

رگرسیون بقا به شیوه بیز ناپارامتری برای داده‌های سانسور شده با استفاده از فرایند دیرخله وابسته

زهرة شاه‌حسینی ، سید مرتضی امینی
بخش آمار دانشکده ریاضی آمار و علوم کامپیوتر، دانشگاه تهران

چکیده

در این مقاله تحلیل بیز ناپارامتری براساس پیشین فرآیند دیرخله وابسته را برای تحلیل بقای داده‌های سانسور شده با در نظر گرفتن عوامل تأثیرگذار در طول عمر معرفی کرده و مورد بررسی و استفاده قرار می‌دهیم. همچنین الگوریتم‌های محاسباتی و نمونه‌گیری زنجیره مارکف را معرفی و تشریح کرده و برای استفاده در این مسأله توسعه می‌دهیم. در نهایت به تحلیل بقای داده‌های واقعی مربوط به طول عمر بیماران سرطانی با استفاده از مدل‌های کاکس، کیپلان مایر و نیز روش تحلیل واریانس دیرخله سلسله مراتبی وابسته می‌پردازیم. نشان داده می‌شود که روش‌های کلاسیک تحلیل بقا در بعضی زمان‌ها مقدار ثابتی را اختیار می‌کنند. با این حال حضور داده‌های سانسور شده در برآورد تابع بقا تحت مدل تحلیل واریانس دیرخله سلسله مراتبی وابسته چنین مشکلی را ایجاد نمی‌کند.

واژه‌های کلیدی: تحلیل بقا، داده‌های سانسور شده، زنجیر مارکف مونت کارلو، فرآیند دیرخله آمیخته، فرآیند دیرخله وابسته، مدل سلسله مراتبی، نمونه‌بردار گیبس.

۱ مقدمه

آماري مختلف، برای بالا بردن قابلیت اطمینان محصولات خود نیاز دارند.

از سویی دیگر، با افزایش آمار مرگ و میر، شیوع بیماری‌های مختلف و پیشرفت‌های پزشکی این سوالات مطرح می‌شود که

۱. طول عمر یک انسان سالم چقدر است؟
۲. چه عواملی بر طول عمر یک انسان اثرگذار هستند؟
۳. در صورتی که فردی دچار بیماری خاصی شود و یا تحت

پیشرفت فناوری و بالا رفتن سطح زندگی مردم با استفاده از دستاوردهای صنعتی، رقابت‌های شدیدی بین تولیدکنندگان برای ارائه محصولات و خدمات مرغوب‌تر پدید آورده است. وجود محصولات مشابه و منابع مختلف اطلاعاتی، به مشتریان این امکان را می‌دهد که خصوصیات محصولات مورد نظر را از نظر قیمت، قابلیت اطمینان و جنبه‌های دیگر مقایسه کرده و بهترین گزینه را انتخاب کنند. با توجه به این مسأله، شرکت‌های معتبر تولیدی در این رقابت به روش‌های

شرایطی درمان شود طول عمر آن چه خواهد بود؟

ناپارامتری فراوانی گرای موجود انعطاف لازم برای مدل‌بندی اثرات متغیرهای کمکی و عوامل تأثیرگذار را ندارند. بنابراین شیوه‌های بیز ناپارامتری برای این منظور به کار گرفته می‌شود. شیوه‌های بیز ناپارامتری وابسته با در نظر گرفتن مدل‌های انعطاف پذیر رگرسیونی در ابر پارامترهای توزیع‌های پیشین خود می‌توانند جایگزین مناسبی برای شیوه‌های پارامتری و نیمه پارامتری وابسته مانند مدل کاکس باشند. در این مقاله به طور خاص به فرآیند دیریکله وابسته و کاربرد آن برای تحلیل رگرسیون بقا برای داده‌های سانسور شده می‌پردازیم. فرآیند دیریکله یکی از متداول‌ترین شیوه‌های بیز ناپارامتری است که دارای خصوصیات بسیار خوب نظری و محاسباتی است.

ساختار کلی این مقاله بدین صورت است که ابتدا مفهوم‌ها و تعریف‌هایی از تحلیل بقا، نمادهای مورد نیاز آن و نیز چند روش خاص تحلیل بقا معرفی خواهد شد. پس از آن به بیان فرآیند دیریکله و نیز رهیافت بیز ناپارامتری پرداخته خواهد شد. همچنین در این مقاله تعمیم‌های این رهیافت از جمله فرآیند دیریکله وابسته و فرآیند دیریکله سلسله مراتبی و نیز نحوه محاسبه توزیع پسین در این مدل با استفاده از روش الگوریتم گیبس بیان می‌شوند. در ادامه روش‌های بیز ناپارامتری بر اساس پیشین فرآیند دیریکله سلسله مراتبی وابسته، در حالتی که داده‌ها سانسور شده از راست باشند و نیز الگوریتم‌های لازم بیان خواهد شد. در پایان مطالب بالا برای بررسی داده‌های مربوط به بیماران سرطانی مورد استفاده قرار می‌گیرد.

۲ تحلیل بقا

تحلیل بقا مجموعه‌ای از روش‌های آماری تحلیل داده‌های طول عمر است. طول عمر یک متغیر تصادفی مثبت است

امروزه تحلیل بقا در اکثر مطالعات علمی که شامل بررسی مدت زمان وقوع یک پیشامد باشد، مورد استفاده قرار می‌گیرد. در کاربرد توزیع داده‌های طول عمر نامعلوم است. بدین سبب به کار بردن روش‌های ناپارامتری در مقایسه با روش‌های پارامتری معمول‌تر است. در مطالعه‌های بقا اگر هدف، توصیف زمان بقا بدون در نظر گرفتن متغیر کمکی باشد، از روش‌های ناپارامتری غیر وابسته همچون جداول عمر و روش کاپلان-مایر استفاده می‌شود. با این حال تعیین عوامل و متغیرهای اثرگذار در زمان رخداد پدیده‌ها از اهمیت زیادی برخوردار است. برای مدل‌سازی چنین عواملی در تحلیل داده‌های بقا و پیش‌بینی از مدل‌های رگرسیون بقا مانند مدل کاکس به عنوان یک روش نیمه ناپارامتری استفاده می‌شود که نخستین بار توسط کالب فلیش [۱۳] با در نظر گرفتن پیشین فرآیند گاما برای تابع توزیع تجمعی هازارد در مدل کاکس به کار گرفته شد. همچنین کریستنسن و جانسون [۳] تحلیل بیز برای مدل‌های زمان شکست شتابیده با در نظر گرفتن پیشین فرآیند دیریکله را مطرح کردند. والکر و مالیک [۱۶] با به کارگیری روش‌های MCMC و با بهره‌گیری از پیشین فرآیند درخت پولیا این شیوه‌های تحلیل بقا را توسعه دادند. همچنین در سال‌های اخیر پیشین فرآیند درخت پولیای آمیخته برای تابع بقای پایه در مدل‌های کاکس، زمان شکست شتابیده و بخت‌های متناسب استفاده شده است. هانسون و جانسون [۱۱] و هانسون، جانسون و لاد [۱۲]، کیو و مالیک [۱۴] و ژلفاند و کوتاس [۱۰] از پیشین فرآیند دیریکله آمیخته برای تابع بقای پایه در مدل‌های زمان شکست شتابیده استفاده کردند. اما این مدل‌ها نیز دارای فرضیات محدود کننده‌ای هستند.

برای حذف چنین محدودیت‌هایی استفاده از شیوه‌های ناپارامتری راه حل مناسبی است. با این حال شیوه‌های

که می‌تواند زمان شکست یک مؤلفه فیزیکی و یا زمان مرگ یک واحد زنده باشد. فرض کنید یک سیستم، یا قطعه الکتریکی، یا هر وسیله یا شیء دیگر، دارای طول عمر T است و مقادیرهای خود را در بازه $(0, \infty)$ اختیار می‌کند. اگر $f(t)$ و $F(t)$ به ترتیب نمایانگر تابع چگالی احتمال و تابع توزیع تجمعی احتمال این متغیر تصادفی باشند، آنگاه $S(t) = 1 - F(t) = P(T > t)$ را تابع بقای سیستم تا زمان t می‌گوییم. معمولاً در کاربردهای غیر پزشکی به تابع بقا، تابع قابلیت اعتماد گفته می‌شود و آن را با $R(t)$ نمایش می‌دهند.

۳ استنباط بیز ناپارامتری بر اساس فرآیند دیرخله

در شیوه استنباط بیزی پارامتری، برای نمونه‌های متناهی یک توزیع معلوم با پارامترهای نامعلوم در نظر گرفته می‌شود و با در نظر گرفتن پارامترهای مدل به عنوان متغیرهای تصادفی و با تعیین توزیع پیشین برای آن‌ها، به استنباط بیزی در خصوص پارامترها می‌پردازیم. در عمل با موارد بسیاری روبرو هستیم که تعیین توزیع معلوم برای داده‌ها دور از واقعیت است. در این شرایط می‌توان از استنباط بیزی ناپارامتری استفاده کرد.

در شیوه‌های بیز ناپارامتری تابع احتمال را به عنوان یک فرآیند تصادفی روی σ -جبر فضای نمونه در نظر می‌گیریم. به عبارت دیگر تابع احتمال P که تمامی تابع‌های مورد استفاده در تحلیل بقا را می‌توان بر اساس آن تعریف کرد را روی σ -جبر A که P روی آن تعریف شده است به عنوان یک فرآیند تصادفی در نظر می‌گیریم به طوری که به ازای هر پیشامد $A \in \mathcal{A}$ ، $P(A)$ یک متغیر تصادفی در نظر گرفته می‌شود. فرآیندهای پیشین متفاوتی برای تحلیل

یکی دیگر از معیارهای مهم در مطالعه قابلیت اطمینان و طول عمر تابع نرخ خطر است. تابع نرخ خطر احتمال شکست یک موجود در یک بازه زمانی کوچک (δ) به شرط آن که تا زمان t از کار نیفتاده باشد را بیان می‌کند و آن را با $h(t)$ نمایش می‌دهند.

روش‌های محاسبه تابع بقا بر اساس دیدگاه‌های پارامتری و ناپارامتری و نیز فراوانی‌گرا یا بیزی دسته‌بندی می‌شوند. به عنوان مثال روش درست‌نمایی ماکزیم دیدگاه پارامتری فراوانی‌گرا است که در آن، برآورد تابع چگالی براساس برآورد پارامترها انجام می‌شود. از دیگر نمونه‌های این دیدگاه می‌توان به روش نیمه پارامتری کاکس اشاره کرد که در آن، متغیرهای کمکی نیز در مدل‌بندی نقش دارند و بر پایه فرض متناسب بودن تابع‌های نرخ خطر مشاهده‌ها و استقلال زمان‌های رخداد پیشامدها تشکیل شده است. دیدگاه بیز پارامتری شیوه‌ای استنباطی پارامتری براساس اطلاعات پیشین در خصوص پارامترهای مجهول است. در این شیوه پارامترهای مجهول توزیع پارامتری معلوم به عنوان متغیرهای تصادفی در نظر گرفته شده و با توجه به آگاهی‌های پیشین برای آن‌ها توزیع‌های پیشین در نظر گرفته می‌شود. در دیدگاه فراوانی‌گرای ناپارامتری همچون روش کاپلان

^۱Kaplan-meier

بیزی ناپارامتری معرفی شده‌اند، که معروف‌ترین آن‌ها فرآیند دیریکله^۲ است، و توسط فرگوسن [۱۹] معرفی شده است.

$$G(x) \sim \text{Beta}(\alpha((-\infty, x]), \alpha((x, \infty))) , \quad (2)$$

امروزه روش‌های بیز ناپارامتری بسیار متنوع و هر یک دارای خصوصیت‌های فراوان و پرکاربرد هستند. برای تعریف فرآیند دیریکله ابتدا لازم است با توزیع دیریکله^۳ آشنا شویم.

$$E(G(x)) = \frac{\alpha(-\infty, x]}{\alpha(\mathbb{R})}$$

$$\text{Var}(G(x)) = \frac{\alpha(-\infty, x]\alpha(x, \infty)}{\alpha(\mathbb{R}) + 1}$$

گاهی اوقات اندازه $\alpha(A)$ را به صورت $\alpha(A) = \alpha_0 Q(A)$ تعریف می‌کنند، که در آن Q یک تابع احتمال روی Θ و $\alpha_0 > 0$ یک پارامتر تمرکز است. در این صورت اگر G_0 تابع توزیع تجمعی متناظر با Q باشد داریم:

$$E(G(x)) = G_0(x), \quad \text{Var}(G(x)) = \frac{G_0(x)(1 - G_0(x))}{\alpha_0 + 1}, \quad (3)$$

توزیع G_0 را می‌توان به عنوان حدسی از G و نیز پارامتر دقت α_0 را به عنوان درجه تمرکز توزیع G و یا شدت قطعیت G حول G_0 در نظر گرفت. برای مقادیر بزرگ α_0 ، یک نمونه از G با احتمال زیاد مقادیر نزدیک به G_0 را می‌گیرد و برای مقادیر کوچک α_0 ، نمونه G با احتمال زیاد جرم احتمال زیادی را تنها در تعداد معین اتم‌ها می‌گذارد. همچنین نشان داده شده است که فرآیند دیریکله یک پیشین مزدوج است. بنابراین توزیع پسین نیز فرآیند دیریکله است. برای مطالب بیشتر در این خصوص می‌توانید به اسکوبار و وست [۴] و وست، مولر و اسکوبار [۱۷] مراجعه کنید.

در صورتیکه در تعریف فوق G نامعلوم باشد، می‌توان براساس دیدگاه بیز ناپارامتری برای G توزیع پیشینی در نظر گرفت. اگر این توزیع پیشین فرآیند دیریکله باشد،

گوئیم (p_1, p_2, \dots, p_k) به ازای $\sum_{i=1}^k p_i = 1$ ، $p_i \geq 0$ دارای توزیع دیریکله با پارامترهای $(\alpha_1, \dots, \alpha_k)$ است و با نماد

$$(p_1, p_2, \dots, p_k) \sim \text{Dir}(\alpha_0, \alpha_1, \dots, \alpha_k)$$

نشان می‌دهیم، اگر $p = (p_1, p_2, \dots, p_k)$ دارای تابع چگالی احتمال زیر باشد

$$f(p_1, p_2, \dots, p_k) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k p_i^{\alpha_i - 1}, \quad (1)$$

که در آن $\Gamma(\cdot)$ تابع گاما را نشان می‌دهد.

حال به تعریف فرآیند دیریکله می‌پردازیم.

تعریف: فرض کنید Θ فضای نمونه و $\alpha(\cdot)$ یک اندازه جمع پذیر روی Θ باشد. می‌گوئیم فرآیند تصادفی $P(\cdot)$ روی فضای نمونه Θ یک فرآیند دیریکله با پارامتر $\alpha(\cdot)$ است، اگر به ازای هر افزایش متناهی اندازه پذیر A_1, \dots, A_r ، $r \geq 1$ ، از فضای نمونه Θ ، بردار تصادفی $(P(A_1), \dots, P(A_r))$ دارای توزیع دیریکله با پارامترهای $(\alpha(A_1), \dots, \alpha(A_r))$ است و آن را با نماد $P \sim DP(\alpha)$ نشان می‌دهیم.

بنابراین اگر افزایش $(-\infty, x] \cup (x, \infty)$ از فضای نمونه $\Theta = \mathbb{R}$ را در نظر بگیریم و تابع توزیع تجمعی متناظر با $P(\cdot)$ را

^۲Dirichlet Process

^۳Dirichlet distribution

مدل حاصل، مدل فرآیند دیریکله آمیخته نامیده می‌شود. می‌آید:

$$G|Y_1, Y_2, \dots, Y_n \sim \int DP(\alpha + \sum_{i=1}^n \delta_{\theta_i}) \quad (7)$$

$$dP(\theta_1, \dots, \theta_n | Y_1, Y_2, \dots, Y_n),$$

می‌شود:

$$(4)$$

$$Y_i | \theta_i \sim H(\cdot | \theta_i), \quad \theta_i | G \sim G(\cdot), \quad G | \alpha \sim DP(\cdot | \alpha).$$

که در آن δ تابع نشانگر و $dP(\theta_1, \dots, \theta_n | Y_1, Y_2, \dots, Y_n)$ توزیع پسین $\theta_1, \theta_2, \dots, \theta_n$ به شرط Y_1, \dots, Y_n است. بنابراین یک برآورد بیز ناپارامتری برای $G(x)$ براساس تابع زیان درجه دوم خطا برابر است با

به عبارت دیگر، اگر Y_1, Y_2, \dots, Y_n هم توزیع با توزیع نامعلوم F باشند، آنگاه

$$E((G(x) | Y_1, \dots, Y_n)) = \quad (8)$$

$$\int \frac{\alpha((-\infty, x]) + \sum_{i=1}^n \delta_{\theta_i}((-\infty, x])}{\alpha(\mathbb{R}) + n}$$

$$dP(\theta_1, \dots, \theta_n | Y_1, Y_2, \dots, Y_n).$$

$$F(y_i) = \int H(y_i | \theta) dG(\theta), \quad G \sim DP(\alpha)$$

که در آن $H(\cdot | \theta)$ یک توزیع پارامتری است، که هسته آمیخته نامیده می‌شود و توسط یک پارامتر با بعد محدود θ اندیس‌گذاری شده است.

با استفاده از نمونه به دست آمده، نمونه بردار گیبس و رابطه (7) می‌توان یک برآورد بیز ناپارامتری از $G(x)$ با تقریب رابطه (8) به صورت زیر به دست آورد:

برای یافتن یک برآورد بیز ناپارامتری برای $G(x)$ ، تعریف آمیخته فرآیند دیریکله ضروری است. آمیخته فرآیند دیریکله تعمیم فرآیند دیریکله است، به صورتی که توزیع نهایی پیشین آمیخته‌ای از فرآیند دیریکله با اندازه‌های پایه گوناگون است. با در نظر گرفتن توزیع آمیزنده^۴ H می‌نویسیم

$$\frac{1}{M} \sum_{j=1}^M \int \frac{\alpha((-\infty, x]) + \sum_{i=1}^n \delta_{\theta_{ij}}((-\infty, x])}{\alpha(\mathbb{R}) + n},$$

$$G | \eta \sim DP(\alpha_\eta), \quad \eta \sim H(\eta) \quad (5)$$

و در نتیجه

که در آن $\theta_j = (\theta_{1j}, \dots, \theta_{nj}), j = 1, \dots, M$ نمونه‌های تولید شده در نمونه‌برداری گیبس پس از طی دوره داغیدن هستند. جهت اطلاعات بیشتر در رابطه با روش‌های شبیه‌سازی مونت کارلو مانند نمونه بردار گیبس به داس و هافر [7] مراجعه شود.

$$G \sim \int DP(\alpha_\eta) dH(\eta). \quad (6)$$

برای مدل (4)، توزیع پسین G به صورت زیر به دست

از طرفی بر اساس مدل بلکول و مک کوپین [2] می‌توان

^۴Mixing distribution

نوشت

$\theta_1, \dots, \theta_n$ از توزیع شرطی هر θ_i به شرط

$$\theta_{-i} = \{\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_n\} \quad G(\theta_{n+1} | \theta_1, \theta_2, \dots, \theta_n) \propto \alpha((-\infty, \theta_{n+1}]) \quad (9)$$

$$+ \sum_{j=1}^n \delta_{\theta_j}((-\infty, \theta_{n+1}])$$

$$= \alpha((-\infty, \theta_{n+1}]) + \sum_{j=1}^n n_j \delta_{\theta_j^*}((-\infty, \theta_{n+1}]),$$

به صورت زیر استفاده می‌شود.

$$G(\theta_i | \theta_{-i}) \propto \alpha_0 G_0(\theta_i) + \sum_{j=1, j \neq i}^n \delta_{\theta_j}((-\infty, \theta_i]). \quad (13)$$

بنابراین می‌توان توزیع پسین شرطی را به صورت زیر نوشت

$$P(\theta_i | \theta_{-i}, Y) \propto H(Y | \theta) \cdot G(\theta_i | \theta_{-i})$$

$$= \alpha_0 H(Y_i | \theta_i) G_0(\theta_i)$$

$$+ \sum_{j=1, j \neq i}^n H(Y_i | \theta_j) \delta_{\theta_j}((-\infty, \theta_i]) \quad (14)$$

که در آن θ_j^* ها مقادیرهای متمایز $\theta_1, \dots, \theta_n$ و n_j تعداد تکرارها θ_j^* در میان $\theta_1, \dots, \theta_n$ است. فرض کنید $\theta = (Y_1, Y_2, \dots, Y_n)$ و $(\theta_1, \theta_2, \dots, \theta_n)$ توزیع پسین $P(\theta | Y)$ در (۷) به صورت زیر به دست می‌آید:

$$P(\theta | Y) \propto H(Y | \theta) \cdot G(\theta), \quad (10)$$

که در آن $G(\theta) = \prod_{i=1}^n G(\theta_i | \theta_1, \theta_2, \dots, \theta_{i-1})$ براساس رابطه (۱۰) به دست می‌آید. بنابراین داریم

$$P(\theta | Y) \propto \prod_{i=1}^n H(Y_i | \theta_i) \cdot \alpha((-\infty, \theta_i]) + \sum_{j=1}^{i-1} \delta_{\theta_j}((-\infty, \theta_i]). \quad (11)$$

از طرفی

$$P(\theta_i | Y_i)$$

$$= \frac{H(Y_i | \theta_i) G_0(\theta_i)}{\int H(Y_i | \theta_i) dG_0(\theta_i)}, \quad (15)$$

برای $\alpha(A) = \alpha_0 Q(A)$ و با در نظر گرفتن G_0 به عنوان تابع توزیع تجمعی متناظر با Q و همچنین مقادیرهای غیر تکراری θ_j^* با فراوانی n_j ، $j = 1, \dots, k_i$ در دنباله $\theta_1, \dots, \theta_n$ می‌توان نوشت

$$\alpha_0 H(Y_i | \theta_i) G_0(\theta_i) = \alpha_0 \int H(Y_i | \theta_i) dG_0(\theta_i) P(\theta_i | Y_i),$$

از طرف دیگر از آن جا که از رابطه (۱۵) داریم

$$P(\theta_i | Y_i) \propto H(Y_i | \theta_i) G_0(\theta_i).$$

می‌توان نوشت

$$P(\theta | Y) = \prod_{i=1}^n H(Y_i | \theta_i) \cdot \mathcal{P}_{\theta_i}, \quad (12)$$

که در آن

$$\mathcal{P}_{\theta_i} = \frac{\alpha_0 G_0(\theta_i) + \sum_{j=1}^{k_i} n_{ij} \delta_{\theta_j^*}((-\infty, (\theta_i]))}{\alpha_0 + i - 1}.$$

$$\alpha_0 H(Y_i | \theta_i) G_0(\theta_i) \propto q_0 t H(Y_i, \theta_i) G_0(\theta_i),$$

در نمونه‌گیری گیبس از توزیع $G(\theta)$ ، برای تولید نمونه

می‌شود. در این بخش ابتدا به معرفی این مدل پرداخته می‌شود و پس از آن روش محاسبه توزیع پسین در حالتی که مدل فرایند دیریکله سلسله مراتبی باشد مختصراً بیان خواهد شد. به منظور معرفی فرایند دیریکله سلسله مراتبی داریم:

که در آن $q_0 \propto \int H(Y_i, \theta_i) dG_0(\theta_i)$ است. بنابراین

$$P(\theta_i | \theta_{-i}, Y) = q_0 H(Y_i | \theta_i) G_0(\theta_i) + \sum_{j=1, j \neq i}^n q_j \delta_{\theta_j}((-\infty, \theta_i]) \quad (17)$$

فرض کنید Y_1, Y_2, \dots, Y_n نمونه‌ای تصادفی باشند به طوری که برای $i = 1, \dots, n$ $Y_i \sim H(\cdot | \theta_i, \zeta_i)$ و Y_i ها به شرط θ_i و ζ_i مستقل هستند. هم چنین قرار دهید

که در آن $q_i \propto H(Y_i | \theta_j)$ و $q_0 + \sum_{j \neq i} q_j = 1$ است.

بنابراین نمونه بردار گیس در سه گام نمونه‌ای از θ را از طریق رابطه (۱۷) به صورت زیر تولید می‌کند:

$$Y = (Y_1, Y_2, \dots, Y_n), \quad \theta = (\theta_1, \theta_2, \dots, \theta_n),$$

گام ۰: یک مقدار اولیه برای $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ انتخاب می‌کنیم.

$$\zeta = (\zeta_1, \zeta_2, \dots, \zeta_n).$$

گام ۱: مقدار θ_1 را از توزیع $P(\theta_1 | \theta_{-1}, Y)$ تولید می‌کنیم، سپس از جایگذاری مقدار θ_1 در θ مقدار θ_2 را از توزیع $P(\theta_2 | \theta_{-2}, Y)$ تولید می‌کنیم و این کار را تا تولید θ_n ادامه می‌دهیم.

توزیع‌های پیشین زیر را در نظر بگیرید.

$$\theta_i \overset{i.i.d}{\sim} G(\cdot | \lambda), \quad \zeta_i \overset{i.i.d}{\sim} H_1(\cdot), \quad \lambda \sim H_2(\cdot).$$

گام ۲: گام ۱ را تا رسیدن به همگرایی الگوریتم تکرار می‌کنیم. بنابراین تعداد نمونه‌هایی که تا این جا تولید کرده‌ایم به عنوان دوره سوخت کنار می‌گذاریم و تعداد نمونه‌های پس از همگرایی را به عنوان نمونه تولید شده از پسین در نظر می‌گیریم.

علاوه بر آن فرض کنید Y و ζ به شرط θ از λ مستقل هستند.

همچنین فرض کنید که θ و G به شرط Y, ζ و λ مستقل از هر پارامتر دیگر در مدل سلسله مراتبی باشند. عبارت سلسله مراتبی در این جا، در گسترده ترین معنی آن، به معنی یک مدل دارای سطوح مختلف است که به وسیله دنباله‌ای از توزیع‌های شرطی آشیانه‌ای مدل بندی شده باشد. یک مدل پارامتری بیز سلسله مراتبی که می‌خواهیم به توصیف آن بپردازیم، مدلی است که عدم قطعیت در خصوص شکل تابعی توزیع پیشین G را توسط یک مدل پیشین فرآیند دیریکله مدل بندی می‌کند. بنابراین می‌توان چنین ادعا کرد که محدودیت معلوم بودن شکل توزیع پیشین که در استنباط بیز پارامتری مفروض است، در این مدل وجود ندارد.

۴ مدل بیز ناپارامتری دیریکله سلسله مراتبی

روش بیز سلسله مراتبی یک روش اساسی در آمار بیزی است. در مدل‌های بیز سلسله مراتبی برای پارامترهای توزیع پیشین (ابر پارامترها) توزیع‌های پیشینی در نظر گرفته می‌شود. به کارگیری روش‌های سلسله مراتبی در بیز ناپارامتری منجر به مدل‌های بیز ناپارامتری سلسله مراتبی

توزیع‌های پسینی هستند که می‌توان از آن‌ها در الگوریتم نمونه بردار گیبس برای نمونه‌گیری از این ابرپارامترها و نهایتاً نمونه‌گیری از θ_i ها استفاده کرد. همان طور که در ابتدای بخش گفتیم معمولاً فرض می‌شود که ζ به شرط Y و θ از λ و α_0 مستقل است. بنابراین

$$\begin{aligned} \Pi(\zeta|Y, \theta, \lambda, \alpha_0) &= \Pi(\zeta|Y, \theta) \\ &\propto \prod_{i=1}^n H(Y_i|\theta_i, \zeta_i) G_0(\theta_i|\lambda) H_1(\zeta). \end{aligned} \quad (21)$$

همچنین همان‌طور که در ابتدای بخش بیان کردیم توزیع شرطی ابر پارامتر λ به شرط α_0 و ζ و θ و Y تنها به θ بستگی دارد. بنابراین

$$\begin{aligned} \Pi(\lambda|Y, \theta, \zeta, \alpha_0) &= \Pi(\lambda|Y, \theta) \\ &\propto G_0(\theta|\lambda) H_2(\lambda). \end{aligned} \quad (22)$$

سپس با توجه به عبارات فوق، با استفاده از الگوریتم نمونه‌گیری گیبس اقدام به نمونه‌گیری از θ می‌کنیم. اما با توجه مطالب بیان شده در بخش‌های فوق، انتخاب یک مقدار ثابت برای α_0 در توزیع پیشین فرآیند دیریکله به طرز قابل توجهی روی نمونه‌گیری از θ_i ها تأثیر می‌گذارد. به منظور کاهش میزان تأثیر این پارامتر می‌توان یک مدل سلسله مراتبی برای این پارامتر با در نظر گرفتن یک توزیع پیشین برای α_0 و یادگیری در مورد آن در نمونه برداری گیبس ساخت که به شرح زیر می‌باشد:

فرض کنید $p(\alpha_0)$ چگالی پیشین پیوسته برای α_0 باشد، که می‌تواند به حجم نمونه، n ، وابسته باشد. آنتونی‌اک [۱] نشان داد که توزیع تعداد خوشه‌ها، k ، به شرط معلوم بودن n و α_0 به صورت زیر است:

$$p(k|\alpha_0, n) = C_n(k) \alpha_0^k f\Gamma(\alpha_0) \Gamma(\alpha_0 + n), \quad k = 1, \dots, n, \quad (23)$$

۱.۴ محاسبه توزیع پسین در مدل فرایند دیریکله سلسله مراتبی

به منظور نمونه‌گیری از پارامتر θ با توجه به بخش سوم و رابطه (۱۷) می‌دانیم که

$$\begin{aligned} P(\theta_i|\theta_{-i}, Y, \zeta, \lambda) &\propto q_0 H(Y_i|\theta_i, \zeta_i) G_0(\theta_i|\lambda) \\ &+ \sum_{j=1, j \neq i}^n q_j \delta_{\theta_j}((-\infty, \theta_i]), \end{aligned} \quad (18)$$

که در آن داریم

$$q_0 \propto \alpha_0 \int H(Y_i|\theta_i, \zeta_i) dG_0(\theta_i|\lambda),$$

همچنین $q_j \propto H(Y_i|\theta_j, \zeta_i)$ و $q_0 + \sum_{j \neq i} q_j = 1$ است. همچنین همان‌طور که در بخش سوم گفتیم رابطه (۱۸) را می‌توان با خوشه‌بندی θ_i ها به صورت θ_j^* , $j = 1, \dots, k$ در k خوشه با تعداد تکرار n_j برای θ_j^* متمایز در مجموعه $\theta_1, \dots, \theta_n$ به صورت زیر بازنویسی کرد:

$$\begin{aligned} P(\theta_i|\theta_{-i}, Y, \zeta, \lambda) &\propto q_0 H(Y_i|\theta_i, \zeta_i) G_0(\theta_i|\lambda) \\ &+ \sum_{j=1, \theta_j^* \neq \theta_i}^k n_j q_j^* \delta_{\theta_j^*}((-\infty, \theta_i]), \end{aligned} \quad (19)$$

که در آن $q_0 + \sum_j n_j q_j^* = 1$ و $q_j^* \propto H(Y_i|\theta_j^*, \zeta_i)$ است. علاوه بر آن برای جلوگیری از پایدار شدن زنجیر مارکف روی تعداد کمی از خوشه‌ها می‌توان از نمونه‌گیری خوشه‌ها به شیوه بازآمیختن با استفاده از توزیع پسین زیر نمونه‌گیری کرد.

$$(20)$$

$$P(\theta_j^*|Y, S, k, \zeta, \lambda) \propto \prod_{i \in J_j} H(Y_i|\theta_j^*, \zeta_i) dG_0(\theta_i|\lambda).$$

در مدل سلسله مراتبی ابرپارامترهای ζ و λ نیز دارای

با توجه به رابطه

$$\frac{\Gamma(\alpha_0)}{\Gamma(\alpha_0 + n)} = \frac{(\alpha_0 + n)}{\Gamma(n)} \int_0^1 \eta^{\alpha_0} (1 - \eta)^{n-1} d\eta,$$

رابطه (۲۳) را می‌توان به صورت زیر نوشت:

$$p(k|\alpha_0, n) = C_n(k) \alpha_0^{k-1} \frac{(\alpha_0 + n)}{\Gamma(n)} \int_0^1 \eta^{\alpha_0} (1 - \eta)^{n-1} d\eta,$$

که در آن $k = 1, \dots, n$. بنابراین

(۲۶)

$$p(\alpha_0|k, n) \propto (\alpha_0 + n) \alpha_0^{k+a-2} e^{-\alpha_0 b} \int_0^1 \eta^{\alpha_0} (1 - \eta)^{n-1} d\eta.$$

حال متغیر پنهان η را به صورت زیر در نظر بگیرید (اسکو بار و وست [۴])

$$(\eta|\alpha_0, n) \propto \text{Beta}(\alpha_0 + 1, n). \quad (27)$$

با دقت در رابطه‌های (۲۵) و (۲۶) می‌توان دریافت که $p(\alpha_0|k, n)$ توزیع حاشیه‌ای (α_0, η) برای α_0 است که از توزیع توأم α_0 و η : به صورت زیر به دست آمده است

$$\begin{aligned} p(\alpha_0, \eta|k, n) &\propto \alpha_0^{a+k-2} (\alpha_0 + n) \exp\{-\alpha_0 (b - \log(\eta))\} \\ &\propto \alpha_0^{a+k-1} \exp\{-\alpha_0 (b - \log(\eta))\} \\ &+ n \alpha_0^{a+k-2} \exp\{-\alpha_0 (b - \log(\eta))\}. \end{aligned} \quad (28)$$

بنابراین توزیع پسین α_0 به شرط متغیر پنهان η ، $p(\alpha_0|k, n, \eta)$ ، به صورت آمیخته‌ای از دو توزیع گاما به

که در آن ثابت $C_n(k)$ به α_0 بستگی ندارد. البته در صورت نیاز ثابت‌های $C_n(k)$ به آسانی و با استفاده از رابطه بازگشتی برای اعداد استرلینگ^۵ محاسبه می‌شوند.

همچنین

$$p(k|n) = E(p(k|\alpha_0, n)) = \int_0^\infty p(k|\alpha_0, n) p(\alpha_0) d\alpha_0.$$

با فرض این که توزیع α_0 به شرط k و n از θ و Y مستقل باشد، توزیع پسین α_0 به شرط معلوم بودن k و n برابر است با

$$\begin{aligned} p(\alpha_0|k, n, \theta, Y) &= p(\alpha_0|k, n) \quad (24) \\ &= \frac{p(k|\alpha_0, n) p(\alpha_0)}{p(k|n)}. \end{aligned}$$

معمولاً توزیع پیشین برای α_0 توزیع گاما یا آمیخته‌ای از چند توزیع گاما در نظر گرفته می‌شود. در این حالت نمونه‌گیری از توزیع پسین α_0 به سادگی انجام می‌پذیرد. در این جا ما توزیع پیشین α_0 را گاما در نظر می‌گیریم. می‌توانید تعمیم این پیشین به آمیخته‌ای از چند پیشین گاما را در وست [۱۸] ببینید.

با در نظر گرفتن $\alpha_0 \sim \Gamma(a, b)$ داریم:

$$p(\alpha_0|k, n) \propto \frac{\alpha_0^{k+a-2} \Gamma(\alpha_0) e^{-\alpha_0 b}}{\Gamma(\alpha_0 + n)}, \quad (25)$$

که توزیع پسین مشخص و سهل‌الوصولی نیست و نمونه‌گیری از آن به راحتی امکان‌پذیر نیست.

یک راه‌کار برای نمونه‌گیری از پسین α_0 در این حالت، که به آن افزایش داده^۶ گفته می‌شود، به صورت زیر است.

^۵Stirling numbers

^۶Data augmentation

صورت زیر است:

با حذف چنین محدودیت‌هایی، انعطاف‌پذیری بسیار زیادی به تحلیل داده‌های بقا بخشیده است. در این بخش ابتدا فرآیند دیریکله وابسته معرفی شده و پس از آن مدل تحلیل واریانس بر پایه این فرآیند دیریکله وابسته بیان می‌شود.

$$p(\alpha_0|k, \eta) \propto \pi_\eta \Gamma(a+k, b-\log(\eta)) + (1-\pi_\eta) \Gamma(a+k-1, b-\log(\eta)), \quad (29)$$

که در آن $\pi_\eta = \frac{a+k-1}{a+k-1+n(b-\log(\eta))}$ است.

بنابراین برای نمونه‌گیری از α_0 به شرط معلوم بودن k کفایت ابتدا η را از توزیع بتا (۲۶) به شرط α_0 مرحله قبل نمونه‌گیری گیبس تولید کنیم و سپس به شرط η و k معلوم α_0 را از پسین آمیخته (۲۹) نمونه‌گیری نماییم.

۱.۵ فرآیند دیریکله وابسته

یکی از جذاب‌ترین شیوه‌های بیز ناپارامتری برای تحلیل مدل‌های رگرسیون بقا مدل فرآیند دیریکله وابسته است که در این بخش به آن می‌پردازیم. مدل فرآیند دیریکله وابسته که توسط مک ایچرن [۱۵] معرفی شد، به صورت زیر تعریف می‌شود.

۵ مدل‌های بیز ناپارامتری دیریکله سلسله مراتبی وابسته برای داده‌های سانسور شده از راست

فرض کنید یک مجموعه از توابع توزیع‌های $\{F_x, x \in X\}$ توسط بردار p بعدی $x = (x_1, \dots, x_p)$ اندیس گذاری شده‌اند. برای مثال در یک آزمایش بالینی $F_{(x_1, x_2)}$ تابع توزیع زمان‌های بهبودی بیماران با دو دوز دارو x_1 و x_2 می‌باشد. اگر در این مجموعه برای هر x ، توزیع پیشین F_x فرآیند دیریکله با پارامتر دقت α_0 و اندازه پایه F_{0x} باشد، $F_x \sim DP(\alpha_0 F_{0x})$ ، آن‌گاه چنین فرآیندی را فرآیند دیریکله وابسته می‌نامیم.

همانطور که در بخش اول بیان شد، مدل کاکس رابطه‌ای برای مدل‌بندی تابع نرخ خطر و دیگر توابع بقا بر حسب متغیرهای پیشگو تأثیرگذار بر روی بقا می‌باشد اما این مدل با فرضیات محدود کننده متناسب بودن مخاطرات و استقلال زمان‌های رخداد پیشامدها در نظر گرفته شده است. مدل‌های نیمه پارامتری دیگر نظیر مدل زمان شکست شتابیده^۷ و یا مدل بخت‌های متناسب^۸ نیز با در نظر گرفتن مدل‌های پارامتری برای توابع بقا دارای چنین فرض‌های محدود کننده‌ای هستند. برای داده‌های بقای پیچیده، پیش فرض‌های این گونه مدل‌ها مشکلات و محدودیت‌هایی در مدل سازی به وجود می‌آورند. استفاده از شیوه‌های بیز ناپارامتری در تحلیل مدل‌های رگرسیون بقا

به بیان دقیق‌تر و کلی‌تر، فرض کنید $P_x(A)$ یک فرآیند تصادفی اندیس گذاری شده توسط x و A باشد، به طوری که x روی X و A روی Θ تغییر می‌کند. همچنین فرض کنید $\alpha_x(\cdot)$ به ازای هر $x \in X$ یک اندازه جمع پذیر روی Θ باشد. آن‌گاه می‌گوییم $\{P_x; x \in X\}$ یک فرآیند دیریکله وابسته است، اگر به ازای هر $x \in X$ و هر افراز A_1, \dots, A_r دلخواه از Θ

$$(P_x(A_1), \dots, P_x(A_r)) \sim \text{Dir}(\alpha_x(A_1), \dots, \alpha_x(A_r)).$$

^۷Accelerated Failure Time

^۸Proportional Odds

۲.۵ تحلیل واریانس فرآیند دیریکله وابسته

پیچش^۹ توابع $\alpha_0 p_m^0$ ، $\alpha_0 p_{A_v}^0$ و $\alpha_0 p_{B_w}^0$ می‌باشد.

میزان وابستگی توزیع‌های تصادفی توسط ساختار کواریانس جرم‌های نقطه‌ای θ_{xh} تعریف می‌شود. از طرفی در مدل تحلیل واریانس استاندارد (پارامتری) سطح وابستگی توسط واریانس‌های p_m^0 ، $p_{A_v}^0$ و $p_{B_w}^0$ مشخص شده است. برای مثال به ازای نقاط ثابت z_1 و z_2 ، دو نقطه $x = (v, z_1)$ و $x' = (v, z_2)$ و نمونه‌های تصادفی $y_x \sim F_x$ و $y_{x'} \sim F_{x'}$ را در نظر بگیرید. همچنین قرار دهید $p_m^0 = N(\mu_m, \sigma_m^2)$ ، $p_{A_v}^0 = N(\mu_{A_v}, \sigma_{A_v}^2)$ و $p_{B_w}^0 = N(\mu_{B_w}, \sigma_{B_w}^2)$ در مدل تحلیل واریانس استاندارد، $\text{Cov}(y_x, y_{x'}) = \sigma_m^2 + \sigma_{A_v}^2$ است. در مدل تحلیل واریانس فرآیند دیریکله وابسته می‌توان نشان داد

$$\text{cov}(y_x, y_{x'}) = \frac{\sigma_m^2 + \sigma_{A_v}^2}{\alpha_0 + 1}.$$

در واقع کواریانس پاسخ‌های مشاهده شده y_x همانند کواریانس تحت مدل تحلیل واریانس استاندارد است که توسط ضریب $\frac{1}{\alpha_0 + 1}$ کاهش یافته است.

یک مزیت مدل تحلیل واریانس فرآیند دیریکله وابسته قابلیت تفسیر ساده آن بر اساس اصطلاحات مربوط به مفاهیم تحلیل واریانس استاندارد می‌باشد. در واقع m_h را می‌توان به عنوان میانگین کلی^{۱۰} و A_{vh} و B_h به عنوان اثرات اصلی^{۱۱} برای سطح v و متغیر پیوسته z تفسیر کرد.

د آوریو و همکاران [۶] مدل‌های ناپارامتری وابسته ای را برای یک مجموعه از تابع‌ها یا توزیع‌های احتمال تصادفی در نوعی از مدل تحلیل واریانس مطرح کردند که به شرح زیر است:

فرض کنید $\{F_x, x \in X\}$ آرایه‌ای از توزیع‌های تصادفی باشد که توسط متغیر کمکی x اندیس گذاری شده است. وابستگی میان توزیع‌های تصادفی توسط وابستگی تحقق‌های آنان (θ_{xh}) ایجاد شده است که مدل تحلیل واریانس برای θ_{xh} بصورت زیر می‌باشد

$$\theta_{xh} = m_h + A_{vh} + B_h z, \quad h = 1, 2, \dots, \quad (3^0)$$

که در آن متغیر کمکی x را بتوان به صورت $x = (v, z)$ افراز کرد به طوری که v یک متغیر رسته‌ای شامل سطوح $1, 2, \dots, V$ و z یک متغیر پیوسته با تکیه‌گاه \mathbb{Z} باشد. برای مثال، متغیر کمکی v می‌تواند سطوح یک تیمار در یک آزمایش بالینی و متغیر کمکی z می‌تواند میزان یک متغیر پیوسته مربوط به هر بیمار مثل فشار خون باشد، که در این حالت F_x می‌تواند توزیع تصادفی زمان بهبودی هر بیمار در سطح تیمار معین v و فشار خون معین z باشد. همچنین $m_h \sim p_m^0$ ، $A_{vh} \sim p_{A_v}^0$ و $B_h \sim p_B^0$ ، m_h ، A_{vh} و B_h ها از هم مستقل هستند. در این حالت مدل احتمالاتی توأم فرآیند $\{F_x, x \in X\}$ را تحلیل واریانس فرآیند دیریکله وابسته می‌نامیم و می‌نویسیم $\{F_x, x \in X\} \sim \text{ANOVADDP}(\alpha_0 p^0)$ که در آن p^0 توزیع توأم (m_h, A_{vh}, B_h) است و اندازه پایه روی اثرات تحلیل واریانس در رابطه (3^0) نامیده می‌شود. از آن جا که برای هر $x = (v, z)$ توزیع تصادفی F_x دارای پیشین فرآیند دیریکله با اندازه $\alpha_0 F_{0x}$ است، بنابراین اندازه پایه F_x در حالت کلی

^۹Convolution

^{۱۰}Overall mean

^{۱۱}Main effects

۳.۵ مدل تحلیل واریانس دیریکله سلسله و مراتبی برای داده‌های سانسور شده وابسته

$$G(\cdot) | \alpha_o, G_o, \lambda \sim DP(\cdot | \alpha_o G_o(\cdot | \lambda))$$

آن‌گاه داریم

$$f_x(y_i^*) = \int [h(Y_i^* | \alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{\delta_i} [1 - H(Y_i^* | \alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{1-\delta_i} dG((\alpha_i, \theta_{oi})) \quad (32)$$

مدل تحلیل واریانس دیریکله سلسله مراتبی وابسته برای داده‌های سانسور شده به صورت زیر بیان می‌شود.

$$G(\cdot) | \alpha_o, G_o, \lambda \sim DP(\cdot | \alpha_o G_o(\cdot | \lambda)),$$

فرض کنید T_1, \dots, T_n مشاهدات حاصله باشند و نیز بعضی از مشاهدات از راست سانسور شده باشند، یعنی به ازای برخی از مقادیر i ، اگر $T_i \leq C_i$ باشد، مشاهده شده باشد و اگر $T_i > C_i$ باشد، T_i سانسور شده باشد.

اگر تابع چگالی متناظر با $H(T_i | \theta_i, \zeta)$ برابر $h(T_i | \theta_i, \zeta)$ باشد و مقادیر مشاهده شده T_1^*, \dots, T_n^* را به صورت $T_i^* = \min(T_i, C_i)$ تعریف کنیم، آن‌گاه با توجه به مطالب بیان شده در بخش سه برای $i = 1, \dots, n$ مدل فرایند دیریکله سلسله مراتبی وابسته بصورت زیر ارائه می‌شود:

$$\begin{aligned} \zeta &\sim H_1(\zeta), \\ \lambda &\sim H_2(\lambda), \\ \alpha_o &\sim \Gamma(a, b). \end{aligned} \quad (33)$$

$$T_i | X_i = x, \theta_{ix}, \zeta, \lambda \sim H(T_i | \alpha_i \eta_x, \theta_{oi}, \zeta_i),$$

در این صورت نمونه‌گیری گیبس برای نمونه‌گیری از $\alpha_i = (m_i, A_{1i}, \dots, A_{Vi}, B_i)$ ، $x = (v, z)$ که در آن $\eta_x = (1, e_v, z)'$ و α_i ها به صورت زیر انجام می‌شود:

گام ۵: پارامترهای α_o ، ζ و λ را به ترتیب از توزیع‌های پیشین $\Gamma(a, b)$ ، H_1 و H_2 تولید می‌کنیم؛

همچنین فرض کنید $\delta_i = I(Y_i \leq C_i)$ و تابع چگالی متناظر با H باشد، آن‌گاه برای $i = 1, \dots, n$:

گام ۱: برای هر $i = 1, 2, \dots, n$ از توزیع پسین زیر نمونه‌گیری می‌کنیم

$$Y_i^* | X_i = x, \theta_{ix}, \zeta_i \sim [h(Y_i^* | \alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{\delta_i} \times [1 - H(Y_i^* | \alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{1-\delta_i} \quad (31)$$

$$((\alpha_i, \theta_{oi}) | (\alpha_{-i}, \theta_{-i}), Y^*, \zeta, \lambda) \sim$$

$$q_o [h(Y_i^* | \alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{\delta_i} \times [1 - H(Y_i^* | \alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{1-\delta_i} dG_o((\alpha_i, \theta_{oi}) | \lambda) + \sum_{j=1, j \neq i}^n q_j \delta_{(\alpha_j, \theta_{oj})}(\alpha_i, \theta_{oi}), \quad (\alpha_i, \theta_{oi}) \sim G$$

از طرف دیگر فرض کنید

که در آن

به صورت زیر برآورد می‌شود:

$$\hat{F}_x(t) = \frac{1}{M} \sum_{j=1}^M H(t|\alpha_i \eta_x, \theta_{oj}, \zeta_j). \quad (34)$$

$$q_o \propto \alpha_o \int [h(Y_i^*|\alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{\delta_i} [1 - H(Y_i^*|\alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{1-\delta_i} G_o(d(\alpha_i, \theta_{oi})|\lambda),$$

۶ تحلیل داده های سرطان

در این بخش به بررسی یک مجموعه از داده‌های مربوط به طول عمر بیماران مبتلا به سرطان معده در استان اردبیل پرداخته خواهد شد. اگرچه عوامل بسیاری بر بقای بیماران اثر می‌گذارند، با این حال در اینجا ما اثر عامل دو سطحی شیمی‌درمانی مدنظر قرار گرفته است که تابع بقای بیماران با استفاده از تحلیل واریانس آمیخته سلسله مراتبی برای داده‌های سانسور شده از راست برآورد خواهد شد.

$$q_j \propto [h(Y_i^*|\alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{\delta_i} [1 - H(Y_i^*|\alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{1-\delta_i}$$

و $1 = q_o + \sum_{j \neq i} q_j$ است.

گام ۲: پارامتر ζ را از توزیع پسین زیر نمونه‌گیری می‌کنیم:

$$\pi(\zeta|Y^*, \alpha, \theta_o, \lambda) \propto \prod_{i=1}^n [h(Y_i^*|\alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{\delta_i} [1 - H(Y_i^*|\alpha_i \eta_x, \theta_{oi}, \zeta_i)]^{1-\delta_i} \times dG_o((\alpha_i, \theta_{oi})|\lambda) h_1(\zeta_i),$$

که در آن h_1 تابع چگالی احتمال متناظر با H_1 است.

گام ۳: پارامتر λ را از توزیع پسین زیر نمونه‌گیری می‌کنیم:

$$\Pi(\lambda|Y, \alpha, \theta_o, \zeta, \alpha_o) \propto G_o(\alpha, \theta_o|\lambda) H_2(\lambda).$$

گام ۴: پارامتر پنهان η را از توزیع $Beta(\alpha_o + 1, n)$ و سپس پارامتر α_o جاری را از توزیع آمیخته ۲۷ نمونه‌گیری می‌کنیم؛

گام ۵: به گام ۱ باز می‌گردیم و تا زمان رسیدن به همگرایی زنجیر الگوریتم را تکرار می‌کنیم. تعداد نمونه‌های تولید شده تا قبل از همگرایی را به عنوان دوره سوخت کنار می‌گذاریم. پس از طی دوره سوخت بر اساس نمونه $j = 1, \dots, M$ تابع توزیع تجمعی $(\alpha_j, \theta_{oj}, \zeta_j, \lambda_j)$

۱.۶ تحلیل بقای داده های سرطان در حضور متغیر کمکی شیمی درمانی

جدول ۱ خلاصه ای از داده‌های مربوط به بیماران سرطانی را به تفکیک شیمی درمانی و نیز سانسور نشان می‌دهد.

جدول ۱: خلاصه داده‌های سرطان

متغیر کمکی شیمی درمانی	سطوح	رخداد	سانسور	جمع
بله	۶۹	۲۲	۹۱	
خیر	۴۲	۱۴	۵۶	

ملاحظه می‌شود که داده‌ها شامل $n = 147$ مشاهده طول عمر سانسور شده از راست T_1^*, \dots, T_n^* هستند. مشاهدات سانسور شده به صورت $T_i^* = \min(T_i, C_i)$ تعریف می‌شوند به طوری که C_1, \dots, C_n زمان‌های سانسور معین هستند. حال فرض کنید برای $i = 1, \dots, n$

$$T_i|X_i = x, \theta_{ix} \sim H(T_i|\beta_i \eta_x, \gamma_i),$$

$$F_x(y_i) = \int H(y_i|\alpha_i \eta_x, \theta_{oi}, \zeta) dG(\alpha_i, \theta_{oi}).$$

که در آن، $x, \eta_x = (1, e_x)'$ متغیر رسته‌ای دو سطحی شیمی درمانی (۱) تحت شیمی درمانی قرار گرفته‌اند، $e_x = 2$ تحت شیمی درمانی قرار نگرفته‌اند)، $\beta_i = (\beta_{0i}, \beta_{1i})$ و e_x به صورت زیر است

$$((\beta_i, \gamma_i) | (\beta_{-i}, \gamma_{-i}), T^*) \sim q_0 [h(T_i^* | \beta_i \eta_x, \gamma_i)]^{\delta_i} \times [1 - H(T_i^* | \beta_i \eta_x, \gamma_i)]^{1-\delta_i} dG_0(\beta_i, \gamma_i | m, B, S) + \sum_{j=1, j \neq i}^n q_j \delta_{(\beta_j, \gamma_j)}(\beta_i, \gamma_i),$$

که در آن

$$e_x = \begin{cases} 1 & , v = 1 \\ -1 & , v = 2 \end{cases}$$

هسته H را توزیع لگ نرمال در نظر گرفته و با توجه به سانسور از راست مشاهدات، داریم

$$q_0 \propto \alpha_0 \int [h(T_i^* | \beta_i \eta_x, \gamma_i)]^{\delta_i} [1 - H(T_i^* | \beta_i \eta_x, \gamma_i)]^{1-\delta_i} dG_0(\beta_i, \gamma_i | m, B, S),$$

$$f_{T_i^*}(t_i^*) = \int [h(T_i^* | \beta_i \eta_x, \gamma_i)]^{\delta_i} [1 - H(T_i^* | \beta_i \eta_x, \gamma_i)]^{1-\delta_i} dG(\beta_i, \gamma_i | \lambda).$$

همچنین

ساختار توزیع‌های پیشین در این تحلیل به صورت زیر است:

$$q_j \propto [h(T_i^* | \beta_i \eta_x, \gamma_i)]^{\delta_i} [1 - H(T_i^* | \beta_i \eta_x, \gamma_i)]^{1-\delta_i}$$

و نیز

$$T_i^* | \beta_i \eta_x, \gamma_i \sim [h(T_i^* | \beta_i \eta_x, \gamma_i)]^{\delta_i} [1 - H(T_i^* | \beta_i \eta_x, \gamma_i)]^{1-\delta_i}$$

$$(\beta_i, \gamma_i) \sim G, \quad i = 1, \dots, n,$$

$$G(\cdot) | \alpha_0, m, B, S \sim DP(\cdot | \alpha_0 G_0(\cdot | m, B, S)),$$

$$G_0(\beta, \gamma) = \Gamma(\gamma^{-1} | \frac{s}{\gamma}, \frac{sS}{\gamma}) N_{\Gamma}(\beta | m, B),$$

$$S \sim \Gamma(\frac{s}{5} q, \frac{s \Delta q}{R}),$$

$$m \sim N_{\Gamma}(a, A),$$

$$B^{-1} \sim Wish_{\mathcal{R}}((cC)^{-1}, c),$$

$$\alpha_0 \sim \Gamma(a_0, b_0),$$

ب. با انجام ۷۰۰ بار الگوریتم متروپولیس-هستینگز به ایستایی می‌رسد. شکل‌های ۱ و ۲ نمودار خودهمبستگی مقادیر $\beta_i^{(j)}$ ها و $\gamma_i^{(j)}$ ها را برای $j = 1, \dots, 700$ نشان می‌دهد.

شکل ۳ برآورد تابع توزیع تجمعی تجربی برای مشاهدات طول عمر مربوط به دو سطح متغیر کمکی را نشان می‌دهد.

همچنین شکل ۴ برآورد تابع بقا با استفاده از سه روش کاپلان-مایر، مدل کاکس و مدل تحلیل واریانس دیرینله سلسله مراتبی را نشان می‌دهد. همانطور که مشاهده می‌شود، برآوردهای مدل کاپلان مایر و کاکس در بعضی از زمان‌ها مقدار ثابتی اختیار می‌کنند. این در حالی است که حضور داده‌های سانسور شده در برآورد تابع بقا تحت مدل

که در آن $R = 1, C = A = 0.00412, q = s = 5$ و $a_0 = b_0 = 1$ است.

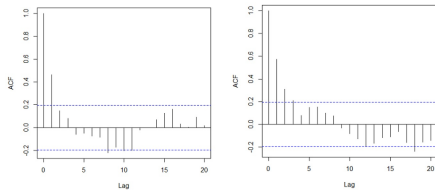
الگوریتم نمونه‌گیری (β_i, γ_i) ها مشابه بخش قبل است و تنها در موارد زیر اختلاف دارد:

الف. مقادیر (β_i, γ_i) را برای $i = 1, 2, \dots, n$ از توزیع

واریانس تاثیری نمی‌گذارند.

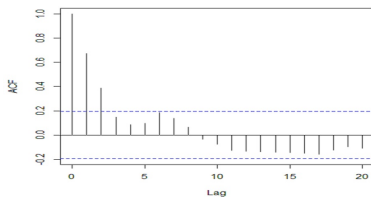
گیس و متروپولیس-هستینگز محاسبات پیچیده در این شیوه را به طرز قابل توجهی ساده می‌سازد. همچنین تحلیل نمودارهای هموار برآورد شده در خروجی بسیار ساده‌تر و مطلوب‌تر از سایر شیوه‌ها است.

طبق برآورد تابع بقا تحت مدل تحلیل واریانس فرایند دیریکله سلسله مراتبی می‌توان دریافت که با احتمال ۴۰ درصد شخصی که تحت شیمی درمانی قرار نگرفته است تا ۲۰ سال اول بعد از شروع بیماری زنده می‌ماند. این در حالی است که با احتمال ۴۰ درصد شخصی که تحت شیمی درمانی قرار گرفته است بیش از ۲۰ ماه بعد از شروع بیماری زنده می‌ماند. بنابراین شیمی درمانی بر روی بقای بیماران سرطانی تاثیرگذار است.



شکل ۱: نمودار خودهمبستگی مربوط به مقادیر $\beta_i^{(j)}$ $j = 1, \dots, 700$.

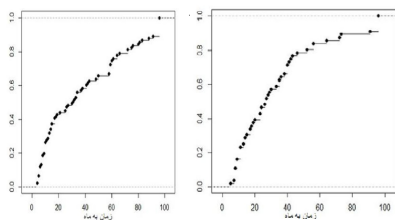
در شکل ۵ برآورد تابع چگالی احتمال، تابع توزیع تجمعی، تابع نرخ خطر و تابع نرخ خطر تجمعی براساس مدل تحلیل واریانس دیریکله سلسله مراتبی را نشان داده می‌شود.



شکل ۲: نمودار خودهمبستگی مربوط به مقادیر $\gamma_i^{(j)}$ $j = 1, \dots, 700$.

۷ بحث و نتیجه‌گیری

یکی از محدودیت‌های مدل‌های ناپارامتری نظیر روش برآورد کاپلان-مایر عدم حضور متغیر کمکی می‌باشد که تنها برای متغیر رسته‌ای می‌توان برای هر سطح تیمار منحنی‌های جداگانه‌ای رسم کرد. این در حالی است که مدل کاکس و تحلیل واریانس سلسله مراتبی برای تمامی مقادیر هر نوع متغیر کمکی به خوبی تابع بقا را برآورد می‌کنند.



همچنین در مدل‌هایی نظیر مدل کاکس فرضیه‌های محدودکننده‌ای برای ارتباط توابع بقا با متغیرهای کمکی وجود دارد که باعث کاهش کارایی این مدل‌ها می‌شود. مدل تحلیل واریانس دیریکله سلسله مراتبی دارای مفروضات اولیه بسیار کمتری است.

شکل ۳: راست؛ برآورد تابع توزیع تجمعی تجربی مشاهدات طول عمر بیمارانی که تحت شیمی‌درمانی قرار نگرفته‌اند. چپ؛ برآورد تابع توزیع تجمعی تجربی مشاهدات طول عمر بیمارانی که تحت شیمی‌درمانی قرار گرفته‌اند.

استفاده از تحلیل بیز سلسله مراتبی به طرز قابل توجهی وابستگی این روش بیزی را به توزیع‌های پیشین کاهش می‌دهد. همچنین استفاده از روش‌های عددی نمونه‌گیری

Annals of Statistics. 2, 1152-1174.

[3] Blackwell, D. and MacQueen, J. B. (1973). Ferguson Distribution via Polya Urn Schemes. Annals of Statistics. 1, 353-355.

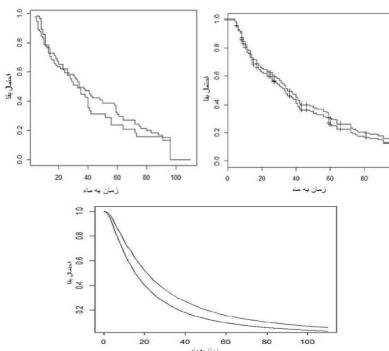
[4] Christensen, R. and Johnson, W. (1988). Modeling Accelerated Failure Time with a Dirichlet Process. Biometrika. 75, 693-704.

[5] Escobar, M. D. and West, M. (1995). Bayesian density estimation and inference using mixture. Journal of the American Statistical Association. 90, 577-588.

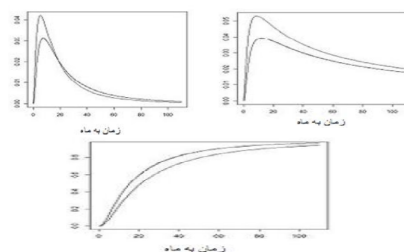
[6] Escobar, M. D. and West, M. (1998). Computing Nonparametric Hierarchical Models. Springer Verlag. New York. USA.

[7] De Iorio, M., Johnson, W. O., Muller, P. and Rosner, G. L. (2009). Bayesian Nonparametric Nonproportional Hazard Survival Modeling. Biometrics. 65, 762-771.

[8] Doss, H. and Huffer, F. W. (1999). Monte Carlo Method for Bayesian Analysis Survival Data using of Dirichlet Process Prior. Journal of Computational and Graphical Statistics. 282-307.



شکل ۴: برآورد تابع بقا؛ سمت چپ مدل کاکس، سمت راست مدل کاپلان مایر، پایین مدل تحلیل واریانس دیریکله سلسله مراتبی



شکل ۵: برآورد توابع عمده مورد استفاده در تحلیل بقا: بالا راست؛ تابع چگالی احتمال، بالا چپ؛ تابع نرخ خطر، پایین؛ تابع توزیع تجمعی.

مراجع

[۱] عاطفه جاویدی، سمیه راه پیما، مجید جعفری خالدی (۱۳۹۲). معرفی پیشین فرایند دیریکله در چارچوب مدل‌های بیزی ناپارامتری. اندیشه آماری، سال هجدهم، شماره ۲، ۶۱-۷۲.

[2] Antoniak, C. E. (1974). Mixture of Dirichlet Process with applications to Bayesian nonparametric problems.

- data. *Journal of the Royal Statistical Society*. 40, 214-241.
- [16] Kuo, L. and Mallick, B. (1997). Bayesian semiparametric inference for the accelerated failure- time model. *Canadian Journal of Statistics*. 25, 457-472.
- [17] MacEachern, S. N. (1999). Dependent nonparametric processes. American Statistical Association. Virginia. USA.
- [18] Walker, S. and Mallick, B. K. (1999). A Bayesian semiparametric accelerated failure time model. *Biometrics*. 55, 477-483.
- [19] West, M. (1992). Hyperparameter estimation in Dirichlet process mixture models. ISDS Discussion Paper # 92-A 03. Duke University. USA.
- [20] West, M., Muller, P. and Escobar, M. D. (1994). Hierarchical Priors and Mixture Models, with Applications in Regression and Density Estimation. John Wiley and Sons. London. UK.
- [21] Zhou, M. (2004). Nonparametric Bayes estimator of survival function for doubly/ interval Censored data. *Statistica Sinica*. 14, 533-546.
- [9] Escobar, M. D. (1988). Estimating the means of several normal populations by estimating the distribution the means. Ph. D. thesis. Yale University. New Haven. USA.
- [10] Escobar, M. D. and West, M. (1998). Computing Nonparametric Hierarchical Models. Springer Verlag. New York. USA.
- [11] Ferguson, T. S. (1973). A Bayesian analysis of some nonparametric problems. *Annals of Statistics*. 1, 209-230.
- [12] Gelfand, A. E. and Kottas, A. (2002). Bayesian semiparametric regression for median residual life. *Scandinavia Journal of Statistics*. 30, 651-665.
- [13] Hanson, T. and Johnson, W. O. (2002). Modeling regression error with a mixture of poly trees. *Journal of the American Statistical Association*. 97, 1020-1033.
- [14] Hanson, T., Johnson, W. O. and Laud, P. (2008). A unified approach to semiparametric inference for survival models with step process covariates. *Canadian Journal of Statistics*. 37, 60-79.
- [15] Kalbfleisch, J. (1978). Nonparametric Bayesian analysis of survival time