

حاشا  
الرحمن الرحيم



دانشکده علوم ریاضی

گروه آمار

پایان نامه کارشناسی ارشد

تحلیل بیزی مدل‌های نرخ شفایافتگی  
غیرآمیخته با تابع بقای نیمه پارامتری کاکس به  
کمک  $P$ -اسپلاین‌ها

مذاشمس

استادان راهنما

دکتر حسین باغی‌شینی و دکتر نگار اقبال

تیر ۱۳۹۷

تقدیم به آن ها که سرآغاز تولد من هستند. از یکی زاده می شوم و از دیگری جاودانه.  
استادی که سپیدی را بر تخته سیاه زندگیم نگاشت و پدر و مادری که تار مویی از آن ها به  
پای من سیاه نامند.

سپاس الطاف بی کران پرودگار متعال را که آغاز و پایان همه چیز است. تشکر  
و سپاس از اساتید فرزانه ام سرکار خانم دکتر نگار اقبال و جناب آقای دکتر حسین  
باغیثی که از محضر پر فیض تدریستان بهره‌های فراوان بردم. از اساتید دلسوز جناب  
آقای دکتر داوود شاهسونی و جناب آقای دکتر محمد رضا ربیعی که زحمات داورمی  
این پایان نامه را مستقبل شدند، کمال تشکر و امتنان را دارم. بر خود وظیفه می‌دانم  
از زحمات بی دریغ خانواده عزیزم و دوستان خوبم خانم هاجر هره فرهادی و محدثه جامعی  
قدر دانی نمایم که بدون کمک آن‌ها این کار به نتیجه نمی‌رسید.

نذات شمس  
تیر ۱۳۹۷

## تعمیر نامه

اینجانب ندا شمس دانشجوی کارشناسی ارشد رشته آمار ریاضی دانشکده علوم ریاضی دانشگاه صنعتی شاهرود، نویسنده پایان نامه با عنوان تحلیل بیزی مدل های نرخ شفایافتگی غیر آمیخته با تابع بقای نیمه پارامتری کاکس به کمک  $P$ -اسپلاین ها، تحت راهنمایی دکتر حسین باغیشنی و دکتر نگار اقبال متعهد می شوم:

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش های دیگر پژوهشگران، به مرجع مورد استفاده استناد شده است.
- مطالب این پایان نامه، تا کنون توسط خود، یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارایه نشده است.
- حقوق معنوی این اثر، به دانشگاه صنعتی شاهرود تعلق دارد، و مقالات مستخرج با نام “دانشگاه صنعتی شاهرود” یا “Shahrood University of Technology” به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آوردن نتایج اصلی پایان نامه تاثیرگذار بوده اند، در مقالات مستخرج از پایان نامه رعایت می گردد.
- در تمام مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت های آنها) استفاده شده است، ضوابط و اصول اخلاقی رعایت شده است.
- در تمام مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته (یا استفاده) شده است، اصل رازداری و اصول اخلاق انسانی رعایت شده است.

ندا شمس  
تیر ۱۳۹۷

## مالکیت نتایج و حق نشر

- تمام حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه های رایانه ای، نرم افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می باشد. این مطلب باید به نحو مقتضی، در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در این پایان نامه بدون ذکر منبع مجاز نمی باشد.

## چکیده

در تحلیل بقا، معمولاً فرض می‌شود که هر واحد، پیشامد مورد نظر را در صورت مشاهده طی مدت طولانی، تجربه می‌کند. با این حال می‌توان به صراحت فرض کرد که نسبت ناشناخته‌ای از جمعیت هرگز پیشامد مورد آزمایش را تجربه نخواهند کرد و به اصطلاح شفایافته محسوب می‌شوند. مدل‌های زمان پیشرفت، زیررده‌ای از مدل‌های بقای شفایافته هستند که این ویژگی را در نظر می‌گیرند. رهیافت‌های متعددی برای تعیین تابع توزیع این مدل ارائه شده‌اند که یکی از آن‌ها رهیافت مدل مخاطرات متناسب کاکس است. در این پایان‌نامه از یک رهیافت نیمه‌پارامتری برای مدل‌بندی تابع بقا استفاده شده و اثرات ناپارامتری در آن با  $P$ -اسپلین‌ها مدل‌بندی شده‌اند. هم‌چنین دیدگاه استنباط بر مبنای مکتب بیزی است. مسأله شناسایی‌پذیری مدل هنگامی که طول زمان مطالعه به اندازه کافی طولانی نباشد، مورد بحث قرار گرفته شده است. دقت مدل پیشنهادی با استفاده از شبیه‌سازی مورد ارزیابی قرار گرفته و کاربرد مدل بر اساس داده‌های حاصل از یک آزمایش بالینی نمایش داده شده است.

کلمات کلیدی: تحلیل بقا، مدل زمان پیشرفت،  $P$ -اسپلین، مدل مخاطرات متناسب کاکس، استنباط بیزی.

# فهرست مطالب

خ فهرست تصاویر

ذ فهرست جداول

۱	مقدمات و تعاریف اولیه	۱
۱	۱.۱ مقدمه	۱.۱
۲	۲.۱ تابع بقا	۲.۱
۴	۳.۱ تابع نرخ مخاطره	۳.۱
۵	۴.۱ سانسور	۴.۱
۸	۵.۱ مدل‌های رگرسیون در بقا	۵.۱
۸	۱.۵.۱ برازش مدل‌های رگرسیونی برای داده‌های بقا	۱.۵.۱
۹	۲.۵.۱ مخاطرات نسبی	۲.۵.۱
۱۱	۶.۱ مدل‌های مخاطره	۶.۱
۱۱	۱.۶.۱ مدل‌های پارامتری	۱.۶.۱
۱۳	۲.۶.۱ مدل‌های نیمه پارامتری	۲.۶.۱
۱۴	۷.۱ روش‌های برآورد در مدل‌های نیمه پارامتری	۷.۱
۱۵	۸.۱ دیدگاه بیزی در استنباط آماری	۸.۱
۱۶	۹.۱ روش‌های نمونه‌گیری MCMC	۹.۱
۱۶	۱.۹.۱ الگوریتم متروپولیس-هستینگز	۱.۹.۱
۱۷	۲.۹.۱ نمونه‌گیر گیبز	۲.۹.۱
۱۹	۲ مدل‌های نرخ شفایافتگی غیر آمیخته	۲
۱۹	۱.۲ مدل‌های نرخ شفایافتگی	۱.۲
۲۰	۲.۲ مدل نرخ شفایافتگی آمیخته	۲.۲
۲۱	۳.۲ مدل نرخ شفایافتگی غیر آمیخته	۳.۲
۲۴	۴.۲ ویژگی‌های مدل نرخ شفایافتگی غیر آمیخته	۴.۲
۲۵	۵.۲ شناسایی پذیری مدل	۵.۲

۲۹	تحلیل بیزی مدل نرخ شفایافتگی غیرآمیخته	۳
۲۹	..... مقدمه	۱.۳
۲۹	..... اسپلین	۲.۳
۳۴	..... $B$ -اسپلین	۱.۲.۳
۳۵	..... $P$ -اسپلین	۲.۲.۳
۳۶	..... برآورد تابع مخاطره پایه با $P$ -اسپلین	۳.۳
۳۷	..... برآورد تابع بقای پایه	۱.۳.۳
۳۸	..... تابع درستنمایی مدل	۴.۳
۳۸	..... تحلیل بیزی	۵.۳
۴۰	..... تولید نمونه از توزیع پسین با الگوریتم‌های $MCMC$	۶.۳
۴۱	ارزیابی عملکرد مدل پیشنهادی	۴
۴۱	..... مطالعه شبیه‌سازی	۱.۴
۴۱	..... پیگیری طولانی مدت	۱.۱.۴
۵۰	..... پیگیری کوتاه مدت	۲.۱.۴
۵۴	..... مثال واقعی	۲.۴
۵۷	..... نتیجه‌گیری و پیشنهادات	۳.۴
۵۹	برآوردگر ناپارامتری کاپلان-مایر	آ
۶۱	ب معیار تشخیص همگرایی جی‌وک برای الگوریتم‌های $MCMC$	
۶۳	پ کد تولید داده‌ها از مدل نرخ شفایافتگی زمان پیشرفت	
۶۵	مراجع	
۶۹	واژه‌نامه فارسی به انگلیسی	
۷۱	واژه‌نامه انگلیسی به فارسی	



# فهرست تصاویر

۳	یک نمودار معمول تابع بقا	۱۰۱
۴	انواع تابع مخاطره	۲۰۱
۶	علامت $X$ به معنی رخداد پیشامد است.	۳۰۱
	برآورد $S_0(t)$ در مثال شبیه‌سازی با پیگیری طولانی مدت برای ۵۰۰ مجموعه داده و $n = ۳۰۰$ . هر منحنی خاکستری برآورد تابع بقا در یک مجموعه داده است. هر ردیف شکل به درصدی از افراد شفایافته (ردیف بالا: ۲۵٪ و ردیف پایین: ۴۰٪) با نرخ کلی سانسور از راست (چپ: سناریو ۱ و راست: سناریو ۲) اختصاص دارد. منحنی خط ممتد نمودار واقعی $S_0(t)$ و منحنی خط چین میانه حاصل از ۵۰۰ تکرار است.	۱۰۴
۴۶	برآورد $S_0(t)$ در مثال شبیه‌سازی با پیگیری طولانی مدت برای ۵۰۰ مجموعه داده و $n = ۶۰۰$ . هر منحنی خاکستری برآورد تابع بقا در یک مجموعه داده است. هر ردیف شکل به درصدی از افراد شفایافته (ردیف بالا: ۲۵٪ و ردیف پایین: ۴۰٪) با نرخ کلی سانسور از راست (چپ: سناریو ۱ و راست: سناریو ۲) اختصاص دارد. منحنی خط ممتد نمودار واقعی $S_0(t)$ و منحنی خط چین میانه حاصل از ۵۰۰ تکرار است.	۲۰۴
۴۷	برآورد نسبت لگاریتم مخاطره جامعه $\log(HR_p(t))$ در مثال شبیه‌سازی با پیگیری طولانی مدت برای ۵۰۰ مجموعه داده و $n = ۳۰۰$ . هر یک منحنی خاکستری برآورد در یک مجموعه داده است. ردیف بالا به ۲۵٪ و ردیف پایین به ۴۰٪ از افراد شفایافته اشاره دارند. سمت چپ سناریو ۱ و سمت راست سناریو ۲ سانسور از راست را نشان می‌دهند.	۳۰۴
۴۸	برآورد نسبت لگاریتم مخاطره جامعه $\log(HR_p(t))$ در مثال شبیه‌سازی با پیگیری طولانی مدت برای ۵۰۰ مجموعه داده و $n = ۶۰۰$ . هر یک منحنی خاکستری برآورد در یک مجموعه داده است. ردیف بالا به ۲۵٪ و ردیف پایین به ۴۰٪ از افراد شفایافته اشاره دارند. سمت چپ سناریو ۱ و سمت راست سناریو ۲ سانسور از راست را نشان می‌دهند.	۴۰۴
۴۹	برآورد نسبت لگاریتم مخاطره افراد غیرمصون در مثال شبیه‌سازی با پیگیری طولانی مدت برای ۵۰۰ مجموعه داده و $n = ۳۰۰$ . هر یک منحنی خاکستری برآورد در یک مجموعه داده است. ردیف بالا به ۲۵٪ و ردیف پایین به ۴۰٪ از افراد شفایافته اشاره دارند. سمت چپ سناریو ۱ و سمت راست سناریو ۲ سانسور از راست را نشان می‌دهند.	۵۰۴
۵۰	چپ سناریو ۱ و سمت راست سناریو ۲ سانسور از راست را نشان می‌دهند.	۵۰

- ۶.۴ برآورد  $S_0(t)$  در مثال شبیه‌سازی با پیگیری کوتاه مدت برای  $500$  مجموعه داده و  $n = 300$ . هر منحنی خاکستری برآورد تابع بقا در یک مجموعه داده است. هر ردیف شکل به درصدی از افراد شفایافته (ردیف بالا: ۲۵٪ و ردیف پایین: ۴۰٪) با نرخ کلی سانسور از راست (چپ: سناریو ۳ و راست: سناریو ۴) اختصاص دارد. منحنی خط ممتد نمودار واقعی
- ۵۲  $S_0(t)$  و منحنی خط‌چین میانه حاصل از  $500$  تکرار است. . . . .
- ۷.۴ برآورد نسبت لگاریتم مخاطره افراد غیرمصون در مثال شبیه‌سازی با پیگیری کوتاه مدت برای  $500$  مجموعه داده و  $n = 300$ . هر منحنی خاکستری برآورد تابع بقا در یک مجموعه داده است. ردیف بالا به ۲۵٪ و ردیف پایین به ۴۰٪ از افراد شفایافته اشاره دارند. سمت چپ سناریو ۳ و سمت راست سناریو ۴ سانسور از راست را نشان می‌دهند. . . . .
- ۵۳ برآورد کاپلان-مایر منحنی بقا در مثال واقعی. نمودار سمت راست: برای همه جامعه با
- ۵۴ علامت صلیب برای افراد سانسور شده و نمودار سمت چپ: برای افراد سانسور شده. . . . .
- ۹.۴ منحنی‌های برآورد شده بقا به تفکیک گروه‌های بیمار. چپ: همه جامعه و راست: بیماران شفایافته. . . . .
- ۵۶ نمودارهای برآورد شده نسبت لگاریتم مخاطره با ناحیه اعتبار ۹۵٪ برای، چپ: همه مردان
- ۵۷ و راست: مردان شفایافته. . . . .

## فهرست جداول

۴۳	نتایج شبیه‌سازی با پیگیری طولانی مدت برای $n = 300$ و $500$ مجموعه داده . . . . .	۱۰۴
۴۴	نتایج شبیه‌سازی با پیگیری طولانی مدت برای $n = 600$ و $500$ مجموعه داده . . . . .	۲۰۴
۵۱	نتایج شبیه‌سازی با پیگیری کوتاه مدت برای $n = 300$ و $500$ مجموعه داده . . . . .	۳۰۴
	برآورد میانه‌های پسین، فاصله اعتبار ۹۵٪ و انحراف استاندارد پسین پارامترها در مثال	۴۰۴
۵۵	واقعی. . . . .	

# فصل ۱

## مقدمات و تعاریف اولیه

### ۱.۱ مقدمه

تحلیل بقا<sup>۱</sup> مجموعه‌ای از روش‌های آماری متنوع برای تحلیل متغیرهای تصادفی نامنفی است. در واقع، متغیر تصادفی  $T$ ، متغیر تصادفی بقا نامیده می‌شود، اگر تکیه‌گاه آن در بازه‌ی  $(0, \infty)$  باشد. مشاهده‌های  $t_1, t_2, \dots, t_n$  از متغیر تصادفی نامنفی  $T$  را داده‌های بقا می‌نامند.

یکی از مهم‌ترین مشخصات داده‌های بقا، زمان از کارافتادگی یک پدیده‌ی فیزیکی، مکانیکی، الکتریکی یا زمان مرگ یک موجود زنده است. در واقع متغیر پاسخ، زمان وقوع (طول عمر) یک پدیده است. در این جا به‌طور خلاصه به چند مورد از مسائل تحلیل بقا اشاره می‌کنیم:

- مطالعه چند هفته‌ای بیماران مبتلا به سرطان خون، برای این‌که مشخص شود چه زمانی طول می‌کشد تا آن‌ها بهبود یابند.
- مطالعه چند ساله افراد سالم به صورت گروهی، برای این‌که مشخص شود که آیا این افراد به بیماری قلبی مبتلا می‌شوند یا خیر.
- مطالعه بر روی زندانی‌هایی که تازه از زندان آزاد شده‌اند و می‌خواهیم ببینیم که چند هفته طول می‌کشد تا برای بار دوم دستگیر شوند.
- مطالعه بر روی بیماران قلبی که در گذشته پیوند قلب انجام داده‌اند و می‌خواهیم ببینیم که چند هفته بعد از عمل زنده می‌مانند.

تمامی مسائل ذکرشده، مسائل مربوط به تحلیل بقا می‌باشند زیرا متغیر پاسخ، زمان لازم تا رخداد پیشامد است. در مثال اول با سرطان خون مواجه هستیم و پیشامد مورد نظر، بهبودی است و متغیر پاسخ، زمان لازم (هفته) تا بهبودی فرد بیمار است. در مثال دوم پیشامد مورد نظر، بهتر شدن بیماری قلبی و پاسخ، زمان لازم تا بهبودی کامل است. در مثال سوم مطالعات روان‌شناسی بدون مطالعات پزشکی داریم که

<sup>۱</sup>Survival Analysis

در آن پیشامد مورد نظر، دستگیری مجدد و پاسخ، زمان لازم (هفته) تا بازداشت دوباره فرد است. نهایتاً در مثال چهارم همان طور که ملاحظه می‌شود پیشامد مورد نظر مرگ و متغیر پاسخ، زمان (هفته) لازم تا وقوع مرگ (تعداد هفته‌های بعد از دریافت پیوند) است. گاهی ممکن است متغیر تصادفی مورد نظر مربوط به زمان نباشد. مانند مقادیر متغیر مربوط به هزینه‌ی پرداختی یک شرکت بیمه به بیمه‌شدگان در یک وضعیت خاص. تحلیل بقا نخستین بار توسط جان گرانت (۱۶۶۲) و ستاره‌شناس معروف، ادموند هالی<sup>۲</sup> به کارهای ابتدایی روی جدول‌های مرگ و میر در قرن هفدهم انجام گرفت، اما دوره‌ی نوین آن از پنجاه سال پیش با کاربردهای مهندسی آغاز شد.

## ۲.۱ تابع بقا

در ادبیات تحلیل بقا واژگان و عبارتهای زیادی وجود دارند که دارای معانی یکسانی هستند. این اصطلاحات در عرصه‌های مختلف مهندسی و پزشکی به کار می‌روند. به عنوان مثال، طول زمان از شروع یک مشاهده تا هنگامی که پیشامد مرگ یا شکست رخ دهد، "طول عمر" یا "زمان بقا" نامیده می‌شود. به پیشامد مورد نظر عنوان "شکست" داده می‌شود زیرا نوع پیشامد مورد نظر معمولاً مرگ یا تجربه‌های منفی افراد است. البته گاهی اوقات آن پیشامد می‌تواند یک تجربه مثبت نیز باشد. مانند بازگشت به کار بعد از انجام یک عمل جراحی قلب.

فرض کنید  $T$  دارای تابع چگالی احتمال  $f$  و تابع توزیع تجمعی  $F$  باشد، در این صورت

$$F(t) = P(T \leq t) = \int_0^t f(u) du.$$

در تحلیل بقا به تابع  $F(t)$ ، تابع شکست<sup>۳</sup> می‌گویند.

تعریف ۱.۲.۱. تابع بقا (تابع قابلیت اعتماد)، برای تمام مقادیر  $t$ ، به صورت  $S(t) = 1 - F(t)$  تعریف می‌شود، به عبارت دیگر

$$P(T > t) = 1 - F(t) \equiv \bar{F}(t) \equiv S(t).$$

بنابراین

$$S(t) = P(T > t) = \int_t^{\infty} f(u) du \quad (1.1)$$

و

$$f(t) = -\frac{d}{dt} S(t).$$

<sup>۲</sup>کاشف ستاره‌ی دنباله‌دار هالی

<sup>۳</sup>Failure Function

بررسی توابع بقا از اصلی‌ترین بخش‌های تحلیل بقاست. با استفاده از رابطه‌ی (۱.۱) و هم‌چنین با توجه به این نکته که انتگرال تابع چگالی احتمال برابر یک است، داریم

$$S(0) = P(T > 0) = \int_0^{\infty} f(u)du = 1.$$

هم‌چنین

$$S(\infty) = \lim_{t \rightarrow \infty} S(t) = \lim_{t \rightarrow \infty} \int_t^{\infty} f(u)du = 0.$$

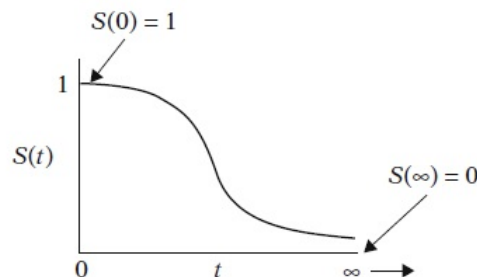
در نهایت، اگر  $a \geq b$ ، در این صورت

$$S(b) - S(a) = \int_b^a f(u)du \geq 0.$$

بنابراین

$$0 \leq S(t) \leq 1.$$

در واقع هر تابع یکنوا و نزولی  $S(\cdot)$ ، با تکیه‌گاه  $[0, \infty)$  که  $S(0) = 1$  و  $S(\infty) = 0$ ، تابع بقای یک متغیر تصادفی بقا را تعریف می‌کند (شکل ۱.۱). متغیر تصادفی مورد نظر متغیری است که دارای تابع چگالی احتمال  $f(t) = -\frac{d}{dt}S(t)$  است.



شکل ۱.۱: یک نمودار معمول تابع بقا

### ویژگی‌های تابع بقا

برای تابع بقای  $S(\cdot)$  می‌توان ویژگی‌های زیر را برشمرد:

$$\frac{dS}{dt} < 0 \quad (1)$$

$$\frac{d^2S}{dt^2} \geq 0 \quad (2)$$

$$f(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T \leq t + \Delta t)}{\Delta t} = \frac{dF(t)}{dt} = -\frac{dS(t)}{dt} \geq 0 \quad (3)$$

### ۳.۱ تابع نرخ مخاطره

همان‌طور که اشاره شد، تابع بقا  $S(t)$ ، بیانگر احتمال زنده بودن در سن  $t$ ، یا قابلیت اعتماد یک واحد در زمان  $t$  است. در این بخش مفهوم دیگری تحت عنوان تابع مخاطره مورد بررسی قرار می‌گیرد. تابع مخاطره<sup>۴</sup> عبارت است از احتمال شکست یا از کار افتادگی واحد مورد نظر در بازه  $(t, t + \Delta t)$ ، وقتی که بدانیم تا زمان  $t$  زنده بوده است. به این تابع، میزان مرگ و میر آنی، میزان شرطی از کار افتادگی یا شدت مرگ و میر نیز گفته می‌شود که به صورت

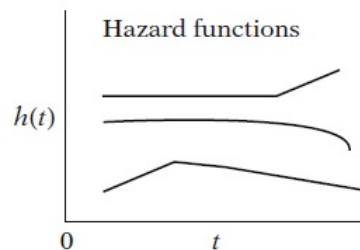
$$\begin{aligned} P(t < T < t + \Delta t | T > t) &= \frac{P(t < T < t + \Delta t, T > t)}{P(T > t)} \\ &= \frac{P(t < T < t + \Delta t)}{P(T > t)} \end{aligned}$$

تعریف می‌شود. کسر فوق را متوسط نرخ شکست در بازه  $(t, t + \Delta t)$  گویند. اگر  $\Delta t$  را به سمت صفر میل دهیم، در این صورت متوسط نرخ شکست به نرخ شکست آنی تبدیل می‌شود. این تابع می‌تواند یکنوا یا غیر یکنوا باشد (شکل ۲.۱).

تعریف ۱.۳.۱. متغیر تصادفی بقای  $T$ ، دارای تابع مخاطره یا نرخ مخاطره یا شدت مرگ و میر  $h(t)$  می‌باشد، که برای مقادیر  $t > 0$ ، به صورت

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

تعریف می‌شود.  $h(t)$  نرخ آنی مرگ در زمان  $t$ ، برای واحدهایی است که تا لحظه  $t$  زنده مانده یا از کار نیفتاده‌اند و  $f(t)$  تابع چگالی احتمال  $T$  است.



شکل ۲.۱: انواع تابع مخاطره

<sup>۴</sup>Hazard Function

روابط بین تابع بقا و تابع مخاطره  
طبق تعریف تابع مخاطره می‌توان نوشت

$$h(t) = \frac{f(t)}{\lambda - F(t)} = \frac{-\frac{d[\lambda - F(t)]}{dt}}{\lambda - F(t)} = \frac{d\{-\ln[\lambda - F(t)]\}}{dt} = \frac{d\{-\ln[S(t)]\}}{dt}.$$

بنابراین

$$\int_{\circ}^t h(u) du = -\ln[\lambda - F(u)]|_{\circ}^t.$$

اما  $F(\circ) = \circ$  و  $\ln(\lambda) = \circ$  پس

$$\ln[\lambda - F(t)] = \ln[S(t)] = -\int_{\circ}^t h(u) du$$

و در نتیجه

$$S(t) = \exp\left\{-\int_{\circ}^t h(u) du\right\}.$$

پس

$$H(t) = \int_{\circ}^t h(u) du = -\ln[S(t)]$$

که  $H(t)$ ، تابع توزیع تجمعی  $h(t)$  است. بنابراین

$$S(t) = \exp[-H(t)].$$

لازم به ذکر است که

$$H(t) \geq \circ$$

و

$$\frac{dH(t)}{dt} = h(t).$$

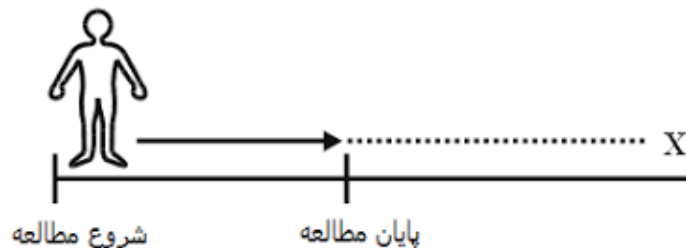
## ۴.۱ سانسور

داده‌های مربوط به زمان، تا لحظه‌ی رخداد یک پیشامد به انواع مختلفی پدیدار می‌شوند که کار با آن‌ها را سخت می‌سازد. یکی از این ویژگی‌ها، که از ویژگی‌های اصلی داده‌های بقا در مقایسه با دیگر داده‌های



آماری است و گاهی اوقات به وقوع می‌پیوندد، سانسور<sup>۵</sup> نام دارد. منظور از داده سانسور شده، داده‌ای است که مقدار دقیق آن مجهول بوده و به جای آن از اطلاعات موجود استفاده شده باشد. سانسور به‌طور کلی زمانی اتفاق می‌افتد که تعدادی از طول عمرها، فقط در فاصله‌های مشخصی اتفاق افتاده‌اند و بقیه طول عمرها به‌طور کامل مشاهده می‌شوند. در واقع سانسور به نداشتن دقیق زمان بقا اطلاق می‌شود (امکان عدم رخداد پیشامد مورد نظر در زمان تعیین شده برای مطالعه).  
به‌طور کلی سه دلیل عمده برای سانسور شدن وجود دارد:

۱. مطالعه به پایان می‌رسد اما پیشامدی رخ نمی‌دهد (زمان بقا بیشتر از دوران مطالعه). شکل ۳.۱ این نوع از سانسور را نشان می‌دهد.
۲. فرد مورد مطالعه در زمان مطالعه گم شود و تحت پیگیری قرار نگیرد.
۳. فرد مورد مطالعه بمیرد (اگر مرگ پیشامد مورد نظر نباشد) یا عوارض جانبی دیگر اتفاق بیافتد.



شکل ۳.۱: علامت  $X$  به معنی رخداد پیشامد است.

تعریف ۱.۴.۱. نمونه‌های سانسور شده آن‌هایی هستند که ممکن است در طول مطالعه به‌طور کامل شرکت نداشته یا تا پایان مطالعه از کار نیفتاده‌اند. زمان سانسور را با  $T_{cens}$  یا با  $C$  نشان می‌دهند. به‌طور مشابه زمان پیشامد بقا را بقا  $T_{surv}$  نشان می‌دهند. اگر  $F$  تابع توزیع تجمعی باشد، آن‌گاه

$$F_{cens} = P(T_{cens} \leq t).$$

هم‌چنین اگر  $G(t)$  تابع سانسور باشد، آن را به‌صورت زیر تعریف می‌کنیم:

$$G(t) = 1 - F_{cens}(t) = P(T_{cens} > t).$$

حال با توجه به  $F_{surv}(t) = P(T_{surv} \leq t)$  و  $S(t) = 1 - F_{surv}(t) = P(T_{surv} > t)$  و این‌که می‌دانیم  $T_{cens}$  و  $T_{surv}$  در عمل با هم اتفاق نمی‌افتند، پس

$$T = \min\{T_{surv}, T_{cens}\}.$$

<sup>۵</sup>Censor

نمونه‌های سانسور شده اغلب نتیجه‌ای از آزمایش‌های طول عمر و زمان واکنش هستند. نخستین بار کریچ استفاده از عنوان سانسور شده را برای این نمونه‌ها پیشنهاد کرد. به‌عنوان مثال برای یک بیمار مبتلا به سرطان خون تا زمان بهبود، اگر مطالعه به پایان برسد بدون این‌که بهبودی حاصل شود (پیشامد رخ ندهد)، یا شخص به دلایل مختلف از مطالعه خارج شود، گوئیم این پیشامد برای بیمار سانسور شده است. در نتیجه می‌بینیم که زمان بقا برای این شخص طولانی‌تر از مدت زمانی است که شخص تحت مطالعه است. حال اگر شخص بعد از پایان مطالعه بهبود یابد ما زمان دقیق بقا را نخواهیم دانست، که به آن سانسور از راست<sup>۶</sup> می‌گویند. در موارد نادری نیز ممکن است، پیشامد مورد مطالعه برای یک فرد، قبل از ورود به مطالعه، رخ داده باشد. اما زمان دقیق رخداد پیشامد نامعلوم باشد که به این وضعیت سانسور از چپ<sup>۷</sup> می‌گویند. سانسور از چپ برای داده‌ها زمانی اتفاق می‌افتد که زمان بقای فرد، کوچکتر یا مساوی زمان بقای مشاهده‌شده‌ی آن فرد باشد. در واقع زمان بقای فرد مورد مطالعه، کوچکتر از زمان سانسور باشد. به‌عنوان مثال، فرض کنید در یک مطالعه هدف پیگیری افرادی است که به ویروس *H.I.V* مبتلا هستند. این مطالعه زمانی آغاز می‌شود که برای اولین بار شخص مورد مطالعه دارای آزمایش *H.I.V* مثبت باشد. به هر حال زمان دقیق در معرض ویروس قرار گرفتن، معلوم نیست. بنابراین، زمان بقا در سمت چپ مطالعه سانسور شده است. همچنین ممکن است در مطالعه‌ای، افراد در نقاطی از زمان که قبلاً در نظر گرفته شده‌اند، مورد معاینه قرار گیرند و پیشامد مورد نظر در فاصله زمانی بین دو معاینه رخ دهد. در این صورت زمان ثبت‌شده، با یک بازه عددی مشخص می‌گردد که به آن سانسور فاصله‌ای<sup>۸</sup> گویند.

## انواع سانسور از راست

باتوجه به این‌که بیش‌تر داده‌ها از سمت راست سانسور می‌شوند، سانسور از راست دارای اهمیت ویژه‌ای است. سانسور از راست انواع مختلفی دارد.

### سانسور نوع یک

در سانسور نوع یک، که ساده‌ترین نوع سانسور از راست می‌باشد، افراد به‌طور هم‌زمان وارد مطالعه شده و مطالعه پس از گذشت زمانی مشخص به پایان می‌رسد.

### سانسور نوع دو

در سانسور نوع دو، افراد هم‌زمان وارد مطالعه شده و پس از وقوع پیشامد برای تعداد یا درصد معینی از افراد مطالعه پایان می‌یابد.

<sup>۶</sup>Right Censoring

<sup>۷</sup>Left Censoring

<sup>۸</sup>Interval Censoring

## سانسور تصادفی

نوع دیگری از انواع سانسور از راست، سانسور تصادفی<sup>۹</sup> نامیده می‌شود. در این حالت افراد در زمان‌های متفاوت وارد مطالعه می‌شوند و تا زمان وقوع پیشامد مورد نظر تحت پیگیری قرار می‌گیرند. برخی افراد به دلیلی غیر از پیشامد مورد مطالعه از مطالعه خارج می‌شوند. در چنین مواردی، پیشامد مورد نظر قابل مشاهده نیست. برخی عواملی که سبب سانسور شدن تصادفی فرد می‌شوند، عبارتند از: مرگ‌های ناشی از تصادف، مهاجرت جوامع انسانی، مرگ به دلیلی غیر از پیشامد مورد مطالعه و کناره‌گیری فرد از آزمایش بالینی.

## سانسور نوع سه

در برخی مطالعات، طرح سانسور، ترکیبی از سانسور تصادفی و سانسور نوع یک است. در چنین مطالعاتی، برخی بیماران به‌طور تصادفی سانسور می‌شوند. به‌عنوان مثال به دلیلی غیر از پیشامد مورد مطالعه، مطالعه را ترک می‌کنند یا می‌میرند. در حالی‌که دیگران هنگامی که دوره ثابت مطالعه پایان می‌یابد، از نوع یک سانسور می‌شوند.

یکی از مفروضات اساسی در تحلیل بقا، این است که سانسور شدن زمان‌های بقا ناآگاهی‌بخش<sup>۱۰</sup> باشد. یعنی علت سانسور شدن زمان‌های بقای یک فرد به پیشامد مورد مطالعه بستگی نداشته باشد. به عبارتی زمان‌های شکست و زمان‌های بقا برای هر فرد از هم مستقل باشند.

## ۵.۱ مدل‌های رگرسیون در بقا

تحلیل رگرسیونی یکی از روش‌های آماری برای تحلیل داده‌های چندعاملی است، که حوزه کاربرد آن وسعت زیادی دارد. نتایج آن از نظر مفهوم، فرآیند ساده به‌کارگیری یک معادله است، که ارتباط بین مجموعه‌ای از متغیرها را بیان می‌کند. در واقع تحلیل رگرسیونی تکنیکی آماری برای بررسی و ارتباط بین متغیرها است.

ساده‌ترین حالت در مسائل رگرسیونی رابطه بین یک متغیر تبیینی رگرسیونی با متغیر پاسخ است. معمولاً  $X$ ، متغیر تبیینی و  $Y$ ، متغیر پاسخ نامیده می‌شود.

### ۱.۵.۱ برآزش مدل‌های رگرسیونی برای داده‌های بقا

شناسایی تابع بقا، تابع مخاطره و برآورد پارامترهای آن‌ها از مهم‌ترین اهداف تحلیل بقا است. تابع بقا رفتار احتمالی زمان بقا را نشان می‌دهد و تابع مخاطره بیانگر قدرت آنی خطر در زمان مشخص می‌باشد. این دو تابع اطلاعات هم‌ارزی در مورد زمان‌های بقا در اختیار محقق قرار می‌دهند. از دیگر اهداف تحلیل بقا می‌توان به مقایسه توابع بقا برای گروه‌های مختلف اشاره کرد. با این حال در بسیاری از کاربردهای

<sup>۹</sup>Random Censoring

<sup>۱۰</sup>Noninformative

بقا، هدف مورد نظر بررسی رابطه متغیرهای تبیینی با زمان بقا است. این کار را می‌توان به وسیله یک رابطه رگرسیونی انجام داد. با توجه به ماهیت داده‌های زمان بقا (حضور داده‌های سانسور شده و چولگی مثبت یا منفی زمان‌های بقا)، مدل‌های رگرسیونی معمولی جوابگو نیست. در این صورت لازم است از مدل‌های رگرسیونی مخاطره استفاده نمود. به کمک یک مدل مخاطره می‌توان علاوه بر تعیین میزان اثر عوامل خطر بر زمان‌های بقا، میزان مخاطره یا بقا را برای افراد و گروه‌های مختلف پیش‌بینی نمود. در این بخش رابطه بین زمان بقای  $T$  و مقادیر یک متغیر تبیینی (پیش‌بین)  $\mathbf{X}$  بررسی خواهد شد. در چارچوب کلی برای رگرسیون  $T$  روی  $\mathbf{X}$ ، اصطلاحات دیگری نیز برای  $\mathbf{X}$  به کار می‌رود که از جمله می‌توان به متغیر کمکی<sup>۱۱</sup> و متغیر رگرسیونی اشاره کرد، که این اصطلاحات می‌توانند به جای یکدیگر استفاده شوند.

در حالت کلی‌تر انتظار می‌رود که زمان بقا وابسته به متغیرهای تبیینی مختلفی باشد که در این صورت تمام این متغیرهای تبیینی به صورت بردار  $\mathbf{X}$  نشان داده می‌شوند. در حالت کلی،  $\mathbf{X}$  برداری شامل  $k$  متغیر تبیینی است که در آن  $\mathbf{X}^T = (x_1, \dots, x_k)$ .

## ۲.۵.۱ مخاطرات نسبی

از آن جایی که متغیر پاسخ به‌گونه‌ای مدل‌بندی می‌شود که وابسته به مقدار بردار متغیرهای تبیینی  $\mathbf{X}$  باشد، واضح است که تابع مخاطره نیز به  $\mathbf{X}$  وابسته خواهد بود. فرض کنید  $h(t)$  نشان‌دهنده‌ی تابع مخاطره  $T$  باشد، که در آن برداری از متغیرهای کمکی است. حال آن‌چه که باید بررسی شود این است که متغیرهای کمکی چه تاثیری بر تابع مخاطره دارند و این تاثیر چگونه باید مدل‌بندی شود. یکی از رایج‌ترین مدل‌های مورد استفاده که در آن متغیرهای کمکی دارای اثر ضربی بر تابع مخاطره پایه<sup>۱۲</sup> می‌باشند، مدل مخاطره پایه نامیده می‌شود. این مدل را مدل مخاطرات نسبی<sup>۱۳</sup> نیز می‌نامند.

**تعریف ۱.۵.۱.** فرض کنید  $T$  نشان‌دهنده متغیر پاسخی باشد که وابسته به بردار متغیرهای کمکی،  $\mathbf{X}$  است. مدل مخاطرات نسبی برای  $T$  به صورت

$$h(t|\mathbf{X}) = h_0(t)g(\mathbf{X})$$

تعریف می‌شود که در آن  $g(\cdot)$  تابعی نامنفی از  $\mathbf{X}$  و  $h_0(t)$  مخاطره پایه نامیده می‌شود.

بر طبق این تعریف، متغیرهای کمکی به واسطه حاصل ضرب  $h_0$  که به  $\mathbf{X}$  وابسته نیست و  $g$  که وابسته به  $\mathbf{X}$  است، دارای اثر ضربی بر مخاطره هستند. تابع  $g(\cdot)$  در تعریف ۱.۵.۱ به‌گونه‌ای است که به ازای هر  $\mathbf{X} \geq \mathbf{0}$  و فقط  $g(\mathbf{0}) = 1$ . این تابع با دو بردار  $\mathbf{X}^T = (x_1, \dots, x_k)$  و  $\beta^T = (\beta_1, \dots, \beta_k)$  معمولاً به صورت

$$g(\mathbf{X}) = \exp(\mathbf{X}^T \beta) \quad (2.1)$$

<sup>۱۱</sup>Covariates

<sup>۱۲</sup>Baseline Hazard Function

<sup>۱۳</sup>Proportional Hazard

ساخته می‌شود. لذا  $g(\mathbf{0}) = \exp(\beta^T \mathbf{0}) = \exp(0) = 1$  نکته قابل توجه این است که در رابطه (۲.۱)،  $\mathbf{x}^T \beta = \beta_1 x_1 + \dots + \beta_k x_k$ ، بردار پارامترهای  $\beta$  گرفته شده‌اند. در نتیجه مدل مخاطرات نسبی استاندارد در ۱.۵.۱ به صورت

$$h(t|\mathbf{x}) = h_0(t) \exp(\mathbf{x}^T \beta) \quad (3.1)$$

می‌باشد. بنابراین مخاطره پایه هنگامی رخ می‌دهد که  $\mathbf{x} = \mathbf{0}$ . هنگامی که از رابطه (۳.۱) استفاده شود، فرآیند برازش یک مدل مخاطرات نسبی شامل برآورد  $k$  پارامتر بردار  $\beta$  با استفاده از مقادیر متغیر پاسخ و متغیرهای کمکی خواهد بود.

### رابطه مخاطرات نسبی و تابع بقا

اگر  $S(t|\mathbf{x})$  تابع بقای یک آزمایش با تابع مخاطره  $h(t|\mathbf{x})$  باشد، در این صورت با توجه به روابط موجود بین تابع بقا و تابع مخاطره، داریم

$$\begin{aligned} S(t|\mathbf{x}) &= \exp\left\{-\int_0^t h(u|\mathbf{x}) du\right\} \\ &= \exp\left\{-\int_0^t h_0(u) g(\mathbf{x}) du\right\} \\ &= \exp\left\{-g(\mathbf{x}) \int_0^t h_0(u) du\right\} \\ &= \left(\exp\left\{-\int_0^t h_0(u) du\right\}\right)^{g(\mathbf{x})} \\ &= (S_0(t))^{g(\mathbf{x})} \end{aligned}$$

که در آن

$$S_0(t) = \exp\left\{-\int_0^t h_0(u) du\right\}$$

تابع بقای پایه است و  $S_0(t) = S_0(t)$  هنگامی رخ می‌دهد که  $\mathbf{x} = \mathbf{0}$ . به طور خلاصه، اثر پذیره مخاطرات نسبی بر  $S(t|\mathbf{x})$  این است که یک تابع مثبت مانند  $g$  و یک تابع بقا پایه مانند  $S_0$  را طوری مشخص می‌کند که

$$S(t|\mathbf{x}) = S_0(t)^{g(\mathbf{x})}. \quad (4.1)$$

رابطه (۴.۱) نشان می‌دهد که در مدل مخاطرات نسبی، بردار متغیرهای کمکی دارای یک اثر توانی بر تابع بقای پایه است.

## ۶.۱ مدل‌های مخاطره

در هنگام تحلیل داده‌های بقا، دو هدف عمده مطرح هستند. یکی مدل‌بندی برای پیدا کردن ترکیبی مناسب از متغیرهای کمکی که طول بقای اعضای مورد مطالعه را تحت تاثیر قرار می‌دهند و هدف بعدی پیدا کردن برآوردهای مطمئن برای تابع مخاطره در زیرگروه‌های مورد بررسی است. به‌طور معمول دو راه‌کار اساسی در تحلیل داده‌های آماری، روش‌های نیمه‌پارامتری<sup>۱۴</sup> و پارامتری<sup>۱۵</sup> هستند و هنگامی در تحلیل بقا با این دیدگاه به بررسی مسأله می‌پردازیم، با توجه به هرکدام از اهداف فوق یکی از راه‌های ذکر شده از اهمیت بیشتری برخوردار خواهند بود.

### ۱.۶.۱ مدل‌های پارامتری

مدل‌های پارامتری، که بر روی تبعیت زمان‌های شکست از خانواده خاصی از توزیع‌های پارامتری پایه‌گذاری شده‌اند. توزیع‌های بسیار زیادی مانند وایبل، نمایی، لگ-نرمال، گامپرتز و لگ-لجستیک را می‌توانند شامل شوند. در ادامه برخی از آن‌ها را معرفی می‌کنیم.

#### مدل مخاطره ثابت

فرض کنید  $\beta > 0$  یک مقدار ثابت و  $\omega = \frac{1}{\beta}$  باشد. یک مدل مخاطره ثابت معمولاً زمانی به‌کار می‌رود که خطر حتمی شکست با زمان تغییر نمی‌کند و به‌صورت

$$h(t) = \omega, \quad t > 0$$

است. تابع بقا برای  $t > 0$  به‌صورت زیر محاسبه می‌شود:

$$S(t) = \exp\left\{-\int_0^t h(u) du\right\} = \exp\left\{-\int_0^t \frac{1}{\beta} du\right\} = \exp\left\{-\frac{1}{\beta} \int_0^t du\right\} = \exp\left\{-\frac{t}{\beta}\right\}.$$

هم‌چنین

$$f(t) = -\frac{dS(t)}{dt} = \frac{1}{\beta} \exp^{-\frac{t}{\beta}}$$

$$\text{و } E(T) = \beta, \text{ Var}(T) = \beta^2.$$

نتیجه ۱.۶.۱. متغیر طول عمر  $T$ ، دارای مدل احتمال نمایی با میانگین  $\beta > 0$  است و به‌صورت  $T \sim \exp(\beta)$  نشان داده می‌شود، هرگاه تابع مخاطره آن ثابت و به‌صورت  $h(t) = \frac{1}{\beta}$  باشد.

مدل نمایی یکی از رایج‌ترین مدل‌های احتمال مورد استفاده برای مدل‌بندی طول عمر قطعات است. این مدل برای داده‌های بقا نیز مناسب است.

<sup>۱۴</sup>Semi-Parametric

<sup>۱۵</sup>Parametric

## مدل مخاطره توانی

فرض کنید  $\alpha > 0$  و  $\beta > 0$  مقادیر ثابتی باشند. یک مدل مخاطره توانی معمولا زمانی به کار می‌رود که خطر حتمی شکست به سرعت با زمان افزایش یابد و به صورت

$$h(t) = \frac{\alpha}{\beta^\alpha} t^{\alpha-1}, \quad t > 0$$

است. تابع بقا برای مدل مخاطره توانی به صورت

$$S(t) = \exp\left\{-\int_0^t h(u)du\right\} = \exp\left\{-\int_0^t \frac{\alpha}{\beta^\alpha} u^{\alpha-1} du\right\} = \exp\left\{-\left(\frac{t}{\beta}\right)^\alpha\right\}$$

محاسبه می‌شود و تابع چگالی احتمال آن به صورت زیر است:

$$f(t) = -\frac{dS(t)}{dt} = \frac{\alpha}{\beta^\alpha} t^{\alpha-1} \exp\left\{-\left(\frac{t}{\beta}\right)^\alpha\right\}, \quad t > 0.$$

نتیجه ۲.۶.۱. متغیر تصادفی طول عمر  $T$ ، دارای مدل احتمال وایبل با پارامترهای  $\alpha > 0$  و  $\beta > 0$  است و  $T \sim W(\alpha, \beta)$ ، هرگاه  $T$  دارای یک مخاطره توانی به شکل  $h(t) = \frac{\alpha}{\beta^\alpha} t^{\alpha-1}$  باشد.

## مدل مخاطره نمایی

مخاطره نمایی با پارامترهای مقیاس  $b$  و مکان  $v$  به صورت

$$h(t) = \frac{1}{b} \exp\left(\frac{t-v}{b}\right), \quad t > 0$$

تعریف می‌شود و اغلب در علوم آماری برای مدل‌بندی طول عمر انسان رخ می‌دهند. تابع بقا برای  $t > 0$  عبارت است از

$$\begin{aligned} S(t) &= \exp\left\{-\int_0^t h(u)du\right\} \\ &= \exp\left\{-\int_0^t \frac{1}{b} \exp\left(\frac{u-v}{b}\right) du\right\} \\ &= \exp\left\{-\exp\left(\frac{u-v}{b}\right)\Big|_0^t\right\} \\ &= \exp\left\{-\exp\left(\frac{t-v}{b}\right) + \exp\left(-\frac{v}{b}\right)\right\}. \end{aligned}$$

این تابع بقا را می‌توان به صورت زیر نیز نشان داد:

$$S(t) = \exp(\theta(1 - \exp(\alpha t)))$$

که در این صورت تابع بقا، تابع بقای گامپرتز<sup>۱۶</sup> نامیده می‌شود و اولین بار توسط موش‌برگر و کلین (۱۹۹۷) معرفی شد.

<sup>۱۶</sup>Gompertz Survival Function

روابط پارامتری بین  $(\alpha, \theta)$  و  $(t, b)$  عبارتند از

$$\alpha = \frac{1}{b}, \quad \theta = \exp\left(-\frac{\nu}{b}\right).$$

تابع چگالی احتمال آن نیز برابر است با

$$f(t) = -\frac{dS(t)}{dt} = \frac{1}{b} \exp\left(\frac{t-\nu}{b}\right) S(t).$$

نتیجه ۳.۶.۱. هرگاه متغیر تصادفی طول عمر  $T$  دارای یک مخاطره نمایی به شکل  $h(t) = \frac{1}{b} \exp\left(\frac{t-\nu}{b}\right)$ ، باشد، در این صورت  $T$  دارای مدل احتمال گامپرتز با پارامترهای  $b > 0$  و  $-\infty < \nu < +\infty$  است و با نماد  $T \sim \text{Gom}(\nu, b)$  نشان داده می‌شود.

مدل‌های مخاطره دیگری نیز وجود دارند که به علت عدم استفاده از آن‌ها در این تحقیق، از بیان آن‌ها اجتناب می‌شود.

## ۲.۶.۱ مدل‌های نیمه پارامتری

مدل‌های نیمه پارامتری هیچ پذیره پارامتری خاصی بر روی توزیع زمان شکست پایه ندارند ولی در عوض فرض می‌کنند که مخاطرات متناسب هستند. یکی از معمول‌ترین مدل‌های نیمه پارامتری، مدل مخاطرات متناسب کاکس<sup>۱۷</sup> است که در ادامه آن را معرفی می‌کنیم.

### مدل مخاطرات متناسب کاکس

یکی از مدل‌های آماری بسیار مهم در تحلیل بقا و تحقیقات پزشکی مدل مخاطرات متناسب کاکس است. در سال ۱۹۷۲ کاکس آماردان انگلیسی مدل رگرسیون کاکس را که با داده‌های زمان شکست سروکار دارد، ارائه نمود و از آن زمان تا کنون به صورت گسترده مورد استفاده قرار می‌گیرد. تحلیل آماری داده‌های زمان شکست عرصه‌ای مهم در تحقیقات است که رشته‌های کاربردی متعددی به آن توجه دارند. تحلیل بقا در آزمایش‌های بالینی و قابلیت اعتماد در سیستم‌های تولیدی دو نمونه از داده‌های زمان شکست هستند. رگرسیون کاکس هنگامی مورد استفاده قرار می‌گیرد که متغیر پاسخ آن زمان بقا یا مدت زمان تا وقوع پیشامد است. معادله رگرسیون با استفاده از یک یا چند متغیر کمکی، احتمال مرگ یا پیشامد مورد نظر را در طول زمان پیش‌بینی می‌کند. از آن‌جا که مدل رگرسیون کاکس می‌تواند هم یک مدل پارامتری و هم یک مدل ناپارامتری محسوب شود، بسیار حائز اهمیت است. ارتباط گسترده‌ای بین تحلیل بقا و رگرسیون کاکس وجود دارد. در واقع این مدل یک نظریه‌ی آماری فرآیندهای شمارشی است که به توسعه‌ی تحلیل بقای ناپارامتری می‌پردازد و مهم‌ترین مدل رگرسیون نیمه پارامتری برای تحلیل داده‌های سانسور شده است.

همان‌طور که در قبل گفته شد، مدل مخاطرات نسبی به صورت زیر در نظر گرفته می‌شود:

$$h(t|\mathbf{X}) = h_0(t) \exp(x_1\beta_1 + \dots + x_k\beta_k) = h_0(t) \exp(\mathbf{X}^T \boldsymbol{\beta}). \quad (5.1)$$

<sup>۱۷</sup>Cox Proportional Hazard Model



نکته قابل توجه در این مدل آن است که عرض از مبدا  $\beta_0$  وجود ندارد. موفقیت قابل توجه کاکس این بود که با استفاده از روش وی می‌توان در رابطه‌ی (۵.۱) ابتدا  $\beta$  و سپس مخاطره پایه را برآورد کرد. بدیهی است که

$$h_0(t) = h_*(t)$$

همان تابع مخاطره پایه می‌باشد.  $h_*(t)$  می‌تواند به‌عنوان هر تابعی از  $t$  تعریف شود. تنها شرط لازم این است که  $h_*(t)$  مثبت باشد. این قسمت ناپارامتری مدل و  $\mathbf{x}^T \beta$  قسمت پارامتری مدل است. لذا مدل مخاطرات کاکس یک مدل نیمه‌پارامتری است.

حال با توجه به رابطه (۵.۱) داریم

$$S(t|\mathbf{x}) = (S_*(t))^{\exp(\mathbf{x}^T \beta)}.$$

بنابراین

$$F(t|\mathbf{x}) = 1 - S_*(t)^{\exp(\mathbf{x}^T \beta)}. \quad (6.1)$$

رابطه (۶.۱) را مدل رگرسیون کاکس می‌نامند.

الزامی نبودن یک توزیع احتمالی برای زمان‌های بقا، یکی از مزایای مدل مخاطرات متناسب کاکس است. مدل کاکس را مدل مخاطرات متناسب می‌نامند زیرا نسبت مخاطره‌ها برای دو دسته یا دو گروه از متغیرهای توضیحی  $\mathbf{x}_1$  و  $\mathbf{x}_2$  به متغیر زمان  $t$  بستگی ندارد. یعنی

$$\frac{h_*(t|\mathbf{x}_1)}{h_*(t|\mathbf{x}_2)} = \frac{h_*(t) \exp(\mathbf{x}_1^T \beta)}{h_*(t) \exp(\mathbf{x}_2^T \beta)} = \frac{\exp(\mathbf{x}_1^T \beta)}{\exp(\mathbf{x}_2^T \beta)} = \exp((\mathbf{x}_1^T - \mathbf{x}_2^T) \beta).$$

به دلیل وجود پذیره‌های کمتر در روش‌های نیمه‌پارامتری نسبت به روش‌های پارامتری، محققین علوم پزشکی اغلب به استفاده از این روش‌ها راغب‌ترند. یک پذیره مهم و اساسی در مدل کاکس وجود دارد و آن پذیره متناسب بودن خطر برای تمامی متغیرهای کمکی موجود در مدل نهایی است. یعنی نسبت خطر برای یک نفر در زمان  $t$  به خطر برای فرد دیگر در همان زمان، مقدار ثابتی دارد. در صورت برقراری این پذیره، تفسیر مدل به‌دست‌آمده، ساده‌تر از مدل‌های پارامتری خواهد بود.

## ۷.۱ روش‌های برآورد در مدل‌های نیمه‌پارامتری

روش‌های مختلفی به‌منظور برآورد پارامترهای مدل نیمه‌پارامتری ارائه شده‌اند. روش هموارسازی اسپلاین<sup>۱۸</sup> (انگل و همکاران، ۱۹۸۶؛ هکمن، ۱۹۸۶؛ رایس، ۱۹۸۶ و اسپیک من، ۱۹۸۸)، روش هسته<sup>۱۹</sup> (رابینسون، ۱۹۸۸)، روش رگرسیونی تطبیقی چندگانه<sup>۲۰</sup> (فریدمن، ۱۹۹۰)، روش مانده‌های جزئی

<sup>۱۸</sup>Smoothing Spline

<sup>۱۹</sup>Kernel

<sup>۲۰</sup>Multivariate Variable Adaptive Regression

(کیوزیک و همکاران، ۱۹۹۲)، روش درست‌نمایی موضعی<sup>۲۱</sup> (کارمن و تاتز، ۲۰۰۱) و روش برآورد موجک<sup>۲۲</sup> (لمینک کیو و چانگ، ۲۰۰۴) از جمله روش‌هایی هستند که توسط محققین پیشنهاد شده‌اند. در این پایان‌نامه از رهیافت اسپلین برای برآورد قسمت ناپارامتری مدل استفاده خواهیم کرد.

## ۸.۱ دیدگاه بیزی در استنباط آماری

به‌طور کلی، روش‌های استنباط آماری متعلق به یکی از دو مکتب کلاسیک و بیزی هستند. با توجه به مشکلات موجود در دیدگاه استنباط کلاسیک، دیدگاه بیزی، که استنباط آن مبتنی بر توزیع پسین است، می‌تواند جایگزین مناسبی باشد. در روش‌های متعلق به دیدگاه کلاسیک، پارامتر  $\theta$  را مقداری ثابت اما نامعلوم فرض می‌کنیم و بر اساس یک نمونه تصادفی  $X_1, \dots, X_n$  که دارای توزیع  $f(x|\theta)$  هستند، در مورد  $\theta$  تصمیم‌گیری (برآورد) می‌کنیم. در دیدگاه بیزی، فرض می‌شود که پارامتر  $\theta$  خود یک متغیر تصادفی است که تغییرات آن توسط یک توزیع احتمالی بیان می‌گردد. این توزیع احتمالی را توزیع پیشین<sup>۲۳</sup> می‌گویند که بیانگر اطلاعات اولیه محقق (قبل از مشاهده داده‌ها) نسبت به پارامتر است. پس از مشاهده نمونه تصادفی، توزیع پیشین تصحیح می‌شود. به عبارت دیگر اطلاعات اولیه و به‌دست آمده ترکیب می‌شوند که منجر به تصحیح توزیع پیشین می‌شود. به این توزیع پیشین تصحیح‌شده، توزیع پسین<sup>۲۴</sup> می‌گویند. نحوه‌ی تشکیل این توزیع برای یک نمونه معمول  $n$  تایی، به صورت زیر است:

فرض کنید  $X_1, \dots, X_n$  یک نمونه تصادفی  $n$  تایی با تابع چگالی احتمال  $f(\cdot|\theta)$  باشد، که در آن پارامتر مورد علاقه می‌باشد و مقادیر مختلف خود را در فضای پارامتر  $\Theta$  اختیار می‌کند. در یک مدل بیزی، فرض می‌شود مقدار واقعی پارامتر  $\theta$  نیز تحقق‌ی از یک توزیع احتمالی به نام توزیع پیشین است که پیشتر معرفی شد. اگر این توزیع را با  $\pi(\theta)$  و تابع درست‌نمایی را با

$$L(\theta|\mathbf{X}) = f_{\mathbf{X}|\theta}(x_1, \dots, x_n|\theta)$$

نشان دهیم، آن‌گاه بنابر قاعده بیز، توزیع پسین  $\pi(\theta|\mathbf{X})$  معادل است با

$$\pi(\theta|\mathbf{X}) = \frac{L(\theta|\mathbf{X})\pi(\theta)}{\int_{\Theta} L(\theta|\mathbf{X})\pi(\theta)d\theta} = c(\mathbf{X})L(\theta|\mathbf{X})\pi(\theta)$$

که در آن  $c(x)$  ثابت نرمال‌ساز نامیده می‌شود و برابر است با

$$c(\mathbf{X}) = \frac{1}{\int_{\Theta} L(\theta|\mathbf{X})\pi(\theta)d\theta}$$

کمیت‌های مورد علاقه در استنباط بیزی معمولاً به صورت یک انتگرال بر حسب تابع چگالی توزیع پسین ظاهر می‌شوند. محاسبه این انتگرال‌ها در صورت پیچیده بودن توزیع پسین، مشکل اساسی بر سر راه

<sup>۲۱</sup>Local Likelihood Method

<sup>۲۲</sup>Wavlate

<sup>۲۳</sup>Prior

<sup>۲۴</sup>Posterior

استنباط بیزی است و باید آن‌ها را به صورت تقریبی محاسبه کرد. رهیافت بیزی برای تقریب انتگرال‌های با بعد بالا از روش‌های مبتنی بر نمونه‌گیری (نمونه‌گیری از توزیع پسین) استفاده می‌کند. از جمله این روش‌ها می‌توان به روش مونت کارلوی زنجیر مارکوفی<sup>۲۵</sup> (MCMC) اشاره کرد.

## ۹.۱ روش‌های نمونه‌گیری MCMC

در روش‌های انتگرال‌گیری مونت کارلویی، هدف، ارائه تقریبی مناسب برای انتگرال‌های به شکل

$$\tau = \int h(x)f(x)dx$$

می‌باشد.

### ۱.۹.۱ الگوریتم متروپولیس-هستینگز

الگوریتم متروپولیس-هستینگز<sup>۲۶</sup> مثالی از مجموعه روش‌های نمونه‌گیری MCMC است. این الگوریتم، با شرط داشتن توزیع پسین توام  $\pi(\theta|\mathbf{X})$ ، مجموعه‌ای از متغیرهای تصادفی  $\{\theta^{(t)}\}$ ، برای  $t = 1, 2, \dots$  را به‌طور پی‌درپی تولید می‌کند. برای تولید این نمونه‌ها، این الگوریتم، ما را ملزم به تعیین چگالی پیشنهادی  $P(\theta^{(t+1)}, \theta^{(t)})$  می‌کند. این تابع، تابع چگالی احتمال  $\theta$  در زمان  $t + 1$ ، با توجه به مقدار آن در زمان  $t$  است. بنابراین عملکرد الگوریتم متروپولیس-هستینگز به صورت زیر بیان می‌شود:

با شرط داشتن  $\theta^{(t)}$

۱. نمونه  $\theta^{(t+1)} \sim P(\theta^{(t+1)}, \theta^{(t)})$  را تولید کن.

۲. قرار بده

$$\theta^{(t)} = \begin{cases} \theta^{(t+1)} & \text{با احتمال } P(\theta^{(t)}, \theta^{(t+1)}) \\ \theta^{(t)} & \text{با احتمال } 1 - P(\theta^{(t)}, \theta^{(t+1)}) \end{cases}$$

به‌طوری که

$$P(\theta^{(t)}, \theta^{(t+1)}) = \min \left\{ 1, \frac{f(\mathbf{X}|\theta^{(t+1)})\pi(\theta^{(t+1)}) P(\theta^{(t)}, \theta^{(t+1)})}{f(\mathbf{X}|\theta^{(t)})\pi(\theta^{(t)}) P(\theta^{(t+1)}, \theta^{(t)})} \right\}$$

که به آن احتمال پذیرش متروپولیس-هستینگز می‌گویند.

<sup>۲۵</sup>Markov Chain Monte Carlo

<sup>۲۶</sup>Metropolis-Hastings Algorithm

## ۲.۹.۱ نمونه‌گیر گیبز

بیش‌ترین کاربرد نمونه‌گیر گیبز<sup>۲۷</sup> در مدل‌های بیزی، زمانی است که توزیع‌های شرطی کامل پارامترهای مورد علاقه به صورت بسته وجود دارند. برای تولید نمونه از توزیع پسین توام  $\pi(\boldsymbol{\theta}|\mathbf{X})$ ، که در آن  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_B)$  بردار پارامترهای مدل است، ابتدا چگالی‌های شرطی کامل را به دست می‌آوریم. این چگالی شرطی عبارت‌اند از

$$\begin{aligned} \pi(\theta_1|\mathbf{X}, \theta_2, \dots, \theta_B), \\ \pi(\theta_2|\mathbf{X}, \theta_1, \theta_3, \dots, \theta_B), \\ \vdots \\ \pi(\theta_B|\mathbf{X}, \theta_1, \dots, \theta_{B-1}). \end{aligned}$$

نمونه‌گیر گیبز با فرض داشتن  $\boldsymbol{\theta}^{(t)} = (\theta_1^{(t)}, \dots, \theta_B^{(t)})$  در مرحله  $t$ ام،  $t = 1, 2, \dots$ ، برای تولید یک زنجیر مارکوف از توزیع پسین توام  $\pi(\boldsymbol{\theta}|\mathbf{X})$ ، دنباله‌ای از نمونه‌ها را به صورت زیر تولید می‌کند:

$$\begin{aligned} 1 : \theta_1^{(t+1)} &\sim \pi(\theta_1|\mathbf{X}, \theta_2^{(t)}, \dots, \theta_B^{(t)}), \\ 2 : \theta_2^{(t+1)} &\sim \pi(\theta_2|\mathbf{X}, \theta_1^{(t+1)}, \theta_3^{(t)}, \dots, \theta_B^{(t)}), \\ \vdots & \\ B : \theta_B^{(t+1)} &\sim \pi(\theta_B|\mathbf{X}, \theta_1^{(t+1)}, \dots, \theta_{B-1}^{(t+1)}). \end{aligned}$$

در صورتی‌که برخی از چگالی‌های شرطی کامل صورت پیچیده داشته باشند و نتوان به صورت مستقیم و ساده از آن‌ها نمونه تولید کرد، برای تولید از آن‌ها می‌توان از الگوریتم متروپولیس-هستینگز استفاده کرد که به روش نمونه‌گیری حاصل الگوریتم نمونه‌گیری متروپولیس-هستینگز درون گیبز<sup>۲۸</sup> می‌گویند (رابت و کسلا، ۲۰۰۴).

<sup>۲۷</sup>Gibbs Sampler

<sup>۲۸</sup>Metropolis-Hastings Within Gibbs



## فصل ۲

# مدل‌های نرخ شفایافتگی غیرآمیخته

در این فصل به معرفی انواع مدل‌های نرخ شفایافتگی می‌پردازیم. در ادامه ویژگی‌های هر مدل مورد بررسی قرار خواهد گرفت و در انتها با یک لم شرایط شناسایی‌پذیری مدل غیرآمیخته را بیان می‌کنیم:

### ۱.۲ مدل‌های نرخ شفایافتگی

همان‌طور که در فصل قبل بیان شد تحلیل بقا مجموعه‌ای از روش‌های تحلیل داده‌ها است که متغیر پاسخ در آن‌ها، زمان وقوع (طول عمر) یک پدیده است. این پدیده‌ها می‌توانند ماهیت‌های متفاوتی داشته باشند، اما معمولاً مرگ یک انسان به سبب یک بیماری خاص مثل سرطان است. پیشرفت‌های چشم‌گیر پزشکی در درمان بیماری (سرطان) موجب شده است، در صورت درمان به‌موقع، برخی از بیماران از طول عمر طبیعی برخوردار شوند. به عبارت دیگر، برخی از بیماران نسبت به بیماری مورد نظر شفایافته (مصون) هستند. بنابراین یکی از پذیره‌های اساسی تحلیل بقا که احتمال رخداد پیشامد را برای تمامی بیماران یک در نظر می‌گیرد، دچار خدشه می‌شود. در این موارد، یک مدل بقای مطلوب، مدلی است که بیماران شفایافته و سایر را به تفکیک مدل‌بندی کند.

در تحلیل داده‌های بقا طولانی مدت، با توجه به پیشرفت‌های حاصل در درمان بیمارها، در انتهای مطالعه با تعداد زیادی از موارد سانسور شده مواجه خواهیم شد. اگر طول مدت مطالعه به اندازه کافی طولانی باشد، استفاده از مدل‌های نرخ شفایافتگی<sup>۱</sup> نتایج بهتری خواهد داشت. به‌طور کلی مدل‌های نرخ شفایافتگی به دو دسته کلی مدل‌های نرخ شفایافتگی آمیخته<sup>۲</sup> و مدل‌های نرخ شفایافتگی غیرآمیخته<sup>۳</sup> تقسیم می‌شوند. در کاربردهای مختلف، مدل نرخ شفایافتگی آمیخته از استقبال بیشتری برخوردار بوده است. اما مدل نرخ شفایافتگی غیرآمیخته بر اساس مکانیسم زیستی توسعه بیماری‌های معمول این نوع از داده‌های بقا (مثل برخی از سرطان‌ها)، تعریف شده است. بنابراین از تفسیرپذیری بالاتری برخوردار است.

<sup>۱</sup>Cure Rate Models

<sup>۲</sup>Mixture Cure Models

<sup>۳</sup>Non-Mixture Cure Models

## ۲.۲ مدل نرخ شفایافتگی آمیخته

اولین خانواده از مدل‌های نرخ شفایافتگی، نخستین بار توسط بوگ (۱۹۴۹) مطرح شد و برکسون و گیچ در سال (۱۹۵۲) آن را بسط و توسعه دادند. این رده از مدل‌ها، مدل‌های نرخ شفایافتگی آمیخته نام دارد. در این مدل‌ها فرض بر این است که جامعه از دو دسته بیمار ناهمگن تشکیل شده است. برای مثال، در بیماری سرطان پستان، به‌عنوان شایع‌ترین بیماری در میان زنان که گاهی اوقات در میان مردان نیز دیده می‌شود، دسته اول بیمارانی هستند که در معرض رخداد مرگ قرار دارند و در مدت زمان معقولی پس از شروع مطالعه مرگ را تجربه خواهند کرد. دسته دوم بیمارانی هستند که در معرض رخداد حادثه قرار ندارند. البته این به معنای تجربه نکردن مرگ تا زمان بی‌نهایت نیست. بلکه تعبیری ریاضی است، به این مفهوم که در زمان معقولی رخداد مرگ را تجربه نمی‌کنند، یا مرگ آن‌ها به علت دیگری جز سرطان پستان است. دسته اول را غیرمصون<sup>۴</sup> (شفانیافته) و دسته دوم را مصون (شفایافته) می‌نامند. اگر  $T$ ، زمان رخداد پیشامد مورد نظر باشد، تابع مشخصه  $Y$  را به صورت زیر در نظر بگیرید:

$$Y = \begin{cases} ۱ & \text{شفانیافته} \\ ۰ & \text{شفایافته} \end{cases}$$

هم‌چنین اگر  $1 - P(\mathbf{x})$  نسبت شفایافتگی افراد تحت متغیر کمکی‌های  $\mathbf{x}$  باشد، آن‌گاه

$$P(Y = ۱) = P(\mathbf{x}).$$

در نتیجه

$$P(Y = ۰) = ۱ - P(\mathbf{x}).$$

تابع بقای افراد غیرمصون به صورت

$$S_u(t|\mathbf{z}) = P(T > t|Y = ۱, \mathbf{z})$$

تعریف می‌شود که در آن  $\mathbf{z}$  بردار متغیرهای کمکی مختص تابع بقای افراد شفایافته است. تابع بقای افراد مصون نیز با  $S_c(t|\mathbf{z})$  نشان داده شده و به صورت

$$S_c(t|\mathbf{z}) = P(T > t|Y = ۰, \mathbf{z})$$

تعریف می‌شود و همواره برابر با ۱ است. بنابراین تابع بقای جامعه که آن را با  $S_p(t|\mathbf{x}, \mathbf{z})$  نشان می‌دهند به صورت

$$S_p(t|\mathbf{x}, \mathbf{z}) = P(\mathbf{x})S_u(t|\mathbf{z}) + (1 - P(\mathbf{x})) \quad (۱.۲)$$

<sup>۴</sup>Uncured

تعریف می‌شود. برکسون و گیج،  $P(\mathbf{x})$  را به‌عنوان یک ثابت نامعلوم و  $S_u(t|\mathbf{z})$  را مدلی پارامتری در نظر گرفتند.

تابع احتمال و تابع نرخ مخاطره در این مدل به‌صورت زیر است:

$$f(t|\mathbf{z}) = P(\mathbf{x})f_u(t|\mathbf{z})$$

و

$$h(t|\mathbf{z}) = \frac{P(\mathbf{x})f_u(t|\mathbf{z})}{(1 - P(\mathbf{x})) + P(\mathbf{x})f_u(t|\mathbf{z})}.$$

فیروز (۱۹۸۶، ۱۹۸۲) نسبت شفایافتگی را با یک مدل رگرسیون لجستیک به‌صورت زیر بسط داد:

$$P(\mathbf{x}) = \frac{\exp(\mathbf{x}^T \boldsymbol{\beta})}{1 + \exp(\mathbf{x}^T \boldsymbol{\beta})}.$$

رگرسیون لجستیک، یکی از رایج‌ترین مدل‌های مورد استفاده برای طبقه‌بندی و پیش‌بینی در داده‌های بقا است.

### تابع درستنمایی مدل شفایافتگی آمیخته

فرض کنید  $t_i$  زمان بقای مشاهده‌شده برای  $i$ -امین بیمار باشد، و تابع مشخصه  $y_i$  را به‌صورت زیر در نظر بگیرید:

$$y_i = \begin{cases} 1 & \text{اگر } t_i \text{ سانسور شده است} \\ 0 & \text{در غیر این صورت} \end{cases} \quad i = 1, \dots, n$$

در این صورت تابع درستنمایی مدل به‌صورت زیر است:

$$L = \prod_{i=1}^n [P(\mathbf{x})f_u(t_i|\mathbf{z})]^{y_i} [(1 - P(\mathbf{x})) + P(\mathbf{x})S_u(t_i|\mathbf{z})]^{1-y_i}.$$

بسیاری از محققین این عرصه از جمله کوک و چن (۱۹۹۲)، تیلور (۱۹۹۵)، پنگ و دیر (۲۰۰۰)، سای و تیلور (۲۰۰۰)، لی و تیلور (۲۰۰۰)، پنگ (۲۰۰۳) و ژانگ، پنگ و لی (۲۰۱۳)، مدل نرخ شفایافتگی آمیخته را در قالب یک مدل نیمه‌پارامتری ارائه کردند و برای اولین بار وانگ و لیانگ (۲۰۱۲) این مدل را به‌صورت کاملاً ناپارامتری معرفی کردند.

## ۳.۲ مدل نرخ شفایافتگی غیرآمیخته

دومین خانواده از مدل‌های نرخ شفایافتگی، مدل‌های نرخ شفایافتگی غیرآمیخته یا مدل‌های زمان پیشرفت<sup>۵</sup> هستند. این مدل‌ها برای جوامعی تعریف شده‌اند که، تابع توزیع زمان بقای کل افراد جامعه

<sup>۵</sup>Promotion Time Cure Model



در آن‌ها کامل نیست. به این معنی که مقدار تابع توزیع تجمعی به یک نمی‌رسد. در این جوامع، میزان اختلاف تابع توزیع تجمعی با مقدار یک، نسبت افراد شفایافته را نشان می‌دهد. اولین مدل نرخ شفایافتگی غیرآمیخته، توسط یاکولف و زودایکف (۱۹۹۶) و زودایکف (۱۹۹۸) ارائه شد و سپس چن و همکاران (۱۹۹۹) مکانیسمی بر اساس توزیع متغیر پنهان<sup>۶</sup> برای آن ارائه کردند. این فرآیند توسط کونر و همکاران (۲۰۰۷) بسط و توسعه یافت.

تابع بقا برای متغیر تصادفی زمان  $T$ ، در این رده از مدل‌ها، به صورت زیر به دست می‌آید:

- فرض کنید  $N$  بیانگر تعداد سلول‌های سرطانی از توزیع پواسون با پارامتر  $\theta$  باشد. یعنی

$$N \sim Pois(\theta)$$

- فرض کنید  $\{Y_1, \dots, Y_N\}$  دنباله‌ای از متغیرهای تصادفی مستقل و هم‌توزیع با تابع توزیع تجمعی  $F(t) = 1 - S(t)$  هستند که به  $N$  نیز وابسته نیستند.

- متغیر تصادفی زمان بقا به صورت زیر تعریف می‌شود:

$$T = \begin{cases} \infty & N = 0 \\ Y_{(i)} & N \geq 1 \end{cases}$$

که در آن  $Y_{(1)} \leq \dots \leq Y_{(i)} \leq \dots \leq Y_{(N)}$  در واقع زمان بازگشت مجدد بیماری به صورت

$$T = \min\{Y_i, \quad 1 \leq i \leq N\}$$

به دست می‌آید.

بنابراین تابع بقای جامعه به صورت زیر قابل محاسبه است:

$$\begin{aligned} S_p(t) &= P(N = 0) + P(Y_1 > t, \dots, Y_N > t, N \geq 1) & (2.2) \\ &= P(N = 0) + P(Y_1 > t, \dots, Y_N > t | N \geq 1)P(N \geq 1) \\ &= \exp(-\theta) + \sum_{k=1}^{\infty} S(t)^k \frac{\theta^k \exp(-\theta)}{k!} \\ &= \sum_{k=0}^{\infty} \frac{(S(t)\theta^k e^{-\theta})}{k!} = \frac{e^{-\theta}}{e^{-\theta S(t)}} = \exp[-\theta(1 - S(t))]. \end{aligned}$$

در این مدل، نسبت افراد شفایافته در داخل مدل بوده و می‌توان آن را به طور مستقیم با گرفتن حد بقا در

<sup>۶</sup>Latent Variable

بینهایت، محاسبه کرد:

$$\begin{aligned} P(N = \circ) &= \lim_{t \rightarrow \infty} S_p(t) \\ &= \lim_{t \rightarrow \infty} \exp(-\theta(1 - S(t))) \\ &= \lim_{t \rightarrow \infty} \exp(-\theta + \theta S(t)) \\ &= \exp(-\theta) \lim_{t \rightarrow \infty} \exp(\theta S(t)) \\ &= \exp(-\theta) \exp(\theta \lim_{t \rightarrow \infty} S(t)) \end{aligned}$$

و با توجه به این که  $S(\infty) = \circ$ ، آنگاه

$$P(N = \circ) = \exp(-\theta) \exp(\circ) = \exp(-\theta).$$

با در نظر گرفتن مدل مخاطره کاکس به صورت

$$F(t|\mathbf{z}) = 1 - S_{\circ}(t)^{\exp(\mathbf{z}^T \boldsymbol{\gamma})}$$

که در آن  $S_{\circ}(t)$  تابع بقا پایه است، رابطه (۲.۲) با دو بردار کمکی  $\mathbf{x}$  و  $\mathbf{z}$  و متغیر پنهان با میانگین  $\exp(\beta_{\circ} + \mathbf{x}^T \boldsymbol{\beta})$  به صورت زیر درمی آید:

$$\begin{aligned} S_p(t|\mathbf{x}, \mathbf{z}) &= \exp(-\theta(\mathbf{x})F(t|\mathbf{z})) \\ &= \exp(-\exp(\beta_{\circ} + \mathbf{x}^T \boldsymbol{\beta})(1 - S_{\circ}(t)^{\exp(\mathbf{z}^T \boldsymbol{\gamma})})). \end{aligned} \quad (۳.۲)$$

در سال‌های اخیر توزیع‌های متفاوتی برای متغیر پنهان  $N$ ، در نظر گرفته شده‌اند. برای مثال، اگر توزیع متغیر پنهان  $N$  برنولی در نظر گرفته شود، مدل شفایافته آمیخته حاصل می‌شود. همچنین کونر و همکاران (۲۰۰۷) توزیع هندسی، برنولی و دوجمله‌ای، رودریگز و همکاران (۲۰۰۹) توزیع پواسون کانوی-ماکسول<sup>۷</sup>، بورگر و همکاران (۲۰۱۲) توزیع سری‌های توانی، باغستانی و همکاران (۲۰۱۳) توزیع تعمیم‌یافته پواسون لیندلی و رحیم‌زاده و همکاران (۲۰۱۳) توزیع فوق هندسی تعمیم‌یافته و دوجمله‌ای منفی را در مدل‌های شفایافته زمان پیشرفت ارائه کردند. شایان ذکر است توزیع متغیر پنهان، به منظور تبیین احتمال شفایافتگی، می‌تواند هر یک از توزیع‌های گسسته که دارای جرم احتمال در نقطه صفر هستند را شامل شود. به نظر می‌رسد رابطه نزدیکی بین توزیع‌هایی که در تحلیل داده‌های شمارشی و مدل‌های شفایافته زمان پیشرفت مورد استفاده قرار می‌گیرد، وجود دارد. مشکل دیگری که علاوه بر بیش‌پراکنشی در تحلیل داده‌های بقا وجود دارد، چوله بودن داده‌ها است که دلیلی دیگر برای ارائه توزیع‌های متفاوت برای حل این مشکل است.

<sup>۷</sup>Conway-Maxwell-Poisson Distribution

## ۴.۲ ویژگی‌های مدل نرخ شفایافتگی غیرآمیخته

مدل نرخ شفایافتگی غیرآمیخته با تابع بقا جامعه در رابطه (۳.۲) معرفی شد را در نظر بگیرید. تابع نرخ مخاطره جامعه با  $h_p(t|\mathbf{X}, \mathbf{Z})$  نشان داده می‌شود و به صورت زیر تعریف می‌شود:

$$\begin{aligned} h_p(t|\mathbf{X}, \mathbf{Z}) &= \frac{-d(\ln(S_p(t|\mathbf{X}, \mathbf{Z})))}{dt} \\ &= \frac{-d(\ln(\exp(-\exp(\beta_0 + \mathbf{X}^T \boldsymbol{\beta}))(\mathbb{1} - S_0(t)^{\exp(\mathbf{Z}^T \boldsymbol{\gamma})})))}{dt} \\ &= \frac{-d(-\exp(\beta_0 + \mathbf{X}^T \boldsymbol{\beta}))(\mathbb{1} - S_0(t)^{\exp(\mathbf{Z}^T \boldsymbol{\gamma})})}{dt}. \end{aligned}$$

با توجه به این نکته که  $\frac{d(g(x)^n)}{dx} = ng'(x)g(x)^{n-1}$  و  $f(t) = -\frac{dS(t)}{dt}$  داریم

$$\begin{aligned} h_p(t|\mathbf{X}, \mathbf{Z}) &= -(-\exp(\beta_0 + \mathbf{X}^T \boldsymbol{\beta}))(\exp(\mathbf{Z}^T \boldsymbol{\gamma})f_0(t)S_0(t)^{\exp(\mathbf{Z}^T \boldsymbol{\gamma})-1}) \\ &= \frac{\exp(\beta_0 + \mathbf{X}^T \boldsymbol{\beta}) \exp(\mathbf{Z}^T \boldsymbol{\gamma})f_0(t)S_0(t)^{\exp(\mathbf{Z}^T \boldsymbol{\gamma})}}{S_0(t)} \\ &= \theta(\mathbf{X})f(t|\mathbf{Z}). \end{aligned}$$

نسبت مخاطره جامعه برای دو گروه ۱ و ۲ به صورت زیر به دست می‌آید:

$$\begin{aligned} HR_p &= \frac{h_p(t|\mathbf{X}_1, \mathbf{Z}_1)}{h_p(t|\mathbf{X}_2, \mathbf{Z}_2)} \\ &= \frac{\frac{\exp(\beta_0 + \mathbf{X}_1^T \boldsymbol{\beta}) \exp(\mathbf{Z}_1^T \boldsymbol{\gamma})f_0(t)S_0(t)^{\exp(\mathbf{Z}_1^T \boldsymbol{\gamma})}}{S_0(t)}}{\frac{\exp(\beta_0 + \mathbf{X}_2^T \boldsymbol{\beta}) \exp(\mathbf{Z}_2^T \boldsymbol{\gamma})f_0(t)S_0(t)^{\exp(\mathbf{Z}_2^T \boldsymbol{\gamma})}}{S_0(t)}} \\ &= \frac{\exp(\beta_0 + \mathbf{X}_1^T \boldsymbol{\beta}) \exp(\mathbf{Z}_1^T \boldsymbol{\gamma})S_0(t)^{\exp(\mathbf{Z}_1^T \boldsymbol{\gamma})}}{\exp(\beta_0 + \mathbf{X}_2^T \boldsymbol{\beta}) \exp(\mathbf{Z}_2^T \boldsymbol{\gamma})S_0(t)^{\exp(\mathbf{Z}_2^T \boldsymbol{\gamma})}} \\ &= \exp((\mathbf{X}_1^T - \mathbf{X}_2^T)\boldsymbol{\beta}) \exp((\mathbf{Z}_1^T - \mathbf{Z}_2^T)\boldsymbol{\gamma})S_0(t)^{(\exp(\mathbf{Z}_1^T \boldsymbol{\gamma}) - \exp(\mathbf{Z}_2^T \boldsymbol{\gamma}))}. \end{aligned}$$

اگر  $\mathbf{Z}_1 = \mathbf{Z}_2$ ، آن‌گاه

$$HR_p = \frac{h_p(t|\mathbf{X}_1, \mathbf{Z}_1)}{h_p(t|\mathbf{X}_2, \mathbf{Z}_2)} = \exp((\mathbf{X}_1^T - \mathbf{X}_2^T)\boldsymbol{\beta})$$

و این مسأله نشان می‌دهد که نسبت مخاطره‌ها به زمان  $t$  بستگی ندارد. تابع بقا برای افراد شفانیافته به صورت زیر به دست می‌آید:

$$\begin{aligned} S_u(t|\mathbf{X}, \mathbf{Z}) &= P(T > t | N \geq 1, \mathbf{X}, \mathbf{Z}) \\ &= \frac{S_p(t|\mathbf{X}, \mathbf{Z}) - P(N = 0)}{\mathbb{1} - P(N = 0)} \\ &= \frac{S_p(t|\mathbf{X}, \mathbf{Z}) - \exp(-\theta(\mathbf{X}))}{\mathbb{1} - \exp(-\theta(\mathbf{X}))}. \end{aligned}$$

هم‌چنین تابع مخاطره برای افراد شفانیافته نیز به صورت زیر است:

$$\begin{aligned} h_u(t|\mathbf{x}, \mathbf{z}) &= \frac{-d(\ln(S_u(t|\mathbf{x}, \mathbf{z})))}{dt} & (۴.۲) \\ &= \frac{S_p(t|\mathbf{x}, \mathbf{z})}{S_p(t|\mathbf{x}, \mathbf{z}) - \exp(-\theta(\mathbf{x}))} h_p(t|\mathbf{x}, \mathbf{z}) \\ &= \frac{h_p(t|\mathbf{x}, \mathbf{z})}{P(T < +\infty | T > t, \mathbf{x}, \mathbf{z})}. \end{aligned}$$

همان‌طور که ملاحظه می‌شود رابطه (۴.۲)، ساختار مخاطرات متناسب را ندارد. به عبارت دیگر ویژگی مخاطرات متناسب برای کل افراد تحت مطالعه برقرار است نه فقط افراد شفانیافته.

## ۵.۲ شناسایی‌پذیری مدل

تعریف ۱.۵.۲. (لیو و چن، ۲۰۰۹). در مدل  $F(x; \beta)$  پارامتر  $\beta$  را شناسایی‌پذیر گویند، هرگاه به ازای هر  $\beta_1$  و  $\beta_2$  از فضای پارامتر که  $F(x; \beta_1) = F(x; \beta_2)$ ، آن‌گاه  $\beta_1 = \beta_2$ .

در نظریه تحلیل بقا، پیگیری یک مطالعه به اندازه کافی طولانی است اگر بزرگترین زمان سانسور بیشتر از بزرگترین زمان شکست باشد؛ یعنی در صورتی که زمان پیگیری هر واحد غیرمصون (شفانیافته) به اندازه کافی طولانی باشد تا شکست آن را ببیند. اما در عمل، اگر منحنی بقای جامعه در دم راست خود صاف باشد، فرض پیگیری کافی منطقی به نظر می‌رسد. شرط‌های لازم برای شناسایی‌پذیر بودن مدل نرخ شفایافتگی غیرآمیخته در لم زیر ارائه شده‌اند.

لم ۲.۵.۲. (برمه‌ورست و لامبرت، ۲۰۱۶). فرض کنید  $A_1$ . بردار متغیرهای کمکی  $\mathbf{z}$  شامل هیچ‌گونه داده‌ی سانسور شده‌ای نیست.  $A_2$ .  $\mathbf{X}^T \mathbf{X}$  و  $\mathbf{Z}^T \mathbf{Z}$  دو ماتریس با رتبه کامل هستند که  $\mathbf{x}$  و  $\mathbf{z}$  به ترتیب متغیرهای کمکی آن‌ها هستند.

$A_3$ .  $F_0(t) = 1 - S_0(t)$  یک تابع توزیع تجمعی سره<sup>۸</sup> است. یعنی

$$\lim_{t \rightarrow \infty} S_0(t) = 0.$$

آن‌گاه تحت شرط‌های  $A_1 - A_3$  می‌توان گفت

(۱) اگر زمان مطالعه به‌طور کافی طولانی باشد، مدل (۳.۲) شناسایی‌پذیر است.

(۲) اگر زمان مطالعه به‌قدر کافی طولانی نباشد و بردارهای  $\mathbf{x}$  و  $\mathbf{z}$  مستقل از هم باشند، آن‌گاه برآوردهای پارامترهای موجود در احتمال درمان و زمان شکست برای یک سلول سرطانی، شناسایی‌پذیر است.

اثبات. برهان قسمت (۱)

<sup>۸</sup>Proper

(a) فرض کنید  $(\beta_0, \beta, \gamma, S_0)$  و  $(\tilde{\beta}_0, \tilde{\beta}, \tilde{\gamma}, \tilde{S}_0)$  دو مجموعه از پارامترها هستند که در رابطه (۳.۲) صدق می‌کنند و  $\hat{\chi}$  مجموعه‌ای از تمام مقادیر بردار  $\mathbf{x}$  است. اگر

$$\begin{aligned} S_p(t|\mathbf{x}) &= \exp \left[ -\exp(\beta_0 + \mathbf{x}^T \beta) \left( 1 - S_0(t)^{\exp(\mathbf{z}^T \gamma)} \right) \right] \\ &= \exp \left[ -\exp(\tilde{\beta}_0 + \mathbf{x}^T \tilde{\beta}) \left( 1 - \tilde{S}_0(t)^{\exp(\mathbf{z}^T \tilde{\gamma})} \right) \right] \\ &= \tilde{S}_p(t|\mathbf{x}) \quad \forall \mathbf{x} \in \chi \quad ; \quad \forall t \in [0, \infty] \end{aligned}$$

آن‌گاه

$$\beta_0 = \tilde{\beta}_0 \quad , \quad \beta = \tilde{\beta} \quad , \quad \gamma = \tilde{\gamma} \quad , \quad S_0(t) = \tilde{S}_0(t) \quad \forall t \in [0, \infty]$$

(b) از آن‌جا که تابع نمایی یکنوا است،

$$\forall \mathbf{x} \in \chi \quad \forall t \in [0, \infty]$$

اگر

$$\exp(\beta_0 + \mathbf{x}^T \beta) \left( 1 - S_0(t)^{\exp(\mathbf{z}^T \gamma)} \right) = \exp(\tilde{\beta}_0 + \mathbf{x}^T \tilde{\beta}) \left( 1 - \tilde{S}_0(t)^{\exp(\mathbf{z}^T \tilde{\gamma})} \right)$$

آن‌گاه

$$\beta_0 = \tilde{\beta}_0 \quad , \quad \beta = \tilde{\beta} \quad , \quad \gamma = \tilde{\gamma} \quad , \quad S_0(t) = \tilde{S}_0(t) \quad \forall t \in [0, \infty].$$

(c) از آن‌جا که  $F_0(t)$  تابع توزیع تجمعی سره است، با جایگذاری در روابط قبل داریم

$$F(t|\mathbf{z}) = 1 - S_0(t)^{\exp(\mathbf{z}^T \gamma)}$$

که یک تابع توزیع تجمعی سره است. ویژگی (b) به‌ازای  $t = \infty$  نیز برقرار است. حال با توجه به این‌که  $F(t|\mathbf{z})$  و  $F(\tilde{t}|\mathbf{z})$  هر دو تابع توزیع تجمعی سره هستند، اگر

$$\beta_0 + \mathbf{x}^T \beta = \tilde{\beta}_0 + \mathbf{x}^T \tilde{\beta} \quad \forall \mathbf{x} \in \chi$$

آن‌گاه

$$\beta_0 = \tilde{\beta}_0, \quad \beta = \tilde{\beta}.$$

تحت  $A_2$  مستقیماً نتیجه می‌شود

$$P_\chi(\mathbf{x}^T \tilde{\beta} = \alpha) = 1 \Rightarrow \alpha = 0, \beta = 0.$$

(d) تا زمانی که بردار  $\mathbf{z}$  شامل هیچ سانسوری نیست، واضح است که تحت رابطه  $A_2$  مدل مخاطره متناسب کاکس شناسایی‌پذیر است.

(e) با ادغام نتایج (c) و (d) قسمت (۱) اثبات می‌شود.

برهان قسمت (۲)

(a) اگر تومور در حال رشد از سلول سرطانی نتواند با اتمام مطالعه در زمان  $t$  تشخیص داده شود، پیگیری به اندازه کافی طولانی نبوده است. این حالت زمانی که  $F_0(t)$  خیلی کوچک (نزدیک به ۰) باشد، اتفاق می‌افتد. در این صورت

$$\begin{aligned} F(t|\mathbf{z}) &= 1 - S_0(t)^{\exp(\mathbf{z}^T \boldsymbol{\gamma})} \\ &= 1 - \left[ (1 - F_0(t))^{\exp(\mathbf{z}^T \boldsymbol{\gamma})} \right] \\ &= 1 - \left[ (1 - \exp(\mathbf{z}^T \boldsymbol{\gamma}) F_0(t)) \right] + o^*(F_0(t)) \\ &\approx 1 - \left[ (1 - \exp(\mathbf{z}^T \boldsymbol{\gamma}) F_0(t)) \right] \\ &= \exp(\mathbf{z}^T \boldsymbol{\gamma}) F_0(t) \end{aligned}$$

بنابراین با استفاده از این تقریب، رابطه (۳.۲) به صورت زیر قابل بازنویسی است:

$$\begin{aligned} S_p(t|\mathbf{x}, \mathbf{z}) &\approx \exp \left[ -\theta(\mathbf{x}) \exp(\mathbf{z}^T \boldsymbol{\gamma}) F_0(t) \right] \\ &= \exp \left[ -\exp(\beta_0 + \mathbf{x}^T \boldsymbol{\beta}) \exp(\mathbf{z}^T \boldsymbol{\gamma}) F_0(t) \right] \\ &= \exp \left[ -\exp(\beta_0 + \mathbf{x}^T \boldsymbol{\beta} + \mathbf{z}^T \boldsymbol{\gamma}) F_0(t) \right]. \end{aligned}$$

بنابراین نتیجه می‌گیریم اگر بردارهای  $\mathbf{x}$  و  $\mathbf{z}$  با هم اشتراک داشته باشند، تقریب تاثیرات متغیرهای کمکی قابل شناسایی نیست، زیرا پیگیری به اندازه کافی طولانی نبوده است.

(c) برآورد اثرات متغیرهای کمکی شناسایی‌پذیر است، هرگاه پیگیری مطالعه به اندازه کافی طولانی نیست ولی بردارهای کمکی واحد مشترک ندارند. در واقع می‌توان نوشت

$$\begin{aligned} \theta(\mathbf{x}) F(t|\mathbf{z})|_{t_{max}} &\approx \exp(\beta_0 + \mathbf{x}^T \boldsymbol{\beta} + \mathbf{z}^T \boldsymbol{\gamma}) F_0(t)_{t_{max}} \\ &= \exp(\beta_0 + \mathbf{x}^T \boldsymbol{\beta} + \mathbf{z}^T \boldsymbol{\gamma} + \text{Log}(1 - \epsilon)) \\ &= \exp(\tilde{\beta}_0 + \mathbf{x}^T \boldsymbol{\beta} + \mathbf{z}^T \boldsymbol{\gamma}) \end{aligned}$$

که در آن

$$\epsilon = F_0(+\infty) - F_0(t_{max}).$$

در نتیجه برآوردی از تاثیر متغیرها در احتمال درمان و زمان شکست برای یک سلول سرطانی، تحت تاثیر قرار نمی‌گیرد. بنابراین اگر بردارهای  $\mathbf{x}$  و  $\mathbf{z}$  اشتراک نداشته باشند، پارامترهای مذکور می‌توانند قابل شناسایی باشند.

□

# فصل ۳

## تحلیل بیزی مدل نرخ شفایافتگی غیرآمیخته

### ۱.۳ مقدمه

تحلیل بیزی مدل‌های نرخ شفایافتگی غیرآمیخته با تابع بقای نیمه‌پارامتری از سال ۱۹۸۰ بسط و توسعه یافت. روش درست‌نمایی جزئی که روش معمول برآورد پارامترها در مدل مخاطره متناسب کاکس است، در این مدل کاربرد ندارد. برای برآورد پارامترها باید از تابع درست‌نمایی کامل استفاده شود که نیازمند برآورد تابع مخاطره پایه است. این رهیافت چالش‌های خاص خود را دارد. بنابراین استفاده از دیدگاه استنباط بیزی در این مدل‌ها می‌تواند جانشین مناسب‌تری باشد. در این فصل دیدگاه مدل‌بندی بیزی در مدل‌های نرخ شفایافتگی غیرآمیخته را بیان می‌کنیم. برای مدل‌بندی ناپارامتری تابع بقای پایه از  $P$ -اسپلاین‌ها استفاده خواهیم کرد. بنابراین ابتدا مقدمه‌ای بر اسپلاین‌ها را مطرح می‌کنیم.

### ۲.۳ اسپلاین

در آمار، هموار کردن<sup>۱</sup> یک نمودار پراکنش، استفاده از الگوهای مهم داده‌ها در ایجاد یک تابع تقریبی است. یک هموارکننده تابعی برای خلاصه‌سازی متغیر پاسخ به‌عنوان تابعی از متغیرهای پیش‌بین  $X_1, \dots, X_k$  است. یکی از خصوصیات مهم هموارکننده‌ها، ماهیت ناپارامتری آن‌ها است. هموارکننده‌ها همواره صورت تابعی ساده‌ای از وابستگی  $Y$  به  $X_1, X_2, \dots, X_k$  را ارائه می‌دهند. به همین دلیل اغلب در مدل رگرسیون، ساختار ناپارامتری دارند. برآورد تولیدشده توسط یک هموارکننده را یک هموار می‌نامیم.

فرض کنید  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  مقادیر متغیر پاسخ  $Y$  و  $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$  مقادیر متغیر پیش‌بین باشند. یک هموارکننده، تابعی از  $\mathbf{x}$  و  $\mathbf{y}$  است که برازش مناسبی در راستای خلاصه‌سازی داده‌ها ایجاد می‌کند. یکی از انواع هموارکننده‌های نمودار پراکنش، اسپلاین‌ها هستند. اسپلاین توسط شخصی به نام شوئنبرگ در سال ۱۹۴۶ معرفی شد.

<sup>۱</sup>Smoothing



اسپلاین یا منحنی منعطف، از یک نوار بلند تشکیل شده است که شامل نقاط ثابت و یک منحنی هموار است که هدف آن عبور از همه‌ی نقاط برای انتقال از یک نقطه به نقطه دیگر است. ریشه‌ی کلمه اسپلاین در کارهای چوبی بوده و از لحاظ واژه‌شناسی به‌خاطر ارتباط کلمه‌ی اسپلاین با کلمه‌ی باریکه‌ی چوب است که بعدها صنعتگران و کارگران ماهر، اسپلاین‌ها را از لاستیک، استیل و هر ماده کشپار ساختند. برای کشتی‌سازی اغلب از چوب‌های باریک، بلند و انعطاف‌پذیر استفاده می‌کردند که به آن‌ها اسپلاین گفته می‌شد. از اسپلاین‌ها برای خم کردن چوب برای پیانو و ویولون استفاده شده است. به دلایل فیزیکی، منحنی‌های اسپلاین تقریباً چند ضابطه‌ای‌هایی مرتبه‌ی سوم با مشتق دوم پیوسته‌اند.

### تابع توانی بریده

تابع توانی بریده<sup>۲</sup> از مرتبه‌ی  $p$  را که در نقطه‌ی  $t$  بریده شده باشد، به‌صورت زیر تعریف می‌کنیم:

$$(x - t)_+^p = (x - t)^p I(x > t)$$

که مقدار حقیقی  $t$  را گره<sup>۳</sup> می‌نامند و  $I$  تابع نشانگر به‌صورت

$$I(x > t) = \begin{cases} 1 & x > t \\ 0 & \text{سایر نقاط} \end{cases}$$

است.

### اسپلاین مرتبه‌ی اول

فاصله‌ی  $[a, b]$  را به‌صورت اجتماع دو بازه‌ی  $[a, t] \cup [t, b]$  در نظر بگیرید. یک تابع اسپلاین از مرتبه‌ی یک، تابعی است که در هر یک از بازه‌های مذکور، یک چندجمله‌ای از مرتبه‌ی یک (خط راست) باشد و در نقطه‌ی  $t$  نیز به‌هم پیوسته باشند. این تابع را می‌توان به‌صورت رابطه‌ی زیر نوشت:

$$f(x) = \begin{cases} \phi_0 + \phi_1 x & x \leq t \\ \phi'_0 + \phi'_1 x & x > t \end{cases} \quad (1.3)$$

با اعمال شرط پیوستگی تابع  $f$  در نقطه‌ی  $x = t$  داریم

$$\phi_0 + \phi_1 t = \phi'_0 + \phi'_1 t$$

$$\phi'_0 = \phi_0 + \phi_1 t - \phi'_1 t.$$

با جایگذاری  $\phi'_0$  در ضابطه‌ی دوم  $f$  داریم

$$\phi'_0 + \phi'_1 x = \phi_0 + \phi_1 t - \phi'_1 t + \phi'_1 x = \phi_0 + \phi_1 t + \phi'_1 (x - t)$$

<sup>۲</sup>Truncated Power Function

<sup>۳</sup>Knot

بنابراین

$$f(x) = \begin{cases} \phi_0 + \phi_1 x & x \leq t \\ \phi_0 + \phi_1 t + \phi'_1(x-t) & x > t \end{cases}$$

که می‌توان این تابع دو ضابطه‌ای را به صورت رابطه‌ی زیر نوشت:

$$f(x) = \phi_0 + \phi_1 x + \phi_2(x-t)_+ \quad (۲.۳)$$

که  $\phi_2 = \phi'_1 - \phi_1$  تفاوت بین شیب خطوط اول و دوم رابطه‌ی (۱.۳) را بیان می‌کند. در واقع می‌توان گفت که تابع  $f$ ، ترکیبی خطی از توابع پایه‌ای  $1$ ،  $x$  و  $(x-t)_+$  است. به همین ترتیب اگر در فاصله‌ی  $[a, b]$  تعداد دو گره  $t_1 < t_2$  وجود داشته باشند، می‌توان فاصله‌ی  $[a, b]$  را به صورت اجتماع  $[a, t_1] \cup [t_1, t_2] \cup [t_2, b]$  در نظر گرفت و تابع اسپالین مرتبه‌ی یک را به طریق مشابه رابطه‌ی (۲.۳) به صورت زیر نوشت:

$$f(x) = \phi_0 + \phi_1 x + \phi_2(x-t_1)_+ + \phi_3(x-t_2)_+.$$

نهایتاً فرض کنید که بتوان فاصله‌ی  $[a, b]$  را به صورت اجتماع فاصله‌های

$$[a, t_1], [t_1, t_2], \dots, [t_k, b]$$

نوشت که برای گره‌های ذکر شده داریم  $t_1 < \dots < t_k$ . در این صورت تابع اسپالین مرتبه‌ی یک را می‌توان به صورت رابطه‌ی

$$f(x) = \phi_0 + \phi_1 x + \sum_{j=1}^k \phi_{j+1}(x-t_j)_+$$

بیان کرد که ترکیبی خطی از توابع پایه  $1$ ،  $x$ ،  $(x-t_1)_+$ ،  $\dots$  و  $(x-t_k)_+$  به شمار می‌آید.

### اسپالین مرتبه‌ی $p$

توابع اسپالین را برای مراتب بالاتر نیز می‌توان به طور مشابه تعریف کرد. تابع  $f(\cdot)$  را یک اسپالین مرتبه‌ی  $p$  نامند اگر در هر یک از زیر بازه‌های  $[a, t_1]$ ،  $[t_1, t_2]$ ،  $\dots$ ،  $[t_k, b]$  یک چندجمله‌ای از مرتبه‌ی  $p$  باشد به طوری که  $f(\cdot)$  و مشتق‌های آن از مراتب  $1$  تا  $p-1$ ، در گره‌ها پیوسته باشند. تابع اسپالین مرتبه‌ی  $p$  را می‌توان به صورت زیر نوشت:

$$f(x) = \phi_0 + \phi_1 x + \phi_2 x^2 + \dots + \phi_p x^p + \sum_{j=1}^k \phi_{p+j}(x-t_j)_+^p. \quad (۳.۳)$$

تابع اسپالین مرتبه‌ی  $p$ ، ترکیبی خطی از یک چندجمله‌ای مرتبه‌ی  $p$  و توابع توانی بریده‌شده است. به عبارت دیگر مجموعه‌ی  $\{1, x, x^2, \dots, x^p, (x-t_1)_+^p, \dots, (x-t_k)_+^p\}$  یک پایه برای فضای اسپالین‌های مرتبه‌ی  $p$  با گره‌های  $t_1, t_2, \dots, t_k$  است.

## اسپلاین مکعبی

اگر  $p = ۳$ ، آن‌گاه تابع ذکرشده در رابطه‌ی (۳.۳) را یک اسپلاین مکعبی<sup>۴</sup> نامند که به صورت

$$f(x) = \phi_0 + \phi_1 x + \phi_2 x^2 + \phi_3 x^3 + \sum_{j=1}^k \phi_{3+j} (x - t_j)_+^3$$

بیان می‌شود.

## خواص اسپلاین‌ها

۱. تابع  $f(\cdot)$  در رابطه‌ی (۳.۳)، یک چندجمله‌ای تکه‌ای از مرتبه‌ی  $p$ ، با تعداد  $p + k + ۱$  ضریب است.

۲. ”همواری“ یکی از خواص مهم اسپلاین‌ها است. تابع اسپلاین (۳.۳) از جنس چندجمله‌ای تکه‌ای است و همان‌طور که گفته شد، این تابع دارای مشتقات پیوسته تا مرتبه‌ی  $p - ۱$  است. برای سهولت، اسپلاین مرتبه‌ی  $p$  را با یک گره‌ی  $t$  به صورت زیر در نظر بگیرید:

$$\begin{aligned} f(x) &= \phi_0 + \phi_1 x + \phi_2 x^2 + \dots + \phi_p x^p + \phi_{p+1} (x - t)_+^p \\ &= P(x) + \phi_{p+1} (x - t)_+^p. \end{aligned}$$

الف- تابع  $f(\cdot)$  در گره‌ی  $t$  پیوسته است، زیرا

$$\lim_{x \rightarrow t} f(x) = P(t), \quad f(t) = P(t).$$

بنابراین در سرتاسر بازه‌ی تعریف‌شده پیوسته خواهد بود.

ب- مشتقات تابع  $f(\cdot)$  تا مرتبه‌ی  $p - ۱$  نیز توابعی پیوسته هستند، زیرا

$$f^{(j)}(t) = P^{(j)}(t) \quad j = ۱, ۲, \dots, p - ۱$$

و از آن‌جا که  $\lim_{x \rightarrow t} \phi_{p+1} p(p-1) \dots (p-j+1)(x-t)^{p-j} = 0$  بنابراین

$$\lim_{x \rightarrow t} f^{(j)}(x) = P^{(j)}(t).$$

ج- مشتق مرتبه‌ی  $p$ -ام تابع  $f(\cdot)$  دارای رفتار متفاوتی است. در واقع

$$f^{(p)}(t) = p! \phi_p$$

$$\lim_{x \rightarrow t} f^{(p)}(x) = p! \phi_p + p! \phi_{p+1}.$$

<sup>۴</sup>Cubic Spline

## رگرسیون و کمترین توان‌های دوم اسپلاین

در مدل رگرسیونی

$$y = f(x) + \epsilon$$

اگر  $f(\cdot)$  را با رابطه‌ی (۳.۳) جایگزین کنیم، نتیجه‌ی حاصل را رگرسیون اسپلاین<sup>۵</sup> می‌نامند. یک برتری رگرسیون اسپلاین، برای روش‌های ناپارامتری و نیمه‌پارامتری این است که متعلق به رده‌ی رگرسیون خطی است. به عبارت دیگر، ضابطه تابع  $f(\cdot)$  نسبت به پارامترها خطی است. برآزش یک اسپلاین به داده‌ها می‌تواند با استفاده از روش کمترین توان‌های دوم خطا انجام پذیرد. استفاده از توابع پایه‌ای بریده‌شده، منجر به برآزش مدلی به صورت زیر می‌گردد:

$$y_j = \phi_0 + \phi_1 x_j + \dots + \phi_p x_j^p + \phi_{p+1} (x_j - t_1)_+^p + \dots + \phi_{p+k} (x_j - t_k)_+^p + \epsilon_j, \quad j = 1, \dots, n$$

که در آن  $\epsilon_j$  ها در پذیره‌های کلاسیک رگرسیون صدق می‌کنند. شکل ماتریسی این عبارت را می‌توان به صورت زیر بیان نمود:

$$y = D\phi + \epsilon$$

به طوری که  $\phi^T = (\phi_0, \phi_1, \dots, \phi_{p+k})^T$  و  $D$  یک ماتریس طرح  $n \times (p+k+1)$  است که درآیه‌های ستون اول آن برابر ۱، درآیه‌های  $p$  ستون بعدی، مربوط به مشاهدات متغیرهای  $x$ ،  $x^2$ ،  $\dots$  و  $x^p$  و درآیه‌های ستون‌های  $p+2$  تا  $(p+k+1)$  مربوط به مشاهده‌های حاصل از متغیرهای  $(x - t_1)_+^p$ ،  $\dots$  و  $(x - t_k)_+^p$  به ازای گره‌های از قبل تعیین شده  $t_1, t_2, \dots, t_k$  هستند. به عبارت دیگر

$$D = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p & (x_1 - t_1)_+^p & (x_1 - t_2)_+^p & \dots & (x_1 - t_k)_+^p \\ 1 & x_2 & x_2^2 & \dots & x_2^p & (x_2 - t_1)_+^p & (x_2 - t_2)_+^p & \dots & (x_2 - t_k)_+^p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^p & (x_n - t_1)_+^p & (x_n - t_2)_+^p & \dots & (x_n - t_k)_+^p \end{bmatrix}.$$

با به کار بردن روش کمترین توان‌های دوم، برآوردگر زیر حاصل می‌شود

$$\hat{\phi} = (D^T D)^{-1} D^T y. \quad (4.3)$$

مشکلی که در برآوردگر (۴.۳) وجود دارد، هم خطی متغیرهای

$$(x - t_1)_+^p, (x - t_2)_+^p, \dots, (x - t_k)_+^p$$

است که باعث عدم وارون‌پذیری یا بدشرطی<sup>۶</sup> ماتریس  $D^T D$  خواهد شد.

<sup>۵</sup>Spline Regression

<sup>۶</sup>Ill-Condition

۱.۲.۳  $B$ -اسپلینبازه  $[a, b]$  و تعداد  $q$  گره

$$a = t_0 < t_1 < \dots < t_q < t_{q+1} = b$$

را در نظر بگیرید. توابع  $B$ -اسپلین، چندجمله‌ای‌های قطعه‌ای در بازه  $[a, b]$  هستند که در گره‌ها به یکدیگر متصل بوده و دارای تکیه‌گاه کوچکی می‌باشند. متداول‌ترین روش محاسبه توابع  $B$ -اسپلین مرتبه  $p$ ، به وسیله توابع  $B$ -اسپلین مرتبه  $p-1$  و بر اساس رابطه بازگشتی زیر به دست می‌آید:

$$b_{i,p}(x) = \left(\frac{x - t_i}{t_{i+p} - t_i}\right)b_{i,p-1}(x) + \left(\frac{t_{i+p+1} - x}{t_{i+p+1} - t_{i+1}}\right)b_{i+1,p-1}(x), \quad i = -p, -p+1, \dots, q \quad (5.3)$$

که  $t_{q+1} = b$  و  $t_{-p} = t_{-p+1} = \dots = t_{-2} = t_{-1} = t_0 = a$  است. بنابراین برای یک  $B$ -اسپلین مرتبه  $p$  می‌توان تعداد  $p+q+1$  تابع  $B$ -اسپلین تعریف کرد.

کمترین توان‌های دوم  $B$ -اسپلین

ابتدا تعداد  $p+q+1$  تابع  $B$ -اسپلین مرتبه  $p$  رابطه  $(5.3)$  را با اندیس‌گذاری مجدد، به صورت

$$b_{0,p}, b_{1,p}, \dots, b_{p+q,p}$$

در نظر می‌گیریم. مدل رگرسیونی  $B$ -اسپلین به صورت زیر است:

$$s(x) = \sum_{j=0}^{p+q} \phi_j b_{j,p}(x) + \epsilon. \quad (6.3)$$

از آنجا که صورت ماتریسی این مدل به صورت

$$s = \mathbf{b}\phi + \epsilon$$

است، لذا برآورد کمترین توان‌های دوم  $\phi$  عبارت است از

$$\hat{\phi} = (\mathbf{b}^T \mathbf{b})^{-1} \mathbf{b}^T \mathbf{y}.$$

متعامد بودن توابع  $B$ -اسپلین موجب می‌شود که ماتریس  $\mathbf{b}^T \mathbf{b}$  وارون‌پذیر باشد و لذا مشکل موجود در رابطه  $(4.3)$  مرتفع می‌شود. لازم به ذکر است که علیرغم به دست آوردن این مزیت، درک توابع  $B$ -اسپلین به آسانی درک توابع بریده توانی نیست و لذا ضرایب مدل مذکور غیر قابل تفسیر است.

## ۲.۲.۳ P-اسپلین

ایلرز و مارکس (۱۹۹۶) اسپلین‌ها را به عنوان ترکیبات خطی از توابع B-اسپلین روی یک بازه از گره‌های هم‌فاصله نشان دادند. طبق مدل رگرسیونی B-اسپلین در رابطه‌ی (۶.۳) برآورد کمترین توان‌های دوم  $\beta$  از مینیم کردن رابطه‌ی زیر به دست می‌آید:

$$SSE = \sum_{i=1}^n \left\{ y_i - \sum_{j=0}^{p+q} \phi_j b_{j,p}(x_i) \right\}^2. \quad (7.3)$$

فرض کنید تعداد گره‌ها نسبتاً زیاد است، به طوری که منحنی برازش شده بیشترین تغییراتی که توسط داده‌ها نتیجه شده را نشان خواهد داد. برای اجتناب از پیچیدگی بیش از حد، سیلون (۱۹۸۶) مجموع توان‌های دوم جریمه‌شده<sup>۷</sup> را به صورت رابطه‌ی

$$S_\lambda(\phi_j) = \sum_{i=1}^n \left\{ y_i - \sum_{j=0}^{p+q} \phi_j b_{j,p}(x_i) \right\}^2 + \lambda \int_{x_{min}}^{x_{max}} \left\{ \sum_{j=0}^{p+q} \phi_j b_j''(x) \right\}^2 dx \quad (8.3)$$

معرفی کرد. ایلرز و مارکس (۱۹۹۶) جریمه مشتق دوم برای هموارسازی اسپلین‌ها را با تفاضل مرتبه‌ی دوم<sup>۸</sup> ضرایب B-اسپلین جایگذاری کردند و دلیل استفاده آن‌ها از تفاضل مرتبه‌ی دوم، مشتق دوم تابع B-اسپلین است.

بردار ضرایب  $\phi = (\phi_1, \phi_2, \dots, \phi_j)$  را در نظر بگیرید. تفاضل مرتبه‌ی اول به صورت زیر تعریف می‌شود:

$$\Delta\phi_j = \phi_j - \phi_{j-1}, \quad j = 2, \dots, p+q.$$

تفاضل مرتبه‌ی دوم، در حالت کلی به شکل زیر است:

$$\begin{aligned} \Delta^2\phi_j &= \Delta(\Delta\phi_j) \\ &= (\phi_j - \phi_{j-1}) - (\phi_{j-1} - \phi_{j-2}) = \phi_j - 2\phi_{j-1} + \phi_{j-2}, \quad j = 3, 4, \dots, p+q. \end{aligned}$$

به ازای هر  $j$  می‌توان این تفاضل را به صورت زیر نوشت:

$$\begin{aligned} \Delta(\Delta\phi_3) &= \phi_1 - 2\phi_2 + \phi_3 \\ \Delta(\Delta\phi_4) &= \phi_2 - 2\phi_3 + \phi_4 \\ \Delta(\Delta\phi_5) &= \phi_3 - 2\phi_4 + \phi_5 \\ &\vdots \end{aligned}$$

<sup>۷</sup>Penalized Sum Of Squares

<sup>۸</sup>2nd Order Differences

با چنین صورت نگارشی که ضرایب مشابه در زیر هم واقع شده‌اند، می‌توان ماتریس  $D_{(K-2) \times K}$  را که  $K = p + q$ ، متناظر با ضرایب  $\Delta^2 \phi_j$  به صورت ماتریس (۹.۳) در زیر تعریف کرد:

$$D = \begin{bmatrix} 1 & -2 & 1 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & -2 & 1 & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & -2 & 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 1 & -2 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & 1 & -2 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -2 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & -2 & 1 \end{bmatrix}. \quad (9.3)$$

حال ماتریس  $\mathcal{G}_{K \times K} = D^T D$  عبارت است از

$$\mathcal{G} = \begin{bmatrix} 1 & -2 & 1 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \\ -2 & 5 & -4 & 1 & 0 & 0 & 0 & 0 & \dots & 0 \\ 1 & -4 & 6 & -4 & 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & -4 & 6 & -4 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & -4 & 6 & -4 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 1 & -4 & 6 & -4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & -4 & 6 & -4 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -4 & 5 & -2 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & -2 & 1 \end{bmatrix}. \quad (10.3)$$

به آسانی می‌توان تحقیق کرد که تساوی زیر برقرار است (الیز و مارکس، ۱۹۹۶):

$$\phi^T \mathcal{G} \phi = \sum_{j=3}^K (\Delta^2 \phi_j)^2. \quad (11.3)$$

### ۳.۳ برآورد تابع مخاطره پایه با $P$ -اسپلین

تعداد  $K$  گره را روی بازه  $[0, t_{Rcens}]$  در نظر بگیرید که در آن  $t_{Rcens}$  کران بالای بازه تحت مطالعه یا همان بزرگترین زمان سانسور از راست است. این گره‌ها از پیش تعیین شده و به فاصله برابر از هم قرار دارند. برای برآورد تابع بقای پایه در مدل (۳.۲) به کمک اسپلین‌ها، لگاریتم تابع مخاطره پایه را به صورت ترکیبی خطی از  $B$ -اسپلین‌های مکعبی می‌نویسیم. فرض کنید  $b_1(t), \dots, b_K(t)$ ، توابع پایه‌ای  $B$ -اسپلین‌های مکعبی باشند. در این صورت می‌توان تابع بقای پایه را به صورت زیر نوشت:

$$h_0(t) = \exp\left(\sum_{k=1}^K b_k(t) \phi_k\right) \quad (12.3)$$

که در آن  $\phi = (\phi_1, \dots, \phi_K)$  ضرایب مجهول  $B$ -اسپلاین هستند و باید برآورد شوند. برای اطمینان از انعطاف لازم همراه با عدم بیش‌برازش<sup>۹</sup>، ایلرز و مارکس (۱۹۹۶) انتخاب تعداد بزرگی از  $B$ -اسپلاین‌ها را با افزودن یک جمله جریمه به لگاریتم تابع درست‌نمایی مدل به شکل

$$\tau \sum_k (\Delta^r \phi_k)^2 = \tau \phi^T \mathbf{D}^T \mathbf{D} \phi$$

پیشنهاد کردند، که در آن پارامتر جریمه و  $\mathbf{D}$  ماتریس تفاضلی جریمه مرتبه  $r$  است. برای مثال، اگر  $r = 2$  ماتریس تفاضلی مرتبه دو  $\mathbf{D}$  به صورت

$$\mathbf{D} = \begin{bmatrix} 1 & -2 & 1 & 0 & \dots & 0 \\ 0 & 1 & -2 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & -2 & 1 \end{bmatrix}$$

و اگر  $r = 3$  به صورت

$$\mathbf{D} = \begin{bmatrix} 1 & -3 & 3 & -1 & 0 & \dots & 0 \\ 0 & 1 & -3 & 3 & -1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & -3 & 3 & -1 \end{bmatrix}$$

تعریف می‌شود.

$P$ -اسپلاین‌ها در موقعیت‌های مختلفی مورد استفاده قرار می‌گیرند. محققین زیادی از جمله ایلرز و مارکس (۱۹۹۶)، لانگ و برزگر (۲۰۰۴) و لامبرت (۲۰۱۳، ۲۰۰۷) از آن‌ها در چارچوب بیزی استفاده کرده‌اند. نکته قابل توجه در همه این مراجع آن است که اگر  $K$  به اندازه کافی بزرگ (بین ۱۰ تا ۲۰) انتخاب شود، نیازی به معرفی معیارهای انتخاب مدل نیست؛ زیرا تمام  $K$ ‌ها تقریباً نتایج مشابهی را نشان می‌دهند.

### ۱.۳.۳ برآورد تابع بقای پایه

اگر  $f$  تابعی پیوسته (یا حتی ناپیوسته) و نامنفی بر بازه  $[a, b]$  باشد، برای محاسبه مساحت محصور به نمودار  $f$  و محور  $x$ ‌ها و خطوط  $x = a$  و  $x = b$ ، دنبال محاسبه مساحت مستطیل‌هایی هستیم که این محدوده (یعنی از  $a$  تا  $b$ ) را پوشش دهند. اگر ارتفاع مستطیل‌ها بالاتر از نمودار تابع باشد، جمع همه مساحت را تقریب اضافی<sup>۱۰</sup> و اگر ارتفاع مستطیل‌ها پایین‌تر از نمودار تابع باشد به آن تقریب نقصانی<sup>۱۱</sup> می‌گویند. در نهایت سطح زیر نمودار تقریب برابر است با مجموع همه مستطیل‌ها.

<sup>۹</sup>Over Fitting

<sup>۱۰</sup>Additional Approximation

<sup>۱۱</sup>Defective Approximation



همان‌طور که در فصل ۲ بیان شد

$$S_{\circ}(t) = \exp\left\{-\int_{\circ}^t h_{\circ}(u)du\right\}.$$

لذا با دانستن این رابطه بین تابع مخاطره پایه و تابع بقای پایه و جایگذاری رابطه (۱۲.۳) در آن، داریم

$$S_{\circ}(t) = \exp\left(-\int_{\circ}^t \exp\left(\sum_{k=1}^K b_k(u)\phi_k\right)du\right). \quad (13.3)$$

واضح است که رابطه‌ی (۱۳.۳) به صورت عددی غیرقابل بررسی است. در نتیجه باید آن را به صورت تقریبی به دست آوریم. بازه  $[\circ, t_{Rcens}]$  را در نظر بگیرید. این بازه را به یک دنباله متناهی از زیربازه‌ها با کران‌های  $\tau_0 = \circ < \tau_1 < \dots < \tau_j = t_{Rcens}$  تقسیم می‌کنیم که هر زیربازه به صورت  $J_j = [\tau_{j-1}, \tau_j]$  قابل نمایش است. انتخاب  $J = 3^{\circ}$  می‌تواند یک انتخاب معقول باشد. فرض کنید  $u_j$  و  $\delta_j$  به ترتیب میانه و عرض  $J_j$  باشند. با استفاده از محاسبه انتگرال به روش ریمان رابطه (۱۳.۳) را می‌توان به صورت زیر تقریب زد:

$$S_{\circ}(t) \approx \exp\left(-\sum_{j=1}^{j(t)} \exp\left(\sum_{k=1}^K b_k(u_j)\phi_k\right)\delta_j\right) \quad (14.3)$$

به طوری که  $j(t)$  نشان‌دهنده فاصله شامل  $t$  است. برای شناسایی پذیر شدن مدل، آخرین پارامتر اسپلین، یعنی  $\phi_K$ ، را یک مقدار به اندازه کافی بزرگ مثل  $\phi_K = 1^{\circ}$  در نظر می‌گیرند. این انتخاب باعث می‌شود تابع بقای پایه برآورد شده  $\hat{S}_{\circ}(\cdot)$  در پایان زمان مطالعه صفر شود.

### ۴.۳ تابع درستنمایی مدل

فرض کنید  $\mathbb{D}_i = (t_i, \nu_i, \mathbf{x}_i, \mathbf{z}_i)$  مجموعه مشاهدات آزمودنی  $i$ ام باشد به طوری که  $t_i$  نشان‌دهنده زمان شکست یا زمان سانسور مشاهده‌شده،  $\nu_i$  تابع مشخصه پیشامد مورد مطالعه و  $\mathbf{x}_i$  و  $\mathbf{z}_i$  دو بردار متغیرهای تبیینی در مدل نرخ شفایافتگی زمان پیشرفت هستند. اگر  $\mathbf{D} = \{\mathbb{D}_i; i = 1, \dots, n\}$  مجموعه همه مشاهدات و  $\Phi$  بردار تمام پارامترهای مدل باشد، تابع درستنمایی مدل به صورت زیر خواهد بود:

$$L(\Phi|\mathbf{D}) = \prod_{i=1}^n h_p(t_i)^{\nu_i} S_p(t_i)$$

که در آن  $S_p(t_i)$  در (۳.۲) نوشته شده است.

### ۵.۳ تحلیل بیزی

استنباط داده‌های بقا نسبت به سایر مدل‌های آماری پیچیده است. بنابراین همان‌طور که ذکر شد، رهیافت بیزی یک انتخاب معمول برای این داده‌ها است. در روش بیزی اگر اطلاعات اضافی در مورد پارامترها

در دست باشد، می‌توان از توزیع‌های پیشین برای وارد کردن آن‌ها در استنباط استفاده کرد. تابع چگالی پسین پارامترهای مدل از حاصل ضرب تابع درستنمایی و توابع چگالی پیشین نتیجه می‌شود. برای معرفی مدل بیزی نرخ شفایافتگی زمان پیشرفت پیشنهادی، ابتدا باید توزیع‌های پیشین پارامترها معرفی شوند.

### مدل بیزی

بر اساس لانگ و برزگر (۲۰۰۴)، تابع چگالی توزیع پیشین بردار ضرایب  $B$ -اسپلاین، یعنی  $\Phi$ ، را به صورت زیر در نظر می‌گیریم:

$$\pi(\phi|\tau) \propto \tau^{\frac{K}{\nu}} \exp\left(-\frac{\tau}{\nu} \phi^T \mathbf{P} \phi\right)$$

که در آن  $\mathbf{P} = \mathbf{D}^T \mathbf{D} + \epsilon \mathbf{I}_k$  یک ماتریس رتبه کامل به ازای یک مقدار کوچک برای  $\epsilon$  مانند  $10^{-6}$  است. به عبارت دیگر، یک توزیع نرمال چندمتغیره با بردار میانگین  $\mathbf{0}$  و ماتریس واریانس  $\mathbf{P}^{-1}$  برای پارامترهای اسپلاین در نظر گرفته شده است. بر اساس پیشنهاد جولین و لامبرت (۲۰۰۷) برای پارامتر جریمه از یک توزیع پیشین سلسله‌مراتبی به صورت زیر استفاده شده است:

$$\tau|\phi \sim G\left(\frac{\nu}{\nu}, \frac{\nu\delta}{\nu}\right)$$

$$\delta \sim G(a_\delta, b_\delta)$$

که در آن  $G(a, b)$  توزیع گاما با میانگین  $\frac{a}{b}$  و واریانس  $\frac{a}{b^2}$  را نشان می‌دهد. آن‌ها نشان دادند اگر یک مقدار کوچک برای  $a_\delta$  و  $b_\delta$  انتخاب شود (مثل  $10^{-4}$ )، انتخاب  $\nu$  (که اینجا برابر ۲ قرار داده شده است) بر شکل منحنی برآوردشده مخاطره پایه تاثیر نمی‌گذارد. اگر اطلاعات قبلی مثل یکنوا بودن در مورد مخاطره پایه در دسترس باشد، می‌توان آن‌را از طریق توزیع پیشین برای پارامترهای اسپلاین وارد کرد. برای تمام پارامترهای رگرسیونی  $\beta$ ،  $\beta_0$  و  $\gamma$  از توزیع پیشین نرمال استفاده شده است.

اکنون با استفاده از قضیه بیز، توزیع پسین مدل به صورت

$$\pi(\Phi|\mathbb{D}) \propto L(\Phi|\mathbf{D})\pi(\phi|\tau)\pi(\delta)\pi(\beta_0, \beta)\pi(\gamma) \quad (15.3)$$

نوشته می‌شود. توزیع پسین (۱۵.۳) شکل بسته ندارد و پیچیده است. بنابراین برای تقریب آن از الگوریتم‌های  $MCMC$  استفاده می‌کنیم.

با توجه به آن که توزیع‌های پیشین انتخابی ذکر شده همگی سره هستند، توزیع پسین نیز حتما سره است. اما تنها توابع چگالی شرطی کامل<sup>۱۲</sup> پارامترهای  $\tau$  و  $\delta$  شکل بسته دارند و برای سایر پارامترها شکل چگالی شناخته شده نیست. در واقع می‌توان نشان داد

$$\tau|\phi, \delta, \mathbb{D} \sim G\left(\frac{\nu + K}{\nu}, \frac{\nu\delta + \phi^T \mathbf{P} \phi}{\nu}\right)$$

<sup>۱۲</sup>Full Conditional Densities

و

$$\delta | \tau, \mathbb{D} \sim G(a_\delta + \frac{\nu}{\nu}, b_\delta + \frac{\nu \tau}{\nu}).$$

بنابراین از یک الگوریتم متروپولیس-هستینگز درون گیز برای تولید نمونه از پسین (۱۵.۳) باید استفاده کرد.

### ۶.۳ تولید نمونه از توزیع پسین با الگوریتم‌های MCMC

زمانی که توزیع پسین به دست آمده دارای ابعاد زیاد و هم‌چنین پیچیدگی زیادی باشد امکان محاسبه توزی‌های پسین پارامترهای مدل به روش مستقیم وجود ندارد. بنابراین می‌توان با استفاده از روش‌های نمونه‌گیری MCMC برای هر حجم نمونه دلخواه، استنباط قابل قبولی از پارامترها ارائه داد. در این روش با نمونه‌گیری‌های متوالی از توزیع شرطی کامل هر یک از پارامترها، با استفاده از یک الگوریتم مناسب مانند نمونه‌گیری گیز و متروپولیس-هستینگز. زنجیرهای مارکوفی تشکیل می‌شود که در صورت همگرایی آن‌ها، می‌توان مقادیر شبیه‌سازی شده را نمونه‌هایی از توزیع پسین پارامترها دانست و بر اساس آن استنباط در خصوص پارامترهای مدل را انجام داد. الگوریتم تولید نمونه از توزیع پسین (۱۵.۳) به صورت زیر است:

فرض کنید  $\tilde{\beta}_0 = (\beta_0, \beta)$  و  $\phi^{(0)}, \tau^{(0)}, \delta^{(0)}, \tilde{\beta}^{(0)}$  و  $\gamma^{(0)}$  مقادیر اولیه برای شروع زنجیر باشد. در تکرار  $m$  مراحل الگوریتم عبارتند از:

(۱) تولید  $\phi^{(m)}$  از  $(\phi | \tau^{(m-1)}, \delta^{(m-1)}, \tilde{\beta}^{(m-1)}, \gamma^{(m-1)})$  با استفاده از یک الگوریتم متروپولیس-هستینگز.

(۲) تولید  $\tau^{(m)}$  از  $G(\frac{\nu+k}{\nu}, \frac{\nu \delta^{(m-1)} + \phi^{(m)T} \mathbf{P} \phi^{(m)}}{\nu})$  با استفاده از یک مرحله نمونه‌گیری گیز.

(۳) تولید  $\delta^{(m)}$  از  $G(a_\delta + \frac{\nu}{\nu}, b_\delta + \frac{\nu \tau^{(m)}}{\nu})$  با استفاده از یک مرحله نمونه‌گیری گیز.

(۴) تولید  $\tilde{\beta}^{(m)}$  از چگالی شرطی کامل  $\pi(\tilde{\beta} | \phi^{(m)}, \tau^{(m)}, \delta^{(m)}, \gamma^{(m-1)})$  با استفاده از یک الگوریتم متروپولیس-هستینگز چندمتغیره.

(۵) تولید  $\gamma^{(m)}$  از چگالی شرطی کامل  $\pi(\tilde{\gamma} | \phi^{(m)}, \tau^{(m)}, \delta^{(m)}, \beta^{(m)})$  با استفاده از یک الگوریتم متروپولیس-هستینگز چندمتغیره.

# فصل ۴

## ارزیابی عملکرد مدل پیشنهادی

در این فصل با استفاده از مثال‌های واقعی و شبیه‌سازی، عملکرد مدل پیشنهادی معرفی شده در فصل سوم را ارزیابی می‌کنیم. برای این منظور، در مطالعه شبیه‌سازی دو سناریوی زمان پیگیری طولانی و کوتاه را در نظر می‌گیریم تا ویژگی شناسایی‌پذیر بودن مدل را نیز بررسی کنیم. مثال واقعی نیز یک مجموعه داده از بیماران مبتلا به ملانوما پستی را شامل می‌شود. نتایج گزارش شده در این فصل برگرفته از مقاله برمهورست و لامبرت (۲۰۱۶) هستند.

### ۱.۴ مطالعه شبیه‌سازی

#### ۱.۱.۴ پیگیری طولانی مدت

همان‌طور که می‌دانیم اگر روند مطالعه به اندازه کافی طولانی باشد، مدل پیشنهادی در (۱۲.۳) شناسایی‌پذیر است و در نتیجه بزرگترین زمان سانسور بیشتر از بزرگترین زمان شکست است. در این شبیه‌سازی در هر سناریوی در نظر گرفته شده، تابع بقای پایه در مدل (۱۲.۳) دارای توزیع وایبل با میانگین ۸ و انحراف معیار ۴/۱۸ است. هم‌چنین متغیرهای تبیینی رگرسیونی، دو متغیر به صورت  $W_1 \sim N(0, 1)$  و  $W_2 \sim Bernoulli(0.5)$  در نظر گرفته شدند. از آنجایی که در این شبیه‌سازی فرض پیگیری کافی برقرار است، بنابراین هر دو متغیر می‌توانند به‌طور هم‌زمان برای مدل نسبت شفایافتگی و زمان لازم برای تشکیل یک تومور قابل تشخیص، بدون ایجاد یک مشکل شناسایی‌پذیری، مورد استفاده قرار گیرند. بنابراین

$$\mathbf{X} = \{W_1, W_2\} = \mathbf{Z}$$

در نظر گرفته شد. ضرایب رگرسیونی مربوط به  $W_1$  و  $W_2$  در مدل مخاطرات متناسب کاکس، به ترتیب ۰/۴ و -۰/۴ در نظر گرفته شدند. کران بالا برای زمان شکست مشاهده‌شده نیز عدد ۲۳ قرار داده شد. با این انتخاب بیش از ۹۹٪ پیشامدهای تولید شده تحت توزیع وایبل انتخاب‌شده قبل از  $t = ۲۳$  اتفاق

می‌افتند. دو درصد ۲۵ و ۴۰ نیز برای نسبت افراد شفایافته در نظر گرفته شدند. این دو نسبت با تنظیم مقادیر واقعی ضرایب رگرسیونی  $(\beta_0, \beta_1, \beta_2)$  در پارامتر نسبت شفایافته به دست می‌آیند. با این توضیحات، هر مجموعه داده با استفاده از مدل زمان پیشرفت پیشنهادی به شرح زیر تولید شد:

$$(1) \text{ تولید تعداد سلول‌های سرطان‌زا با استفاده از } N \sim Pois(\theta(\mathbf{x})) \text{ که در آن } \theta(\mathbf{x}) = \exp(\beta_0 + \mathbf{x}^T \beta).$$

(۲) اگر  $N \neq 0$ ، متغیرهای پنهان  $Y_1, \dots, Y_N$  با استفاده از مدل مخاطرات متناسب کاکس تولید می‌شوند. سپس زمان شکست مشاهده‌شده به صورت  $T = \min\{Y_1, \dots, Y_N\}$  تعریف می‌شود. این مرحله تا زمانی که  $T < 23$  تکرار می‌شود. توجه داشته باشید که این فرآیند در کمتر از ۱٪ موارد باید بیش از یک بار تکرار شود. اگر  $N = 0$ ، آنگاه زمان شکست برای فرد شفایافته مقدار دلخواه بزرگی مثل ۹۹۹ قرار داده می‌شود.

(۳) نرخ کلی سانسور از راست تحت یکی از دو توزیع سانسور زیر کنترل می‌شود:

(a) سناریو (۱): یک توزیع یکنواخت در فاصله  $[20, 25]$ . این توزیع سانسور تضمین می‌کند که تقریباً همه افراد از راست سانسور شده، شفایافته و شناسایی‌پذیر هستند.

(b) سناریو (۲): یک توزیع وایبل با میانگین ۲۲/۲۸ و انحراف معیار ۸/۰۸ که در ۲۵ بریده<sup>۱</sup> شده‌است. با استفاده از این توزیع سانسور، زمان‌های سانسور تنها ۲۵٪ افراد شفایافته در دم راست برآوردگر کاپلان‌مایر منحنی بقا قرار می‌گیرند و ۴٪ افراد شفایافته، از راست سانسور می‌شوند.

حال با استفاده از رابطه (۱۲.۳) و در نظر گرفتن ۱۲ گره با فاصله‌های برابر در  $[0, t_{Rcens}]$  که  $t_{Rcens}$  در آن برابر با ۲۵ است و ماتریس جریمه تفاضلی مرتبه سوم اقدامات زیر را انجام می‌دهیم:

- تولید تعداد ۵۰۰ مجموعه داده با حجم‌های نمونه  $n = 300$  و  $n = 600$ .
- تولید یک زنجیر  $MCMC$  با حجم مونت‌کارلوی ۲۳۰۰۰ از توزیع پسین مدل. دوره داغیدن<sup>۲</sup> نیز برابر ۳۰۰۰ انتخاب می‌شود.
- بررسی همگرایی زنجیر تولید شده توسط معیار تشخیص جی‌وک (جی‌وک، ۱۹۹۲). برای جزئیات این معیار پیوست **ب** را مشاهده کنید.
- محاسبه برآورد پارامترها (میانه توزیع پسین) و سایر کمیت‌های مورد نظر.

نتایج حاصل برای پارامترهای رگرسیونی  $\beta$ ،  $\tilde{\beta}$  و  $\tilde{\gamma}$  برای دو حجم نمونه  $n = 300$  و  $n = 600$  به ترتیب در جدول‌های ۱.۴ و ۲.۴ گزارش شده‌اند. در جدول‌ها، ستون اول نسبت افراد شفایافته، ستون دوم

<sup>۱</sup>Truncated

<sup>۲</sup>Burn-In

سناریوهای کنترل نرخ کلی سانسور، ستون سوم مقادیر واقعی پارامترها، ستون چهارم اربیی برآوردها، دو ستون پنجم و ششم احتمال پوشش متناظر فواصل اعتبار ۹۰ و ۹۵ درصد، ستون هفتم خطای استاندارد تجربی<sup>۳</sup> ( $ESE$ ) و ستون آخر مجذور میانگین توان دوم خطا<sup>۴</sup> ( $RMSE$ ) را نشان می‌دهند. هنگامی که درصد افراد شفایافته ۲۵٪ و ۴۰٪ باشد، جدول‌های ۱.۴ و ۲.۴ نتایج شبیه‌سازی برای پارامترهای رگرسیون را نشان می‌دهند.

جدول ۱.۴: نتایج شبیه‌سازی با پیگیری طولانی مدت برای  $n = 300$  و  $500$  مجموعه داده

شفایافته	سناریو	پارامترها	اربیی	$CV_{/90}$	$CV_{/95}$	$ESE$	$RMSE$
۲۵٪	۱	$\beta_0 = 0.75$	0.028	88.4	93.8	0.130	0.018
		$\beta_1 = 0.80$	0.016	90.8	94.8	0.112	0.013
		$\beta_2 = -0.50$	-0.013	91.6	96.6	0.172	0.030
		$\gamma_1 = 0.40$	-0.035	91.0	95.8	0.138	0.142
	۲	$\gamma_2 = -0.40$	0.008	87.4	94.2	0.221	0.221
		$\beta_0 = 0.75$	0.052	85.8	92.6	0.162	0.030
		$\beta_1 = 0.80$	-0.003	93.4	97.8	0.129	0.017
		$\beta_2 = -0.50$	0.017	88.6	94.0	0.209	0.044
		$\gamma_1 = 0.40$	0.001	92.8	95.8	0.171	0.171
		$\gamma_2 = -0.40$	-0.029	88.2	93.6	0.283	0.282
۴۰٪	۱	$\beta_0 = 0.30$	0.006	90.6	94.0	0.128	0.017
		$\beta_1 = 1.00$	0.013	90.6	95.0	0.119	0.014
		$\beta_2 = -0.75$	-0.008	90.6	94.4	0.182	0.033
		$\gamma_1 = 0.40$	-0.021	90.0	95.6	0.150	0.151
	۲	$\gamma_2 = -0.40$	0.003	92.0	96.8	0.215	0.215
		$\beta_0 = 0.30$	0.043	86.2	91.8	0.155	0.027
		$\beta_1 = 1.00$	-0.014	93.6	97.0	0.137	0.019
		$\beta_2 = -0.75$	0.024	88.0	93.4	0.222	0.049
		$\gamma_1 = 0.40$	0.005	91.8	96.4	0.183	0.183
		$\gamma_2 = -0.40$	-0.025	90.4	94.4	0.276	0.276

<sup>۳</sup>Empirical Standard Error

<sup>۴</sup>Root Mean Squared Error

جدول ۲۰۴: نتایج شبیه‌سازی با پیگیری طولانی مدت برای  $n = ۶۰۰$  و  $۵۰۰$  مجموعه داده

شفایافته	سناریو	پارامترها	اریبی	$CV_{\%۹۰}$	$CV_{\%۹۵}$	$ESE$	$RMSE$
%۲۵	۱	$\beta_0 = ۰٫۷۵$	۰٫۰۲۰	۹۰٫۲	۹۴٫۰	۰٫۰۸۷	۰٫۰۰۸
		$\beta_1 = ۰٫۸۰$	۰٫۰۱۷	۹۰٫۴	۹۵٫۲	۰٫۰۷۷	۰٫۰۰۶
		$\beta_2 = -۰٫۵۰$	-۰٫۰۲۰	۹۱٫۶	۹۵٫۸	۰٫۱۱۵	۰٫۰۱۴
		$\gamma_1 = ۰٫۴۰$	-۰٫۰۳۸	۸۷٫۶	۹۳٫۶	۰٫۰۹۷	۰٫۱۰۶
		$\gamma_2 = -۰٫۴۰$	۰٫۰۲۲	۸۹٫۸	۹۶٫۴	۰٫۱۳۹	۰٫۱۴۱
	۲	$\beta_0 = ۰٫۷۵$	۰٫۰۳۶	۸۷٫۲	۹۲٫۸	۰٫۱۰۷	۰٫۰۱۳
		$\beta_1 = ۰٫۸۰$	۰٫۰۰۹	۹۳٫۲	۹۷٫۲	۰٫۰۹۰	۰٫۰۰۸
		$\beta_2 = -۰٫۵۰$	-۰٫۰۱۱	۸۸٫۶	۹۳٫۶	۰٫۱۵۰	۰٫۰۲۳
		$\gamma_1 = ۰٫۴۰$	-۰٫۰۲۳	۹۱٫۰	۹۶٫۰	۰٫۱۱۵	۰٫۱۱۸
		$\gamma_2 = -۰٫۴۰$	۰٫۰۱۸	۹۱٫۶	۹۵٫۰	۰٫۱۸۲	۰٫۱۸۳
%۴۰	۱	$\beta_0 = ۰٫۳۰$	۰٫۰۰۵	۹۱٫۲	۹۵٫۲	۰٫۰۸۷	۰٫۰۰۸
		$\beta_1 = ۱٫۰۰$	۰٫۰۱۶	۹۱٫۲	۹۵٫۶	۰٫۰۷۹	۰٫۰۰۶
		$\beta_2 = -۰٫۷۵$	-۰٫۰۰۳	۹۱٫۶	۹۵٫۶	۰٫۱۱۹	۰٫۰۱۴
		$\gamma_1 = ۰٫۴۰$	-۰٫۰۲۷	۹۱٫۲	۹۵٫۰	۰٫۰۹۸	۰٫۱۰۱
		$\gamma_2 = -۰٫۴۰$	۰٫۰۳۶	۹۲٫۴	۹۵٫۲	۰٫۱۴۲	۰٫۱۴۶
	۲	$\beta_0 = ۰٫۳۰$	۰٫۰۲۱	۸۹٫۴	۹۴٫۸	۰٫۰۹۴	۰٫۰۰۹
		$\beta_1 = ۱٫۰۰$	-۰٫۰۰۴	۹۲٫۸	۹۵٫۴	۰٫۰۹۳	۰٫۰۰۹
		$\beta_2 = -۰٫۷۵$	۰٫۰۱۰	۹۱٫۲	۹۵٫۸	۰٫۱۴۲	۰٫۰۲۰
		$\gamma_1 = ۰٫۴۰$	-۰٫۰۰۲	۹۳٫۶	۹۷٫۸	۰٫۱۱۵	۰٫۱۱۵
		$\gamma_2 = -۰٫۴۰$	۰٫۰۲۵	۹۲٫۶	۹۶٫۴	۰٫۱۸۴	۰٫۱۸۴

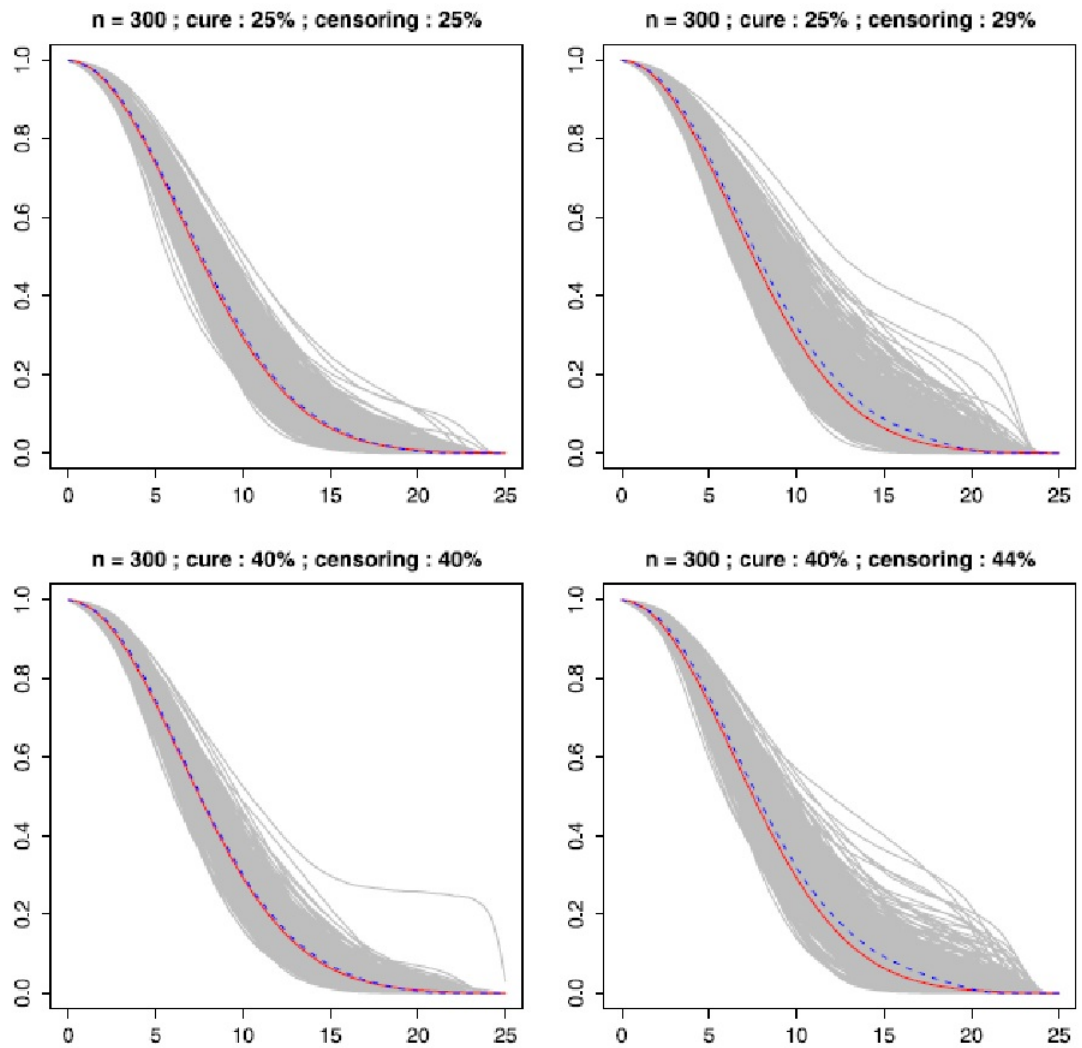
با توجه به اعداد جدول‌ها می‌توان نتایج زیر را خلاصه کرد:

- (۱) برآوردهای ضرایب رگرسیونی برای هر دو حجم نمونه در نظر گرفته شده اریبی ناچیزی دارند.
- (۲) خطای تجربی و  $RMSE$  برآوردهای پارامترهای رگرسیونی با افزایش حجم نمونه کوچک می‌شوند. یعنی دقت برآوردها افزایش می‌یابند.  
برعکس این دو معیار در دو حالت افزایش می‌یابند:
- (الف) زمانی که نسبت افراد شفایافته با زمان سانسور بزرگتر از ماکسیمم زمان شکست مشاهده شده کاهش یابد.
- (ب) درصد سانسور از راست در بین افراد شفانیافته افزایش یابد.
- (۳) احتمال‌های پوشش فواصل اعتبار ۹۰ و ۹۵ درصد به مقادیر اسمی خود نزدیک هستند.

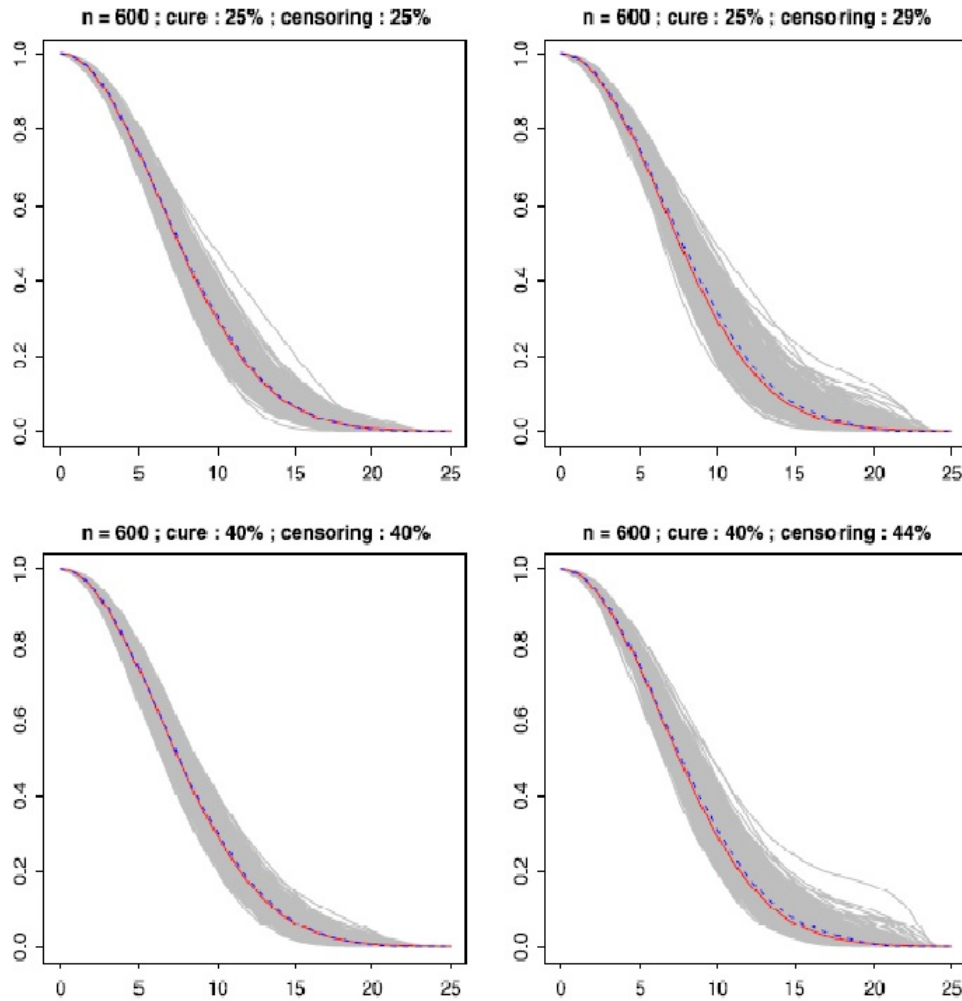
(۴) نسبت افراد شفایافته بر دقت برآوردها تاثیر چندانی ندارند.

برآوردهای تابع بقای پایه برای حجم‌های نمونه  $n = 300$  و  $n = 600$  به ترتیب در شکل‌های ۱.۴ و ۲.۴ نمایش داده شده‌اند. از روی منحنی‌های برآورد شده می‌توان دریافت که زمانی که نسبت افراد شفایافته با زمان سانسور بزرگتر از ماکسیمم زمان شکست مشاهده شده، کاهش یابد و هم‌چنین زمانی که درصد سانسور از راست در بین افراد شفانیافته افزایش یابد، تغییرپذیری توزیع پایه برآورد شده تا حدودی افزایش خواهد یافت. از طرفی این تغییرپذیری با افزایش حجم نمونه کاهش می‌یابد. از طرف دیگر می‌توان گفت تا حدودی اریبی برآورد دم سمت راست  $S_0(t)$  با افزایش حجم نمونه زمانی که نسبت افراد شفایافته با زمان سانسور بزرگتر از ماکسیمم زمان شکست مشاهده شده کاهش یابد و نسبت افراد از راست سانسور شده غیرمصون افزایش یابد، کاهش یافته است.



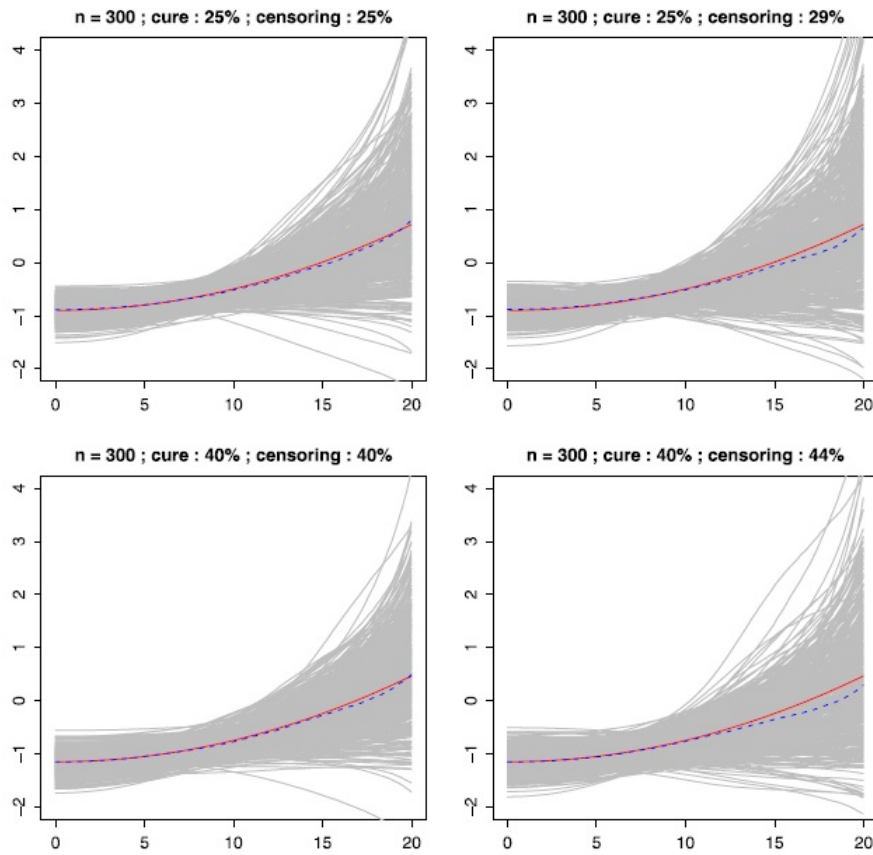


شکل ۱.۴: برآورد  $S_0(t)$  در مثال شبیه‌سازی با پیگیری طولانی مدت برای ۵۰۰ مجموعه داده و  $n = 300$ . هر منحنی خاکستری برآورد تابع بقا در یک مجموعه داده است. هر ردیف شکل به درصدی از افراد شفایافته (ردیف بالا: ۲۵٪ و ردیف پایین: ۴۰٪) با نرخ کلی سانسور از راست (چپ: سناریو ۱ و راست: سناریو ۲) اختصاص دارد. منحنی خط ممتد نمودار واقعی  $S_0(t)$  و منحنی خط چین میانه حاصل از ۵۰۰ تکرار است.

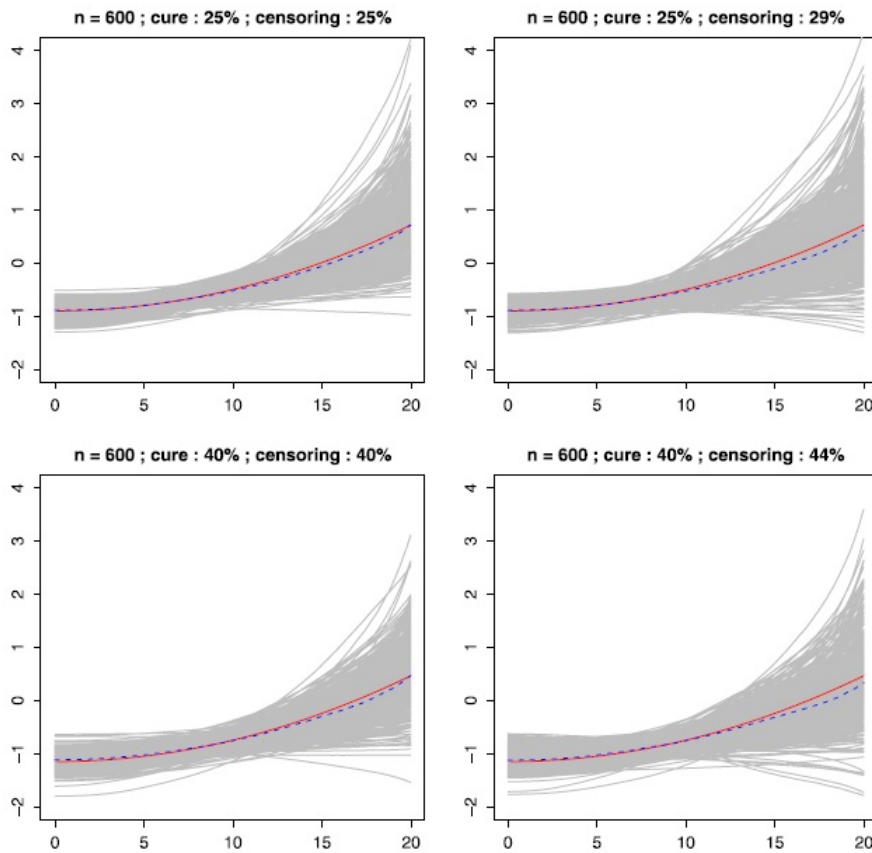


شکل ۲.۴: برآورد  $S_0(t)$  در مثال شبیه‌سازی با پیگیری طولانی مدت برای  $500$  مجموعه داده و  $n = 600$ . هر منحنی خاکستری برآورد تابع بقا در یک مجموعه داده است. هر ردیف شکل به درصدی از افراد شفایافته (ردیف بالا: ۲۵٪ و ردیف پایین: ۴۰٪) با نرخ کلی سانسور از راست (چپ: سناریو ۱ و راست: سناریو ۲) اختصاص دارد. منحنی خط ممتد نمودار واقعی  $S_0(t)$  و منحنی خط چین میانه حاصل از  $500$  تکرار است.

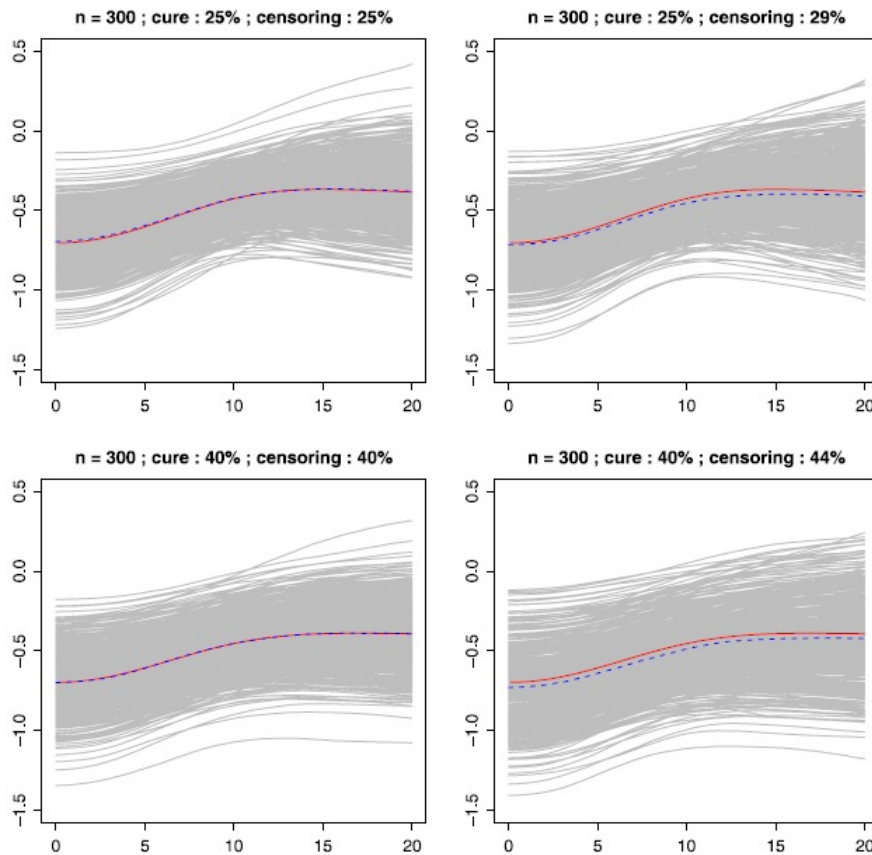
شکل‌های ۳.۴ و ۴.۴، برآورد نسبت لگاریتم مخاطره کل جامعه را برای به ترتیب  $n = 300$  و  $n = 600$  نمایش می‌دهند. نتایج ذکر شده برای برآوردهای  $S_0(t)$  در نمودارهای ۱.۴ و ۲.۴، برای نسبت لگاریتم مخاطره نیز برقرار هستند و قابل درک است که نسبت افراد شفایافته تاثیری بر دقت این برآوردها نداشته‌اند. البته نسبت لگاریتم مخاطره برای افراد غیرمصون (۴.۴ برای  $n = 300$ ) تحت هر سناریو به خوبی برآورد شده است.



شکل ۳.۴: برآورد نسبت لگاریتم مخاطره جامعه  $\log(HR_p(t))$  در مثال شبیه‌سازی با پیگیری طولانی مدت برای ۵۰۰ مجموعه داده و  $n = 300$ . هر یک منحنی خاکستری برآورد در یک مجموعه داده است. ردیف بالا به ۲۵٪ و ردیف پایین به ۴۰٪ از افراد شفایافته اشاره دارند. سمت چپ سناریو ۱ و سمت راست سناریو ۲ سانسور از راست را نشان می‌دهند.



شکل ۴.۴: برآورد نسبت لگاریتم مخاطره جامعه  $\log(HR_p(t))$  در مثال شبیه‌سازی با پیگیری طولانی مدت برای  $500$  مجموعه داده و  $n = 600$ . هر یک منحنی خاکستری برآورد در یک مجموعه داده است. ردیف بالا به  $25\%$  و ردیف پایین به  $40\%$  از افراد شفایافته اشاره دارند. سمت چپ سناریو ۱ و سمت راست سناریو ۲ سانسور از راست را نشان می‌دهند.



شکل ۵.۴: برآورد نسبت لگاریتم مخاطره افراد غیرمضون در مثال شیه‌سازی با پیگیری طولانی مدت برای  $500$  مجموعه داده و  $n = 300$ . هر یک منحنی خاکستری برآورد در یک مجموعه داده است. ردیف بالا به  $25\%$  و ردیف پایین به  $40\%$  از افراد شفایافته اشاره دارند. سمت چپ سناریو ۱ و سمت راست سناریو ۲ سانسور از راست را نشان می‌دهند.

## ۲.۱.۴ پیگیری کوتاه مدت

هدف از این مثال، استفاده محدود از مدل نرخ شفایافتگی پیشرفت زمان هنگامی که پیگیری مطالعه به اندازه کافی طولانی نیست، می‌باشد. مجموعه داده‌ها با استفاده از روش شرح داده‌شده در بخش ۱.۱.۴ تولید شدند. از آنجا که پیگیری مطالعه به اندازه کافی بلند نیست، برای شناسایی پذیر بودن مدل، متغیر تبیینی  $W_1$  را تنها برای بخش تولید یک تومور قابل تشخیص ( $Z = W_1$ ) و متغیر تبیینی  $W_2$  را تنها برای بخش نسبت شفایافتگی ( $X = W_2$ ) در نظر گرفتیم. دو مقدار  $25\%$  و  $40\%$  نیز برای نسبت افراد شفایافته در نظر گرفته شدند. برای اطمینان از اینکه بزرگترین زمان سانسور کوچکتر از بزرگترین زمان شکست مشاهده‌شده (ویژگی پیگیری کوتاه مدت) است، یک توزیع وایبل با میانگین  $17/8$  و انحراف معیار  $6/5$  که در  $13/7$  (سناریوی ۳) و  $10/6$  (سناریوی ۴) بریده شده‌اند، در نظر گرفتیم. مشابه ۱.۱.۴

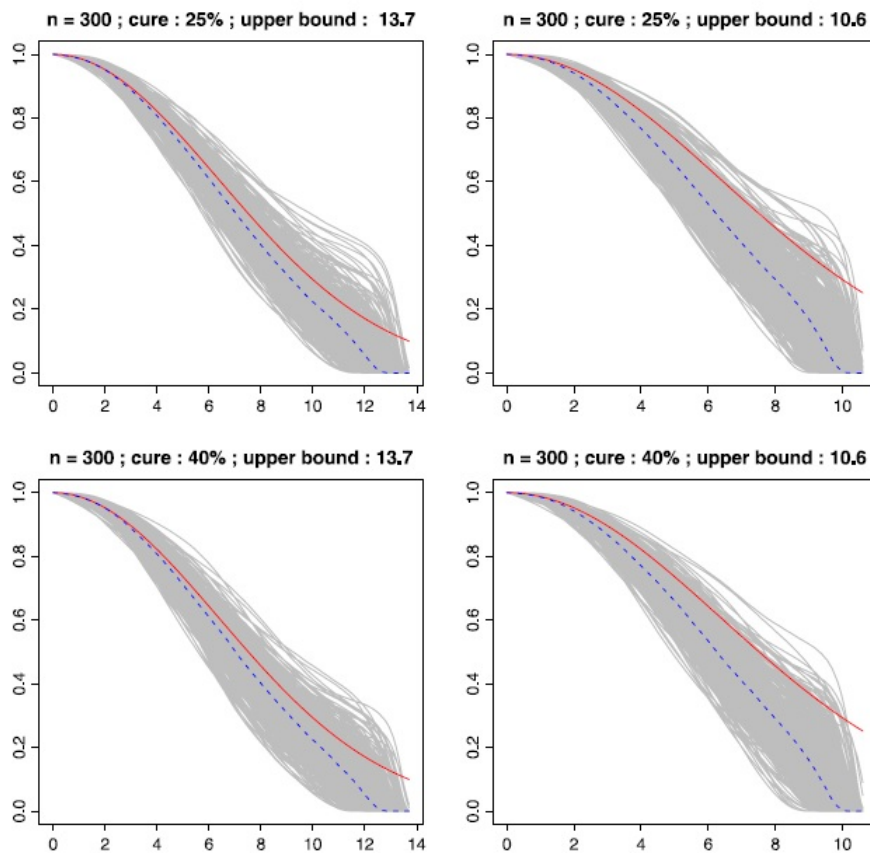
و با در نظر گرفتن ۱۲ گره با فاصله‌های برابر روی  $[0, t_{Rcens}]$  که در آن  $t_{Rcens}$  برابر با  $۱۳/۷$  یا  $۱۰/۶$  است، شبیه‌سازی برای نمونه‌های به حجم  $n = 300$  و  $n = 600$  با  $500$  تکرار انجام داده شد. مشابه مثال قبلی زنجیری به طول  $23000$  با دوره داغیدن  $3000$  تولید شد که آماره  $Z$  معیار جی‌وک تمام پارامترهای مدل برای بررسی همگرایی الگوریتم  $MCMC$  بررسی و همگرایی زنجیر تایید شد. نتایج حاصل برای پارامترهای رگرسیونی  $\beta$ ،  $\tilde{\beta}$ ،  $\gamma$  و  $\tilde{\gamma}$  برای حجم نمونه  $n = 300$  در جدول ۳.۴ گزارش شده‌اند.

جدول ۳.۴: نتایج شبیه‌سازی با پیگیری کوتاه مدت برای  $n = 300$  و  $500$  مجموعه داده

شفایافته	سناریو	پارامترها	اریبی	$CV_{\%90}$	$CV_{\%95}$	$ESE$	$RMSE$
%۲۵	۳	$\beta_0 = 0.70$	$-0.093$	۸۰/۶	۸۷/۸	۰/۱۲۹	۰/۰۲۵
		$\beta_1 = -0.70$	$0.007$	۸۸/۴	۹۵/۶	۰/۱۵۶	۰/۰۲۴
		$\gamma_1 = 0.40$	$-0.029$	۸۶/۶	۹۲/۲	۰/۱۱۱	۰/۱۱۴
	۴	$\beta_0 = 0.70$	$-0.266$	۴۶/۶	۶۰/۶	۰/۱۵۷	۰/۰۹۶
		$\beta_1 = -0.70$	$0.007$	۹۱/۰	۹۳/۸	۰/۱۶۹	۰/۰۲۹
		$\gamma_1 = 0.40$	$-0.046$	۸۵/۶	۹۱/۰	۰/۱۲۴	۰/۱۳۰
%۴۰	۳	$\beta_0 = 0.30$	$-0.112$	۷۷/۴	۸۵/۰	۰/۱۴۲	۰/۰۳۳
		$\beta_1 = -0.80$	$0.024$	۹۱/۰	۹۶/۰	۰/۱۷۸	۰/۰۳۲
		$\gamma_1 = 0.40$	$-0.039$	۸۶/۸	۹۳/۰	۰/۱۲۷	۰/۱۳۳
	۴	$\beta_0 = 0.30$	$-0.292$	۴۵/۴	۵۶/۶	۰/۱۵۹	۰/۱۰۴
		$\beta_1 = -0.80$	$0.024$	۹۲/۰	۹۶/۲	۰/۱۹۵	۰/۰۳۸
		$\gamma_1 = 0.40$	$-0.049$	۸۶/۲	۹۲/۲	۰/۱۳۷	۰/۱۴۵

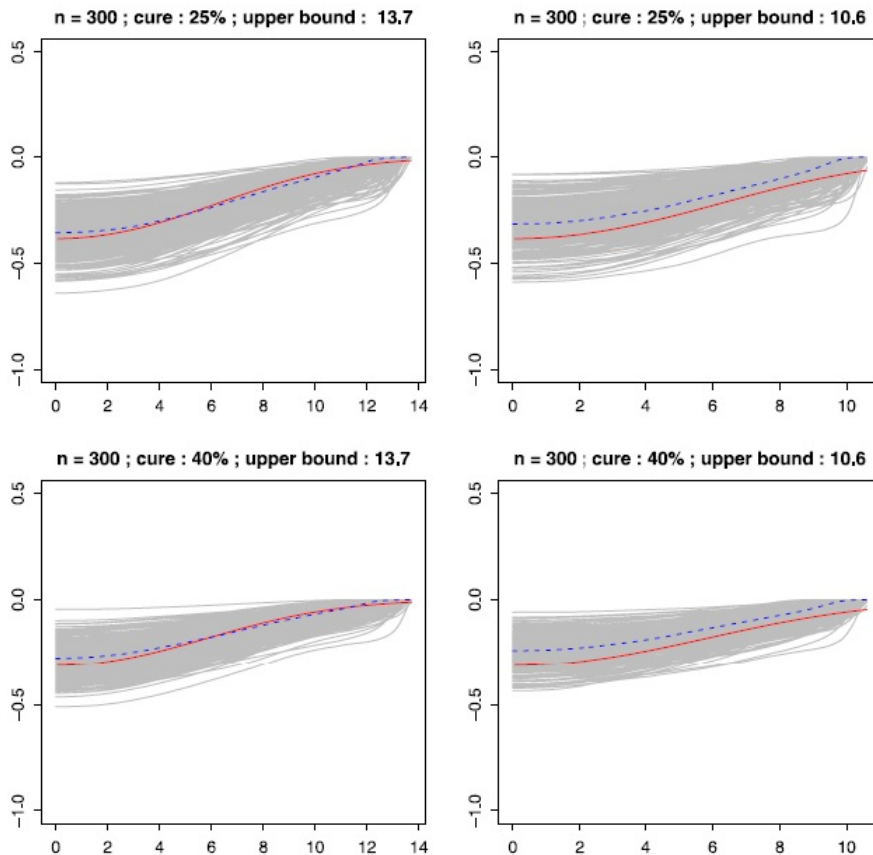
با توجه به اعداد جدول، نتایج زیر قابل طرح هستند:

- (۱) احتمال‌های پوشش فواصل اعتبار ۹۰ و ۹۵ درصد به مقادیر اسمی خود نزدیک هستند.
- (۲) برآوردهای ضرایب رگرسیونی برای حجم نمونه در نظر گرفته‌شده اریبی ناچیزی دارند.



شکل ۶.۴: برآورد  $S_0(t)$  در مثال شبیه‌سازی با پیگیری کوتاه مدت برای  $500$  مجموعه داده و  $n = 300$ . هر منحنی خاکستری برآورد تابع بقا در یک مجموعه داده است. هر ردیف شکل به درصدی از افراد شفایافته (ردیف بالا: ۲۵٪ و ردیف پایین: ۴۰٪) با نرخ کلی سانسور از راست (چپ: سناریو ۳ و راست: سناریو ۴) اختصاص دارد. منحنی خط ممتد نمودار واقعی  $S_0(t)$  و منحنی خط چین میانه حاصل از  $500$  تکرار است.

شکل ۶.۴ نشان می‌دهد که توزیع پایه  $S_0(t)$  به خاطر محدودیت در دم صفر برآورد نمی‌شود.



شکل ۷.۴: برآورد نسبت لگاریتم مخاطره افراد غیرمصون در مثال شبیه‌سازی با پیگیری کوتاه مدت برای  $n = 300$  و مجموعه داده و  $n = 300$ . هر منحنی خاکستری برآورد تابع بقا در یک مجموعه داده است. ردیف بالا به ۲۵٪ و ردیف پایین به ۴۰٪ از افراد شفایافته اشاره دارند. سمت چپ سناریو ۳ و سمت راست سناریو ۴ سانسور از راست را نشان می‌دهند.

نسبت لگاریتم مخاطره برای افراد غیرمصون (شکل ۷.۴ برای  $n = 300$ ) تحت هر سناریو به خوبی برآورد شده است. همان‌طور که در شکل ۷.۴ مشاهده می‌شود، نسبت مخاطره در افراد شفایافته زمانی که کران بالای پیگیری کوچک است، بیش‌برازش می‌شود. این نتیجه زمانی حاصل می‌شود که توزیع بقای پایه در رابطه (۳.۲) در ماکسیمم زمان سانسور ممکن کوچک‌تر از یک باشد. در واقع در کوتاه‌ترین پیگیری (سناریو ۴)  $F_0(10/6) = 0.75$  و در مطلوب‌ترین پیگیری (سناریو ۳)  $F_0(13/7) = 0.9$ .

نتایج مشابهی را می‌توان برای نمونه به حجم  $n = 600$  نیز بیان کرد. به‌طور کلی نتایج به‌دست آمده از این مثال شبیه‌سازی، لم ۲.۵.۲ را تایید می‌کنند.



## ۲.۴ مثال واقعی

این مطالعه بر روی ۲۸۴ بیمار مبتلا به ملانوما پوستی توسط گروه *ECOG*<sup>۵</sup> انجام شده است. در این مطالعه که شامل یک آزمایش بالینی تصادفی دومارحله‌ای است، تاثیر داروی اینترفرون آلفا-۲ بی ارزیابی شده است (کرک‌وود و همکاران، ۱۹۹۶). مطالعه آن‌ها نشان داد که اینترفرون آلفا-۲ بی تاثیر مثبتی بر زمان بقا بدون بازگشت بیماری دارد. از میان بیماران حاضر در مطالعه ۱۴۴ نفر (۵۱٪) درمان *IFN* دریافت کرده‌اند، که ۱۷۱ نفر (۶۰٪) از آن‌ها مرد هستند و عود سرطان برای ۱۹۶ نفر (۶۹٪) آن‌ها اتفاق افتاده است. سایر افراد در این مطالعه از راست سانسور شده‌اند. میانگین سن بیماران ۴۷ سال و انحراف معیار آن‌ها ۱۳ سال است. شکل ۸.۴ منحنی برآوردشده کاپلان مایر را برای تابع بقای بیماران نشان می‌دهد. با توجه به این شکل معلوم است که پیگیری مطالعه به اندازه کافی طولانی است. بنابراین فرض می‌کنیم جنس، سن و درمان تصادفی به‌طور هم‌زمان بر احتمال شفایافتگی و زمان لازم برای تولید یک تومور قابل تشخیص، تاثیر دارند.

داده‌های بقای این آزمایش بالینی قبلاً توسط چن و همکاران (۱۹۹۹) و کونور و همکاران (۲۰۰۷) مورد مطالعه قرار گرفته‌اند. با این وجود، در تحلیل آن‌ها فرض بر این بود که این متغیرهای کمکی فقط بر احتمال شفایافتگی تاثیر دارند. اما می‌توان گفت که متغیرهای کمکی قادر هستند تا زمان لازم برای تولید یک تومور قابل تشخیص را نیز تحت تاثیر قرار دهند.

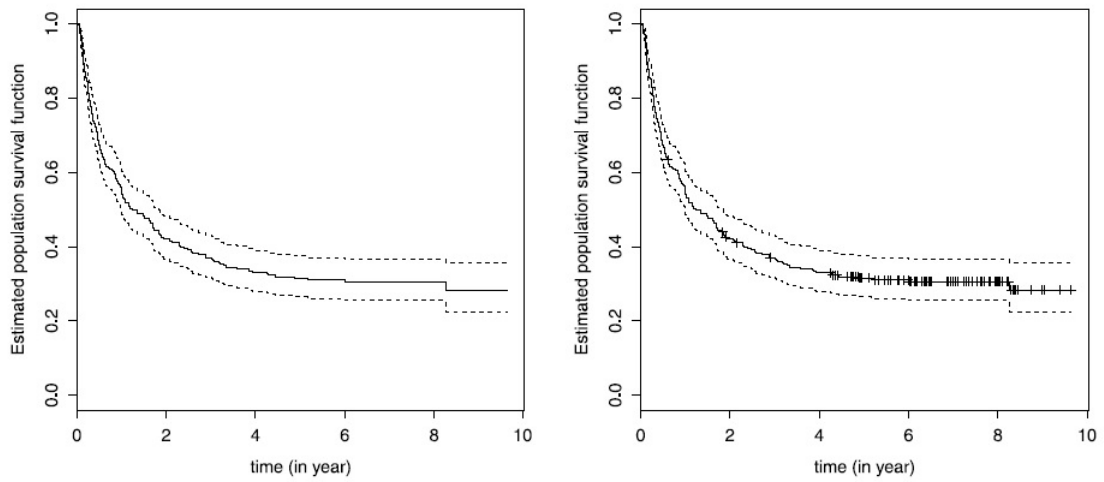
با در نظر گرفتن این داده‌ها و الگوریتم *MCMC* پیشنهادی برای مدل بیزی ذکر شده در فصل ۳، زنجیری به طول ۲۳۰۰۰ از توزیع پسین تولید شد. برای بررسی همگرایی الگوریتم *MCMC*، ردیابی تمام پارامترهای مدل مورد بررسی قرار گرفت. آماره  $Z$  جی‌وک برای تمام پارامترهای مدل بین ۰.۸۶ و ۱.۸۶ یافت شد که نشانگر همگرایی الگوریتم *MCMC* است. جدول ۴.۴ برآوردهای پارامترها را نشان می‌دهد.

می‌توان نتیجه گرفت که تنها درمان، تاثیر قابل توجهی بر احتمال شفایافتگی دارد. به عبارت دیگر، این مدل نشان می‌دهد که اینترفرون آلفا-۲ بی به‌طور قابل توجهی سلول‌های سرطان‌زا را کاهش می‌دهد اما بر زمان انکوباسیون سلول تاثیر نمی‌گذارد. این نتیجه‌گیری در شکل‌های ۹.۴ و ۱۰.۴ نیز نشان داده شده است. منحنی سمت راست شکل ۹.۴ تابع بقا جمعیت و منحنی سمت چپ تابع بقا افراد شفانیافته را نشان می‌دهد. تنها اختلاف بین گروه‌های درمان در تابع بقا جمعیت، مربوط به متوسط سن افراد است. شکل ۱۰.۴ نیز برآورد نسبت لگاریتم مخاطره جمعیت (منحنی چپ) و نسبت لگاریتم مخاطره افراد شفانیافته (منحنی راست) را برای میانگین سن مردان نشان می‌دهد. نتایج مشابهی را می‌توان برای زنان نیز به‌دست آورد.

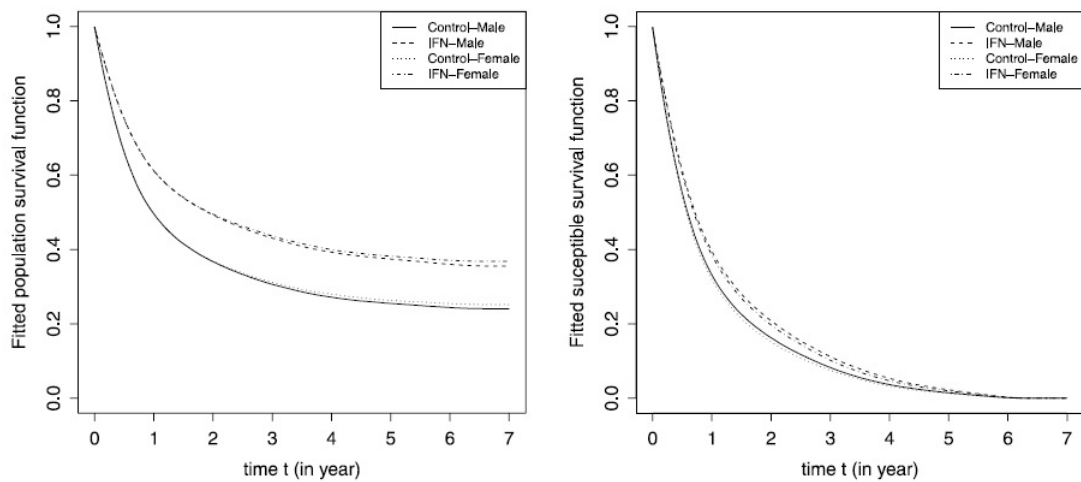
<sup>۵</sup>Eastern Cooperative Oncology Group

جدول ۴.۴: برآورد میانه‌های پسین، فاصله اعتبار ۹۵٪ و انحراف استاندارد پسین پارامترها در مثال واقعی.

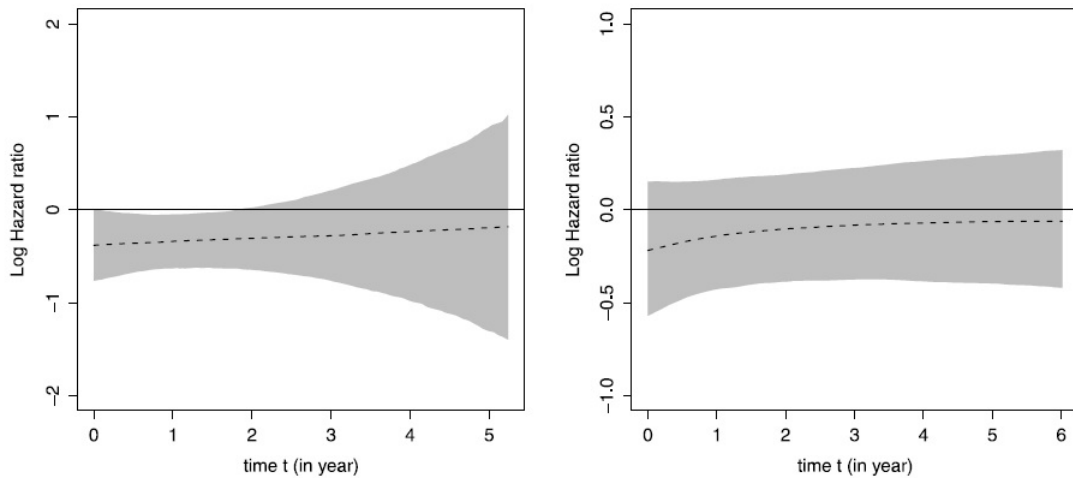
	پارامترها	برآورد	$HPD_{\%95}$	$sd_{post}$
$\theta(\mathbf{x})$	برش	۰/۳۵۱	[۰/۰۸۸, ۰/۵۶۴]	۰/۱۲۳
	سن	۰/۱۰۰	[-۰/۰۴۶, ۰/۲۴۸]	۰/۰۷۴
	$IFN$	-۰/۳۲۱	[-۰/۰۶۳۵, -۰/۰۳۲]	۰/۱۵۵
$F(t \mathbf{z})$	جنس	-۰/۰۳۱	[-۰/۲۳۷, ۰/۲۷۶]	۰/۱۶۱
	سن	-۰/۱۳۶	[-۰/۳۰۳, ۰/۰۳۳]	۰/۰۸۵
	$IFN$	-۰/۰۶۰	[-۰/۰۴۱۷, ۰/۳۴۱]	۰/۱۹۰
	جنس	۰/۰۵۳	[-۰/۳۳۳, ۰/۴۳۰]	۰/۱۹۵



شکل ۸.۴: برآورد کاپلان-مایر منحنی بقا در مثال واقعی. نمودار سمت راست: برای همه جامعه با علامت صلیب برای افراد سانسور شده و نمودار سمت چپ: برای افراد سانسور شده.



شکل ۹.۴: منحنی‌های برآورد شده بقا به تفکیک گروه‌های بیمار. چپ: همه جامعه و راست: بیماران شفایافته.



شکل ۳.۴: نمودارهای برآورد شده نسبت لگاریتم مخاطره با ناحیه اعتبار ۹۵٪ برای، چپ: همه مردان و راست: مردان شفانیافته.

### ۳.۴ نتیجه‌گیری و پیشنهادات

یک نسخه انعطاف‌پذیر از مدل نرخ شفایافتگی پیشرفت زمان هنگامی که متغیرها به‌طور هم‌زمان بر احتمال شفایافتگی و زمان لازم برای سلول به منظور تولید یک تومور قابل تشخیص تأثیر می‌گذارند، پیشنهاد شد. اگر چه مدل پیشنهادی ساختار مخاطرات متناسب در بیمار را ندارد، اما یک برآورد ساده از نسبت‌های خطر در طول زمان را فراهم می‌کند. هنگامی که پیگیری مطالعه به اندازه کافی طولانی نیست، باز هم می‌توان از مدل زمان پیشرفت با برخی محدودیت‌ها استفاده کرد. در این زمینه، ثابت شده است که اگر از متغیرهای تبیینی مشترکی به‌طور هم‌زمان برای مدل احتمال شفایافتگی و زمان لازم برای تشخیص رشد تومور سلول سرطانی استفاده نشود، اثرات متغیرها شناسایی‌پذیر هستند. استفاده از یک پیوند لجیت برای مدل‌سازی احتمال شفایافتگی مورد بررسی قرار گرفت تا سعی کند مشکل شناسایی‌پذیری را حل کند.

با توجه به محتوی پایان‌نامه می‌توان پیشنهادات زیر را برای بررسی و انجام ارائه داد:

- (۱) مدل‌بندی ناپارامتری تابع بقای پایه در دیدگاه بیزی با سایر روش‌ها مانند روش هسته یا چندجمله‌ای موضعی.
- (۲) در نظر گرفتن اثرات ناپارامتری برای متغیرهای رگرسیونی به‌جای اثر خطی آن‌ها در هر دو قسمت نسبت شفایافتگی و تابع بقای افراد غیرمصون.
- (۳) وارد کردن وابستگی داده‌ها به مدل پیشنهادی، مانند وابستگی طولی یا فضایی داده‌های بقا.



# پیوست آ

## برآوردگر ناپارامتری کاپلان-مایر

برای برآورد تابع مخاطره پایه می‌توان از روش کاپلان-مایر استفاده کرد. در این روش برآورد تابع مخاطره پایه در زمان  $t(i)$  از  $i$  امین پیشامد به صورت

$$\hat{h}_{\circ}(t(i)) = \frac{d(i)}{\sum_{j \in R(t(i))} \exp(\hat{\beta}^T x_j)}$$

تعریف می‌شود، که در آن  $d(i)$  تعداد مرگ‌ها در زمان  $t(i)$  است و  $R(t(i))$  مجموعه‌ای از افرادی است که می‌توانند در آن زمان بمیرند. برآورد تابع مخاطره پایه تخمینی مجزا برای یک تابع پیوسته است. بنابراین با توجه به

$$S(t) = \exp\left\{-\int_0^t h(\tau) d\tau\right\}$$

می‌توان از  $\hat{h}_{\circ}(t(i))$  برای برآورد  $\int_{t(i-1)}^{t(i)} h_{\circ}(\tau) d\tau$  استفاده کرد. در این صورت برآورد تابع بقای پایه به صورت زیر به دست می‌آید:

$$\hat{S}_{\circ}(t(i)) = \exp\left[-\sum_{j \leq i} \hat{h}_{\circ}(t(j))\right].$$



# پیوست ب

## معیار تشخیص همگرایی جیوک برای الگوریتم‌های MCMC

جیوک در سال ۱۹۹۲ یک آماره آزمون تشخیص همگرایی که برابری میانگین‌های بخش اول و آخر زنجیره‌های مارکوف را اندازه‌گیری می‌کند، ارائه داد. این دو بخش به‌طور معمول، ۱۰ درصد اول و ۵۰ درصد آخر زنجیر در نظر گرفته می‌شوند. آزمون تشخیص همگرایی برای این‌که چه تعداد از تکرارهای آغازی باید نادیده گرفته شوند، مورد استفاده قرار می‌گیرد. به عبارتی این آزمون به‌عنوان یک برآوردگر برای دوره داغیدن شناخته می‌شود.

اگر  $n$  تعداد تکرارها و  $n_0$  برابر با اولین تکراری باشد که بخواهیم با انجام این آزمون بررسی کنیم زنجیر همگرا می‌شود یا نه، باید آماره آزمون جیوک که یک آماره  $Z$  استاندارد است، محاسبه شود. فرض کنید  $\theta(X)$  کمیت مورد دلخواه باشد و  $\theta^t = \theta(X^{(t+n_0)})$ . اگر  $n_A$  تکرار اول، تعداد تکرارهای بخش  $A$  و  $n_B$  تکرار آخر، تعداد تکرارهای بخش  $B$ ، که در آن  $n^*$  مقداری است که  $n - n^* + 1 = n_B$  باشد آن‌گاه بخش‌های  $A$  و  $B$  و میانگین‌های آن‌ها به صورت زیر تعریف می‌شوند:

$$\begin{aligned}A &= \{t : 1 \leq t \leq n_A\}, \\B &= \{t : n^* \leq t \leq n\}, \\1 &< n_A < n^* < n, \quad \frac{n_A + n_B}{n} < 1, \\ \bar{\theta}_A &= \frac{1}{n_A} \sum_{t \in A} \theta^t, \\ \bar{\theta}_B &= \frac{1}{n - n^* + 1} \sum_{t \in B} \theta^t.\end{aligned}$$

آزمون جیوک بیان‌گر این مطلب است که اگر زنجیر در  $n_0$  همگرا شود، آن‌گاه نمونه‌های تولیدشده از توزیع مانای زنجیر گرفته شده‌اند و در نتیجه میانگین‌های محاسبه‌شده از اولین و آخرین بخش زنجیر، معادل هستند. بنابراین آماره جیوک که به‌طور مجانبی دارای توزیع نرمال استاندارد است، عبارت

<sup>۱</sup>Z-score



است از

$$Z_n = \frac{\bar{\theta}_A - \bar{\theta}_B}{\sqrt{\frac{\hat{S}_\theta^A}{n_A} + \frac{\hat{S}_\theta^B}{n - n^* + 1}}} \rightarrow N(0, 1) \quad n \rightarrow \infty,$$

که در آن  $\frac{\hat{S}_\theta^B}{n - n^* + 1}$  و  $\frac{\hat{S}_\theta^A}{n_A}$  واریانس‌های مجانبی  $\bar{\theta}_A$  و  $\bar{\theta}_B$  هستند. حال با توجه به مطالب ذکر شده می‌توان فرضیه صفر، فرضیه برابری میانگین‌ها، را روی فضای حالت  $\theta$  آزمون کرد. در این آزمون، اگر مقدار  $|Z_n|$  بزرگ باشد، آنگاه فرضیه صفر رد می‌شود و این یعنی زنجیر هنوز همگرا نشده است.

# پیوست پ

## کد تولید داده‌ها از مدل نرخ شفایافتگی زمان پیشرفت

```
library(mvtnorm)
library(survival)
library(mixdist)
library(matlib )
library(FAmle)

weibull.par=weibullpar(8,4.18)

K=12
n=K
W11=rnorm(n,0,1);W12=rbinom(n,1,0.5)
W21=rnorm(n,0,1);W22=rbinom(n,1,0.5)
W1=rnorm(n,0,1);W2=rbinom(n,1,0.5)
x=matrix(c(W1,W2),ncol=2)
z=matrix(c(W1,W2),ncol=2)

x1=matrix(c(W11,W12),ncol=2)
z1=matrix(c(W11,W12),ncol=2)
x2=matrix(c(W21,W22),ncol=2)
z2=matrix(c(W21,W22),ncol=2)

beta0=0.75 ; beta1=0.8 ; beta2=-0.5
Beta=c(beta1,beta2)
Beta_real=c(beta0,beta1,beta2)
```

```
gamma1=0.4 ; gamma2=-0.4
Gamma=c(gamma1,gamma2)
theta=exp(beta0+x%*%Beta)

T=c()
N=rpois(n,theta)
for(i in 1:n){
  if(N[i]==0){T[i]=999}
  else
  {y=rweibull(N[i],weibull.par$shape,weibull.par$scale)
  T[i]=min(y)}
}
T
r_cens=runif(n,20,25)
r_cens
nou=as.numeric(r_cens>=23)

D=matrix(c(T,nou,x[,1],x[,2],z[,1],z[,2]),ncol=6)
I <- matrix(0, K, K);diag(I) <- 1
epsilon=10^-6
P=(D%*%t(D))+(epsilon*I)
solve(P)

t_Rcens=25
B=3
n=300
S=500
```

## مراجع

- [1] Bremhorst, V., Lambert, P. (2016). Flexible estimation in cure survival models using P-splines, *Computational Statistics and Data Analysis*, **93**, 270-284.
- [2] Berkson, J., Gage, R.P. (1952). Survival curve for cancer patients following treatment, *Journal of American Statistical Association*, **47**, 501-515.
- [3] Cetinyurek, A., Lambert, P. (2011). Smooth estimation of survival functions and hazard ratios from interval-censored data using Bayesian penalized B-splines, *Statistics in Medicine*, **30**, 75-90.
- [4] Chen, M.H., Ibrahim, J.G., Sinha, D. (1999). A new Bayesian model for survival data with a surviving fraction, *Journal of American Statistical Association*, **94**, 909-919.
- [5] Cooner, F., Banerjee, S., Carlin, B.P., Sinha, D. (2007). Flexible cure rate modelling under latent activation schemes, *Journal of American Statistical Association*, **102**, 560-572.
- [6] Cox, D.R.(1972). *Regression Models and Life-Tables*, *Journal of the Royal Statistical Society: series B*, **34**, 187-220.
- [7] Cuzick, J. (1992). Semiparametric additive regression, *Journal of the Royal Statistical Society: series B*, **54**, 831-843.
- [8] Eilers, P.H.C. (2007). Ill-posed problems with counts, the composite link model and penalized likelihood, *Statistical Modeling*, **7**, 239-254 .
- [9] Eilers, P.H.C., Marx, B.D. (1996). Flexible smoothing with B-splines and penalties (with discussion), *Statistical Society*. **11**, 89-121.
- [10] Engle, R. F., Granger, C. W. J., Rice, J. and Weiss, A. (1986 ). Semiparametric estimates of the relation between weather and electricity sales, *Journals of American Statistical Association*, **81**, 310- 320.
- [11] Farewell, V.T. (1982). The use of mixture models for the analysis of survival data with long-term survivors, *Biometrics*, **38**, 1041-1046.

- [12] Farewell, V.T. (1986). Mixture models in survival analysis: are they worth the risk?, *Canadian Journal of Statistics*, **14**, 257–262.
- [13] Friedman, J. H. (1990). Multivariate adaptive regression splines. *The Annals of Statistics*, **19**, 11- 41.
- [14] Gelman, A., Roberts, G.O., Gilks, W.R. (1996). *Efficient Metropolis jumping rules*. In: Bernardo, J.M., Berger, J.O., Dawid, A.P., Smith, A.F.M. (Eds.), Bayesian Statistics 5. Oxford University Press, UK.
- [15] Geweke, J. (1992). Evaluating the accuracy of sampling-based approaches to calculating posterior moments, In: Bernardo, J.M., Berger, J.O., Dawid, A.P., Smith, A.F.M. (Eds.), Bayesian Statistics 4. Clarendon Press, Oxford, UK.
- [16] Haario, H., Saksman, E., Tamminen, J. (2001). An adaptive Metropolis algorithm. *Bernoulli*, **7**, 223–242.
- [17] Hekman, N. E. (1986). Spline smoothing in partly linear models. *Journal of the Royal Statistical Society: series B*, **48**, 244-248 .
- [18] Ibrahim, J.G., Chen, M.-H., Sinha, D. (2001). Bayesian semiparametric models for survival data with a cure fraction, *Biometrics*, **57**, 383–388.
- [19] Jullion, A., Lambert, P. (2007). Robust specification of the roughness penalty prior distribution in spatially adaptive Bayesian P-splines models, *Computational Statistics and Data Analysis*, **51**, 2542–2558.
- [20] Kirkwood, J.M., Strawderman, M.H., Ernstoff, M.S., Smith, T.J., Borden, E.C., Blum, R.H. (1996). Interferon alpha-2b adjuvant therapy of high-risk resected cutaneous melanoma: the eastern cooperative oncology group trial e1684, *Journal of Clinical Oncology*, **14**, 7–17.
- [21] Kuk, A.Y.C, Chen, M.H. (1992). A mixture model combining logistic regression with proportional hazards regression, *Biometrika*, **79**, 531–541.
- [22] Lambert, P. (2007). Archimedean copula estimation using Bayesian splines smoothing techniques, *Computational Statistics Data and Analysis*, **51**, 6307–6320.
- [23] Lambert, P. (2013). Nonparametric additive location-scale models for interval censored data, *Statistical Computation*, **23**, 75–90.
- [24] Lambert, P., Eilers, P. (2005). Bayesian proportional hazards model with time varying regression coefficients: a penalized Poisson regression approach, *Statistics in Medicine*, **24**, 3977–3989.

- [25] Lang, S., Brezger, A. (2004). Bayesian P-splines, *Journal of Computational and Graphical Statistics*, **13**, 183–212.
- [26] Li, C.S., Taylor, J.M.G. (2002). A semi-parametric accelerated failure time cure model, *Statistics in Medicine*, **21**, 3235–3247.
- [27] Liu, H., Shen, Y. (2009). A semiparametric parametric regression cure model for interval-censored data, *Journal of American Statistical Association*, **104**, 1168–1178.
- [28] Lopes, C.C.M., Bolfarine, H. (2012). Random effects in the promotion times cure rate models, *Computational Statistics and Data Analysis*, **56**, 75–87.
- [29] Lu, W. (2010). Efficient estimation for an accelerated failure time model with a cure fraction, *Statistica Sinica*, **20**, 661–674.
- [30] Ma, Y., Yin, G. (2008). Cure rate model with mismeasured covariates under transformation, *Journal of American Statistical Association*, **103**, 743–756.
- [31] Peng, Y. (2003). Estimating baseline distribution in proportional hazards cure models, *Computational Statistics and Data Analysis*, **42**, 187–201.
- [32] Peng, Y., Dear, K.B.G. (2000). A nonparametric mixture model for cure rate estimation *Biometrics*, **56**, 237–243.
- [33] Roberts, G.O., Rosenthal, J.S. (2001). Optimal scaling for various Metropolis Hastings algorithms, *Statistical Society*, **16**, 351–367.
- [34] Schoenberg I.J. (1946). *Contribution to the problem of approximation of equidistant data by analytic functions, Parts A and B*, Quarterly Applied Mathematics 4 , 45–99, 112–141.
- [35] Speckman, P. (1988). Kernel smoothing in partial linear models. *Journal of the Royal Statistical Society, Series B*, **50**, 413–436.
- [36] Sy, J.P., Taylor, J.M.G. (2000). Estimation in a Cox proportional hazards cure model, *Biometrics*, **56**, 227–236.
- [37] Taylor, J.M.G. (1995). Semi-parametric estimation in failure time mixture models, *Biometrics*, **51**, 899–907.
- [38] Tsodikov, A.D. (1998). A proportional hazard model taking account of long-term survivors, *Biometrics*, **54**, 1508–1516.
- [39] Tsodikov, A.D. (2002). Semi-parametric model of long- and short-term survival: an application to the analysis of breast cancer survival in Utah by age and stage, *Statistics in Medicine*, **21**, 895–920.

- 
- [40] Wang, L., Du, P., Liang, H. (2012). Two-component mixture cure rate model with spline estimated nonparametric components, *Biometrics*, **68**, 726–735.
- [41] Yakovlev, A.Y., Tsodikov, A.D. (1996). *Stochastic Models for Tumor of Latency and their Biostatistical Applications*, World Scientific Publication, Singapore.
- [42] Yin, G., Ibrahim, J.G. (2005). Cure rate models: a unified approach, *Canadian Journal of Statistics*, **33**, 559–570.
- [43] Zeng, D., Yin, G., Ibrahim, J.G. (2006). Semiparametric transformation models for survival data with a cure fraction, *Journal of American Statistics Association*, **101**, 670–684.
- [44] Zhang, J., Peng, Y., Li, H. (2013). A new semiparametric estimation method for accelerated hazards mixture cure model, *Computational Statistics and Data Analysis*, **59**, 95–102.

# واژه‌نامه فارسی به انگلیسی

Cubic Spline	اسپلاین مکعبی
Survival Function	تابع بقا
Failure Function	تابع شکست
Hazard Function	تابع مخاطره
Survival Analysis	تحلیل بقا
Identification	شناسایی پذیری
Cure Rate Models	مدل‌های نرخ شفایافتگی
Non-Mixture Cure Models	مدل‌های نرخ شفایافتگی غیرآمیخته
Cox Proportional Hazard Model	مدل مخاطرات متناسب کاکس
Covariates	متغیرهای کمکی
Semi-Parametric	نیمه‌پارامتری





# واژه‌نامه انگلیسی به فارسی

Covariates	متغیرهای کمکی
Cox Proportional Hazard Model	مدل مخاطرات متناسب کاکس
Cubic Spline	اسپلاین مکعبی
Cure Rate Models	مدل‌های نرخ شفایافتگی
Failure Function	تابع شکست
Hazard Function	تابع مخاطره
Identification	شناسایی‌پذیری
Non-Mixture Cure Models	مدل‌های نرخ شفایافتگی غیرآمیخته
Semi-Parametric	نیمه پارامتری
Survival Analysis	تحلیل بقا
Survival Function	تابع بقا

## **Abstract**

In the analysis of survival data, it is usually assumed that any unit will experience the event of interest if it is observed for a sufficiently long time. However, it can be assumed that an unknown proportion of the population under study will never experience the interest event. The Promotion time model, which has a biological motivation, is one of the survival models taking this feature in to account. The survival function of uncured people is estimated by using a Cox proportional hazard model where the logarithm of the baseline hazard function is specified using Bayesian P-splines. We use a Bayesian framework for implementing Statistical inference in the proposed model. The identification issues of the model are discussed and a restricted use of the model when the follow-up of the study is not sufficiently long is suggested. The accuracy of our methodology is assumed through a simulation study and the model is illustrated on a dataset from a Melanoma clinical trial.

Keywords: Survival analysis, Promotion time model, P-spline, Cox proportional hazard model, Bayesian inference.



**Shahrood University of Technology**

**Faculty of Mathematical Sciences**

**Bayesian analysis of non-mixture cure rate  
models with a Cox semi-parametric survival  
function by using P-splines**

**Neda Shams**

**Supervisors**

**Dr. Hossein Baghishani and Dr. Negar Eghbal**

**Jul 2018**