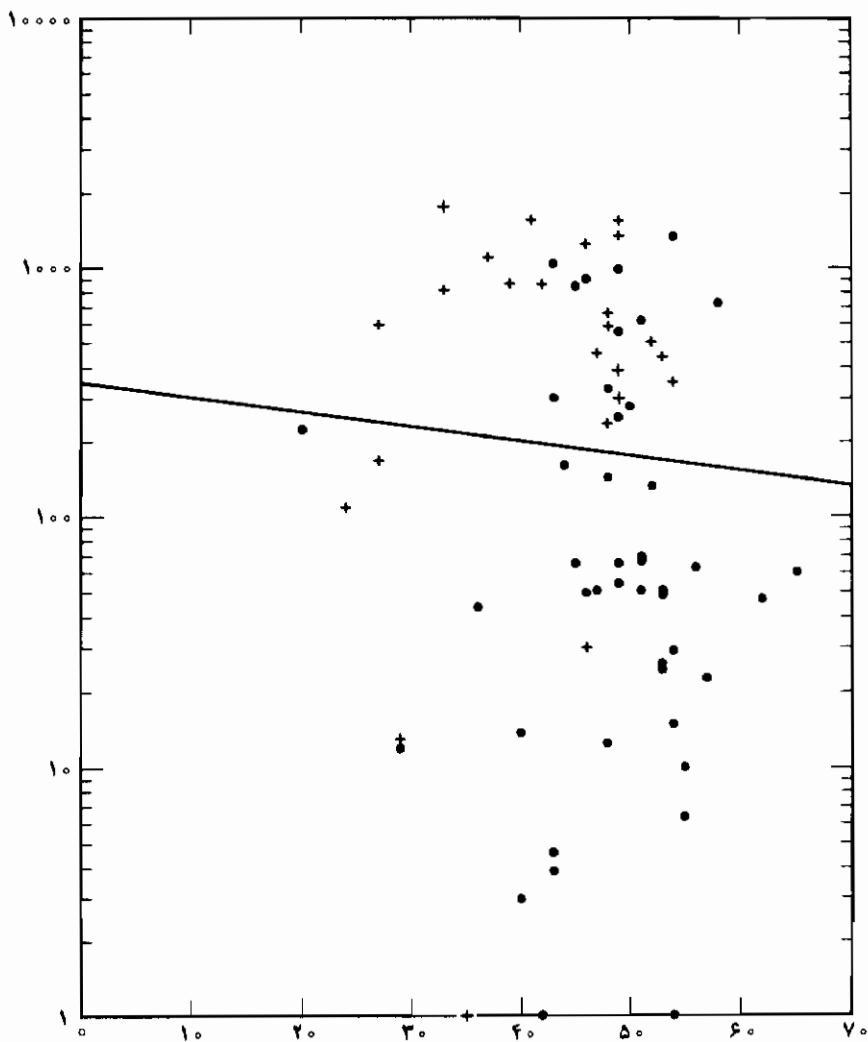


نمودار ۶. بقاء در مقابل سن.

"+" = زننده

"•" = مرده

"-" = خط کمترین مربعات کاپلان-مایر



جدول ۶. رگرسیون لگاریتم زمان بقاء بر امتیاز ناهم‌خوان

روش	$\hat{\alpha}$	$\hat{\beta}$	$\hat{SD}(\hat{\beta})$
کاکس	—	۱,۰۷۶	۰,۳۶۸
میلر	۳,۰۳۶	-۰,۳۹۴	—
تعمیم میلر	۳,۱۲۰	-۰,۴۵۲	۰,۲۳۶
باکلی-جیمز	۳,۱۴۵	-۰,۴۷۱	۰,۲۳۴
—	—	—	—

جدول ۷. رگرسیون لگاریتم زمان بقاء بر سن

روش	$\hat{\alpha}$	$\hat{\beta}$	$\hat{SD}(\hat{\beta})$
کاکس	—	۰,۰۵۷۵	۰,۰۲۳۳
میلر	۲,۵۳۷	-۰,۰۰۵۸	—
تعمیم میلر	۲,۱۱۱	۰,۰۰۳۶	۰,۰۱۶۶
باکلی-جیمز	۲,۱۷۱	۰,۰۰۲۴	۰,۰۱۶۳
—	۳,۵۸۲	-۰,۰۲۷۸	۰,۰۱۴۹

فصل هفتم

نیکویی برازش

۱ روشهای ترسیمی

به راحتی می توان با نگاه کردن، خط را از منحنی تمیز داد. بنابراین برای روش رسم نمودار از اصل زیر استفاده می کنیم:

اصل اساسی. مقیاس محورهای مختصات را به گونه ای اختیار می کنیم، که اگر الگو برقرار باشد، نمودار به شکل خط مستقیم و در غیر این صورت به شکل منحنی در آید.

معمولاً، دو نوع نمودار بقاء و نرخ شکست مورد استفاده قرار می گیرند، این دو مورد بسیار نزدیک هم اند. و در هر حالت مناسبترین انتخاب می شود.

(الف) نمودارهای بقاء

در این نمودارها یا $\hat{S}(y_{u(i)})$ را در مقابل $y_{u(i)}$ و یا $\hat{S}(t)$ را در مقابل t رسم می کنند. این نوع، حالت خاصی از نمودار $Q-Q$ یا نمودارهای احتمالی است.

REFERENCE

Wilk and Gnanadesikan, *Biometrika* (1968).

(ب) نمودارهای نرخ شکست

در این نمودارها یا $\hat{\Lambda}(y_{u(i)})$ در مقابل $y_{u(i)}$ و یا $\hat{\Lambda}(t)$ در مقابل t رسم می شود. برای این کار از رابطه نلسون (فصل سوم، بخش سوم را ببینید) به صورت:

$$\hat{\Lambda}_r(t) = \sum_{y(i) \leq t} \frac{\delta(i)}{n-i+1} \quad \text{و یا رابطه: } \hat{\Lambda}_r(t) = -\log \hat{S}(t), \text{ استفاده می شود.}$$

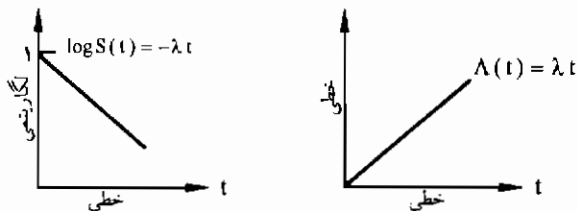
REFERENCES

- Nelson, J. Qual. Tech. (1969).
 _____, Technometrics (1972).

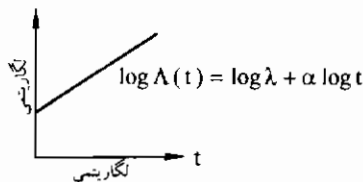
۱.۱ یک نمونه

نمودارهای چند توزیع در زیر رسم شده‌اند:

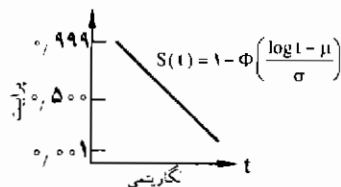
(الف) نمایی



(ب) وایبل



(ج) لوگ نرمال



(د) گاما و سایر توزیعها

بدون استفاده از کاغذهای نموداری، کمیتها بر مبنای فرضهای پارامتری در مقابل کمیتهای مبتنی بر برآوردگر PL، رسم می‌شود.

REFERENCES

- Wilk, Gnanadesikan, and Huyett, Technometrics (1962), for the gamma distribution without censoring.

۲.۱ دو نمونه تا K نمونه

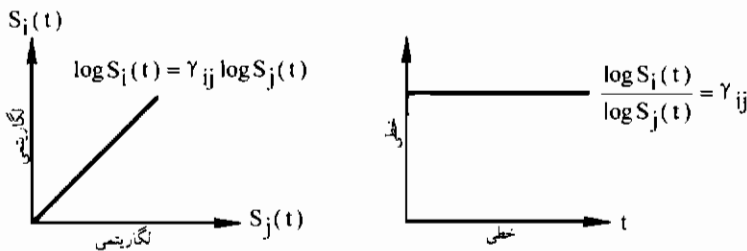
برای الگوهای پارامتری، موارد (الف) تا (د) را روی هر نمونه تکرار می‌کنیم. فرض کنید می‌خواهیم اعتبار الگوی نرخ شکست متناسب کاکس را بررسی کنیم. تحت الگوی $S_i(t) = S_j(t)^{\gamma_{ij}}$ برای چند مقدار از γ_{ij} ، داریم:

$$\log S_i(t) = \gamma_{ij} \log S_j(t)$$

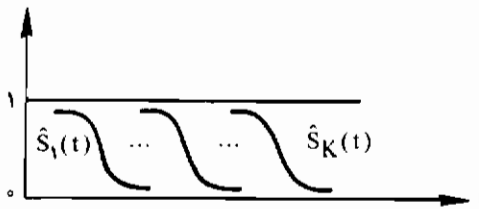
یا

$$\frac{\log S_i(t)}{\log S_j(t)} = \gamma_{ij}$$

برآوردهای جدای PL هر یک از $\hat{S}_1(t)$ تا $\hat{S}_K(t)$ را محاسبه و یکی از نمودارهای زیر را می‌سازیم:



برای واری الگوی خطی، برآورد PL را به طور جداگانه برای هر یک از K نمونه محاسبه و آن را رسم می‌کنیم. تغییر مکانها را توسط انتقال واری می‌کنیم.



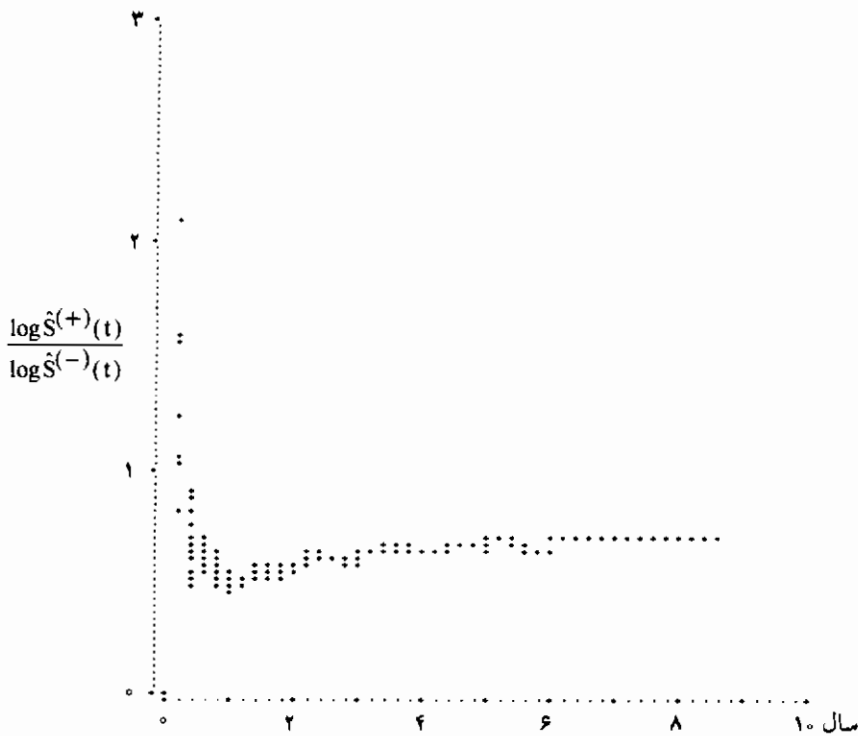
مثال: مطالعه DNCB. بیماران مرض هادکین، مورد حساسیت قرار گرفته و سپس به طور مستمر در معرض شیمی درمانی دینی تروکلروبنزن (DNCB) قرار گرفته‌اند. جامعه (+) شامل آن بیمارانی است که واکنش مثبت به DNCB نشان داده و جامعه (-) شامل آن بیمارانی است که واکنش نشان نداده‌اند. بیماران می‌توانند در میان جمعیتها نقل

مکان کنند، زمان بقاء را تا زمان جایگذاری در نظر می گیریم.

آیا بیماران جامعه (+) بیش از بیماران جامعه (-) عمر می کنند. الگوی نرخ شکست متناسب کاکس به کار می رود. رسم $\log \hat{S}^{(+)} / \log \hat{S}^{(-)}$ در نمودار شماره ۷، نشان می دهد که برای زمان t نزدیک به صفر، نسبت لگاریتم به طور مستدلی ثابت است، که نشان دهنده برقراری الگوست.

REFERENCE

Gong, Stanford Univ. Tech. Report No. 57 (1980).



نمودار ۷. نمایش $\log \hat{S}^{(+)}(t) / \log \hat{S}^{(-)}(t)$ توسط رایانه

۳.۱ رگرسیون

فرض کنید می خواهیم الگوی نرخ شکست متناسب را بررسی نماییم. در حالت یک بعدی، می توان محور x ها را به K بازه افزایش کرد. برآوردگر PL به طور جداگانه برای هر بازه محاسبه نموده و روش K نمونه را به کار برد. اگر x چند بعدی باشد،

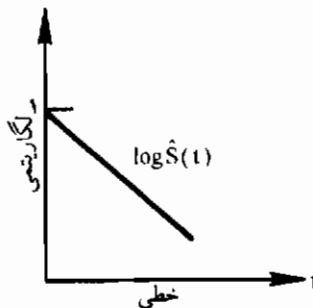
می توان فضای x را به K ناحیه افراز کرد. با این وجود، لازمه طبقه بندی داده ها، زیاد بودن آنهاست و تعداد لازم به سرعت با بعد x افزایش می یابد. یک راه دیگر برای طبقه بندی به شرح زیر است. ابتدا تعریف زیر ارائه می شود:

$$\Lambda_{\underline{x}_i}(T_i) = e^{-\underline{\beta}' \underline{x}_i} \int_0^{T_i} \lambda_0(u) du$$

در این صورت، تحت الگوی نرخ شکست متناسب، رابطه زیر نشان می دهد که $\Lambda_{\underline{x}_i}(T_i)$ یک متغیر نمایی واحد است.

$$P\{\Lambda_{\underline{x}_i}(T_i) > t\} = P\{T_i > \Lambda_{\underline{x}_i}^{-1}(t)\} = \exp\{-\Lambda_{\underline{x}_i}(\Lambda_{\underline{x}_i}^{-1}(t))\} = e^{-t}$$

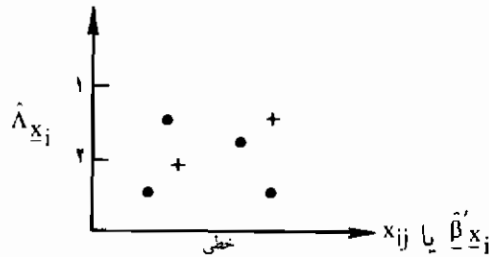
بنابراین، $(\Lambda_{\underline{x}_1}(Y_1), \delta_1)$ تا $(\Lambda_{\underline{x}_n}(Y_n), \delta_n)$ یک نمونه از توزیع نمایی واحد با برش است. چون $\Lambda_{\underline{x}_i}(Y_i)$ به پارامتر مجهول $\underline{\beta}$ و $\lambda_0(t)$ بستگی دارد، برآوردها را جانشین می کنیم. با تعریف $\hat{\Lambda}_i = \hat{\Lambda}_{\underline{x}_i}(Y_i) = e^{-\hat{\underline{\beta}}' \underline{x}_i} \int_0^{Y_i} \hat{\lambda}_0(u) du$ ، فرض کنید \hat{S} بر آوردگر PL بر مبنای $(\hat{\Lambda}_1, \delta_1)$ تا $(\hat{\Lambda}_n, \delta_n)$ ، تحت الگوی نرخ شکست متناسب باشد، باید $\log \hat{S}(t)$ تقریباً تابع خطی از t شود.



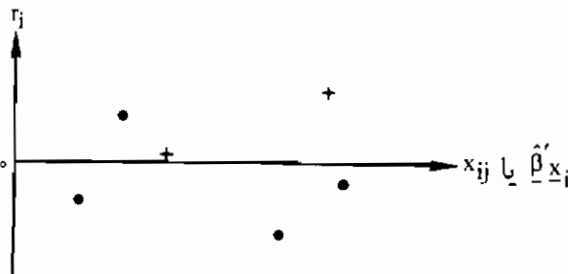
اگر نمودار $\log \hat{S}(t)$ در مقابل t خطی نباشد، امکان دارد تصمیم گیری در مورد الگوی مناسب دیگری مشکل باشد.

تحت الگوی نرخ شکست متناسب، $(\hat{\Lambda}_1, \delta_1)$ تا $(\hat{\Lambda}_n, \delta_n)$ ، مقادیر بریده شده (بریده نشده)، تقریباً متغیرهای تصادفی iid هستند. بنابراین، نباید رسم $\hat{\Lambda}_i(t)$ در مقابل یک متغیر کمکی خاص همانند: x_{ij} یا در مقابل $\hat{\underline{\beta}}' \underline{x}_i$ هیچ گونه طرح مجانبی را آشکار

سازد. برآوردهای $\hat{\Lambda}_1$ تا $\hat{\Lambda}_n$ را مانده‌های تعمیم یافته می‌نامند.



در واریس الگوی خطی، در صورت زیاد بودن تعداد مشاهدات، افزایش ناحیه x به K زیر ناحیه و استفاده از روش K نمونه پیشنهاد می‌شود. از طرف دیگر، می‌توان مانده‌های: $r_i = y_i - \hat{\beta}'_i x_i$ را در مقابل متغیر کمکی x_{ij} یا $\hat{\beta}'_i x_i$ رسم کرد.



برای هر دو مورد نرخ شکست متناسب و الگوی خطی، حساسیت نمودار مانده‌ها را رسم تا اثر یک متغیر وابسته x_{ij} را کشف کنیم. امکان دارد این عمل با محاسبه باقی‌مانده و حذف این متغیر وابسته در الگو انجام شود.

REFERENCES

Cox and Snell, JRSS B (1968), discuss generalized residuals.

Crowley and Hu, JASA (1977), plot generalized residuals for the Stanford heart transplant data.

Kay, Appl. Stat. (JRSS C) (1977), discusses plotting generalized residuals.

۲ آزمونها

۱.۲ یک نمونه

می‌خواهیم فرض: $H_0: F = F_0$ را با معلوم بودن F_0 ، آزمون کنیم.

(الف) تعمیم آزمون کلموگروف-اسمیرنوف: فرض H_0 را می‌پذیریم، اگر

$$\sqrt{n} |\hat{F}(t) - F_0(t)| \leq \hat{C}_n(t) \quad t \geq 0$$

که در آن، $\hat{F}(t)$ برآوردگر PL و $\hat{C}_n(t)$ را از جدول محاسبه کنیم. این آزمون را می‌توان برای ساختن فاصله اطمینان برای $F_0(t)$ نیز به کار برد، داریم:

$$P \left\{ \hat{F}(t) - \frac{\hat{C}_n(t)}{\sqrt{n}} \leq F_0(t) \leq \hat{F}(t) + \frac{\hat{C}_n(t)}{\sqrt{n}}, \forall t \geq 0 \right\} = 1 - \alpha$$

REFERENCES

- Barr and Davidson, *Technometrics* (1973), and
 Koziol and Byar, *Technometrics* (1975), and
 Dufour and Maag, *Technometrics* (1978), consider Type I and
 Type II censoring.
 Gillespie and Fisher, *Ann. Stat.* (1979), and
 Hall and Wellner, *Biometrika* (1980), consider the PL estimator
 and random censoring.

(ب) تعمیم آزمون کرامر-وون مایسز: بعد از انجام یک تبدیل احتمالی انتگرال به گونه‌ای که $F_0(t) = t$ باشد، تابع توزیع، یکنواخت شود. آزمون کرامر-وون مایسز از آماره زیر استفاده می‌کند.

$$n \int_0^1 (\hat{F}(t) - t)^2 dt$$

که در آن \hat{F} برآوردگر PL است.

REFERENCES

- Koziol and Green, *Biometrika* (1976), consider the PL estimator
 and random censoring.
 Pettit and Stephens, *Biometrika* (1976), consider Type I and
 Type II censoring. Pettit specializes to the normal
 and exponential distributions in
 Pettit, *Biometrika* (1976), and
 _____, *Biometrika* (1977), respectively.

(ج) آزمون نوع مانتل - هانزل

REFERENCE

Hyde, *Biometrika* (1977).

(د) حد آزمون افرون

REFERENCE

Hollander and Proschan, *Biometrics* (1979).

(ه) خانواده‌های پارامتری

فرض کنید می‌خواهیم فرض: $\theta \in \Theta_0$ ، $H_0: F = F_{\theta}$ را آزمون کنیم. به روش معمول یک برآورد $\hat{\theta}$ را به دست می‌آوریم. سپس، واری می‌کنیم که آیا \hat{F} به اندازه کافی به $F_{\hat{\theta}}$ نزدیک است یا خیر.

REFERENCE

Mihalko and Moore, *Ann. Stat.* (1980), consider χ^2 - tests for Type II censoring with estimates that are asymptotically equivalent to linear combinations of order statistics.

اگر $\Theta_0 \subset \Theta$ و بخواهیم $H_0: \theta \in \Theta_0$ را آزمون کنیم، باید آزمون نسبت درست‌نمایی را به کار ببریم.

REFERENCE

Turnbull and Weiss, *Biometrics* (1978), consider likelihood ratio tests for discrete or grouped data.

۲.۲ رگرسیون

(الف) خانواده‌های پارامتری

الگو را در یک الگوی بزرگتر می‌نشانیم (مثلاً، الگویی که دارای اثرات درجه ۲ یا ۳ اثرات متقابل است). آزمون می‌کنیم که آیا الگوی کوچکتر برقرار است؟ در واقع فرض زیر را آزمون می‌کنیم

$$H_0: \theta \in \Theta_0 \subset \Theta$$

(ب) آزمونهای χ^2

REFERENCES

Schoenfeld, *Biometrika* (1980), considers proportional hazards models with regions in the time \times covariate space.

Lamborn, Stanford Univ. Tech. Report No. 21 (1969), looks at χ^2 – tests for exponential regression.

فصل هشتم

مباحث مختلف

۱ برآوردگر دو متغیری کاپلان-مایر

فرض کنید، $\underline{T}_j = (T_{j1}, T_{j2})$ ، یک زوج از زمانهای شکست باشد. برای مثال، ممکن است، زمانهای از کار افتادن کلیه‌های راست و چپ یا زمانهای تشخیص سرطان در شش راست و چپ باشد. امکان دارد، هر کدام یا هر دو زمان شکست به علت متغیر تصادفی بریده شده بعدی C_j ، قابل مشاهده و بردار نشانگر به شرح زیر است:

$$\underline{Y}_j = (Y_{j1}, Y_{j2}) = (T_{j1} \wedge C_j, T_{j2} \wedge C_j)$$

$$\underline{\delta}_j = (\delta_{j1}, \delta_{j2}) = (I(T_{j1} \leq C_j), I(T_{j2} \leq C_j))$$

می‌نویز نشان داده است که چگونه می‌توان برآوردگر تعمیم یافته کاپلان-مایر دوبعدی را به کمک الگوریتم خودسازگاری و تجدید نظر در توزیع راست، محاسبه کرد. همچنین، او ثابت نمود که این برآوردگر، تعمیم برآورد حداکثر درست‌نمایی و یک برآوردگر سازگار توزیع دو متغیری $\{T_{j1} \leq t_1, T_{j2} \leq t_2\}$ است. $F(t_1, t_2)$ است. کامل الگو را با زمانهای برش دو متغیری و رفتار شرایط داده‌های طبقه‌بندی شده در نظر گرفته است، همچنین، کوروار داده‌های طبقه‌بندی شده دو متغیره با هر دو نوع برش راست و چپ را بررسی کرده است.

REFERENCES

Campbell, Purdue univ. Mimeoseries #79 – 25 (1979),

and

Korwar, unpublished manuscript (1980), treat bivariate grouped

data with censoring.

Muñoz, Stanford Univ. Tech. Report No. 60 (1980), defines the two – dimensional KM estimator through algorithms and proves it is the GMLE.

_____, Stanford Univ. Tech. Report No. 61 (1980), proves consistency of the two – dimensional estimator.

۲ نرخ شکست رقیب

فرض کنید $T_i = (T_{i1}, \dots, T_{ip})$ ، یک بردار p بعدی از زمانهای شکست باشد. هر مختص، زمان شکست حاصل از علت خاصی مانند: از کار افتادن قلب، سرطان، خرابی کبد و غیره است. موضوع فقط هنگامی قابل مشاهده است که اولین شکست رخ دهد: زمانهای شکست تمام عوامل دیگر، توسط از کار افتادگی دستگاه در اولین زمان شکست، بریده می‌شوند. کمتهای قابل مشاهده به صورت زیر است:

$$T_i = \min \{T_{i1}, \dots, T_{ip}\}$$

و

$$\delta_i = (\delta_{i1}, \dots, \delta_{ip}) = (I(T_{i1} \leq T_i), \dots, I(T_{ip} \leq T_i))$$

بردار نشانگر δ_i ، علت خاص شکست را نشان می‌دهد.

احتمال $\{T_{ij} \leq t, \delta_{ij} = 1\}$ را احتمال خام مردن به علت z در زمان t می‌نامند. این احتمال، مستقیماً توسط نسبت مشاهده‌ای، برآورد می‌شود.

$$\frac{1}{n} \sum_{i=1}^n I(T_i \leq t, \delta_{ij} = 1)$$

عَلت خام، عبارت از $P(T_{ij} \leq t)$ و اگر علتها مستقل باشند، آن را می‌توان به طور سازگار به وسیله روش PL برآورد کرد، که در آن تمام زمانهای شکست را به علت z ، به علت یکی از زیر مجموعه‌های ممکن علتها در نظر می‌گیرد.

یک نتیجه اصلی در نظریه نرخهای شکست رقیب این است که بر مبنای نمونه T_i و δ_i ، $i=1, \dots, n$ ، نمی‌توان مدعی شد که این رابطه

$$P\{T_{i1} \leq t_1, \dots, T_{ip} \leq t_p\} = \prod_{j=1}^p P\{T_{ij} \leq t_j\}$$

برقرار است یا زمانهای شکست T_{i1}, \dots, T_{ip} وابسته‌اند. اثباتهای مختلف این نتیجه با شرایط متفاوت، سالهای مورد بحث بوده است. به مقالات برمن، آلت‌شولر، تسیاتیس، پترسن، لانگبرگ - پروچان - کوینزی نگاه کنید.

REFERENCES

- Chiang, Introduction to Stochastic Processes in Biostatistics (1968), discusses the relationships between crude, net, and partial crude probabilities in Chapter 11.
- Moeschberger and David, Biometrics (1971), consider parametric likelihood methods.
- Gail, Biometrics (1975), is a review article.
- Prentice et al. , Biometrics (1978), review competing risks from the hazard rate point of view.
- Berman, Ann. Math. Stat. (1963),
- Altshuler, Mathematical Biosciences (1970),
- Tsiatis, Proc. Natl. Acad. Sci. (1975),
- Peterson, Stanford Univ. Tech. Report No. 13 (1975),
- _____, Proc. Natl. Acad. Sci. (1976), and
- Langberg Proshan, and Quinzi, Ann. Stat. (1981), examine the identifiability question.

۳ برش وابسته

برای حالتی که زمانهای شکست و زمانهای برش وابسته‌اند، کار زیادی صورت نگرفته است. بعضی از کارهای انجام شده در نرخهای شکست رقیب وابسته در این خصوص انجام شده است. در مقالات لاگاکوس و ویلیامز بحث‌های کلی و نتایج ارائه شده است.

REFERENCES

- Williams and Lagakos, Biometrika (1977).
- Lagakos and Williams, Biometrika (1978).
- Lagakos, Biometrics (1979).

۴ روش جک‌نایف و بوت‌استرپ

فرض کنید پارامتر θ یک تابعی $T(F)$ از تابع توزیع F باشد. در به یاری از موارد θ توسط تابع توزیع نمونه F_n ، به جای F در T ، به صورت $\hat{\theta} = T(F_n)$ ، برآورد می‌شود. برآورد جک‌نایف θ برای یک نمونه Y_1 تا Y_n ، با توزیع F به صورت زیر تعریف می‌شود:

$$\tilde{\theta}_i = n\hat{\theta} - (n-1)\hat{\theta}_{-i}, \quad i=1, \dots, n$$

$$\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \tilde{\theta}_i = n\hat{\theta} - \frac{(n-1)}{n} \sum_{i=1}^n \hat{\theta}_{-i}$$

که در آن $\tilde{\theta}_{-i} = T(F_{n-1,-i})$ برآورد θ با این فرض که نمونه i ام یعنی Y_i در نمونه حذف شده، می‌باشد.

در صورت نبودن برش، رابطه زیر برای T به اندازه کافی هموار قابل اثبات است.

$$\frac{\tilde{\theta} - \theta}{\sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (\tilde{\theta}_i - \tilde{\theta})^2}} \stackrel{a}{\sim} N(0, 1) \quad (18)$$

برای بررسی (18) در شرایط مختلف به مقاله تجدید نظر شده میلر (1974) مراجعه شود. میلر نشان داده که (18) برای داده‌هایی که به طور تصادفی بریده شده - با این فرض که $\hat{\theta} = T(\hat{F})$ و \hat{F} برآوردگر PL است - برقرار است.

هموار بودن T که در (18) استفاده می‌شود، با هموار بودن تابع تأثیر زیر در ارتباط است.

$$IC(y; F) = \lim_{\varepsilon \rightarrow 0} \frac{T((1-\varepsilon)F + \varepsilon\delta_y) - T(F)}{\varepsilon}$$

که در آن δ_y ، تابع توزیعی است که به نقطه y جرم واحد را نسبت می‌دهد. برای داده‌های بریده نشده، تابع جک‌نایف و تابع تأثیر به شکل زیر با هم در ارتباط اند.

$$(n-1)(\hat{\theta} - \hat{\theta}_{-i}) = \frac{T((1-\varepsilon)F + \varepsilon\delta_y) - T(F)}{\varepsilon} \Big|_{\varepsilon = -\frac{1}{n-1}, F = \hat{F}_n, y = y_i}$$

راید برای داده‌های بریده شده، توابع تأثیر (یعنی مشتقات جزئی نسبت به F_u و F_c) را برای تابعی از برآوردگر PL، پیدا نموده است.

بوت استرپ افرون، به صورت زیر انجام می‌شود. فرض کنید Y_1^* تا Y_n^* ، یک نمونه با جایگذاری از y_1 تا y_n ، در حالت بدون برش باشد. در صورت وجود برش، فرض کنید (Y_1^*, δ_1^*) تا (Y_n^*, δ_n^*) ، نمونه‌ای با جایگذاری از (y_1, δ_1) تا (y_n, δ_n) باشد. در این صورت F_n^* و \hat{F}^* ، به ترتیب: توزیع نمونه بوت استرپ و برآوردگر PL هستند. همچنین، $\hat{\theta}^* = T(F_n^*)$ یا $T(\hat{F}^*)$ ، بستگی به این دارد که داده‌ها بریده شده یا بریده نشده باشند. این روش نمونه‌گیری N بار تکرار می‌شود تا $\hat{\theta}_1^*, \dots, \hat{\theta}_N^*$ ، حاصل شود. توزیع تجربی $\hat{\theta}_1^*$ تا $\hat{\theta}_N^*$ ، برای محاسبه توزیع تقریبی $\hat{\theta}$ به کار می‌رود. به ویژه کمیتهای محوری توزیع تجربی $\hat{\theta}^* - \hat{\theta}$ برای تقریب توزیع $\hat{\theta} - \theta$ ، مورد استفاده قرار می‌گیرد.

REFERENCES

- Miller, *Biometrika* (1974), reviews the jackknife for uncensored data problems.
 _____, Stanford Univ. Tech. Report No. ۱۴ (۱۹۷۵), establishes the validity of
 jackknifing the PL estimator.
 Reid, *Ann. Stat.* (۱۹۸۱), derives the influence functions for the PL estimator.
 Efron, *Ann. Stat.* (۱۹۷۹), introduces bootstrapping for uncensored data problems.
 _____, Stanford Univ. Tech. Report No. ۵۳ (۱۹۸۰), studies bootstrapping

فصل نهم

مسائل

مسألة ۱.

ثابت کنید توزیع گاما در ازای $\alpha > 1$ دارای IFR و در ازای $\alpha < 1$ دارای DFR است.

حل:

$$\frac{1}{\lambda(t)} = \frac{\int_t^{\infty} x^{\alpha-1} e^{-\lambda x} dx}{t^{\alpha-1} e^{-\lambda t}} = \int_t^{\infty} \left(\frac{x}{t}\right)^{\alpha-1} e^{-\lambda(x-t)} dx$$

$$= \int_0^{\infty} \left(1 + \frac{u}{t}\right)^{\alpha-1} e^{-\lambda u} du \quad \text{با تغییر متغیر } u = x - t$$

اگر $\alpha > 1$ باشد، عبارت: $\left(1 + \frac{u}{t}\right)^{\alpha-1}$ بر حسب t نزولی و $\lambda(t)$ صعودی خواهد بود.
اگر $\alpha < 1$ باشد، عبارت بالا بر حسب t صعودی و $\lambda(t)$ نزولی خواهد بود.

مسألة ۲.

تابع اطلاع فیشر را برای یک مشاهده در توزیع نمایی با برش نوع اول به دست آورید.

حل:

فرض کنید t_c زمان برش، ثابت باشد. لگاریتم درست نمایی به شرح زیر است:

$$\delta \log \lambda - \delta \lambda y - (1 - \delta) \lambda t_c$$

اگر نسبت به λ دوبار مشتق بگیریم، $-\frac{\delta}{\lambda^2}$ به دست می‌آید. در نتیجه اطلاع فیشر به شرح زیر خواهد بود:

$$I(\lambda) = \frac{1}{\lambda^2} E(\delta) = \frac{1}{\lambda^2} P\{T \leq t_c\} = \frac{1}{\lambda^2} (1 - e^{-\lambda t_c})$$

مسألة ۳.

ماتریس اطلاع نمونه را برای توزیع وایبل تحت برش تصادفی به دست آورید.

حل:

از معادله شماره ۲ فصل دوم، داریم:

$$\frac{\partial}{\partial \gamma} \log L = \frac{n_u}{\gamma} - \sum_{i=1}^n y_i^\alpha$$

$$\frac{\partial}{\partial \alpha} \log L = \frac{n_u}{\alpha} + \sum_u \log t_i - \gamma \sum_{i=1}^n y_i^\alpha \log y_i$$

ماتریس اطلاع نمونه در (γ, α) ، به صورت زیر است:

$$-\begin{pmatrix} \frac{\partial^2}{\partial \gamma^2} \log L & \frac{\partial^2}{\partial \gamma \partial \alpha} \log L \\ \frac{\partial^2}{\partial \alpha^2} \log L \end{pmatrix}$$

مسألة ۴.

از فوریه ۱۹۷۲ تا فوریه ۱۹۷۵، تعداد ۲۹ بیمار در دو بیمارستان تحت معالجه قرار گرفته و به تصادف به دو گروه تیمار و شاهد تقسیم شده‌اند. زمانهای بقاء (برحسب هفته) چهارده بیمار در گروه تیمار، به شرح زیر است:

16^+ و 16^+ و 16^+ و 12^+ و 10^+ و 10^+ و 8 و 7 و 5 و 4^+ و 1^+ و 1 و 1

توزیع نمایی $S(t) = \exp(-\lambda t)$ را فرض کنید و به سؤالات زیر پاسخ دهید.

(الف) مقدار λ را به روش حداکثر درست‌نمایی برآورد کنید و یک فاصله اطمینان ۹۵٪ را برای آن پیدا کنید.

(ب) مقدار $S(16)$ را برآورد کنید و فاصله اطمینان ۹۵٪ را برای آن پیدا کنید.

(پ) میانهٔ زمان بقاء را برآورد کنید. فاصلهٔ اطمینان ۹۵٪ را برای آن پیدا کنید.

حل:

(الف) از مثال شماره ۱ فصل دوم بخش (۱.۲)، داریم:

$$\hat{\lambda} = \frac{n_u}{\sum_{i=1}^n y_i} = \frac{7}{10.8} = 0.65 \quad \text{و} \quad \log \hat{\lambda} \stackrel{a}{\sim} N\left(\log \lambda, \frac{1}{n_u}\right)$$

در نتیجه فاصلهٔ اطمینان ۹۵٪ برای λ ، به شرح زیر است:

$$\left(\hat{\lambda} \exp\left(\frac{-Z_{0.025}}{\sqrt{n_u}}\right), \hat{\lambda} \exp\left(\frac{Z_{0.025}}{\sqrt{n_u}}\right) \right) = (0.31, 0.136)$$

(ب) داریم: $\hat{S}(16) = \exp(-\hat{\lambda} \times 16) = 0.355$. یک فاصلهٔ اطمینان ۹۵٪ برای $S(16)$ به

شرح زیر است:

$$\left(e^{-0.136 \times 16}, e^{-0.31 \times 16} \right) = (0.113, 0.609)$$

(پ) داریم: $\hat{t}_{\text{med}} = \log(2)/\hat{\lambda} = 1.069$. یک فاصلهٔ اطمینان ۹۵٪ برای میانه به شرح زیر

است:

$$\left(\frac{\log 2}{0.136}, \frac{\log 2}{0.31} \right) = (5.097, 22.36)$$

مسألة ۵.

برای داده‌های بیمارستانی مسألة شماره ۴، برآورد حدّ حاصل ضرب کاپلان-مایر، تابع بقاء را به دست آورید. نمودار آن و نمودار تابع بقاء را تحت فرض نمایی روی همان محور مختصات خطی-لگاریتمی رسم کنید. آیا فکر می‌کنید فرض نمایی در مدت شانزده هفته صادق است.

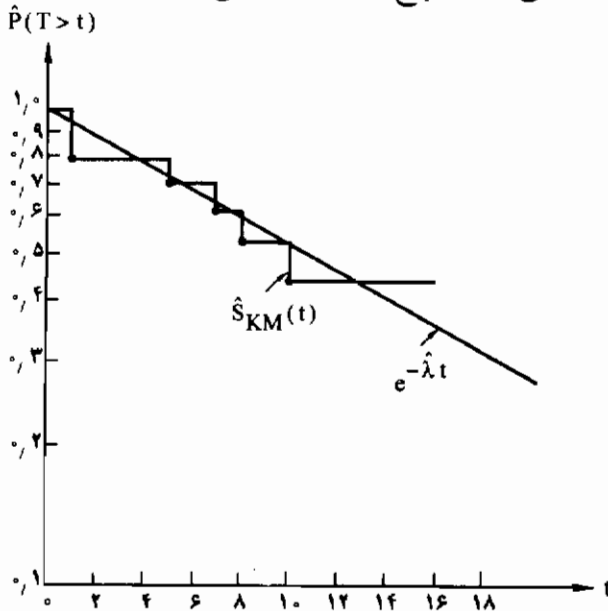
حل:

داریم: (شکل زیر را نیز ببینید)

$$\hat{S}(t) = \begin{cases} 1 & 0 \leq t < 1 \\ 11/14 = 0.786 & 1 \leq t < 5 \\ 11 \times 8/14 \times 9 = 0.698 & 5 \leq t < 7 \end{cases}$$

$$\hat{S}(t) = \begin{cases} 11 \times 7 / 14 \times 9 = 0,611 & 10 \leq t < 16 \\ 11 \times 6 / 14 \times 9 = 0,524 & 8 \leq t < 10 \\ 11 \times 5 / 14 \times 9 = 0,437 & 7 \leq t < 8 \end{cases}$$

نتایج عادلانه نیکویی برازش عبارت است از: پانزده نمونه از بیست و یک نمونه رسم شده توزیع نمایی را تأیید می کنند و پنج نمونه تأیید نمی کنند و یک مورد جواب نمی دهد.



مسأله ۶.

از جدول طول عمر (جدول شماره ۱) خطای استاندارد $S(5)$ را حساب کنید.

حل:

با استفاده از رابطه گزین وود، داریم:

$$\begin{aligned} \widehat{\text{Var}}(\hat{S}(5)) \cong (0,44)^2 & \left[\frac{47}{116,5(116,5-47)} + \frac{5}{51,5(51,5-5)} + \frac{2}{30,5(30,5-5)} + \right. \\ & \left. + \frac{2}{16,5(16,5-2)} + \frac{0}{7(7-0)} \right] = 0,003608 \end{aligned}$$

بنابراین $\widehat{SE}(\hat{S}(5)) \cong 0,06$

مسئله ۷.

به کمک داده‌های AML (بخش دوم از فصل سوم) خطای استاندارد $\hat{S}(۲۴)$ را در گروه تیمار حساب کنید.

حل:

با استفاده از رابطه گرین‌وود، داریم:

$$\widehat{\text{Var}}(\hat{S}(۲۴)) = \left(\frac{۶ \times ۹}{۱۱ \times ۸}\right)^2 \left(\frac{۱}{۱۰ \times ۱۱} + \frac{۱}{۹ \times ۱۰} + \frac{۱}{۷ \times ۸} + \frac{۱}{۶ \times ۷}\right) = ۰٫۰۲۳۲۹$$

$$\widehat{\text{SE}}(\hat{S}(۲۴)) = ۰٫۱۵۲۶ \text{ در نتیجه}$$

مسئله ۸.

در اثبات GMLE برآوردگر PL، نشان دهید حداکثر عبارت:

$$\prod_{i=1}^n p_i^{\delta(i)} \left(\sum_{j=i}^n p_j \right)^{1-\delta(i)}$$

به ازای مقدار زیر به دست می‌آید:

$$p_i = \frac{\delta(i)}{n-i+1} \prod_{j=1}^{i-1} \left(1 - \frac{\delta(j)}{n-j+1} \right)$$

حل:

فرض کنید: $\lambda_i = \frac{p_i}{\sum_{j=i}^n p_j}$ ، $i=1, \dots, n$ در این صورت چون $\sum_{j=1}^n p_j = 1$

و $1 - \lambda_i = \frac{\sum_{j=i+1}^n p_j}{\sum_{j=i}^n p_j}$ داریم: $\sum_{j=i}^n p_j = \prod_{j=1}^{i-1} (1 - \lambda_j)$ و چون $\lambda_n = 1$ است، در

نتیجه داریم:

$$\prod_{i=1}^n p_i^{\delta(i)} \left(\sum_{j=i}^n p_j \right)^{1-\delta(i)} = \prod_{i=1}^n \lambda_i^{\delta(i)} \prod_{j=1}^{i-1} (1 - \lambda_j) = \prod_{i=1}^{n-1} \lambda_i^{\delta(i)} (1 - \lambda_i)^{n-i}$$

با توجه به نظریه نمونه‌گیری دو جمله‌ای، هر حاصل ضرب به

ازای: $\hat{\lambda}_i = \frac{\delta(i)}{n-i+\delta(i)} = \frac{\delta(i)}{n-i+1}$ ، حداکثر می‌شود. بنابراین، داریم:

$$\hat{p}_i = \hat{\lambda}_i \left(\sum_{j=i}^n \hat{p}_j \right) = \hat{\lambda}_i \prod_{j=1}^{i-1} (1 - \hat{\lambda}_j) = \frac{\delta(i)}{n-i+1} \prod_{j=1}^{i-1} \left(1 - \frac{\delta(j)}{n-j+1} \right)$$

مسأله ۹.

ثابت کنید که الگوریتم تجدید نظر در توزیع به راست برآوردگر حد حاصل ضرب کاپلان-مایر را بدون تکرار می‌دهد.

حل:

دو راه اصلی برای اثبات این نتیجه وجود دارد:

(الف) با توجه به الگوریتم بالا، کلیه نقاط $y(i)$ در هر دو حالت بریده شده و بریده نشده، در ابتدا دارای جرم $\frac{1}{n}$ هستند. این الگوریتم از چپ به راست در میان آماره‌های مرتب حرکت می‌کند. هنگامی که به $y(i)$ می‌رسد، تمام نقاط باقی‌مانده، $y(i)$ ، $y(i+1)$ و ... و $y(n)$ ، دارای جرم مساوی، با توجه به روش انجام الگوریتم، هستند. فرض کنید جرم کل باقی‌مانده $\tilde{S}(y(i)-)$ باشد. با توجه به مساوی بودن جرمها، $y(i)$ دارای جرم $\tilde{S}(y(i)-)/(n-i+1)$ است. اگر بریده نشده باشد این جرم حفظ می‌شود. اگر بریده شده باشد، این جرم در سمت راست توزیع، توزیع می‌شود.

چون برآوردگر \hat{S} ، \hat{S} ، مانند \tilde{S} از یک شروع می‌شود و اندازه جهشهای آن در مشاهدات بریده نشده به صورت: $\hat{S}(y(i)-)/(n-i+1)$ و در حالت بریده شده صفر است، لذا، دو برآوردگر برابرند.

(ب) برای برآوردگر کاپلان-مایر، داریم:

$$\begin{aligned} \hat{\Delta}(i) &= \hat{S}(y(i)-) - \hat{S}(y(i)) = \prod_{j=1}^{i-1} \left(\frac{n-j}{n-j+1} \right)^{\delta(j)} - \prod_{j=1}^i \left(\frac{n-j}{n-j+1} \right)^{\delta(j)} \\ &= \prod_{j=1}^{i-1} \left(\frac{n-j}{n-j+1} \right)^{\delta(j)} \frac{\delta(i)}{n-i+1} = \prod_{j=1}^{i-1} \left(\frac{n-j+1}{n-j} \right)^{-\delta(j)} \times \frac{1}{n} \times \frac{n}{n-1} \times \dots \times \\ &\quad \times \frac{n-i+2}{1} \times \frac{\delta(i)}{n-i+1} = \frac{\delta(i)}{n} \prod_{j=1}^{i-1} \left(\frac{n-j+1}{n-j} \right)^{1-\delta(j)} \end{aligned}$$

فرض کنید، $j_1 < \dots < j_r$ ، زیرنویسهای مشاهدات بریده شده قبل از $y(i)$ باشند. برای الگوریتم تجدید نظر در توزیع به راست، جرم نسبت داده شده به $y(i)$ با شرط $\delta(i) = 1$ برابر زیر است:

$$\tilde{\Delta}(i) = \frac{1}{n} \left(1 + \frac{1}{n-j_1}\right) \left(1 + \frac{1}{n-j_2}\right) \dots \left(1 + \frac{1}{n-j_r}\right) = \frac{1}{n} \prod_{j=1}^{i-1} \left(\frac{n-j+1}{n-j}\right)^{1-\delta(j)}$$

و اگر $\delta(i) = 0$ باشد، $\tilde{\Delta}(i) = 0$ می شود. این معادل $\hat{\Delta}(i)$ است. بنابراین، الگوریتم تجدید نظر در توزیع به راست، برآوردگر PL را می دهد.

مسألة ۱۰.

اگر برای برآوردگر PL، $\hat{S}(t)$ ، رابطه:

$$ACov(\hat{S}(t_1), \hat{S}(t_2)) = \frac{S(t_1) S(t_2)}{n} \times \int_0^{t_1 \wedge t_2} \frac{dF_U(s)}{(1-H(s))^2}$$

را داشته باشیم، نشان دهید که برای $\hat{\mu} = \int_0^\infty \hat{S}(t) dt$ رابطه زیر برقرار است:

$$AVar(\hat{\mu}) = \frac{1}{n} \int_0^\infty \frac{1}{(1-H(s))^2} \left(\int_s^\infty S(t) dt \right)^2 dF_U(s)$$

که در آن $ACov$ و $AVar$ ، به معنی واریانس و کوواریانس مجانبی است.

حل:

$$Var(\hat{\mu}) = E(\hat{\mu}^2) - (E(\hat{\mu}))^2$$

$$= E \left(\int_0^\infty \int_0^\infty \hat{S}(t_1) \cdot \hat{S}(t_2) dt_1 dt_2 \right) - \left(E \left(\int_0^\infty \hat{S}(t) dt \right) \right)^2$$

$$= \int_0^\infty \int_0^\infty Cov(\hat{S}(t_1), \hat{S}(t_2)) dt_1 dt_2$$

بنابراین، داریم:

$$AVar(\hat{\mu}) = \frac{1}{n} \int_0^\infty \int_0^\infty S(t_1) \cdot S(t_2) \times \int_0^{t_1 \wedge t_2} \frac{dF_U(s)}{(1-H(s))^2} dt_1 dt_2$$

با تعویض انتگرال (قضیه فوینی)، داریم:

$$\begin{aligned} A\text{Var}(\hat{\mu}) &= \frac{1}{n} \int_0^{\infty} \frac{1}{(1-H(s))^2} \int_s^{\infty} S(t_1) dt_1 \times \int_s^{\infty} S(t_2) dt_2 dF_U(s) \\ &= \frac{1}{n} \int_0^{\infty} \frac{1}{(1-H(s))^2} \left(\int_s^{\infty} S(t) dt \right)^2 dF_U(s) \end{aligned}$$

مسأله ۱۱.

برای داده‌های AML در گروه شاهد، که به شکل زیر هستند، (الف) $\hat{\mu}$ و (ب) $\text{Var}(\hat{\mu})$ را محاسبه کنید.

برحسب هفته ۴۵ و ۴۳ و ۳۳ و ۳۰ و ۲۷ و ۲۳ و ۱۶⁺ و ۱۲ و ۸ و ۸ و ۵ و ۵

حل:

(الف) برآوردگر کاپلان-مایر $\hat{S}(t)$ در جدول زیر داده شده است:

$t \in [0, 5) [5, 8) [8, 12) [12, 23) [23, 27) [27, 30) [30, 33) [33, 43) [43, 45) [45, \infty)$

$\hat{S}(t) =$	۱	$\frac{10}{12}$	$\frac{8}{12}$	$\frac{7}{12}$	$\frac{7 \times 5}{12 \times 6}$	$\frac{7 \times 4}{12 \times 6}$	$\frac{7 \times 3}{12 \times 6}$	$\frac{7 \times 2}{12 \times 6}$	$\frac{7 \times 1}{12 \times 6}$.
----------------	---	-----------------	----------------	----------------	----------------------------------	----------------------------------	----------------------------------	----------------------------------	----------------------------------	---

در این صورت، داریم:

$$\begin{aligned} \hat{\mu} &= \int_0^{\infty} \hat{S}(t) dt = 1 \times 5 + \frac{10}{12} \times 3 + \frac{8}{12} \times 4 + \frac{7}{12} \times 11 + \frac{7}{12} \times \frac{5}{6} \times 4 + \frac{7}{12} \times \frac{4}{6} \times 3 + \\ &+ \frac{7}{12} \times \frac{3}{6} \times 3 + \frac{7}{12} \times \frac{2}{6} \times 10 + \frac{7}{12} \times \frac{1}{6} \times 2 = 22,71 \end{aligned}$$

(ب)

$$\begin{aligned} \hat{\text{Var}}(\hat{\mu}) &= \sum_u \left(\int_{y(i)}^{\infty} \hat{S}(t) dt \right)^2 \frac{d_i}{n_i(n_i - d_i)} \\ &= (17,71)^2 \frac{2}{12 \times 10} + (15,21)^2 \frac{2}{10 \times 8} + (12,54)^2 \frac{1}{8 \times 7} + (6,125)^2 \frac{1}{6 \times 5} + \\ &+ (4,18)^2 \frac{1}{5 \times 4} + (3,01)^2 \frac{1}{4 \times 3} + (2,14)^2 \frac{1}{3 \times 2} + (0,19)^2 \frac{1}{2 \times 1} = 17,47 \end{aligned}$$

مسألة ۱۲.

برای داده‌های هیاتیت (مسألة ۴)، تیمار I و شاهد II، زمانهای بقاء گروهها برحسب هفته به شرح زیراند. با $m=14$ و $n=15$ ، مطلوب است محاسبه (الف) آماره گهان، (ب) جایگشت آن و (پ) آماره استاندارد و مقدار P .

حل:

امتیازهای U^* در جدول زیر محاسبه شده‌اند:

Z	گروه	$\# < Z$	$\# > Z$	U^*
۱(۳)	I	۰	۲۶	-۲۶(۳)
۱ ⁺	I	۳	۰	۳
۱ ⁺	II	۳	۰	۳
۲ ⁺	II	۳	۰	۳
۳(۲)	II	۳	۲۱	-۱۸(۲)
۳ ⁺	II	۵	۰	۵
۴ ⁺	I	۵	۰	۵
۵	I	۵	۱۸	-۱۳
۵+(۲)	II	۶	۰	۶(۲)
۷	I	۶	۱۵	-۹
۸	I	۷	۱۴	-۷
۱۰	I	۸	۱۳	-۵
۱۰ ⁺	I	۹	۰	۹
۱۲ ⁺	I	۹	۰	۹
۱۶+(۳)	I	۹	۰	۹(۳)
۱۶+(۸)	II	۹	۰	۹(۸)

(الف) آماره گهان برابر $\sum_{II} U^* = 59$ می شود.

(ب) واریانس جایگشت به شرح زیر است:

$$\frac{14 \times 15}{29 \times 28} \sum_{I, II} (U^*)^2 = 1086,72$$

(پ) آماره استاندارد برابر $\frac{59}{\sqrt{1086,72}} = 1,79$ به دست می آید، که متناظر با P یک طرفه ۰,۰۳۷۵ است.

REFERENCE

Gregory et al. , New England Journal of Medicine (1976).

مسأله ۱۳.

برای داده های بیمارستانی (مسأله ۱۲)، موارد زیر را به دست آورید:

(الف) آماره MH و مقدار P متناظر آن.

(ب) صورت تارون-وایر آماره گهان و مقدار متناظر P.

حل:

محاسبات مانند جدول شماره ۴ از فصل چهارم، بخش دوم، انجام شده است.

z	n	m ₁	n ₁	a	E _o (A)	n(a - E _o (A))	$\frac{m_1(n - m_1)}{n - 1}$	$\frac{n_1}{n}(1 - \frac{n_1}{n})$	n ₁ (n - n ₁)
۱	۲۹	۳	۱۴	۳	۱,۴۴۸	۴۵	۲,۷۸۶	۰,۲۴۹۷	۲۱۰
۳	۲۳	۲	۱۰	۰	۰,۸۶۹	-۲۰	۱,۹۰۹	۰,۲۴۵۷	۱۳۰
۵	۱۹	۱	۹	۱	۰,۴۷۴	۱۰	۱	۰,۲۴۹۳	۹۰
۷	۱۶	۱	۸	۱	۰,۵۰۰	۸	۱	۰,۲۵۰۰	۶۴
۸	۱۵	۱	۷	۱	۰,۴۶۷	۸	۱	۰,۲۴۸۹	۵۶
۱۰	۱۴	۱	۶	۱	۰,۴۲۸	۸	۱	۰,۲۴۴۹	۴۸

(الف)

$$MH = \frac{\text{مجموع ستون } \{a - E_o(A)\}}{\sqrt{\left\{ \frac{m_1(n-m_1)}{n-1} \text{ ستون} \times \frac{n_1}{n} \left(1 - \frac{n_1}{n}\right) \text{ ستون} \right\} \text{مجموع}}} = \frac{2,814}{\sqrt{2,1578}} = 1,916$$

بنابراین، مقدار P یک طرفه متناظر برابر ۰٫۲۷ است.

(ب) فرض کنید U_{TW} صورت تارون-وایر آماره گهان باشد، داریم:

$$U_{TW} = \frac{\text{مجموع ستون } \{n(a - E_o(A))\}}{\sqrt{\left\{ \frac{m_1(n-m_1)}{n-1} \text{ ستون} \times n_1(n-n_1) \text{ ستون} \right\} \text{مجموع}}} = \frac{59}{\sqrt{1091,23}} = 1,786$$

در نتیجه مقدار P یک طرفه برابر ۰٫۳۷ است.

مسأله ۱۴.

پنج نقطه داده‌های مربوط به تعویض قلب بیمارستان استانفورد را در نظر بگیرید (مثال فصل ششم، بخش چهارم)

(الف) فرض $H_o: \beta = 1$ را در الگوی نرخ شکست متناسب با محاسبه مقدار P به وسیله آماره کاکس را، که به صورت زیر است، آزمون کنید.

$$\left(\frac{\partial}{\partial \beta} \log L_c(t) \right)^2 \Bigg/ - \frac{\partial^2}{\partial \beta^2} \log L_c(t)$$

(ب) برآورد تسیاتیس-لینک از $S(t; x)$ را برای $x = 1,5$ و $0 \leq t \leq 297$ ، با استفاده از $\beta = 1$ ، به دست آورید.

شماره مریض	نمره‌های عمل (X)	زمان بقاء (Y)
۳۶	۲,۰۹	۵۴
۱۱	۰,۳۶	۱۲۷ ⁺
۸۴	۰,۶۰	۲۹۷
۸۶	۱,۴۴	۳۸۹ ⁺
۳۳	۰,۹۱	۱۵۳۶ ⁺

حل:

(الف) فرض کنید:

$$i : \quad 1 \quad 2 \quad 3 \quad 4 \quad 5$$

$$x_i : \quad 2,09 \quad 0,36 \quad 0,60 \quad 1,44 \quad 0,91$$

با توجه به عبارت فصل ششم، بخش دوم، داریم:

$$\frac{\partial}{\partial \beta} \log L_c(1) = x_1 + x_3 - \frac{\sum_{j=1}^5 x_j e^{x_j}}{\sum_{j=1}^5 e^{x_j}} - \frac{\sum_{j=3}^5 x_j e^{x_j}}{\sum_{j=3}^5 e^{x_j}} = 0,965$$

$$-\frac{\partial^2}{\partial \beta^2} \log L_c(1) = \frac{\sum_{j=1}^5 x_j^2 e^{x_j}}{\sum_{j=1}^5 e^{x_j}} - \left(\frac{\sum_{j=1}^5 x_j e^{x_j}}{\sum_{j=1}^5 e^{x_j}} \right)^2 +$$

$$+ \frac{\sum_{j=3}^5 x_j^2 e^{x_j}}{\sum_{j=3}^5 e^{x_j}} - \left(\frac{\sum_{j=3}^5 x_j e^{x_j}}{\sum_{j=3}^5 e^{x_j}} \right)^2 = 0,5110$$

بنابراین، داریم:

$$\frac{\left(\frac{\partial}{\partial \beta} \log L_c(1) \right)^2}{-\frac{\partial^2}{\partial \beta^2} \log L_c(1)} = 0,18$$

و مقدار P، که از χ^2 به دست می‌آید، تقریباً ۰,۹ می‌شود.

(ب) برآورد تسیاتیس (بخش اول از فصل ششم) عبارت است از:

$$\hat{\Lambda}_{o,T}(t) = \begin{cases} 1 & 0 \leq t < 54 \\ \frac{1}{\sum_{j=1}^5 e^{x_j}} = 0,0554 & 54 \leq t < 297 \\ \frac{1}{\sum_{j=1}^5 e^{x_j}} + \frac{1}{\sum_{j=3}^5 e^{x_j}} = 0,1727 & t = 297 \end{cases}$$

بنابراین، داریم:

$$\hat{S}_T(t; 1/5) = e^{-\hat{\Lambda}_{o,T}(t)} e^{1/5} = \begin{cases} 1 & 0 \leq t < 54 \\ 0,78 & 54 \leq t < 297 \\ 0,46 & t = 297 \end{cases}$$

برآوردگر مطلوب (بخش اول از فصل ششم) به شرح زیر است:

$$\hat{\Lambda}_{o,T}(t) = \begin{cases} \frac{0,0554}{54} t & 0 \leq t < 54 \\ \frac{0,1727 - 0,0554}{297 - 54} (t - 54) + 0,0554 & 54 \leq t \leq 297 \end{cases}$$

بنابراین، داریم:

$$\hat{S}_L(t; 1/5) = e^{-\hat{\Lambda}_{o,L}(t)} e^{1/5} = \begin{cases} e^{-0,0046 t} & 0 \leq t < 54 \\ 0,877 e^{-0,0022 t} & 54 \leq t \leq 297 \end{cases}$$

مسألة ۱۵.

برای داده‌های AML (مثال بخش دوم از فصل سوم) دو گروه تیمار و شاهد به

شرح زیرند:

۹, ۱۳, ۱۳⁺, ۱۸, ۲۳, ۲۸⁺, ۳۱, ۳۴, ۴۵⁺, ۴۸, ۱۶۱⁺

گروه تیمار

۵, ۵, ۸, ۸, ۱۲, ۱۶⁺, ۲۳, ۲۷, ۳۰, ۳۳, ۴۳, ۴۵

گروه شاهد

دو گروه را به شکل زیر مقایسه کنید:

(الف) با آماره گهان و واریانس جایگشت آن

(ب) با آماره ML

(پ) صورت تارون-وایر آماره گهان

در هر حالت آماره استاندارد و مقدار P متناظر را به دست آورید.

حل:

(الف) در محاسبه نمرات مورد نیاز برای انجام آماره گهان، از جدول صفحه بعد استفاده می‌کنیم:

آماره گهان برابر $-50 = \sum_{NM} U^*$ و واریانس جایگشت ML، به شرح زیر است:

$$\frac{11 \times 12}{23 \times 22} \sum_{N, NM} (U^*)^2 = 912$$

بنابراین آماره استاندارد برابر با: $-1,656 = \frac{-50}{\sqrt{912}}$ بوده، که متناظر مقدار P دوطرفه ۰,۰۹۹ است.

Z	گروه	# < Z	# > Z	U*
۵(۲)	NM	۰	۲۱	-۲۱(۲)
۸(۲)	NM	۲	۱۹	-۱۷(۲)
۹	M	۴	۱۸	-۱۴
۱۲	NM	۵	۱۷	-۱۲
۱۳	M	۶	۱۶	-۱۰
۱۳ ⁺	M	۷	۰	۷
۱۶ ⁺	NM	۷	۰	۷
۱۸	M	۷	۱۳	-۶
۲۳	M	۸	۱۱	-۳
۲۵	NM	۸	۱۱	-۳
۲۵	NM	۱۰	۱۰	۰
۲۸	M	۱۱		۱۱
۳	NM	۱۱	۸	۳
۳۱	M	۱۲	۷	۵
۳۳	NM	۱۳	۶	۷
۳۴	M	۱۲	۵	۹
۳۳	NM	۱۵	۴	۱۱
۴۵	NM	۱۶	۳	۱۳
۴۵ ⁻	M	۱۷	۰	۱۷
۴۸ ⁻	M	۱۷	۱	۱۶
۴۹ ⁺	M	۱۸	۰	۱۸

(ب) و (ج) محاسبات در جدول زیر مشابه مسأله ۱۳، انجام پذیرفته است. آماره استاندارد MH (به جواب مسأله ۱۳ توجه کنید)، برابر با: $-\frac{3,69}{\sqrt{4,0072}} = -1,84$ است. بنابراین، مقدار P دوطرفه برابر $0,066$ است. صورت تارون آماره گهان برابر: $-\frac{50}{\sqrt{917,97}} = -1,65$ است (جواب مسأله ۱۳ را ببینید). بنابراین، مقدار P دوطرفه برابر $0,099$ به دست می آید.

z	n	m_1	n_1	a	$E_0(A)$	$n(a - E_0(A))$	$\frac{m_1(n - m_1)}{n - 1}$	$\frac{n_1}{n}(1 - \frac{n_1}{n})$	$n_1(n - n_1)$
۵	۲۳	۲	۱۱	۰	۰,۹۵۶۵	-۲۲	۱,۹۰۹	۰,۲۴۹۵	۱۳۲
۸	۲۱	۲	۱۱	۰	۱,۰۴۸	-۲۲	۱,۹	۰,۲۴۹۴	۱۱۰
۹	۱۹	۱	۱۱	۱	۰,۵۷۹	۸	۱	۰,۲۴۳۷	۸۸
۱۲	۱۸	۱	۱۰	۰	۰,۵۵۵	-۱۰	۱	۰,۲۴۶۹	۸۰
۱۳	۱۷	۱	۱۰	۱	۰,۵۸۸	۷	۱	۰,۲۴۲۲	۷۰
۱۸	۱۴	۱	۸	۱	۰,۵۷۱	۶	۱	۰,۲۴۴۹	۴۸
۲۳	۱۳	۲	۷	۱	۱,۰۷۷	-۱	۱,۸۳	۰,۲۴۸۵	۴۲
۲۷	۱۱	۱	۶	۰	۰,۵۴۵	-۶	۱	۰,۲۴۷۹	۳۰
۳۰	۹	۱	۵	۰	۰,۵۵۵	-۵	۱	۰,۲۴۶۹	۲۰
۳۱	۸	۱	۵	۱	۰,۶۲۵	۳	۱	۰,۲۳۴۴	۱۵
۳۳	۷	۱	۴	۰	۰,۵۷۱	-۴	۱	۰,۲۴۴۹	۱۲
۳۴	۶	۱	۴	۱	۰,۶۶۷	۲	۱	۰,۲۲۲۲	۸
۴۳	۵	۱	۳	۰	۰,۶	-۳	۱	۰,۲۴۰۰	۶
۴۵	۴	۱	۳	۰	۰,۷۵	-۳	۱	۰,۱۸۷۵	۳
۴۸	۲	۱	۲	۱	۱,۰	۰	۱	۰	۰

مسألة ۱۶.

ثابت کنید که صورت نسخه تارون-وایسر آماره گهان (یعنی: $(\sum n_i(a_i - E_0(A_i)))$) با صورت آماره گهان (یعنی: $U = \sum \sum U_{ij}$) برابر است. بجز احتمالاً برای عامل -1 ، که به خاطر تکرارهاست.

حل:

از بخش اول فصل چهارم نتیجه می شود که:

$$U = \sum_{k=1}^{m+n} U_k^* I(k \in I_1)$$

به گونه ای که:

$$U_k^* = \sum_{\substack{\ell=1 \\ \ell \neq k}}^{m+n} U_{k\ell}$$

است. یعنی اگر مشاهده ای با در نمونه ۱ دارای زیرنویس k باشد، بریده می شود، آن گاه U_k^* برابر تعداد مشاهدات بریده نشده قبل از آن منهای تعداد مشاهدات بعد از آن است.

$$U_k^* = \sum_{j=1}^{k-1} m_{j1} - (n_k - m_{k1}) = \sum_{j=1}^k m_{j1} - n_k \quad (\text{بخش دوم از فصل چهارم را ببینید})$$

از طرف دیگر، اگر مشاهده k ام بریده شده باشد، آن گاه U_k^* برابر تعداد مشاهدات

$$\text{بریده نشده قبل از آن است، یعنی: } U_k^* = \sum_{j=1}^k m_{j1}. \text{ بنابراین، داریم:}$$

$$U = \sum_{k=1}^{m+n} \sum_{j=1}^k m_{j1} I(k \in I_1) + \sum_{k=1}^{m+n} \left(\sum_{j=1}^k m_{j1} - n_k \right) I(k \in I_1)$$

که در این رابطه، c و u به ترتیب به معنی این است که، مجموعها روی مشاهدات بریده شده و نشده انجام می شود. پس، داریم:

$$\begin{aligned}
 U &= \sum_{k=1}^{m+n} \sum_{j=1}^k m_{j1} I(k \in I_1) - \sum_{k=1}^{m+n} n_k I(k \in I_1, \delta_k = 1) \\
 &= \sum_{j=1}^{m+n} m_{j1} \sum_{k=j}^{m+n} I(k \in I_1) - \sum_{k=1}^{m+n} n_k a_k \\
 &= \sum_{j=1}^{m+n} (m_{j1} n_{j1} - n_j a_j) = \sum_u (m_{j1} n_{j1} - n_j a_j) \\
 &= \sum_u n_j (a_j - E_0(A_j))
 \end{aligned}$$

دو تساوی آخر از این حقیقت ناشی می‌شود که m_j (در نتیجه a_j) - تعداد مشاهدات بریده نشده در z برابر صفر است. به شرطی که z یک مشاهده بریده شده باشد. (این قرارداد را به خاطر بیاورید که تکرار بین مشاهدات بریده شده و بریده نشده توسط این ملاحظه که مشاهدات بریده شده بزرگتراند، شکسته شده است).

مسأله ۱۷.

نشان دهید که واریانس جایگشت مانند برای آماره گهان بر $N^3 = (m+n)^3$ ، تقسیم می‌شود. یعنی:

$$\frac{1}{N^3} \times \frac{mn}{(m+n)(m+n-1)} \sum_{j=1}^{m+n} (U_j^*)^2$$

به عبارت دیگر $\int_0^\infty (1-H(t))^2 dH_U(t)$ همگراست.

همچنین اگر $N \rightarrow \infty$ میل کند، آن گاه $\frac{m}{N} \rightarrow \infty$ تحت فرض $H_0^*: F_1 = F_2; G_1 = G_2$ میل خواهد کرد، به گونه‌ای که F و G پیوسته بوده و داشته باشیم:

$$H(t) = P\{Z \leq t\} = \int_0^t (1-G(u)) dF_U(u) + \int_0^t (1-F(u)) dG_U(u)$$

$$H_U(t) = P\{Z \leq t, \xi = 1\} = \int_0^t (1-G(u)) dF_U(u)$$

حل:

فرض کنید:

$$\hat{H}(t) = \frac{1}{N} \sum_{i=1}^N I(Z_i \leq t)$$

$$\hat{H}_u(t) = \frac{1}{N} \sum_{i=1}^N I(Z_i \leq t, \xi_i = 1)$$

باشد. آن گاه،

$$U_i^* = \begin{cases} (Z_i > \text{تعداد مشاهده}) - (Z_i < \text{تعداد بریده شده ها}) & \xi_i = 1 \\ (Z_i < \text{تعداد بریده شده ها}) - & \xi_i = 0 \end{cases}$$

$$= \begin{cases} N \hat{H}_u(Z_i^-) - N(1 - \hat{H}(Z_i)) & \xi_i = 1 \\ N \hat{H}_u(Z_i^-) & \xi_i = 0 \end{cases}$$

$$= N[\hat{H}_u(Z_i^-) - \xi_i(1 - \hat{H}(Z_i))]$$

در نتیجه، داریم:

$$\begin{aligned} \frac{1}{N^2} \sum_{i=1}^N (U_i^*)^2 &= \frac{1}{N^2} \sum_{i=1}^N N^2 [\hat{H}_u(Z_i^-) - \xi_i(1 - \hat{H}(Z_i))]^2 \\ &= \frac{1}{N} \sum_{i=1}^N [\hat{H}_u^2(Z_i^-) - \frac{2}{N} \sum_{i=1}^N \xi_i \hat{H}_u(Z_i^-)(1 - \hat{H}(Z_i)) + \\ &\quad + \frac{1}{N} \sum_{i=1}^N \xi_i (1 - \hat{H}(Z_i))^2] \\ &= \int_0^\infty \hat{H}_u^2(t-) d\hat{H}(t) - \frac{2}{N} \int_0^\infty \hat{H}_u(t-)(1 - \hat{H}(t)) d\hat{H}_u(t) + \\ &\quad + \int_0^\infty (1 - \hat{H}(t))^2 d\hat{H}_u(t) \end{aligned}$$

چون $\hat{H}(t) \xrightarrow{\text{a.s.}} H(t)$ و $\hat{H}_u(t) \xrightarrow{\text{a.s.}} H_u(t)$ به طور یکنواخت نسبت به t

وقتی که $N \rightarrow \infty$ میل کند، لذا، بنا بر قضیه گلیونکو-کانتلی، اگر $N \rightarrow \infty$ ، آن گاه داریم:

$$\frac{1}{N^2} \sum_{i=1}^N (U_i^*)^2 \xrightarrow{\text{a.s.}} \int_0^{\infty} H_U^2(t) dH(t) - \\ - 2 \int_0^{\infty} H_U(t)(1-H(t)) dH_U(t) + \int_0^{\infty} (1-H(t))^2 dH_U(t)$$

با انتگرال گیری جزء به جزء، داریم:

$$2 \int_0^{\infty} H_U(t)(1-H(t)) dH_U(t) \\ = H_U^2(t)(1-H(t)) \Big|_0^{\infty} + \int_0^{\infty} H_U^2(t) dH(t) = \int_0^{\infty} H_U^2(t) dH(t)$$

در نتیجه، دو جمله اول حد بالا حذف شده و همراه با:

$$\frac{mn}{(m+n)(m+n-1)} \rightarrow \lambda(1-\lambda)$$

نتیجه کامل می شود.

پایان

واژه‌نامه

A

Accelerated time model

الگوی زمانی شتاب داده شده

Actuarial method

روش بیمه‌گری

Asymptotic normality

نرمال مجانبی

B

Bias–corrected estimator

برآوردگر تصحیح‌اریبی

Bootstrap

بوت‌استرپ

Bayesian estimate

برآوردگر بیزی

C

Censoring

برش

Chi–square test

آزمون کی دو

Competing risk

نرخ شکست رقیب

Conditional likelihood

درست‌مایی شرطی

Confidence interval

فاصله اطمینان

Consistency

سازگاری

Continuity correction

تصحیح پیوستگی

Cox model

الگوی کاکس

Cumulative hazard function

تابع نرخ شکست تجمعی

Cohort life table

جدول طول عمر گروهی

Current life table	جدول طول عمر جاری
Crude probability	احتمال خام
D	
Delta method	روش دلتا
Density function	تابع چگالی
Dirichlet process	فرایند دریکله
Discret data	داده‌های جدا-گسسته
Distribution function	تابع توزیع
Drop out	قطع ادامه کار
Draw back	باجایگذاری - عیب
Diagnosis	تشخیص بیماری
E	
Exponential function	تابع نمایی
Exponential distribution	توزیع نمایی
Exponential density	چگالی نمایی
Extreme value function	تابع مقادیر فرین
Extreme value distribution	توزیع مقادیر فرین
Extreme value density	چگالی مقادیر فرین
Efron's test	آزمون افرون
Effective sample size	طول مؤثر نمونه
F	
Force of mortality	نرخ مرگ و میر
Full likelihood	درست‌نمایی کامل
G	
Gamma distribution	توزیع گاما

Gamma density

چگالی گاما

Gehan test

آزمون گهان

H

Hypergeometric

فوق هندسی

Hazard function

تابع نرخ شکست

Hazard rate

نرخ شکست

I

IFR

تابع نرخ شکست صعودی

IFRA

تابع متوسط نرخ شکست صعودی

Insurance

بیمه‌گری

Indentially distribution

هم‌توزیع

iid

مستقل و هم‌توزیع

Interval censoring

برش فاصله‌ای (نوع ۳) یا تصادفی

Influence function

تابع تأثیر

Iterative method

روش تکراری

J

Jackknifed method

روش جک‌نایف

K

Kaplan–Meier estimator

برآوردگر کاپلان-مایر

L

Least squares

کمترین مربعات

Linear Rank test

آزمون رتبه خطی

Log Rank test

آزمون رتبه لگاریتمی

Lose to follow-up

عدم بازگشت

Left censoring

چپ برش (برش از چپ)

Logistic function	تابع لوجستیک
Log-linear model	الگوی خطی لگاریتمی
M	
Maximum likelihood (ML)	حداکثر درست‌نمایی
Maximum likelihood estimator (MLE)	برآوردگر درست‌نمایی
Mean	میانگین (حسابی)، معدل
Median	میانه
Miller modified	تعمیم میلر
Mortality table	جدول طول عمر
Mariginal	حاشیه‌ای
Method of scoring	روش امتیازی (چوب خطی)
N	
Newton-Raphson method	روش نیوتن روفسون
Neyman-Pearson	نیمن-پیرسن
Naive estimator	برآوردگر ساده
O	
Order statistics	آماره‌های مرتب
Order populations	جامعه‌های مرتب
P	
Partial likelihood	درست‌نمایی نسبی
Permutation	جایگشت
Permutation theory	نظریه جایگشت
Peterson's representation	نمایش پترسن
Plots	رسم به کمک نقطه‌یابی
Probability plots	رسم احتمالی
Probability paper	کاغذ احتمالی

Ploting position	احتمال تجربی
Poisson process	فرایند پواسن
Product – limit estimator	حدّ حاصل ضرب برآوردگر
Proportional hazards model	الگوهای نرخ شکست متناسب
Pivotal statistic	آماره محوری
Prior distribution	توزیع پیشین
Posterior	توزیع پسین
partial probability	احتمال جزئی
Q	
Q–Q plots	رسم Q–Q
Quadratic	درجه دوم
R	
Rank	رتبه
Redistribute – to – the – right	تجدید نظر در توزیع به راست
Reduced sample method	روش کاهش نمونه
Regression linear model	الگوی خطی رگرسیون
Rao – Blackwell theorem	قضیه راثو – بلکول
Restricted mean	میانگین محدود شده
Randome censoring	برش تصادفی
Right censoring	راست برش (برش از راست)
Risk	شکست – نرخ شکست
Reliability function	تابع اعتماد
Robust estimator	برآوردگر تنومند
S	
Sample information matrix	ماتریس اطلاع نمونه
Score vector	بردار امتیاز
Score method	روش امتیازی (چوب خطی)

Self – consistency	خودسازگاری
Subdistribution function	تابع توزیع جزئی
Survival	بقاء
Survival function	تابع بقاء
Survival time	زمان بقاء
Surviving fraction	کسر بقاء
Single sample	نمونه‌ای به حجم واحد
Subsurvival function	تابع بقاء جزئی
T	
Test	آزمون
Ties	تکرار
Time dependent covarites	متغیرهای وابسته به زمان
Trend	روند
Truncation	قطع
Two by two tables	جدولهای ۲×۲
Trimed mean	پیرایش کردن میانگین
Tukey biweight estimator	برآوردگر دو وزنی توکی
Tarone – Ware class	رده تارون – وایر
Tarone – Ware generalized	تعمیم تارون – وایر
W	
Weak convergence	همگرایی ضعیف
With draw	باجایگذاری
Wilks likelihood ratio	نسبت درست‌نمایی ویلکس
Weibul distribution	توزیع وایبل
U	
Unbias	نااریب

منابع

The numbers in brackets after the references are the numbers of the pages on which the references are cited.

Aalen, O. (1976). Nonparametric inference in connection with multiple decrement models. Scandinavian Journal of Statistics 3, 15-27. [69]

_____ (1978). Nonparametric inference for a family of counting processes. Annals of Statistics 6, 701-726. [69]

Abelson, R. P. and Tukey, J. W. (1963). Efficient utilization of non-numerical information in quantitative analysis: General theory and the case of simple order. Annals of Mathematical Statistics 34, 1347-1369. [112]

Altshuler, B. (1970). Theory for the measurement of competing risks in animal experiments. Mathematical Biosciences 6, 1-11. [179]

- Bailey, K. R. (1979). The general maximum likelihood approach to the Cox regression model. Ph.D. dissertation, University of Chicago, Chicago, Illinois. [133]
- Barlow, R. E. and Proschan, F. (1975). Statistical Theory of Reliability and Life Testing. Holt, Rinehart, and Winston, New York. [15]
- Barr, D. R. and Davidson, T. (1973). A Kolmogorov-Smirnov test for censored samples. Technometrics 15, 739-757. [173]
- Basu, A. P. (1964). Estimates of reliability for some distributions useful in life testing. Technometrics 6, 215-219. [35]
- Berkson, J. and Gage, R. P. (1950). Calculation of survival rates for cancer. Proceedings of the Staff Meetings of the Mayo Clinic 25, 270-286. [46]
- Berman, S. M. (1963). Note on extreme values, competing risks and semi-Markov processes. Annals of Mathematical Statistics 34, 1104-1106. [179]
- Billingsley, P. (1968). Convergence of Probability Measures. Wiley, New York. [65]
- Breslow, N. (1970). A generalized Kruskal-Wallis test for comparing K samples subject to unequal patterns of censorship. Biometrika 57, 579-594. [109, 114]
- _____. (1972). Discussion on Professor Cox's paper. Journal of the Royal Statistical Society, Series B 34, 216-217. [136]

- _____ (1974). Covariance analysis of censored survival data. Biometrics 30, 89-99. [136, 139]
- _____ and Crowley, J. (1974). A large sample study of the life table and product limit estimates under random censorship. Annals of Statistics 2, 437-453. [46, 65, 69]
- Buckley, J. and James, I. (1979). Linear regression with censored data. Biometrika 66, 429-436. [153, 163]
- Campbell, G. (1979). Nonparametric bivariate estimation with randomly censored data. Mimeoseries #79-25, Department of Statistics, Purdue University, West Lafayette, Indiana. [177]
- Chiang, C. L. (1968). Introduction to Stochastic Processes in Biostatistics. Wiley, New York. [46, 179]
- Cohen, A. C. (1965). Maximum likelihood estimation in the Weibull distribution based on complete and on censored samples. Technometrics 7, 579-588. [29]
- Cox, D. R. (1972). Regression models and life-tables. Journal of the Royal Statistical Society, Series B 34, 187-202. [127, 139]
- _____ (1975). Partial likelihood. Biometrika 62, 269-276. [132]
- _____ and Snell, E. J. (1968). A general definition of residuals. Journal of the Royal Statistical Society, Series B 30, 248-275. [172]
- Crowley, J. (1974). Asymptotic normality of a new nonparametric statistic for use in organ transplant studies. Journal of the American Statistical Association 69, 1006-1011. [103]

- _____ and Hu, M. (1977). Covariance analysis of heart transplant survival data. Journal of the American Statistical Association 72, 27-36. [141, 172]
- Cutler, S. J. and Ederer, F. (1958). Maximum utilization of the life table method in analyzing survival. Journal of Chronic Diseases 8, 699-712. [42, 46]
- Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. Journal of the Royal Statistical Society, Series B 39, 1-22. [154]
- Dufour, R. and Maag, U. R. (1978). Distribution results for modified Kolmogorov-Smirnov statistics for truncated or censored samples. Technometrics 20, 29-32. [173]
- Efron, B. (1967). The two sample problem with censored data. Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Vol. IV. University of California Press, Berkeley, California. 831-853. [52, 57, 106]
- _____ (1977). The efficiency of Cox's likelihood function for censored data. Journal of the American Statistical Association 72, 557-565. [132]
- _____ (1979). Bootstrap methods: Another look at the jackknife. Annals of Statistics 7, 1-26. [182]
- _____ (1980). Censored data and the bootstrap. Technical Report No. 53 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [76, 182]

- _____ and Hinkley, D. V. (1978). Assessing the accuracy of the maximum likelihood estimator: Observed versus expected Fisher information. Biometrika 65, 457-487. [21]
- Elveback, L. (1958). Estimation of survivorship in chronic disease: The "actuarial" method. Journal of the American Statistical Association 53, 420-440. [46]
- Embury, S. H., Elias, L., Heller, P. H., Hood, C. E., Greenberg, P. L. and Schrier, S. L. (1977). Remission maintenance therapy in acute myelogenous leukemia. Western Journal of Medicine 126, 267-272. [50]
- Epstein, B. and Sobel, M. (1953). Life testing. Journal of the American Statistical Association 48, 486-502. [25]
- Farewell, V. T. (1977). A model for a binary variable with time-censored observations. Biometrika 64, 43-46. [38]
- Feigl, P. and Zelen, M. (1965). Estimation of exponential survival probabilities with concomitant information. Biometrics 21, 826-838. [36]
- Ferguson, T. S. (1973). A Bayesian analysis of some nonparametric problems. Annals of Statistics 1, 209-230. [79]
- _____ and Phadia, E. G. (1979). Bayesian nonparametric estimation based on censored data. Annals of Statistics 7, 163-186. [79]
- Fleming, T. R. and Harrington, D. P. (1979). Nonparametric estimation of the survival distribution in censored data. Unpublished manuscript. [67]

- Földes, A., Rejtő, L. and Winter, B. B. (1978). Strong consistency properties of nonparametric estimators for randomly censored data. Part II: Estimation of density and failure rate. Unpublished manuscript. [76]
- Gail, M. (1975). A review and critique of some models used in competing risk analysis. Biometrics 31, 209-222. [179]
- Gaver, D. P., Jr. and Hoel, D. G. (1970). Comparison of certain small-sample Poisson probability estimates. Technometrics 12, 835-850. [35]
- Gehan, E. A. (1965). A generalized Wilcoxon test for comparing arbitrarily singly-censored samples. Biometrika 52, 203-223. [89]
- Gilbert, J. P. (1962). Random censorship. Ph.D. dissertation, University of Chicago, Chicago, Illinois. [94]
- Gillespie, M. J. and Fisher, L. (1979). Confidence bands for the Kaplan-Meier survival curve estimate. Annals of Statistics 7, 920-924. [173]
- Glasser, M. (1967). Exponential survival with covariance. Journal of the American Statistical Association 62, 561-568. [36]
- Gong, G. (1980). Do Hodgkin's disease patients with DNCB sensitivity survive longer? Biostatistics Casebook, Vol. III, Technical Report No. 57 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [168]
- Gregory, P. B., Knauer, C. M., Kempson, R. L. and Miller, R. (1976). Steroid therapy in severe viral hepatitis. New England Journal of Medicine 294, 681-686. [186, 196]

- Gross, A. J. and Clark, V. A. (1975). Survival Distributions: Reliability Applications in the Biomedical Sciences. Wiley, New York. [20]
- Hall, W. J. and Wellner, J. A. (1980). Confidence bands for a survival curve from censored data. Biometrika 67, 133-143. [173]
- Hollander, M. and Proschan, F. (1979). Testing to determine the underlying distribution using randomly censored data. Biometrics 35, 393-401. [174]
- Hyde, J. (1977). Testing survival under right censoring and left truncation. Biometrika 64, 225-230. [174]
- _____ (1977). Life testing with incomplete observations. Technical Report No. 30 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [94]
- Johansen, S. (1978). The product limit estimator as maximum likelihood estimator. Scandinavian Journal of Statistics 5, 195-199. [59]
- Johns, M. V., Jr. and Lieberman, G. J. (1966). An exact asymptotically efficient confidence bound for reliability in the case of the Weibull distribution. Technometrics 8, 135-175. [32]
- Kalbfleisch, J. and Prentice, R. L. (1972). Discussion on Professor Cox's paper. Journal of the Royal Statistical Society, Series B 34, 215-216. [139]
- _____ and _____ (1973). Marginal likelihoods based on Cox's regression and life model. Biometrika 60, 267-278. [130]

- _____ and _____ (1980). The Statistical Analysis of Failure Time Data. Wiley, New York. [20, 127, 143, 146]
- Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete observations. Journal of the American Statistical Association 53, 457-481. [48, 59, 72]
- Kay, R. (1977). Proportional hazard regression models and the analysis of censored survival data. Applied Statistics (Journal of the Royal Statistical Society, Series C) 26, 227-237. [172]
- Kiefer, J. and Wolfowitz, J. (1956). Consistency of the maximum likelihood estimator in the presence of infinitely many incidental parameters. Annals of Mathematical Statistics 27, 887-906. [59]
- Korwar, R. M. (1980). Nonparametric estimation of a bivariate survivorship function with doubly censored data. Unpublished manuscript. [177]
- Koul, H., Susarla, V. and Van Ryzin, J. (1979). Regression analysis with randomly right censored data. Unpublished manuscript. [155]
- Koziol, J. A. and Byar, D. P. (1975). Percentage points of the asymptotic distributions of one and two sample K-S statistics for truncated or censored data. Technometrics 17, 507-510. [173]
- _____ and Green, S. B. (1976). A Cramér-von Mises statistic for randomly censored data. Biometrika 63, 465-474. [173]
- Lagakos, S. W. (1979). General right censoring and its impact on the analysis of survival data. Biometrics 35, 139-156. [179]

- and Williams, J. S. (1978). Models for censored survival analysis: A cone class of variable-sum models. Biometrika 65, 181-189. [179]
- Lamb, E. J. and Leurgans, S. (1979). Does adoption affect subsequent fertility? American Journal of Obstetrics and Gynecology 134, 138-144. [141]
- Lamborn, K. (1969). On chi-squared goodness of fit tests for sampling from more than one population with possibly censored data. Technical Report No. 21 (T01 GM00025), Department of Statistics, Stanford University, Stanford, California. [175]
- Langberg, N. A., Proschan, F. and Quinzi, A. J. (1981). Estimating dependent life lengths, with applications to the theory of competing risks. Annals of Statistics 9, 157-167. [179]
- Latta, R. B. (1977). Generalized Wilcoxon statistics for the two-sample problem with censored data. Biometrika 64, 633-635. [146]
- Leavitt, S. S. and Olshen, R. A. (1974). The insurance claims adjuster as patients' advocate: Quantitative impact. Report for Insurance Technology Company, Berkeley, California. [2]
- Leiderman, P. H., Babu, D., Kagia, J., Kraemer, H. C. and Leiderman, G. F. (1973). African infant precocity and some social influences during the first year. Nature 242, 247-249. [7]
- Leurgans, S. (1980). Does adoption affect fertility? A proportional hazards model. Biostatistics Casebook, Vol. III, Technical Report No. 57 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [141]

- Lininger, L., Gail, M. H., Green, S. B. and Byar, D. P. (1979). Comparison of four tests for equality of survival curves in the presence of stratification and censoring. Biometrika 66, 419-428. [103]
- Link, C. L. (1979). Confidence intervals for the survival function using Cox's proportional hazard model with covariates. Technical Report No. 45 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [136]
- Mantel, N. (1967). Ranking procedures for arbitrarily restricted observation. Biometrics 23, 65-78. [89]
- _____ and Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. Journal of the National Cancer Institute 22, 719-748. [103]
- _____ and Myers, M. (1971). Problems of convergence of maximum likelihood iterative procedures in multiparameter situations. Journal of the American Statistical Association 66, 484-491. [36]
- Marcuson, R. and Nordbrock, E. (1981). A K-sample generalization of the Gehan-Gilbert procedure for the analysis of arbitrarily censored survival data. Biometrische Zeitschrift/Biometrical Journal. [113]
- Meier, P. (1975). Estimation of a distribution function from incomplete observations. Perspectives in Probability and Statistics. Papers in Honour of M. S. Bartlett (Ed. J. Gani). Academic Press, New York. 67-82. [72]

- Mihalko, D. P. and Moore, D. S. (1980). Chi-square tests of fit for Type II censored data. Annals of Statistics 8, 625-644. [174]
- Miller, R. G. (1974). The jackknife - a review. Biometrika 61, 1-15. [182]
- _____ (1975). Jackknifing censored data. Technical Report No. 14 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [182]
- _____ (1976). Least squares regression with censored data. Biometrika 63, 449-464. [150, 163]
- Moeschberger, M. L. and David, H. A. (1971). Life tests under competing causes of failure and the theory of competing risks. Biometrics 27, 909-923. [179]
- Morton, R. (1978). Regression analysis of life tables and related nonparametric tests. Biometrika 65, 329-333. [146]
- Muñoz, A. (1980). Nonparametric estimation from censored bivariate observations. Technical Report No. 60 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [177]
- _____ (1980). Consistency of the self-consistent estimator of the distribution function from censored observations. Technical Report No. 61 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [177]
- Nelson, W. (1969). Hazard plotting for incomplete failure data. Journal of Quality Technology 1, 27-52. [67, 165]

- (1972). Theory and applications of hazard plotting for censored failure data. Technometrics 14, 945-966. [67, 165]
- Oakes, D. (1977). The asymptotic information in censored survival data. Biometrika 64, 441-448. [132]
- Peterson, A. V., Jr. (1975). Nonparametric estimation in the competing risks problem. Technical Report No. 13 (RO1 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [179]
- (1976). Bounds for a joint distribution function with fixed sub-distribution functions: Application to competing risks. Proceedings of the National Academy of Sciences 73, 11-13. [179]
- (1977). Expressing the Kaplan-Meier estimator as a function of empirical subsurvival functions. Journal of the American Statistical Association 72, 854-858. [63, 67]
- Peto, R. (1972). Discussion on Professor Cox's paper. Journal of the Royal Statistical Society, Series B 34, 205-207. [139]
- and Peto, J. (1972). Asymptotically efficient rank invariant test procedures. Journal of the Royal Statistical Society, Series A 135, 185-198. [146]
- and Pike, M. C. (1973). Conservatism of the approximation $\Sigma(O-E)^2/E$ in the logrank test for survival data or tumor incidence data. Biometrics 29, 579-584. [117]

, Pike, M. C., Armitage, P., Breslow, N. E., Cox, D. R., Howard, S. V., Mantel, N., McPherson, K., Peto, J. and Smith, P. G. (1976). Design and analysis of randomized clinical trials requiring prolonged observation of each patient. I. Introduction and design. British Journal of Cancer 34, 585-612. [117]

, Pike, M. C., Armitage, P., Breslow, N. E., Cox, D. R., Howard, S. V., Mantel, N., McPherson, K., Peto, J. and Smith, P. G. (1977). Design and analysis of randomized clinical trials requiring prolonged observation of each patient. II. Analysis and examples. British Journal of Cancer 35, 1-39. [117]

Pettit, A. N. (1976). Cramér-von Mises statistics for testing normality with censored samples. Biometrika 63, 475-481. [173]

_____ (1977). Tests for the exponential distribution with censored data using Cramér-von Mises statistics. Biometrika 64, 629-632. [173]

_____ and Stephens, M. A. (1976). Modified Cramér-von Mises statistics for censored data. Biometrika 63, 291-298. [173]

Phadia, E. G. (1980). A note on empirical Bayes estimation of a distribution function based on censored data. Annals of Statistics 8, 226-229. [80]

Prentice, R. L. (1978). Linear rank tests with right censored data. Biometrika 65, 167-179. [146]

_____ and Gloeckler, L. A. (1978). Regression analysis of grouped survival data with application to breast cancer data. Biometrics 34, 57-67. [139]

- _____ and Kalbfleisch, J. D. (1979). Hazard rate models with covariates. Biometrics 35, 25-39. [127, 143]
- _____, Kalbfleisch, J. D., Peterson, A. V., Jr., Flournoy, N., Farewell, V. T. and Breslow, N. E. (1978). The analysis of failure times in the presence of competing risks. Biometrics 34, 541-554. [179]
- Rai, K., Susarla, V. and Van Ryzin, J. (1980). Shrinkage estimation in nonparametric Bayesian survival analysis: A simulation study. Communications in Statistics, Simulation and Computation B9, 271-298. [79]
- Rao, C. R. (1965). Linear Statistical Inference. Wiley, New York. [20, 21]
- Reid, N. M. (1981). Influence functions for censored data. Annals of Statistics 9, 78-92. [73, 74, 76, 182]
- _____ and Iyengar, S. (1979). Estimating the variance of the median. Unpublished notes. [76]
- Sander, J. M. (1975). The weak convergence of quantiles of the product-limit estimator. Technical Report No. 5 (R01 CM21215), Division of Biostatistics, Stanford University, Stanford, California. [76]
- _____ (1975). Asymptotic normality of linear combinations of functions of order statistics with censored data. Technical Report No. 8 (R01 GM21215), Division of Biostatistics, Stanford University, Stanford, California. [72, 73]
- Schmee, J. and Hahn, G. J. (1979). A simple method for regression analysis with censored data. Technometrics 21, 417-432. [154]

- Schoenfeld, D. (1980). Chi-squared goodness-of-fit tests for the proportional hazards regression model. Biometrika 67, 145-153. [175]
- Susarla, V. and Van Ryzin, J. (1976). Nonparametric Bayesian estimation of survival curves from incomplete observations. Journal of the American Statistical Association 71, 897-902. [79]
- _____ and _____ (1978a). Empirical Bayes estimation of a distribution (survival) function from right censored observations. Annals of Statistics 6, 740-754. [80]
- _____ and _____ (1978b). Large sample theory for a Bayesian nonparametric survival curve estimator based on censored samples. Annals of Statistics 6, 755-768. [79]
- _____ and _____ (1980). Large sample theory for an estimator of the mean survival time from censored samples. Annals of Statistics 8, 1001-1016. [72]
- Tarone, R. E. (1975). Tests for trend in life table analysis. Biometrika 62, 679-682. [118]
- _____ and Ware, J. (1977). On distribution-free tests for equality of survival distributions. Biometrika 64, 156-160. [105, 116]
- Thomas, D. R. and Grunkemeier, G. L. (1975). Confidence interval estimation of survival probabilities for censored data. Journal of the American Statistical Association 70, 865-871. [52]
- Tsiatis, A. (1975). A nonidentifiability aspect of the problem of competing risks. Proceedings of the National Academy of Sciences 72, 20-22. [179]

- _____ (1978). A heuristic estimate of the asymptotic variance of the survival probability in Cox's regression model. Technical Report No. 524, Department of Statistics, University of Wisconsin, Madison, Wisconsin. [136]
- _____ (1981). A large sample study of Cox's regression model. Annals of Statistics 9, 93-108. [133, 136]
- Turnbull, B. W. (1974). Nonparametric estimation of a survivorship function with doubly censored data. Journal of the American Statistical Association 69, 169-173. [7, 57]
- _____ (1976). The empirical distribution function with arbitrarily grouped, censored and truncated data. Journal of the Royal Statistical Society, Series B 38, 290-295. [57]
- _____, Brown, B. W., Jr. and Hu, M. (1974). Survivorship analysis of heart transplant data. Journal of the American Statistical Association 69, 74-80. [141]
- _____ and Weiss, L. (1978). A likelihood ratio statistic for testing goodness of fit with randomly censored data. Biometrics 34, 367-375. [174]
- Wilk, M. B. and Gnanadesikan, R. (1968). Probability plotting methods for the analysis of data. Biometrika 55, 1-17. [165]
- _____, Gnanadesikan, R. and Huyett, M. J. (1962). Probability plots for the gamma distribution. Technometrics 4, 1-20. [166]
- Williams, J. S. and Lagakos, S. W. (1977). Models for censored survival analysis: Constant-sum and variable-sum models. Biometrika 64, 215-224. [179]

- Zacks, S. and Even, M. (1966). The efficiencies in small samples of the maximum likelihood and best unbiased estimators of reliability functions. Journal of the American Statistical Association 61, 1033-1051. [35]
- Zipin, C. and Armitage, P. (1966). Use of concomitant variables and incomplete survival information in the estimation of an exponential survival parameter. Biometrics 22, 665-672. [36]
- _____ and Lamborn, K. (1969). Concomitant variables and censored survival data in estimation of an exponential survival parameter, Part II. Technical Report No. 20 (T01 GM00025), Department of Statistics, Stanford University, Stanford, California. [36]



FERDOWSI UNIVERSITY OF MASHHAD

Publication No. 312

Survival Analysis

by

RUPERT G. MILLER, JR.

Translated by

ABOLGHASEM BOZORGNIA - HOJJAT REZAAE PAZHAND

FERDOWSI UNIVERSITY PRESS

2001