

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

۳۳۰۲۰



دانشگاه صنعتی اصفهان
دانشکده علوم ریاضی

برآوردگر رگرسیونی داده‌های بقا
سانسور شده بر روی فواصل با مدل شکست متناسب

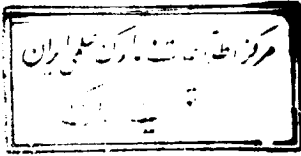
پایان نامه کارشناسی ارشد آمار
محسن ابوالقاسمی

۱۰۱۲۷

استاد راهنما
دکتر سروش غلیمرادی

۱۳۷۹

۳۳۰۳۰



۱۳۸۰ / ۱ / ۲۴

دانشگاه صنعتی اصفهان

دانشکده علوم ریاضی

پایان نامه کارشناسی ارشد رشته آمار آقای محسن ابوالقاسمی

تحت عنوان

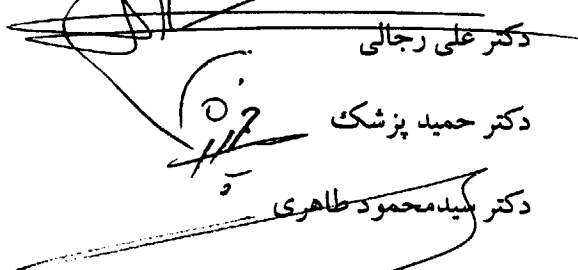
برآوردگر رگرسیونی داده‌های بقا

سانسور شده بر روی فواصل با مدل شکست متناسب

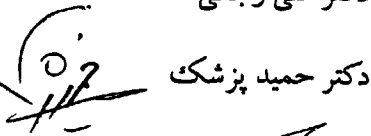
در تاریخ ۷۹/۷/۳۰ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهائی قرار گرفت.


دکتر سروش علی‌رادی

۱- استاد راهنمای پایان‌نامه


دکتر علی رجالی

۲- استاد مشاور پایان‌نامه


دکتر حمید پزشکی

۳- استاد داور ۱


دکتر بکید محمود طاهری

۴- استاد داور ۲


۱۳۸۰

دکتر امیر نادری

سرپرست تحصیلات تکمیلی دانشکده

کلیه حقوق مادی مترتب بر نتایج مطالعات،
ابتکارات و نوآوریهای ناشی از تحقیق موضوع
این پایان نامه (رساله) متعلق به دانشگاه صنعتی
اصفهان است.

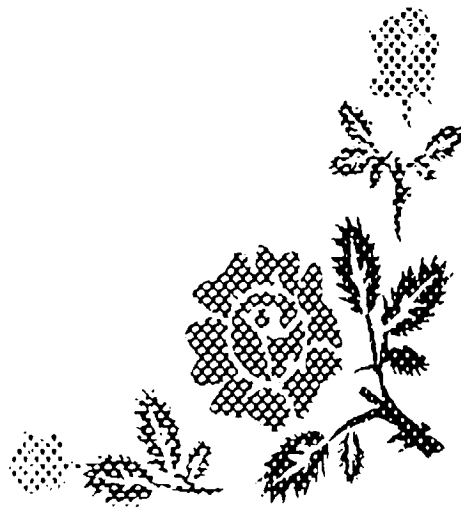
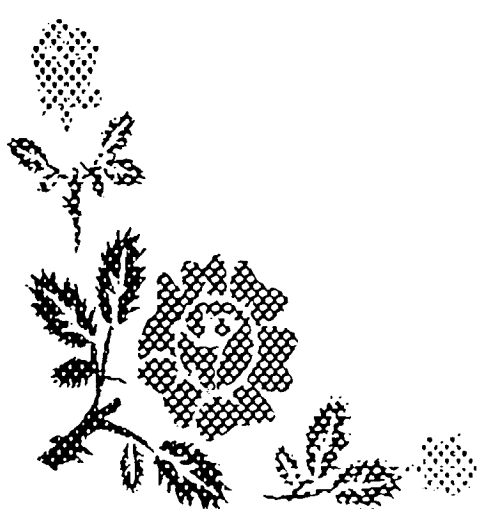


تقدیم به پیشگاه :

پدر و مادر عزیز و زحمتکشم

همسر مهربان و فداکارم

دوستان گرانقدر دوران تحصیلم



فهرست مطالب

۲	فصل اول : مقدمه
۶	فصل دوم : تعاریف و مفاهیم اولیه
۶	۱.۲ مقدمه
۷	۲.۲ تعاریف و نمادها
۷	۱.۲.۲ تعریف قابلیت اعتماد
۷	۲.۲.۲ تعریف تابع بقا
۸	۳.۲.۲ تعریف نرخ شکست
۹	۲.۳.۲ خصوصیات تابع نرخ شکست
۱۰	۲.۳.۲ خصوصیات تابع نرخ شکست تجمعی
۱۰	۳.۳.۲ شکلهای مختلف تابع نرخ شکست و تعبیر آنها
۱۲	۳.۲ معرفی دو مدل احتمالی از زمان شکست
۱۲	۱.۳.۲ مدل شکست نمایی
۱۲	۲.۳.۲ مدل شکست وایبال
۱۳	۴.۲ مدلهایی با نرخ شکست تصادفی
۱۳	۵.۲ بررسی زمان شکست موجود با زمان شکست تشدید شده
۱۴	۶.۲ مدل های خطی
۱۴	۷.۲ داده های طول عمر و سانسور
۱۴	۱.۷.۲ داده های کامل
۱۵	۲.۷.۲ داده های سانسور شده
۱۷	فصل سوم : مدل های زمان شکست
۱۷	۱.۳ مقدمه
۱۸	۲.۳ مدل شکست نمایی
۱۹	۳.۳ مدل شکست وایبال
۲۰	۴.۳ مدل های شکست رگرسیونی
۲۲	۵.۳ مدل شکست متناسب
۲۴	۶.۳ عملکرد سانسور در تابع درستمایی
۲۷	۷.۳ برآورد پارامترها در مدل رگرسیون نمایی

۳۰.....	۸.۳. برآورد پارامترها در مدل رگرسیون لگاریتم خطی
۳۴.....	۹.۳. برآورد پارامترها در مدل شکست متناسب، با داده‌های سانسور شده بر روی فواصل
۴۳.....	فصل چهارم: الگوریتم EM در برآورد پارامترها، به روش حداکثر درستنمایی با داده‌های گم شده
۴۳.....	۱.۴. مقدمه
۴۴.....	۲.۴. برآورد پارامترها به روش حداکثر درستنمایی بدون در نظر گرفتن محدودیت خطی
۴۴.....	۳.۴. برآورد حداکثر درستنمایی با در نظر گرفتن محدودیت خطی روی پارامترها
۴۶.....	۴.۴. الگوریتم EM در برآورد پارامترها با داده‌های گم شده
۴۸.....	۵.۴. الگوریتم محدود EM در برآورد پارامترها با داده‌های گم شده
۴۹.....	۶.۴. کاربرد الگوریتم محدود EM در آزمون فرضها و فواصل اطمینان با داده‌های گم شده
۵۳.....	۷.۴. طرح یک مثال
	فصل پنجم: تحلیل رگرسیونی داده‌های بقا سانسور شده بر روی فواصل، با استفاده از مدل‌های خطی
۵۶.....	۱.۵. مقدمه
۵۷.....	۲.۵. طرح یک مدل
۶۰.....	۳.۵. برآورد حداکثر درستنمایی
۶۳.....	۱.۳.۵. گام E از الگوریتم EM
۶۴.....	۱.۱.۳.۵. محاسبه امید شرطی $\log(\lambda_0(T_i))$
۶۸.....	۲.۱.۳.۵. محاسبه امید شرطی $\lambda_0(T_i)$
۷۲.....	۲.۳.۵. گام M از الگوریتم EM
۷۷.....	۴.۵. آزمون فرض و ساختن فواصل اطمینان
۷۷.....	۱.۴.۵. آزمون فرض روی پارامتر β_h
۸۵.....	۲.۴.۵. ساختن فاصله اطمینان جهت پارامتر β_h
۹۷.....	۵.۵. طرح یک مثال
۹۳.....	پیوست الف: برآورد پارامترها در مدل‌های خطی تعمیم یافته
۹۴.....	الف. ۱) خانواده توزیع‌های نمایی و الگوهای خطی تعمیم یافته
۹۴.....	الف. ۲) خانواده توزیع‌های نمایی
۹۸.....	الف. ۳) برآورد الگوهای خطی تعمیم یافته

پیوست ب : برنامه کامپیوتری ۱۰۳

کتابنامه ۱۰۲

چکیده

ما تحلیل رگرسیون داده های بقا را با مشاهدات سانسور از راست - چپ یا فاصله ای بررسی می کنیم. همچنین با به کارگیری مدل های لگاریتم خطی برای شکستهای سانسور شده فاصله ای روشهای جدول عمر را برای داده های بقای سانسور شده گسترش می دهیم. الگوریتم EM برای محاسبه برآوردهای حداکثر درستنمایی پارامترها استفاده شده است. ما فرض می کنیم تابع شکست یک تابع پله ای روی فواصل مجزای زمان است ، بنابراین مدل نمایی پارامتری و مدل شکست متناسب کاکس به سادگی ابزاری برای حالات خاص هستند. ما الگوریتم محدود شده EM را برای آزمون فرض ها و بنا کردن فواصل اطمینان برای پارامترها سازگار می کنیم. این روشها در یک تحلیل زمان عود بیماری در معالجه مبتلایان ملانوما به کار برده شده اند.

فصل اول

مقدمه

تحلیل داده‌های بقاء^۱ به مفهوم بررسی اطلاعاتی است که با طول عمر موجودات سرو کار دارد. فایده اصلی این کار، شناسایی روابط بین زمان بقاء و مجموعه‌ای از متغیرهای توضیحی^۲ مانند: سن، جنس، رفتار و دیگر مشخصه‌ها است. [۱]

یکی از مهمترین کاربردهای تحلیل بقاء، در پزشکی است. در این نوع تحلیل‌ها ارتباط بین طول مدت بقاء بیمار و مجموعه‌ای از متغیرها مانند: سن، جنس، میزان پیشرفت بیماری، شیوه درمان و سایر خصوصیات بالینی بیمار مورد بررسی قرار می‌گیرد. [۲]

اما مسئله پیچیده در مورد اطلاعات پزشکی بیماران، این است که اطلاعات همیشه در دسترس نیست و ممکن است در قسمتی از زمان درمان در دسترس نباشد. اینجاست که مسئله داده‌های سانسور شده^۱ مطرح می‌گردد. اگر شکست را زمان مرگ بیمار در نظر بگیریم، در صورتی که زمان دقیق شکست مشخص نباشد و یا اطلاعات جزئی در مورد زمان شکست وجود داشته باشد، آن‌گاه داده‌ها سانسور شده‌اند.

در مورد اطلاعات پزشکی سه نوع سانسور وجود دارد. [۲]

۱- **سانسور از راست**^۲: این سانسور زمانی اتفاق می‌افتد که اطلاعات پزشکی بیماران در انتهای طول دوره درمان از دست رفته است.

۲- **سانسور از چپ**^۳: این نوع سانسور به این مفهوم است که، اطلاعات پزشکی بیماران در ابتدای دوره درمان وجود ندارد.

۳- **سانسور بر روی فواصل**^۴: این نوع سانسور به این مفهوم است که، قسمتی از اطلاعات پزشکی بیمار در طول دوره درمان (یک یا چند فاصله زمانی جدا از هم) از دست رفته باشد. در صورتی که اطلاعات شامل زمان‌های دقیق شکست^۵ و یا سانسور باشد، برای تحلیل زمان بقاء، چندین روش پارامتری و نیمه پارامتری وجود دارد.

سانسور بر روی فواصل موقعی اتفاق می‌افتد که زمان دقیق شکست مشخص نباشد و شکست در یک فاصله زمانی پیوسته اتفاق افتاده باشد. یکی از هدفهای مهم مطالعات پزشکی، کوتاه کردن طول فاصله سانسور است. در صورتی که اطلاعات بر روی فواصل سانسور شده باشند، چندین روش آماری برای تحلیل بقاء، وجود دارد.

ترنبال^۶ در سال ۱۹۷۶، یک الگوریتم تکرار شونده، برای تحلیل بقاء با داده‌های سانسور شده بر روی فواصل را ابداع نمود [۳] و در ادامه فین کلشتاین^۷ و ولف^۸ در سال ۱۹۸۵ یک مدل رگرسیون نیمه پارامتری^۹ را بر اساس تجزیه توزیع توأم زمان بقاء و متغیرها فرمول بندی نمودند. [۴]

1-Censored Data

2-Right-Censored

3-Left-Censored

4-Interval-Censored

5-Faiure Time

6- Turnbull

7- Finkelstein

8- Wolfe

9- Semiparametric Regression

همچنین فین کلشتاین در سال ۱۹۸۶ روشی را برای تحلیل بقاء در مدل شکست متناسب با داده‌های سانسور شده بر روی فواصل، ابداع کرد. [5]

در ادامه راخر^۱ و مسرر^۲ در سال ۱۹۸۸ با استفاده از الگوریتم تکرار معرفی شده توسط ترنبال، زمان بهبود بیماران لوکمیا^۳ را برآورد نمودند. [6]

ادل اندرسون^۴ و آگوستینو^۵ در سال ۱۹۹۲، برای تحلیل بقاء با داده‌های سانسور شده بر روی فواصل، از مدل زمان شکست وایبال^۶ استفاده نمودند. [7]

همچنین کیم^۷ و تیلور^۸ در سال ۱۹۹۴ سعی کردند که زمان تشخیص بیماری ایدز^۹ را با استفاده از الگوریتم تکرار معرفی شده توسط ترنبال، دنبال نمایند. [8]

با وجود این که اکثر اطلاعات پزشکی بیماران بر روی فواصل، سانسور شده می‌باشد، اما با توجه به مقالات موجود واضح است که تاکنون تحقیقات بسیار کمی در مورد تحلیل بقاء بر روی اطلاعات سانسور شده روی فواصل، انجام شده است.

دلیل اصلی کمی پیشرفت در تحلیل بقاء اطلاعات سانسور شده روی فواصل را می‌توان پیچیدگی محاسباتی و همچنین عدم پشتیبانی کامل نرم افزارهای آماری این نوع تحلیل‌ها دانست.

هدف اصلی این رساله، معرفی و فرمول بندی یک مدل رگرسیونی در تحلیل بقاء با داده‌های سانسور شده بر روی فواصل و همچنین یک روش نسبتاً ساده برای برآورد پارامترهای مدل است.

در فصل دوم ابتدا به معرفی برخی تعاریف و اصطلاحات مورد نیاز می‌پردازیم و در ادامه، پیش‌نیازها و روابطی را که در طول رساله مورد نیاز است معرفی می‌کنیم.

در فصل سوم، مدل‌های رگرسیونی زمان شکست را بر اساس توزیع‌های وایبال و نمایی، معرفی و در ادامه، پارامترهای مدل را بر اساس داده‌های سانسور شده، برآورد می‌کنیم [9] و در پایان، مدل شکست متناسب^{۱۰} با داده‌های سانسور شده بر روی فواصل را معرفی و پارامترهای مدل را با استفاده از الگوریتم نیوتن رافسون^{۱۱} برآورد می‌کنیم. [۱۰]

در فصل چهارم، ابتدا الگوریتم EM^{۱۲} را معرفی و در ادامه از این الگوریتم برای برآورد پارامترهای مدل (به روش حداکثر درست‌نمایی) با داده‌های گم شده، استفاده می‌کنیم. در ادامه از الگوریتم محدود EM در برآورد پارامترها با در نظر گرفتن محدودیت خطی $\underline{A}\theta = \underline{a}$ استفاده می‌کنیم. (A یک ماتریس معلوم

1-Rücker
2-Messerer
3-Levkemia patients
4-Odell
5-Agostino
6-Weibull failure time

7-Kim
8-Taylor
9-AIDS
10-Proportional hazard model
11-Newton-Raphson algorithm
12- EM-algorithm

$q \times p$ بعدی و a بردار معلوم $q \times 1$ بعدی و θ بردار پارامتری $P \times 1$ بعدی است. اگر θ_j عنصر j ام از بردار پارامتری θ باشد، در آن صورت از الگوریتم محدود EM برای آزمون $H_0: \theta_j = \theta_0$ (θ_0 مقدار ثابت) و ساختن فواصل اطمینان برای پارامترهای θ ($j=1,2,\dots,p$) استفاده می‌کنیم. [۱۱]

در فصل پنجم، یک روش تقریباً ساده برای تخمین زمان بقاء با داده‌های سانسور شده روی فواصل مطرح و یک مدل رگرسیونی برای داده‌های بقاء سانسور شده بر روی فواصل را پایه ریزی می‌کنیم. در ادامه از الگوریتم EM که از مشخصه‌های آن سادگی، ثبات و خاصیت همگرایی مطلوب است، برای برآورد پارامترهای مدل رگرسیونی استفاده می‌کنیم. سپس از الگوریتم محدود EM برای آزمون فرض و ساختن فواصل اطمینان روی پارامترهای مدل استفاده می‌کنیم. در انتهای فصل نیز به عنوان مثال کاربردی، به تحلیل بقاء در بیماری ملانوما^۱ می‌پردازیم. [۱۲]

در پایان رساله نیز به معرفی اجمالی از الگوی خطی تعمیم یافته پرداخته و برنامه‌های کامپیوتری مورد استفاده در رساله را می‌آوریم.

فصل دوم

تعاریف و مفاهیم اولیه

۱.۲. مقدمه

در این فصل ابتدا به معرفی برخی اصطلاحات آماری و پیش‌نیازهایی که در طول رساله از آنها استفاده شده، می‌پردازیم، سپس روابطی را که در تجزیه و تحلیل بقا از آنها استفاده شده، به دست می‌آوریم.

۲.۲. تعاریف و نمادها

همان طور که در فصل اول ذکر شد، تجزیه و تحلیل داده‌های بقا به مفهوم بررسی داده‌هایی است که با طول عمر موجودات جاندار سرو کار دارد. اما قبل از پرداختن به تحلیل رگرسیونی بر روی داده‌های بقا ابتدا به ذکر تعاریف لازم می‌پردازیم.

۱.۲.۲. تعریف قابلیت اعتماد

قابلیت اعتماد^۱ به عنوان مشخصه‌ای از یک موجود، عبارت از احتمال کارکرد (زندگی) موجود تحت شرایط معین، در فاصله زمانی معین و بدون خرابی برای انجام یک ماموریت خاص است. [۱۳]

اگر قابلیت اعتماد را با $R(t)$ نشان دهیم، داریم:

$$R(t) = P(\text{موجود در فاصله } (0, t) \text{ تحت شرایط لازم از انجام کار معین باز نماند}) \quad (1-2)$$

۲.۲.۲. تعریف تابع بقا

تابع بقا^۲ به مفهوم احتمال زنده بودن یک موجود زنده در زمان معین است و آن را با $S(t)$ نشان می‌دهیم. [۳]

$$S(t) = P(\text{زنده باشد}) \quad (\text{موجود در فاصله } (0, t)) \quad (2-2)$$

در رابطه با اطلاعات پزشکی بیماران، از اصطلاح شکست برای بیمارانی که شیوه درمان در بهبود وضعیت بالینی آنها مؤثر نبوده، استفاده می‌گردد.

عدم شکست نیز به این مفهوم است که یا شیوه درمان در بهبود وضعیت بالینی مؤثر بوده و یا بیمار تا زمان t زنده می‌باشد و اطلاعات پزشکی بیمار از زمان t به بعد در دسترس نیست. یعنی این اطلاعات در زمان t از راست سانسور شده است.