

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه اصفهان

دانشکده علوم

گروه آمار

پایان نامه‌ی کارشناسی ارشد رشته‌ی آمار گرایش ریاضی

رهیافت بیزی نیمه پارامتری و کاربرد آن در مدل‌های خطی با اثرهای آمیخته

استاد راهنما:

دکتر ایرج کاظمی

پژوهشگر:

کیوان عسلی صاف

۱۳۸۸/۱۰/۲۷

اصوات بزرگ میانی
تسبیح دارک

شهریور ماه ۱۳۸۸

۱۲۹۹۸۱

کلیه حقوق مادی مترتب بر نتایج مطالعات،
ابتکارات و نوآوری های ناشی از تحقیق
موضوع این پایان نامه متعلق به دانشگاه
اصفهان است.



دانشگاه اصفهان

دانشکده علوم

گروه آمار

پایان نامه
در رشته آمار
تصویب شده است
تصویبات تکمیلی دانشگاه اصفهان

پایان نامه‌ی کارشناسی ارشد رشته‌ی آمار گرایش ریاضی

آقای کیوان عسلی صاف

تحت عنوان

رهیافت بیزی نیمه پارامتری و کاربرد آن در مدل‌های خطی با اثرهای آمیخته

در تاریخ ۸۸/۶/۲۵ توسط هیأت داوران زیر بررسی با نمره ۱۸/۱۹ و با درجه عالی به تصویب نهایی رسید.

۱- استاد راهنمای پایان‌نامه دکتر ایرج کاظمی با مرتبه‌ی علمی استادیار

۲- استاد داور داخل گروه دکتر هوشنگ طالبی با مرتبه‌ی علمی استادیار

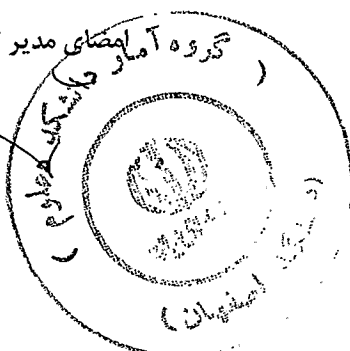
۳- استاد داور خارج از گروه دکتر حمید پزشک با مرتبه‌ی علمی استاد

امضاء

امضاء

امضاء

گروه آمار (امضای مدیر گروه)



سپاسگزاری

سپاس خدای را مرا شوق دانستن و توان ادراک آموخت.

مراتب سپاس و تشکر خود را از استاد گرانقدرم جناب آقای دکتر ایرج کاظمی ابراز می‌دارم که در تمام مراحل تدوین این رساله با تلاش‌های بی‌وقفه و راهنمایی‌های ارزنده خود مرا یاری نمودند. همچنین از زحمات بی‌دریغ اساتید داور جناب آقای دکتر هوشنگ طالبی و جناب آقای دکتر حمید پزشک کمال تشکر را دارم.

از پدرم که سربلند زیستن و از مادرم که صبر و ایثار را به من آموختند، نهایت تشکر و قدردانی را دارم.

تقدیم به

پدر و مادرم

به پاس تعبیر عظیم و انسانی‌شان از واژه ایثار،
به پاس عاطفه سرشار و گرمای امیدبخش وجودشان
و به پاس محبت‌های بی‌دریغشان که هرگز فروکش نمی‌کند.

چکیده

متداول ترین مدل‌ها برای تحلیل داده‌های وابسته در علوم مختلف، از جمله مطالعات زیستی، اقتصادسنجی و پزشکی و غیره، مدل‌های خطی تعمیم‌یافته با اثرهای آمیخته است. از آنجایی که استفاده از روش‌های مناسب برآوردیابی پارامترهای این مدل‌ها حائز اهمیت فراوان در حصول نتایج می‌باشد و با توجه به استفاده چشم‌گیر از آمار بیز در برآزش مدل‌های پیچیده، این پایان‌نامه با به‌کارگیری روش‌های بیز نیمه‌پارامتری به مطالعه این موضوع می‌پردازد. یک مسأله اساسی در استفاده از این مدل‌ها بکارگیری فرآیند مناسبی برای تولید توزیع اثرهای آمیخته است. این پایان‌نامه، با فرض اینکه توزیع این اثرها از فرآیند دیریکله با توزیع پایه و پارامتر دقت مشخص پیروی می‌کنند، مدل‌های بیزی سلسله مراتب نیمه‌پارامتری را برای تحلیل داده‌های وابسته در رگرسیون خطی تعمیم‌یافته مورد بررسی قرار می‌دهد. بدین منظور، ابتدا فرآیند دیریکله معرفی و چگونگی تولید آن توسط روش‌های مختلف شرح داده می‌شود. سپس روش بکارگیری آن به عنوان فرآیند تولید توزیع‌های پیشین در مدل‌های مختلف توضیح داده خواهد شد. در ادامه، توزیع‌های پسین شرطی کامل پارامترهای مدل‌های رگرسیونی با اثرهای آمیخته را یافته و توسط رهیافت بیز، استنباط آماری را با الگوریتم نمونه‌گیری گیبز انجام خواهیم داد. همچنین مدل‌های متفاوتی به داده‌های واقعی از علوم مختلف برآزش داده خواهد شد تا اهمیت کاربرد روش‌های بیز نیمه‌پارامتری در مقایسه با روش‌های متداول مشخص شود.

کلید واژه‌ها : مدل‌های رگرسیون خطی تعمیم‌یافته، مدل‌های بیز سلسله مراتب، اثرهای آمیخته، فرآیند دیریکله، نمونه‌گیری گیبز، توزیع‌های پسین شرطی کامل.

فهرست مطالب

صفحه

عنوان

فصل اول : کلیات

۱	۱ - ۱ مقدمه	
۱	۲ - ۱ موضوع و هدف تحقیق	
۲	۳ - ۱ اهمیت و کاربرد موضوع	
۳	۴ - ۱ زمینه و تاریخچه موضوع تحقیق	

فصل دوم : تعاریف و مفاهیم مقدماتی

۵	۱ - ۲ مقدمه	
۶	۲ - ۲ برآوردگرهای پسین بیز	
۷	۱ - ۲ - ۲ توزیع‌های پیشین و پسین	
۱۹	۳ - ۲ دنباله‌های مونت کارلو مارکوفی (MCMC)	
۱۹	۱ - ۳ - ۲ توزیع‌های پسین شرطی کامل	
۲۰	۲ - ۳ - ۲ نمونه‌گیری گیبز	
۲۶	۴ - ۲ رهیافت بیزی سلسله مراتبی	
۳۱	۵ - ۲ نتیجه‌گیری	

فصل سوم : مدل‌های رگرسیون خطی در قالب بیز و سنجش اعتبار آنها

۳۲	۱ - ۳ مقدمه	
۳۳	۲ - ۳ مدل‌های رگرسیون خطی ساده	
۳۳	۳ - ۳ مدل رگرسیون خطی چندگانه	
۳۴	۴ - ۳ مدل خطی	
۳۵	۱ - ۴ - ۲ برآوردگرهای ماکزیمم درست‌نمایی β و σ^2 و خواص آنها	
۳۶	۵ - ۳ مدل رگرسیون خطی در قالب بیزی	
۳۹	۱ - ۵ - ۳ چگالی‌های پسین حاشیه‌ای β و τ	
۴۱	۲ - ۵ - ۳ استنباط در رگرسیون خطی بیزی برای β	

۴۲	۳-۵-۳ استنباط در رگرسیون خطی بیزی برای τ
۴۳	۳-۶ استنباط بیز بر اساس شبیه‌سازی دنباله‌های مونت کارلوی مارکوفی
۴۷	۳-۷ مدل‌های رگرسیون خطی با اثرات آمیخته
۵۰	۳-۷-۱ بهترین پیش‌بینی خطی (BLP) متغیر تصادفی α
۵۲	۳-۸ مدل‌های رگرسیون خطی بیزی با اثرات آمیخته
۶۰	۳-۹ مدل‌های خطی تعمیم یافته
۶۰	۳-۹-۱ خانواده توزیع‌های نمایی
۶۱	۳-۹-۲ شکل قراردادی برای مدل‌های خطی تعمیم یافته
۶۵	۳-۹-۳ معادلات درست‌نمایی برای مدل‌های خطی تعمیم یافته
۶۷	۳-۱۰ معیارهای مقایسه مدل‌ها
۶۷	۳-۱۰-۱ DIC برای مقایسه مدل‌های بیزی
۶۸	۳-۱۰-۲ AIC و BIC برای مقایسه مدل‌های بیز
۷۲	۳-۱۱ خلاصه و نتیجه‌گیری

فصل چهارم: برازش مدل آمیخته فرآیند دیریکله

۷۳	۴-۱ مقدمه
۷۴	۴-۲ توزیع آمیخته
۷۴	۴-۲-۱ کاربردی از توزیع‌های آمیخته
۷۵	۴-۳ توزیع دیریکله
۷۶	۴-۳-۱ توزیع دیریکله با شکل ظاهری متفاوت
۷۷	۴-۳-۲ خواص مهم کلاس توزیع‌های دیریکله
۷۸	۴-۴ فرآیند دیریکله
۷۸	۴-۴-۱ خواص مهم فرآیند دیریکله
۸۲	۴-۵ روش‌های شبیه‌سازی فرآیند دیریکله
۸۲	۴-۵-۱ روش شکستن میله
۸۴	۴-۵-۲ روش کیسه پولیا
۸۵	۴-۶ مدل آمیخته فرآیند دیریکله (DPM)
۸۸	۴-۷ نمونه‌گیری پسین تحت DPM

عنوان	صفحه
۴ - ۸ افزایش کردن در استنباط برای مدل آمیخته	۹۴
۴ - ۹ نتیجه‌گیری	۹۸
فصل پنجم : کاربرد فرآیند دیریکله در رگرسیون خطی با اثرات آمیخته	
۵ - ۱ مقدمه	۹۹
۵ - ۲ MDP در مدل رگرسیون خطی چندگانه با اثرهای آمیخته	۱۰۰
۵ - ۲ - ۱ توزیع پسین شرطی کامل اثرات تصادفی در رگرسیون خطی چندگانه با	
اثرات آمیخته	۱۰۲
۵ - ۳ MDP در مدل رگرسیون خطی تعمیم یافته با اثرات آمیخته	۱۱۳
۵ - ۳ - ۱ توزیع پسین شرطی کامل اثرات تصادفی رگرسیون خطی تعمیم یافته	
با اثرات آمیخته	۱۱۶
۵ - ۴ نتیجه‌گیری	۱۲۰
پیوست ها	۱۲۱
منابع و مآخذ	۱۳۳

فهرست جدول‌ها

صفحه	عنوان
۲۵	جدول (۱ - ۲) برآورد پارامترهای b و λ مثال ۲ - ۶
۳۱	جدول (۲ - ۲) برآورد پارامترهای b و λ مثال ۲ - ۸ (مدل پواسن سلسله مراتبی)
۴۶	جدول (۱ - ۳) برآورد $\beta_1, \beta_2, \beta_3, \sigma^2$ مثال ۳ - ۱
۴۷	جدول (۲ - ۳) داده‌های چربی بدن مثال ۳ - ۱
۵۹	جدول (۳ - ۳) برآورد پارامترهای مجهول تحت مدل A برای مثال ۳ - ۲
۶۳	جدول (۳ - ۴) توابع پیوند کانونی برای توزیع‌های پرکاربرد در مدل‌های خطی تعمیم‌یافته
۷۱	جدول (۳ - ۵) برآورد β_0 و β_1 و Deviance در مثال ۳ - ۳ (دستگاه‌های ریسندگی)
۷۱	جدول (۳ - ۶) مقایسه مدل رگرسیونی لجیت و پروبیت در مثال ۳ - ۳ (دستگاه‌های ریسندگی)
۹۶	جدول (۱ - ۴) برآورد پارامترهای مجهول در مدل B برای مثال ۴ - ۳
	جدول (۲ - ۴) برخی از برآوردهای ضرایب اثرات نمرات اکتسابی، به همراه برآورد فواصل
۹۷	اطمینان هر کدام و انحراف معیار آنها به تفکیک سه مدل A, B و C در مثال (۴ - ۳)
۹۷	جدول (۳ - ۴) معیار AIC برای مقایسه مدل‌های A, B و C در مثال ۴ - ۳
۱۰۷	جدول (۱ - ۵) برآورد پارامترهای مجهول در مثال ۵ - ۱ (قسمت اول)
۱۰۷	جدول (۲ - ۵) برآورد اثرات تصادفی مدارس در مثال ۵ - ۱ (قسمت اول)
۱۰۸	جدول (۳ - ۵) برآورد Deviance در مثال ۵ - ۱ (قسمت اول)
۱۱۰	جدول (۴ - ۵) برآورد پارامترهای مجهول در مثال ۵ - ۱ (قسمت دوم)
	جدول (۵ - ۵) برآورد احتمالات متعلق بودن اثرات تصادفی به خوشه ۱ و ۲
۱۱۰	در مثال ۵ - ۱ (قسمت دوم)
۱۱۱	جدول (۵ - ۶) برآورد اثرات تصادفی مدارس در مثال ۵ - ۱ (قسمت دوم)
۱۱۲	جدول (۵ - ۷) برآورد Deviance در مثال ۵ - ۱ (قسمت دوم)
۱۱۸	جدول (۵ - ۸) برآورد میانگین ضرایب رگرسیون و اثرات تصادفی در مثال ۵ - ۲
۱۱۹	جدول (۵ - ۹) برآورد احتمالات متعلق بودن اثرات تصادفی به خوشه ۱ و ۲ مثال ۵ - ۲
۱۲۰	جدول (۵ - ۱۰) برآورد اثرات تصادفی بخش‌های کارخانه ریسندگی در مثال ۵ - ۲

فهرست شکل‌ها

صفحه	عنوان
۲۵	شکل (۱ - ۲) توزیع‌های کناری μ و σ^2 مثال ۲ - ۶
۲۶	شکل (۲ - ۲) تکرارهای نمونه‌های استخراجی برای μ و σ^2 مثال ۲ - ۶
۲۶	شکل (۳ - ۲) همبستگی نمونه‌ها برای μ و σ^2 مثال ۲ - ۶
۴۶	شکل (۱ - ۳) توزیع‌های $\sigma^2, \beta_3, \beta_2, \beta_1$ مثال ۳ - ۱
۸۳	شکل (۱ - ۴) شکستن قسمتهایی از یک میله به طول واحد برای شبیه‌سازی فرآیند دیریکله
۱۰۸	شکل (۱ - ۵) نمودار اثرات تصادفی ۱۶۰ دبیرستان مثال ۵ - ۱ (قسمت اول)
۱۱۱	شکل (۲ - ۵) نمودار اثرات تصادفی ۱۶۰ دبیرستان با پیشین دیریکله مثال ۵ - ۱ (قسمت دوم)
۱۱۹	شکل (۳ - ۵) اثرات تصادفی در بخش‌های کارخانه ریسندگی مثال (۵ - ۲)

پیوست‌ها

صفحه

عنوان

۱۲۳	پیوست ۱ : برنامه‌های کامپیوتری
۱۳۰	پیوست ۲ : توزیع‌ها
۱۳۲	پیوست ۳ : مخفف‌ها و نمادها

فصل اول

کلیات

۱-۱ مقدمه

در این فصل ابتدا موضوع تحقیق را بیان کرده و سپس هدف اصلی از انجام این تحقیق را ذکر می‌کنیم. در ادامه به بررسی اهمیت و کاربرد موضوع مورد نظر خواهیم پرداخت و در این راستا به بیان زمینه و تاریخچه موضوع تحقیق نیز اشاره کامل خواهد شد. لازم به ذکر است که در تمام فصول این پایان‌نامه از حرف بزرگ لاتین برای نمایش متغیر تصادفی و از پررنگ کردن حروف بزرگ لاتین برای نمایش بردار و ماتریس استفاده کرده‌ایم.

۲-۱ موضوع و هدف تحقیق

موضوع این پایان‌نامه رهیافت بیزی نیمه پارامتری و کاربرد آن در مدل‌های با اثرهای آمیخته می‌باشد. در این پایان‌نامه سعی شده است که روش‌های جدید مرتبط با موضوع تحقیق مورد بحث قرار گیرند. از جمله به بررسی مسئله برآوردیابی پارامترها و مدل‌سازی در حالت‌های مختلف خطی خواهیم پرداخت. هدف از انجام این تحقیق، معرفی راهکار جدیدی برای برآورد پارامترهای مجهول مدل‌های مورد بررسی می‌باشد که با توجه به مفهوم توزیع‌های آمیخته با رویکرد بیزی و اتخاذ فرآیند دیریکله به عنوان پیشین برای اثرات تصادفی، به بحث در زمینه تحلیل داده‌های وابسته و بخصوص موضوع تغییرپذیری بین واحدها خواهیم پرداخت.

۳-۱ اهمیت و کاربرد موضوع

از آنجایی که برای تحلیل داده‌های وابسته در علوم مختلف، به طور متداول از مدل‌های خطی تعمیم‌یافته با اثرهای آمیخته استفاده می‌شود، لذا استفاده از روش‌های مناسب برازش مدل و برآوردیابی پارامترهای این مدل‌ها حائز اهمیت فراوان در حصول نتایج می‌باشد. با توجه به گسترش زیاد آمار بیز در برازش مدل‌ها، در این پایان‌نامه با به کارگیری روش‌های بیز نیمه‌پارامتری به مطالعه برآورد پارامترهای مدل‌های مذکور می‌پردازیم.

واقف بودن محقق به پاره‌ای از اطلاعات قبلی در مورد مدل‌های با اثرات آمیخته و پارامترهای آن اهمیت ویژه‌ای در حصول نتایج برآوردیابی دارد بنابراین یک مسأله اساسی در استفاده از این مدل‌ها بکارگیری فرآیند مناسبی برای تولید توزیع پیشین اثرهای آمیخته است. در این پایان‌نامه، فرض می‌کنیم که توزیع این اثرها از فرآیند دیریکله با توزیع پایه و پارامتر دقت مشخص پیروی می‌کنند. در نهایت برای مقایسه روش ارائه شده با روش‌های متداول بیز به برآورد پارامترهای مجهول و اثرات آمیخته در مدل‌های خطی می‌پردازیم و سپس با توجه به معیارهای اعتبار مدل‌ها از قبیل معیار اطلاع پراکندگی (DIC)،^۱ AIC^۲ و پراکندگی^۳ و غیره که در فصل سوم به طور مفصل مورد بحث قرار خواهند گرفت، به مقایسه و سنجش روش ارائه شده با روش‌های معمول بیز خواهیم پرداخت.

۴-۱ زمینه و تاریخچه موضوع تحقیق

کارهای تحقیقاتی بر روی توزیع‌های آمیخته دارای قدمت زیادی در علم آمار می‌باشد. برای مثال باتاچاریا در سال ۱۹۶۷ توزیع طول ماهی‌ها را به صورت یک جامعه آمیخته در نظر گرفته و از آن در تعیین سن گروه‌ها استفاده کرد. دقت کنید که هر گروه سنی یک توزیع متفاوت از نظر طولی دارد به طوری که طول ماهی‌ها، برخی اطلاعات را در مورد سن فراهم می‌آورد.

جهت یک آشنایی کلی با موضوع مورد بحث، فرض کنید که برای $Y_i, i = 1, 2, \dots, k$ ها به شرط پارامترهای θ_i به صورت آمیخته‌ای از توزیع‌های $f(\cdot)$ که در پارامترهایشان متفاوت هستند، پیروی کنند. همچنین فرض کنید که θ_i ها دارای توزیع یکسان G_0 باشند. فرگوسن^۴ (۱۹۷۳) این مدل را به صورت سلسله‌مراتبی بیزی به صورت زیر بیان نمود

^۱ Deviance Informatin Criterion
^۲ Akaike Information Criterion
^۳ Deviance
^۴ Ferguson (1973)

$$\begin{aligned}
 Y_i | c_i, \theta &\sim f(Y | \theta_{c_i}) \\
 \mathbf{c} | \pi_{1:k} &\sim \text{discrete}(\pi_1, \pi_2, \dots, \pi_k) \\
 \theta_i &\sim G_0(\cdot) \\
 (\pi_1, \pi_2, \dots, \pi_k) &\sim \text{Dir}(\alpha, M)
 \end{aligned}$$

که در آن c_i ها نشانگرهایی هستند که مقادیر Y_i را به یک مقدار θ_{c_i} نسبت می‌دهند و π_i ها ضرایب آمیختگی هستند که نامنفی بوده و مجموعشان یک است. $\text{Dir}(\alpha, M)$ نیز نشانگر توزیع دیریکله با پارامترهای α و M می‌باشد. برای بیان مفهوم توزیع دیریکله فرض کنید $\Theta = \{\theta_1, \dots, \theta_n\}$ یک توزیع احتمال روی فضای گسسته $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$ باشد، به طوری که $P(X = X_i) = \theta_i$ که در آن X یک متغیر تصادفی در فضای \mathbf{X} می‌باشد. توزیع دیریکله روی Θ به صورت زیر می‌باشد

$$P(\Theta | \alpha, M) = \frac{\Gamma(\alpha)}{\prod_{i=1}^n \Gamma(\alpha m_i)} \prod_{i=1}^n \theta_i^{\alpha m_i - 1}$$

که در آن $\mathbf{M} = \{m_1, \dots, m_n\}$ و $i = 1, \dots, n$ اندازه پایه تعریف شده روی \mathbf{X} و بردار میانگین Θ و α پارامتر دقی است که چگونگی تمرکز توزیع حول \mathbf{M} را نشان می‌دهد. بزرگ بودن α اطمینان ما را در مورد انتخاب مناسب اندازه پایه \mathbf{M} بیشتر می‌کند و از این رو توزیع بیشتر حول \mathbf{M} متمرکز می‌شود. برای مقادیر مفروض ضرایب آمیختگی فرض معمول بر آن است که نشانگرهای c_i دارای توزیع چندجمله‌ای باشند. حال اگر تعداد مولفه‌های جامعه آمیخته زیاد باشد و به عبارت دیگر وقتی که $k \rightarrow \infty$ آنتونیاک^۱ (۱۹۷۴) مدل آمیخته فوق را به صورت

$$\begin{aligned}
 Y_i | \theta &\sim f(Y | \theta_i) \\
 \theta_i &\sim G(\cdot) \\
 G &\sim \text{DP}(\alpha G_0(\theta))
 \end{aligned}$$

^۱ Antoniak (1974)

نشان داد که در آن G یک توزیع عمومی و آزاد می‌باشد و از فرآیند دیریکله $DP(\alpha G_0(\theta))$ با اندازه پایه G_0 و پارامتر دقت α پیروی می‌کند.

کاربردهای اساسی ساختار فوق را که به مدل بیزی سلسله‌مراتبی نیمه پارامتری معروف است، اسکویار و وست (۱۹۹۸) و همچنین کلینمن و ابراهیم^۱ (۱۹۹۸) در مدل‌های خطی تعمیم‌یافته با اثرات آمیخته بیان کردند. واژه نیمه پارامتری از آن جهت به این مدل‌ها اطلاق می‌گردد که شکل توزیع $f(Y | \theta_i)$ برای محقق معلوم می‌باشد و به عبارت دیگر این توزیع بخش پارامتری ساختار را شرح می‌دهد و از طرفی $\theta_i \sim G(\cdot)$ ، که در آن G یک توزیع عمومی و آزاد می‌باشد، بخش ناپارامتری ساختار را بیان می‌کند. لذا با تلفیق دو بخش ذکر شده مدل فوق را مدل بیزی سلسله‌مراتبی نیمه پارامتری می‌گوییم.

ایشوارن و جیمز^۲ (۲۰۰۱) توزیع‌های شرطی کامل پارامترهای مدل با اثرات آمیخته را برای انجام استنباط بیزی به دست آوردند. اسپیگل‌هالتر^۳ (۲۰۰۲) معیارهای سنجش اعتبار مدل‌ها از قبیل DIC و پراکندگی را برای مقایسه چند مدل برازش داده‌شده به داده‌ها معرفی و به کار برد. به طوری که مدل با DIC کمتر به بقیه مدل‌ها ارجحیت دارد. اولسن و همکاران^۴ (۲۰۰۶) یک مدل رگرسیون خطی تعمیم‌یافته با اثرات آمیخته را برای بررسی اثرات بیمارستان‌های تحت بررسی بر روی بهبود بیماران به داده‌ها برازش دادند و با تخصیص فرآیند دیریکله به عنوان پیشین به اثرات بیمارستان‌ها با استفاده از روش‌های بیز نیمه پارامتری و نیز با بکارگیری طرح نمونه‌گیرگیز با نرم‌افزار WinBUGS به برآورد و تحقیق در مورد اثرات بیمارستان‌ها پرداختند.

^۱ Kleinma. and Ibrahim (1998)

^۲ Ishwaran and James (2001)

^۳ Spiegelhalter et al (2002)

^۴ Ohlssen et al (2006)

فصل دوم

تعاریف و مفاهیم مقدماتی

۱-۲ مقدمه

در این فصل ابتدا مطالعه‌ای مقدماتی بر مفهوم آمار بیز و استفاده از آن در برآوردیابی پارامترها خواهیم داشت. از این رو برای آشنایی بیشتر با آمار بیز، نخست به بررسی مفاهیم توزیع پیشین و پسین خواهیم پرداخت که به نوعی اساس و بنیاد مبحث آمار بیز می‌باشند، سپس با بررسی توابع زیان، مخاطره و مخاطره بیز به بیان مفهوم برآوردگرهای پسین بیزی و همچنین برآوردگرهای بیزی خواهیم پرداخت. در ادامه به دلیل اینکه روش‌های عددی و تحلیلی برای استنباط و استخراج نمونه به جهت برآورد پارامترهای مجهول، در بسیاری از توزیع‌های پسین پیچیده میسر نمی‌باشد، از این رو یکی از روش‌های متداول مبتنی بر دنباله‌های مونت کالو را که به نمونه‌گیر گیبز معروف است، بیان می‌کنیم. در انتهای این فصل نیز به بیان تعاریف، مفاهیم و قضایایی خواهیم پرداخت که در فصل‌های بعدی از آنها استفاده می‌کنیم. لازم به ذکر است که تمام توابع چگالی احتمال را با نماد $P(\cdot)$ یا $f(\cdot)$ متغیرهای تصادفی را با حروف بزرگ لاتین و بردارها و ماتریس‌ها را با پررنگ کردن حروف نمایش داده‌ایم.

۲-۲ برآوردگرهای پسین نیز^۱

در بررسی مسائل برآورد، اغلب فرض می‌کنیم که نمونه تصادفی از یک تابع چگالی احتمال $P(\cdot; \theta)$ که در آن $P(\cdot; \theta)$ معلوم فرض می‌شود، بدست آمده است. علاوه بر این، فرض می‌کنیم که θ هر چند برای ما نامعلوم است، ثابت می‌باشد. در بسیاری از کاربردهای واقعی که از چگالی $P(\cdot; \theta)$ استفاده می‌شود، اغلب اطلاعات بیشتری درباره θ وجود دارد. برای مثال، پژوهشگر ممکن است شواهدی داشته باشد که خود θ در نقش یک متغیر تصادفی عمل کند و در نتیجه بتواند برای آن یک تابع چگالی احتمال منطقی را در نظر بگیرد. برای مثال، فرض کنید ماشینی که قطعات اتومبیل را تولید می‌کند قرار است از نقطه نظر نسبت معیوب‌های تولید بررسی شود. در یک روز معین، ۱۰ قطعه از خروجی ماشین مورد آزمایش قرار گرفته و مشاهدات را با X_1, X_2, \dots, X_{10} نشان می‌دهیم که در آن اگر قطعه نام معیوب باشد، $X_i = 1$ و اگر سالم باشد $X_i = 0$ که آنها را می‌توان به عنوان یک نمونه تصادفی به حجم ۱۰ از چگالی احتمال برنولی

$$P(X; \theta) = \theta^x (1-\theta)^{1-x}$$

در نظر گرفت که در آن $x = 0, 1$ و $0 \leq \theta \leq 1$. این چگالی احتمال نشان می‌دهد که احتمال اینکه قطعه مفروضی معیوب باشد برابر با عدد نامعلوم θ است. چگالی احتمال توأم ۱۰ متغیر تصادفی X_1, X_2, \dots, X_{10} عبارت است از

$$\theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}$$

می‌دانیم که برآوردگر ماکزیمم درستمایی^۲ θ ، برابر $\hat{\theta} = \bar{X}$ است. روش برآورد گشتاوری^۳ هم همین برآوردگر را نتیجه می‌دهد. با وجود این، فرض کنید که پژوهشگر دارای اطلاعات بیشتری درباره θ است و او مشاهده کرده است که در روزهای مختلف، مقدار θ تغییر می‌کند و چنین می‌نماید که این تغییر را می‌توان به عنوان یک متغیر تصادفی با چگالی احتمال $g(\theta)$ نشان داد. یک سؤال مهم این است که، چگونه می‌توان از این اطلاعات اضافی درباره θ برای برآورد θ_0 استفاده کرد، که θ_0 مقداری است که θ در روز استخراج نمونه با آن برابر بوده است.

¹ Bayesian Posterior Estimators

² Maximum Likelihood Estimator

³ Moment Estimation

برای بررسی این مسأله، علاوه بر این فرض که نمونه تصادفی از چگالی احتمال $P(\cdot; \theta)$ استخراج شده است فرض می‌کنیم که پارامتر نامعلوم θ مقداری از یک متغیر تصادفی، مانند Θ است. باز هم علاقه‌مند به برآورد کردن تابعی از θ ، مانند $\tau(\theta)$ خواهیم بود. اگر Θ متغیری تصادفی باشد پس دارای یک توزیع است. فرض می‌کنیم $G(\cdot) = G_\theta(\cdot)$ تابع توزیع تجمعی Θ و $g(\cdot) = g_\theta(\cdot)$ تابع چگالی احتمال Θ را نشان می‌دهد و این توابع دارای پارامترهای مجهول نیستند.

اگر فرض کنیم توزیع Θ معلوم است، دارای اطلاعات اضافی می‌باشیم. بنابراین یک سؤال مهم این است که چگونه این اطلاعات اضافی را در برآورد به کار برد؟ این مسئله‌ای است که در ادامه به آن می‌پردازیم. در بسیاری از مسائل، این فرض که θ مقداری از یک متغیر تصادفی است، ممکن است غیر واقعی باشد. در سایر مسائل، گرچه فرض اینکه θ مقداری از یک متغیر تصادفی است مناسب به نظر می‌رسد، ممکن است توزیع Θ معلوم نباشد، یا حتی اگر معلوم باشد، ممکن است شامل پارامترهای نامعلوم دیگری باشد. با این وجود، در برخی مسائل فرض معلوم بودن توزیع Θ واقع بینانه است.

۲-۱-۲ توزیع‌های پیشین و پسین^۱

اکثراً برای هر θ در فضایی که به آن متعلق است، نماد $P(X; \theta)$ را برای نشان دادن تابع چگالی احتمال متغیر X به کار می‌بریم. هرگاه بخواهیم نشان دهیم که پارامتر θ مقداری از متغیر تصادفی Θ است، چگالی X را به جای $P(X; \theta)$ به صورت $P(X|\theta)$ خواهیم نوشت و ممکن است به جای P از f برای نشان دادن تابع چگالی احتمال استفاده کنیم. باید توجه کنیم که $P(X|\theta)$ یک چگالی شرطی است، یعنی چگالی X به شرط $\Theta = \theta$ است.

فرض کنید X_1, X_2, \dots, X_n یک نمونه تصادفی به حجم n از چگالی $P(\cdot|\theta)$ باشد، که در آن θ مقدار متغیر تصادفی Θ می‌باشد. فرض کنید چگالی Θ ، $g(\cdot) = g_\theta(\cdot)$ ، معلوم بوده و دارای هیچ پارامتر مجهولی نیست و فرض کنید می‌خواهیم $\tau(\theta)$ را برآورد کنیم. چگونه باید اطلاعات اضافی معلوم بودن $g_\theta(\cdot)$ را در روش‌های برآورد وارد کنیم؟ تابع درستمایی را به عنوان عبارتی که همه اطلاعات را در بر دارد، می‌شناسیم. عبارت تابع درستمایی، نمونه x_1, x_2, \dots, x_n و همچنین شکل چگالی $P(X; \theta)$ که از آن نمونه گرفتیم را در برداشت. اکنون عبارتی لازم است که علاوه بر دارا بودن تمام اطلاعات موجود در تابع درستمایی، اطلاعات اضافی معلوم بودن چگالی $g_\theta(\cdot)$ را نیز دارا باشد. $g_\theta(\cdot)$ توزیع پیشین Θ نامیده می‌شود. این چگالی آنچه را که ما پیش از

^۱ Prior and Posterior Distribution