



پایان نامه

برای دریافت درجه کارشناسی ارشد در رشته آمار ریاضی

عنوان

تأثیر نمونه‌گیری در طول اربب بر متغیرهای گلی

در برآورد تابع بقا

استاد راهنما

جناب آقای دکتر وحید فلور بافنده

استاد مشاور

جناب آقای دکتر حسنعلی آذرنوش

نگارنده

پطروس اصغری

خرداد ماه ۱۳۸۹

این پایان نامه حاصل گوشه نظر پروردگار مهربان و همراهی و تشویق بی دریغ پدر و مادر عزیزم می باشد.

و نیز بدون کمک‌های فراوان آقایان دکتر فکور و دکتر آذرنوش هیچگاه به سرانجام نمی‌رسیدم. امیدوارم توانسته باشم گوشه‌ای از زحمات آن عزیزان را جبران کرده باشم.

این پایان نامه را تقدیم می‌کنم به برادرم حمیدرضا:

خون کدام عاشق بیدل

نفرین کدام پیر بلاکش

ترا به باد داد ای باغ سبز

چنین گفت سیب سرخ

پطروس اصغری

خرداد ماه یکهزار و سیصد و هشتاد و نه

## اطلاعات مربوط به چکیده پایان نامه



شماره دانشجو : ۸۶۱۳۳۹۵۰۲۳ نام و نام خانوادگی : پطروس اصغری  
دانشکده : علوم ریاضی رشته : آمار ریاضی  
مقطع : کارشناسی ارشد دوره : روزانه

### عنوان پایان نامه : تأثیر نمونه گیری در طول اریب بر متغیرهای کمکی در برآورد تابع بقا

تاریخ دفاع : ۱۳۸۹-۳-۳۱

تعداد صفحات : ۱۰۰

کلمات کلیدی : در طول اریب ; برآورد تابع بقا ; برش چپ ; سانسور از راست ; سانسور حاصلضربی

استاد راهنما : وحید فکور بافنده

استاد مشاور : حسنعلی آذرنوش

### چکیده :

تحلیل داده های بقای سانسور راست و بریده شده از چپ در بسیاری از تالیفات آماری به چشم میخورد . اهمیت تحلیل این نوع داده ها از وسعت کاربرد آنها در عمل سرچشمه می گیرد که در بررسی های پزشکی در رابطه با بیماری هایی از جمله ایدز و دمانس کاربرد فراوان دارد. نکته ای که در مورد این چنین بیماری هایی پیش می آید، این است که زمان دقیقی یا حتی تخمینی شروع بیماری در دست نیست و بیماران معمولاً در اثر حادثه یا کاملاً اتفاقی متوجه بیماری خود می شوند چرا که بیماری هایی مانند ایدز یا تومورهای خوش خیم، مدت نسبتاً زیادی در بدن فرد باقی می مانند تا به مرحله نمود خود برسند و نشانه های خود را آشکار سازند. حتی بیماری هایی نظیر دمانس که ممکن است اطرافیان زود به آن پی ببرند فرد از زمان شروع بیماری اطلاعی ندارد . موضوعی که کمتر به آن توجه می شود، اریب شدن نتیجه حاصل از تحلیل این گونه داده ها می باشد، چرا که مشاهده می شود که افرادی که وارد نمونه می شوند، نسبت به بقیه افراد جامعه دارای طول عمر بیشتری هستند . این اریبی در اصطلاح، در طول اریبی نامیده می شود. حال در کنار داده های اصلی، متغیرهای کمکی نیز می توانند وارد مدل شوند. در این پایان نامه به بررسی چگونگی ورود و نیز نحوه تأثیر در طول اریبی بر متغیرهای کمکی پرداخته می شود . اما هدف اصلی، چگونگی رفتار با این متغیرها در برآورد پارامتری تابع بقا می باشد . در این راستا دو نوع تابع درستنمایی معرفی می شود و از روش درستنمایی ماکسیمم، پارامترها برآورد می شوند . این دو نوع تابع درستنمایی بر پایه دو دیدگاه معرفی می شوند. دیدگاه اول از نادیده گرفتن اطلاعات موجود در متغیرهای کمکی تأثیر می گیرد و تابع درستنمایی را شرطی بر روی آنها معرفی می کند. در دیدگاه دوم تابع درستنمایی به صورت توأم با متغیرهای کمکی تولید می شود. در پایان نتیجه می شود که این دو تابع درستنمایی در حجم نمونه زیاد، تفاوت چندانی با هم نداشته، حال آنکه در حجم نمونه کم، برآوردگرهای حاصل از تابع درستنمایی توأم کارایی بیشتری نسبت به حالت شرطی دارند

## **Thesis Abstract Information**

**Name:** petros asghari

**Student No.:** 8613395023

**Field:** Mathematical Statistics

**Faculty:** Mathematics

**Section:** M.A.



**Thesis title :** The effect of Length Biased sampling on Covariates in estimating the survival function

**Date of thesis defense :** 31-3-1389

**Page No. :** 100

**Keywords :** length bias ; estimating the survival function ; left truncation ; right censored ; multiplicative censoring

---

**Thesis Supervisor :** Vahid Fakoor Bafandeh

**Thesis Advisor :** Hasan Ali Azarnoosh

---

**Abstract :**

Analysis of the right censored and left truncated survival data is seen in many statistical publications. The importance of the analysis of these data is a result of the wide usage of them in practice especially in medical surveys such as studies on some diseases such as AIDS and Dementia. The point about these kinds of data is this that the initiating time of the disease is not known and the patients would accidentally determine that they have the disease. Because these kinds of diseases stay potential in the patient's body for almost a long time until they show their signs and infect the person. Even diseases such as Dementia, which the relatives may recognize it soon, there is no accurate information about the time when the disease has started. The subject that is ignored mostly is this, that these kinds of data are biased. As a matter of fact, it can be seen that the people whom they are in the sample, have a longer life time in comparison to whom they are not or cannot be in the sample. This biasness in statistics is called Length Bias. Beside these data, covariates can be entered into the model. In this thesis, the most concern is on how to enter covariates into the model, and how to make behave with these covariates in estimating the survival function. In this sequel, two kinds of likelihood functions are introduced on the basis of two viewpoints. In the end, it is concluded that these two types of likelihood functions, do not differ a lot when dealing with huge sample volumes. Although in low sample volumes the estimate of the survival function of one of the likelihood functions is more efficient. And that function is the function that takes into account the covariates.

---

## فهرست مندرجات

صفحه	عنوان
۱	پیشگفتار
۴	<b>فصل ۱ تعاریف و مقدمات</b>
۵	۱-۱ مقدمه ای در مفاهیم بقا
۱۱	۲-۱ مدل های پارامتری رایج برای داده های بقا
۱۷	۳-۱ مقدمه ای بر سانسور و برش
۱۸	۴-۱ سانسور راست و انواع آن
۲۱	۵-۱ سانسور تصادفی
۲۲	۶-۱ سانسور چپ و انواع آن
۲۳	۷-۱ سانسور فاصله ای
۲۳	۸-۱ برش و انواع آن
۲۵	۹-۱ ساختار درستی برای داده های سانسور شده و بریده شده
۲۷	<b>فصل ۲ در طول اریبی</b>
۲۸	۱-۲ مفهوم در طول اریبی
۳۲	۲-۲ فرضیه ی مانایی
۳۴	<b>فصل ۳ متغیرهای کمکی</b>
۳۵	۱-۳ روش های ورود متغیر های کمکی به مدل
۳۸	۲-۳ تأثیر در طول اریبی بر متغیرهای کمکی

۴۰	فصل ۴ برآورد پارامتری تابع بقا در طول اریب با وجود متغیرهای کمکی
۴۱	۱-۴ مقدمه
۴۱	۲-۴ مدل حاصلضربی
۴۲	۳-۴ سانسور آگاهی بخش
۴۲	۴-۴ برآورد در فرآیندهای تجدید
۴۵	۵-۴ معرفی درستنمایی مناسب مدل
۴۷	۶-۴ رابطه بین $L_I$ و $L_r$
۴۸	۷-۴ مقایسه بین $L_I$ و $L_r$
۵۳	۸-۴ کارآیی $\hat{\theta}_{r_n}$ نسبت به $\hat{\theta}_{I_n}$
۵۸	۹-۴ خواص مجانبی برآوردگرهای ماکسیمم درستنمایی
۷۷	۱۰-۴ مثالی عملی از تحلیل بقا در بیماران دمانس
۸۱	۱۱-۴ بوت استرپ نیمه پارامتری و کارآیی ها
۸۳	۱۲-۴ نتیجه گیری
۸۶	واژه نامه فارسی به انگلیسی
۸۸	واژه نامه انگلیسی به فارسی
۹۰	کتابنامه

## نمادها

$F_{LB}$	تابع توزیع در طول اریب
$\wedge$	کمینه بین دو مقدار
$\propto$	متناسب بودن با
$\xrightarrow{a.s.}$	همگرایی قریب به یقین
$\xrightarrow{D}$	همگرایی در توزیع

## پیشگفتار

تحلیل داده های بقای سانسور راست و بریده شده از چپ در بسیاری از تا لیفات آماری به چشم می خورد. اهمیت تحلیل این نوع داده ها از وسعت کاربرد آنها در عمل سرچشمه می گیرد که در بررسی های پزشکی در رابطه با بیماری هایی از جمله ایدز و دمانس کاربرد فراوان دارد. نکته ای که در م ورد این چنین بیماری هایی پیش می آید، این است که زمان دقیق یا حتی تخمینی شروع بیماری در دست نیست و بیماران معمولاً در اثر حادثه یا کاملاً اتفاقی متوجه بیماری خود می شوند چرا که بیماری هایی مانند ایدز یا تومور های خوش خیم، مدت نسبتاً زیادی در بدن فرد باقی می ماند تا به مرحله نمود خود برسند و نشانه های خود را آشکار سازند. حتی بیماری هایی نظیر دمانس که ممکن است اطرافیان زود به آن پی ببرند فرد از زمان شروع بیماری اطلاعی ندارد.

مثال دیگر، در مورد تحلیل اطلاعات مربوط به معتادین می باشد. مشکلی که در اینجا پیش می آید، این است که افراد مورد بررسی تمایل زیادی به بیان نحوه و چگونگی ابتلا و زمان شروع اعتیاد، ندارند. معمولاً برای جمع آوری اطلاعات مربوط به این موضوع ها، به آسایشگاه ها و یا بیمارستان های محل نگهداری آنها مراجعه می شود. مشکل بعدی در اینجا پیش می آید که زمان های شروع بیماری ها یا به طور کلی زمان های تجربه پیشامد مورد بررسی ما توسط افراد، با هم متفاوت است و نیز اینکه عده بسیاری قبل از ورود به آسایشگاه یا بیمارستان دچار مرگ شده یا پیشامد پای ان دهنده مطلوب ما که می تواند مرگ یا بهبود بیماری باشد را تجربه کرده اند.

این مشکلات و بسیاری دیگر، اهمیت تحلیل و کار بر روی اینگونه داده ها را روشن می کند و باعث می شود که روشی خاص برای جمع آوری و تحلیل اینگونه داده ها به کار بریم.

یک روش منطقی در جمع آوری اینگونه داده، تقسیم بندی جامعه به سه قسمت می باشد.

**۱- جامعه هدف:** این جامعه شامل تمام افرادی است که ویژگی مورد نظر ما را دارا می باشند و

هدف تحقیق و تعمیم نتایج به این جامعه است.

**۲- جامعه نمونه گیری:** چنانچه نمونه گیری از جامعه هدف ممکن نباشد، از جامعه ای که قابل

نمونه گیری است (و مسلماً کوچکتر از جامعه هدف می باشد) نمونه گیری می کنیم. این جامعه قابل دسترسی برای نمونه گیری را جامعه نمونه گیری گویند.

**۳- نمونه:** که از جامعه نمونه گیری به صورت زیر انتخاب می شود.



از جامعه نمونه گیری ، به صورت مقطعی در زمان کوتاهی که هیچ فرد جدیدی در آن زمان وارد جامعه نمونه گیری نشود، نمونه می گیریم. لازم به ذکر است که این روش ، تحت عنوان روش خاصی در آمار مطرح نمی شود و تنها یک راه منطقی در جمع آوری داده ها است .

بدیهی است فقط افرادی وارد نمونه می شوند که تا زمان نمونه گیری زنده مانده اند و کسانی که قبلا فوت کرده یا بهبود یافته اند وارد نمونه نمی شوند. ( دچار برش چپ می شوند ) به طور شهودی، تا اینجا دیده می شود شانس افرادی که طول عمر بیشتری دارند، برای ورود به نمونه بیشتر است و این باعث اریب شدن نتایج حاصل از تحلیل اینگونه داده ها می باشد .

اریبی که از آن یاد شد ، تحت عنوان در طول اریبی<sup>۱</sup> در مقالات آماری بیان می شود و کمتر مورد توجه قرار می گیرد، حال آنکه عدم در نظر گرفتن آن باعث تغییر بسیار زیادی در نتایج می شود و بیشتر منجر به یک بیش برآوردی می گردد به عنوان مثال، ولفسن و همکاران (۲۰۰۱) نشان داده اند که برآورد واقعی میانه زمان بقا در یک سری داده عملی (با در نظر گرفتن در طول اریبی ) ۳.۳ می باشد حال آنکه بدون در نظر گرفتن در طول اریبی این مقدار ۶.۶ برآورد شده بود. این مطلب، اهمیت در نظر گرفتن در طول اریبی را می رساند .

در زمینه تحلیل این نوع داده ها ، افرادی مانند واردی(۱۹۸۲)، واردی (۱۹۸۵) و گیل، واردی و ولنر (۱۹۸۸) و واردی (۱۹۸۹) در برآورد نا پارامتری تابع بقا با داده های در طول اریب به بررسی خواص حدی آن پرداخته اند .

همچون پاتیل و راثو (۱۹۸۷) و پاتیل، راثو و زلن (۱۹۸۸) حالات کلی تری از اریبی را در نظر گرفته اند .

اصغریان و ولفسن (۲۰۰۱) و اصغریان و ولفسن (۲۰۰۵) به خواص مجانبی برآوردگر ناپارامتری تابع بقای داده های از راست سانسور شده و در طول اریب و نااریب پرداخته اند .

همینطور، ونگک و همکاران (۱۹۹۳)، آلیوم و کامن جس (۱۹۹۶)، نان و رایان (۱۹۸۹) و گیل (۱۹۸۱) متغیر های کمکی در مدل خود وارد کرده اند.

کریستوبال و آلکلا (۲۰۰۰) و کریستوبال و همکاران (۲۰۰۴) نیز روشهایی ناپارامتری برای رگرسیون داده های در طول اریب ارائه کرده اند.

---

<sup>۱</sup> Length bias

اما موضوعی که تا کنون کمتر مورد بررسی قرار گرفته است تأثیر در طول اریبی بر متغیر های کمکی با وجود سانسور راست و برش چپ می باشد. که هدف اصلی این پایان نامه می باشد.

این پایان نامه مشتمل بر ۴ فصل می باشد.

در فصل اول تعاریف اولیه مورد نیاز را ارائه می کنیم.

در فصل دوم به مفهوم در طول اریبی و چگونگی ورود آن به مدل را مورد بررسی قرار می دهیم . در

فصل سوم چگونگی ورود متغیرهای کمکی به مدل های گوناگون داده های بقا را بیان می کنیم. همچنین نحوه

تأثیر در طول اریبی بر متغیر های کمکی مطرح می شود.

در فصل چهارم نحوه بدست آوردن برآوردگر های تابع بقا در حالت پارامتری با معرفی توابع

درست‌نمایی شرطی بر روی متغیر های کمکی و توأم با آنها را شرح می دهیم و چند خاصیت مجانبی توابع

درست‌نمایی و برآوردگرهای آنها بیان می شود کارایی دو دسته برآوردگر را در راستای یک تحقیق عملی مورد

بررسی قنار می گیرد.

# فصل اول

## تعاریف و مقدمات

- ۱-۱ مقدمه ای در مفاهیم بقا
- ۲-۱ مدل های پارامتری رایج برای داده های بقا
- ۳-۱ مقدمه ای بر سانسور و برش
- ۴-۱ سانسور راست و انواع آن
- ۵-۱ سانسور تصادفی
- ۶-۱ سانسور چپ و انواع آن
- ۷-۱ سانسور فاصله ای
- ۸-۱ برش و انواع آن
- ۹-۱ ساختار درستمایی برای داده های سانسور شده و بریده شده

### ۱-۱ مقدمه ای در مفاهیم بقا

در این بخش، تعاریف و مقدماتی در مدل بندی داده های بقا ارائه می کنیم که در فصل های آتی به آنها احتیاج خواهیم داشت.

فرض کنید  $X$  زمان تا پیشامد مشخصی باشد. این پیشامد می تواند مرگ، بروز تومور، پیشرفت یک بیماری در یک بیمار، پایان مهمات در جنگ و یا بسیاری از اتفاقات دیگری که برای یک محقق مهم است، باشد. بعلاوه این اتفاق می تواند اتفاقی مثبت مانند درمان یک بیمار یا ترک سیگار نیز باشد. به طور دقیق تر در این فصل  $X$  یک متغیر تصادفی نامنفی از یک جامعه همگن می باشد. چهار تابع می توانند بیان کننده نوع توزیع  $X$  باشند. مهم تر از همه تابع بقا می باشد که احتمال بقای یک فرد تا زمان  $x$  است، تابع نرخ خطر یا تابع خطر که شانس این است که یک فرد با سن  $x$  در یک لحظه کوتاه بعدی پیشامد مورد بررسی را تجربه کند. تابع چگالی احتمال (یا تابع جرم احتمال) در  $x$  که احتمال اتفاق افتادن پیشامد مورد بررسی در زمان  $x$  می باشد. میانگین باقیمانده طول عمر در زمان  $x$ ، که میانگین زمان تا پیشامد مورد بررسی می باشد به شرط اینکه پیشامد تا زمان  $x$  اتفاق نیفتاده باشد.

اگر هر کدام از این چهار تابع برای ما معلوم باشد، بقیه آنها را به طور یکتا می توان بدست آورد. در عمل این چهار تابع به همراه کمیت دیگری به نام تابع خطر تجمعی، برای بیان مشخصات مختلف توزیع  $X$  به کار می روند.

### ۱-۱-۱ تابع بقا

کمیت اصلی که برای بررسی زمان تا پیشامد به کار می رود تابع بقا می باشد. احتمال اینکه یک فرد (عنصر) بیش از زمان  $x$ ، عمر کند و باقی بماند را تابع بقا نامند و به صورت زیر بدست می آید:

$$S(x) = \Pr(X > x)$$

در مواردی که زمان تا نقص فنی قطعات تولیدی یا پایان ملزومات یا مهمات انجام کاری، مورد بررسی قرار می گیرد،  $S(x)$  را به عنوان یک تابع اطمینان در نظر می گیرند. اگر  $X$  یک متغیر تصادفی پیوسته باشد، آنگاه  $S(x)$  نیز یک تابع پیوسته و اکیداً نزولی است. تابع بقا، متمم تابع توزیع تجمعی می باشد. یعنی:

$$S(x) = 1 - F(x)$$

$$F(x) = \Pr(X \leq x)$$

بعلاوه تابع بقا، از انتگرال گیری تابع چگالی احتمال  $f(x)$  نیز بدست می آید. یعنی:

$$S(x) = \Pr(X > x) = \int_x^{\infty} f(t) dt$$

$$f(x) = -\frac{dS(x)}{dx}$$

نکته اینکه  $f(x)dx$  می تواند به عنوان احتمال تقریبی اینکه پیشامد در زمان  $x$  اتفاق بیفتد، در نظر گرفته شود.

نمودار تابع  $S(x)$ ، اشکال مختلفی را می تواند به خود بگیرد اما نکته مهم اینجا است که همه این اشکال، خواص اساسی یکسانی دارند. همه آنها یکنوا و ناصعودی می باشند و نیز در نقطه صفر برابر با یک بوده و با پیشرفت زمان به بی نهایت، مقدار تابع به صفر می رود.

اگر  $X$  یک متغیر تصادفی گسسته باشد، موضوع به کلی فرق کرده و از فنون متفاوتی برای رفتار با متغیرهای تصادفی گسسته استفاده می شود. متغیرهای تصادفی گسسته معمولاً در مواردی مانند وقتی که داده های حاصل از یک نمونه گیری را گرد می کنیم یا وقتی که زمان های از کار افتادن را گرد می کنیم یا هنگامی که زمان های بقا از اعداد صحیح مربوط به اجزا، بدست می آیند، نمایان می شوند.

فرض کنید  $X$  یک متغیر تصادفی گسسته بلتابع جرم احتمال زیر باشد :

$$P(x_j) = \Pr(X = x_j) \quad j = 1, 2, \dots$$

تابع بقا برای این متغیر تصادفی گسسته بدین صورت بدست می آید :

$$S(x) = \Pr(X > x) = \sum_{x_j > x} p(x_j)$$

## ۱-۱-۲ تابع خطر

یک کمیت اصلی و زیر بنایی در تحلیل بقا، تابع خطر می باشد. این تابع به نرخ شکست شرطی نیز در مباحث قابلیت اعتماد، معروف می باشد و یا در جمعیت شناسی به شدت میرایی و در فرآیندهای تصادفی به تابع شدت و در علم همه گیرشناسی به نرخ شکست و معکوس نسبت میل در اقتصاد معروف می باشد. همه اینها در حقیقت همان نرخ شکست می باشد که به صورت زیر تعریف می شود :

$$h(x) = \lim_{\Delta x \rightarrow 0} \frac{\Pr(x \leq X < x + \Delta x \mid X \geq x)}{\Delta x}$$

اگر  $X$  یک متغیر تصادفی پیوسته باشد، در این صورت :

$$h(x) = \frac{f(x)}{S(x)} = -\frac{d \ln[S(x)]}{dx}$$

یک کمیت وابسته به نرخ خطر، تابع خطر تجمعی می باشد که به صورت زیر تعریف می شود :

$$H(x) = \int_0^x h(u) du = -\ln[S(x)]$$

پس برای طول عمرهای پیوسته داریم:

$$S(x) = \exp[-H(x)] = \exp\left[-\int_0^x h(u) du\right]$$

از تعریف تابع نرخ خطر، می‌توان مشاهده کرد که  $h(x)\Delta x$ ، احتمال تقریبی این است که یک فرد با سن  $x$  پیشامد مورد بررسی را در لحظه ای بعد، تجربه کند. این تابع به خصوص در مشخص کردن توزیع شکست مناسب با استفاده از اطلاعات کیفی در مورد نوع شکست علی‌مرگ مورد استفاده قرار می‌گیرد.

بر خلاف تابع بقا، برای تابع نرخ خطر اشکال مختلفی ممکن است بدست آید. نرخ خطر می‌تواند صعودی، نزولی و یا حتی ثابت و یا غیر یکنوا باشد و نیز می‌تواند مشخصات دیگری را که نحوه مرگ یا شکست را بیان می‌کنند، به خود بگیرد.

مدل‌هایی با نرخ خطر صعودی، ممکن است بیان‌کننده این باشد که با یک سیربقای طبیعی سر و کار داریم. مانند وقتی پیشامد مورد بررسی ما بررسی مرگ در انژیپیری در بین جوامع انسانی باشد.

توابع خطر نزولی، در عمل کمتر پیش می‌آیند و ممکن است موقعی پیش آید که احتمال شکست یا مرگ یا به طور کلی احتمال اتفاق پیشا مد مورد بررسی، بسیار زیاد باشد. به عنوان مثال در بعضی لوازم الکتریکی بسیار حساس که با کوچکترین نوسانی در برق می‌سوزد یا در بعضی موارد از عمل جراحی پیوند اعضا در بیماران، که احتمال عفونت بعد از عمل زیاد باشد، پیش می‌آید. تابع نرخ خطر ثابت در مواردی مشابه بعضی لوازم تولیدی که ممکن است در اثر نقص در بعضی از قسمت‌هایشان، به پلیمان عمر خود برسند، پیش می‌آید.

در انتها نرخ خطری که زود صعود کند و یک باره شروع به نزول کند، یعنی جهت نقر آن همیشه رو به پایین باشد، بیشتر در مواردی مانند مدل بندی داده‌های بقا مربوط به اعمال جراحی که بعد از مدتی از عمل، فرد بیمار دچار ریسک زیادی می‌شود که می‌تواند در اثر عفونت، یا یک عمل جراحی دیگر باشد، به وجود می‌آید.

**مثال ۱-۱-۱** توزیعی که به اندازه کفای انعطاف پذیر است و می‌تواند دارای نرخ خطر صعودی، نزولی یا ثابت باشد، توزیع وایبل می‌باشد. توزیع وایبل زیر را در نظر بگیرید:

$$f_x(x) = \lambda \alpha x^{\alpha-1} e^{-\lambda x^\alpha} \quad \alpha > 0, \lambda > 0, x > 0$$

در این صورت تابع بقا بدین صورت می‌شود:

$$S(x) = \exp(-\lambda x^\alpha)$$

و نیز تابع نرخ خطر نیز بدین صورت می‌شود:

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

مشاهده می‌شود که تابع برای مقادیر مختلف پارامتره‌ایش، بسیار انعطاف پذیر می‌باشد. □  
 وقتی  $X$  یک متغیر تصادفی گسسته باشد، تابع نرخ خطر به صورت زیر تعریف می‌شود:

$$P(x_j) = \Pr(X = x_j) \quad j = 1, 2, \dots$$

$$h(x_j) = \Pr(X = x_j | X > x_j) = \frac{p(x_j)}{S(x_{j-1})} \quad j = 1, 2, \dots$$

که در آن،  $S(x_0) = 1$ .

رابطه دیگری برای تابع نرخ خطر متغیرهای تصادفی گسسته بدین صورت بدست می‌آید:

$$h(x_j) = 1 - \frac{S(x_j)}{S(x_{j-1})} \quad j = 1, 2, \dots$$

نکته اینکه تابع بقای متغیر تصادفی گسسته را می‌توان به صورت ضرب احتمالات بقای شرطی نیز

نوشت:

$$S(x) = \prod_{x_j \leq x} \frac{S(x_j)}{S(x_{j-1})}$$

پس تابع بقای مربوط به متغیر تصادفی گسسته را می‌توان بدین صورت نیز نوشت:

$$S(x) = \prod_{x_j \leq x} [1 - h(x_j)]$$

### ۳-۱-۱ تابع میانگین باقیمانده عمر و میانه طول عمر

چهارمین کمیت مهم و اصلی در تحلیل داده‌های بقا، تابع میانگین باقیمانده عمر در زمان  $x$  است. برای

هر فرد با طول عمر  $x$ ، این تابع میزان انتظار از باقیمانده طول عمر او را اندازه‌گیری می‌کند.

$$mrl(x) = E(X - x | X > x)$$

به راحتی می‌توان بدست آورد که میانگین باقیمانده طول عمر، ناحیه زیر منحنی بقا در سمت راست نقطه  $x$

تقسیم بر  $S(x)$  می‌باشد.

در این راستا، میانگین طول عمر را با نماد  $\mu$  نشان داده و برابر است با:

$$\mu = mrl(0)$$

که این مقدار تمام ناحیه زیر منحنی بقا می باشد.

برای یک متغیر تصادفی پیوسته داریم:

$$mrl(x) = \frac{\int_x^{\infty} (t-x)f(t)dt}{S(x)} = \frac{\int_x^{\infty} S(t)dt}{S(x)}$$

$$\mu = E(X) = \int_0^{\infty} tf(t)dt = \int_0^{\infty} S(t)dt$$

و همین طور واریانس  $X$  بدین صورت بدست می آید:

$$Var(X) = 2 \int_0^{\infty} tS(t)dt - \left[ \int_0^{\infty} S(t)dt \right]^2$$

و در همین راستا، چندک مرتبه  $p$  ام توزیع  $X$ ، کوچکترین  $x_p$  می باشد که در رابطه زیر صدق کند.

$$S(x_p) \leq 1-p$$

یعنی:

$$x_p = \inf \{t : S(t) \leq 1-p\}$$

اگر  $X$  یک متغیر تصادفی پیوسته باشد، چندک مرتبه  $p$  ام از حل معادله زیر بدست می آید:

$$S(x_p) = 1-p$$

و میانه طول عمر، در حقیقت همان چندک مرتبه ۰.۵ ام می باشد. یعنی میانه طول عمر یک متغیر

تصادفی پیوسته باید در معادله زیر صدق کند:

$$S(x_{0.5}) = 0.5$$

## ۲-۱ مدل های پارامتری رایج برای داده های بقا

با اینکه روشهای ناپارامتری و نیمه پارامتری در تحلیل داده های بقه بسیار پر کاربرد هستند، مدل های پارامتری نیز در خور توجه می باشند. در این بخش نگاهی به مدل ها و توزیع هایی که بیشتر برای داده های بقا مورد استفاده می گیرند و دلایل استفاده از آنها، می اندازیم. این مدل ها نه به خاطر محبوبیت و شهرت در بین محققان بلکه به دلیل خواصی که به توابع مربوط به بقا و نیز پارامترهای مختلف جامعه، القا می کنند، مورد استفاده زیادی دارند. که مهمترین این توابع تابع نرخ خطر می باشد. بعضی از مدل های مهمی که بیشتر مورد



استفاده قرار می گیرند عبارتند از توزیع های نمایی، وایبل، گاما، لگ نرمال، لگ لجستیک، نرمال، توانی نمایی، گمپرتز، گوسی معکوس، پارتو و گامای تعمیم یافته.

ابتدا به دلیل سادگی کار توزیع نمایی را مورد بررسی قرار می دهیم. تابع بقای توزیع نمایی را مورد بررسی قرار می دهیم.

## ۱-۲-۱ توزیع نمایی

تابع چگالی و نرخ خطر آن بدین صورت می باشد:

$$f_X(x) = \lambda \exp(-\lambda x) \quad \lambda, x > 0$$

$$h(x) = \lambda$$

مهمترین خاصیت توزیع نمایی، خاصیت بی حافظه بودن آن می باشد. یعنی:

$$\Pr(X \geq x + z | X \geq x) = \Pr(X \geq z)$$

این خاصیت در بسیاری از موارد خاصیت مثبت و مطلوبی به حساب می آید حال آنکه در بسیاری از موارد عملی، کاربرد ندارد. مانند طول عمر قطعات صنعتی یا حتی طول عمر انسان ها. به دلیل خاصیت بی حافظه بودن توزیع نمایی، میانگین باقی مانده عمر در این توزیع مقداری ثابت می شود:

$$E(X - x | X > x) = E(X) = \frac{1}{\lambda}$$

این اتفاق به این دلیل است که در توزیع نمایی، زمان تا اتفاق افتادن یک پیشامد، به طول عمر قبل از این زمان بستگی ندارد. این خاصیت توزیع نمایی، در تابع نرخ خطر آن نیز مشهود می باشد. در اینجا احتمال شرطی شکست در هر زمان، به شرط اینکه پیشامد تا زمان  $t$  اتفاق نیفتاده باشد، به زمان  $t$  بستگی ندارد. با اینکه توزیع نمایی توزیع محبوبی می باشد، نرخ خطر ثابت این توزیع، محدودیت بسیاری را در کارهای عملی در صنعت و بهداشت تحمیل می کند.

توزیع نمایی حالت خاصی از دو توزیع گاما و وایبل می باشد.

## ۲-۲-۱ توزیع وایبل

توزیع وایبل اولین بار توسط وایبل در سال ۱۹۳۹ معرفی شد. وی این توزیع را برای طول عمر مواد شیمیایی به کار می برده است. در توزیع وایبل با شکل زیر،  $\lambda > 0$  یک پارامتر مقیاس و  $\alpha > 0$  یک پارامتر شکل می باشد.

$$f_X(x) = \alpha \lambda x^{\alpha-1} \exp(-\lambda x^\alpha) \quad \alpha, \lambda, x > 0$$

یکی از خواص مهم توزیع وایبل، تابع نرخ خطر انعطاف پذیر آن می باشد.

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

همانطور که دیده می شود این تابع می تواند صعودی ( $\alpha > 1$ )، نزولی ( $\alpha < 1$ ) و یا ثابت ( $\alpha = 1$ ) باشد. این خاصیت تابع نرخ خطر در تابع بقا نیز دیده می شود.

$$S(x) = \exp(-\lambda x^\alpha)$$

گاهی اوقات به جای استفاده از خود طول عمرها، از لگاریتم این طول عمرها استفاده می شود. این کار در بعضی موارد به دلیل سادگی کار انجام می شود. اگر قرار دهیم  $Y = \ln(X)$  که  $X$  یک متغیر تصادفی وایبل است، در این صورت  $Y$  دارای تابع چگالی زیر می شود:

$$f_Y(y) = \alpha \exp \left\{ \alpha \left[ y - \left( \frac{-\ln(\lambda)}{\alpha} \right) \right] - \exp \left\{ \alpha \left[ y - \left( \frac{-\ln(\lambda)}{\alpha} \right) \right] \right\} \right\}, -\infty < y < \infty$$

### ۱-۲-۳ توزیع لگ نرمال

متغیر تصادفی  $X$  را دارای توزیع لگ نرمال گوئیم هرگاه لگاریتم طبیعی آن  $Y = \ln(X)$  دارای توزیع نرمال باشد. برای داده های زمان تا یک پیشامد مطلوب، این توزیع به دلیل رابطه ای که با توزیع نرمال دارد، مورد استفاده زیاد قرار می گیرد. دلیل دیگر کاربرد زیاد این توزیع این است که بسیاری از محققین، مشاهده نموده اند که این توزیع زمان های بقا یا طول عمر در زمان شروع چند بیماری مشخص را تقریب می زند.

$$f_X(x) = \frac{\exp \left[ -\frac{1}{2} \left( \frac{\ln x - \mu}{\sigma} \right)^2 \right]}{x (2\pi)^{1/2} \sigma}, \mu, \sigma, x > 0$$

همانند توزیع نرمال، توزیع لگ نرمال هم کاملاً با دو پارامتر  $\mu, \sigma$  مشخص می شود.

$$S(x) = 1 - \Phi \left( \frac{\ln x - \mu}{\sigma} \right)$$

$$h(x) = \frac{f_X(x)}{S(x)}$$

که  $\Phi$  تابع توزیع چگالی نرمال می باشد.

مقدار تابع نرخ خطر در صفر، صفر می باشد و در ابتدا صعود می کند و بعد با میل کردن  $x$  به بینهایت، به صفر نزول می کند. به دلیل اینکه تابع نرخ خطر این توزیع برای طول عمرهای زیاد کوچک می شود، این مدل در بسیاری از موارد نامناسب عمل می کند. چرا که این خاصیت تابع نرخ خطر، در عمل خیلی کم اتفاق می افتد. این مدل بیشتر در مواردی مورد استفاده دارد که مقادیر بسیار بزرگ  $X$  از اهمیت کمی برخوردار می باشند.

## ۴-۲-۱ توزیع لگ لجستیک

متغیر تصادفی  $X$ ، از توزیع لگ لجستیک تبعیت می کند، هرگاه لگاریتم طبیعی آن  $Y = \ln(X)$  دارای توزیع لجستیک باشد. این توزیع شباهت بسیاری به توزیع نرمال دارد، با این تفاوت که تابع بقای منعطف تری دارد.

$$f_X(x) = \frac{\alpha x^{\alpha-1} \lambda}{[1 + \lambda x^\alpha]^2} \quad \lambda, \alpha > 0, x \geq 0$$

$$f_Y(y) = \frac{\exp\left(\frac{y - \mu}{\sigma}\right)}{\sigma \left[1 + \exp\left(\frac{y - \mu}{\sigma}\right)\right]^2} \quad -\infty < y < \infty, \mu, \sigma > 0$$

$$h(x) = \frac{\alpha x^{\alpha-1} \lambda}{1 + \lambda x^\alpha}$$

$$S(x) = \frac{1}{1 + \lambda x^\alpha}$$

صورت کسر تابع نرخ خطر توزیع لگ لجستیک شبیه تابع خطر وایبل است و مخرج کسر آن تابع، باعث می شود که تابع دارای این مشخصات شود که، برای  $\alpha \leq 1$  نزولی و برای  $\alpha > 1$  تابع نرخ خطر تا مقدار ماکسیمم اش در زمان  $\left(\frac{\alpha-1}{\lambda}\right)^{1/\alpha}$  صعود کند و بعد به صورت نزولی هنگامی که زمان به بینهایت میل می کند، به سمت صفر رود.

این توزیع شریبه به مدل های نمایی و وایبل می باشد که این شباهت در تابع بقا و تابع نرخ خطر این توزیع ها، مشاهده می شود.

## ۵-۲-۱ توزیع گاما

توزیع گاما دارای مشخصاتی شبیه به توزیع وایبل می باشد با اینکه به اندازه توزیع وایبل انعطاف پذیر نمی باشد. این توزیع همانند وایبل، تعمیمی از توزیع نمایی می باشد. اگر  $X$  دارای توزیع گاما باشد، دارای تابع چگالی احتمال زیر می باشد:

$$f_X(x) = \frac{\lambda^\beta x^{\beta-1} \exp(-\lambda x)}{\Gamma(\beta)} \quad \beta, \lambda > 0, x \geq 0$$

وقتی  $\beta$ ، پارامتر شکل توزیع گاما به سمت بینهایت میل کند، این توزیع به توزیع نرمال میل می کند.

$$h(x) = \frac{f_X(x)}{S(x)} = \frac{\lambda^\beta x^{\beta-1} \exp(-\lambda x)}{\Gamma(\beta)} \cdot \frac{\Gamma(\beta)}{\int_0^{\lambda x} u^{\beta-1} \exp(-u) du}$$

این تابع نرخ خطر برای  $\beta > 1$  صعودی می باشد با:

$$h(0) = 0$$

$$h(x) \xrightarrow{x \rightarrow \infty} \lambda$$

و نیز برای  $\beta < 1$  نزولی می باشد با:

$$h(0) = \infty$$

$$h(x) \xrightarrow{x \rightarrow \infty} \lambda$$

وقتی  $\beta < 1$  باشد، نمای تابع خطر در  $x = \frac{\beta-1}{\lambda}$  اتفاق می افتد.

### ۳-۱ سانسور و برش

داده‌های مربوط به زمان تا رخ دادن یک پیشامد، به راه‌ها و انواع گوناگونی نمود پیدا می کنند که باعث مشکلی کار با آنها می شود. یکی از مشخصه‌های این نوع داده‌ها که بعضی اوقات اتفاق می افتد، به سانسور معروف می باشد. که به طور کلی زمانی اتفاق می افتد که می دانیم بعضی از طول عمرها، فقط در بازه‌های مشخصی اتفاق افتاده اند. بقیه طول عمرها را به طور کامل می بینیم. انواع مختلفی از سانسور را داریم، مانند سانسور راست، سانسور چپ و سانسور فاصله‌ای.

برای بررسی جامع پدیده سانسور در تحلیل داده‌های بقا، باید طرح و روشی را که برای بدست آوردن داده‌ها به کار می بریم در نظر بگیریم.

مشخصه دیگری که در مطالعات بقا بسیار پیش می آید، و با سانسور اشتباه گرفته می شود، برش می باشد. برش چپ زمانی اتفاق می افتد که افراد در یک سن خاص وارد مطالعه ما می شوند و نه الزاماً در زمان شروع طول عمرشان یا لحظه تجربه پیشامد آغازین (مانند شروع بیماری)، و از این به اصطلاح زمان تأخیر، تا