

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه اصفهان

دانشکده علوم

گروه آمار

پایان نامه‌ی کارشناسی ارشد رشته‌ی آمار گرایش آمار ریاضی

مباحثی در استفاده از توزیع لاپلاس - چوله چندمتغیره در مدل‌های رگرسیونی

استاد راهنما:

دکتر ایرج کاظمی

پژوهشگر:

راضیه محمدی

مهرماه ۱۳۹۰

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

پایان‌نامه‌ی کارشناسی ارشد رشته‌ی آمار گرایش آمار ریاضی

خانم راضیه محمدی

تحت عنوان

مباحثی در استفاده از توزیع لاپلاس - چوله چند متغیره در مدل های رگرسیونی

در تاریخ ۹۰/۷/۲۷ توسط هیأت داوران زیر با درجه عالی به تصویب نهایی رسید.

- ۱- استاد راهنمای پایان‌نامه دکتر ایرج کاظمی با مرتبه‌ی علمی استادیار
- ۲- استاد داور داخل گروه دکتر محمد بهرامی با مرتبه‌ی علمی استادیار
- ۳- استاد داور خارج از گروه دکتر مجید جعفری خالدی با مرتبه‌ی علمی استادیار
- امضاء
- امضاء
- امضاء

امضای مدیر گروه

تشر و قدردانی

شکر و سپاس یزدان پاک را که در لحظه لحظه‌ی زندگی یگانه یار و راهنمای ماست. چه زیباست سخن
بزرگ‌مرد تاریخ، سعدی شیرازی که

از دست و زبان که برآید کز عهده‌ی شکرش به در آید

و سپاس فروان از استادم، جناب آقای دکتر ایرج کاظمی که مرا در انجام این پایان‌نامه بسیار یاری رساندند. سپاس
برای تمام اوقاتی که به من اختصاص دادند و مرا از راهنمایی‌های ارزشمند خود بهره‌مند کردند. امید آن دارم که
کیفیت این مجموعه‌ی گردآوری‌شده، رضایت ایشان را به‌دنبال داشته باشد و همچنین آرزو می‌کنم در تمامی
مراحل زندگی، موفق و سربلند باشند.

چکیده

در تحلیل مدل‌های رگرسیونی، فرض معمول این است که مؤلفه‌های خطا دارای توزیع نرمال هستند. اما از آن جا که امکان نقض این فرض در تحلیل داده‌های واقعی امکان‌پذیر است، لذا مطالعه بر روی جایگزینی توزیع‌های منعطف‌تر از نرمال موضوع اصلی بسیاری از تحقیقات کاربردی می‌باشد. یکی از این توزیع‌ها لاپلاس - چوله‌ی چندمتغیره است که در دهه‌ی اخیر از سوی محققان زیادی مورد توجه قرار گرفته است. در تحقیقات مختلف، شکل‌های متفاوتی برای نشان دادن تابع چگالی آن ارائه شده است. در این پایان‌نامه، یکی از این شکل‌ها که نسبت به بقیه دارای ساختار ساده‌تر و کاربرد بهتر در مدل‌سازی است بررسی می‌شود. این توزیع به دلیل وجود انعطاف‌پذیری بیشتر که شامل سنگینی دم‌ها، کشیدگی زیاد و چولگی است می‌تواند جایگزین مناسبی در مباحث مدل‌سازی رگرسیونی باشد. در این پایان‌نامه ابتدا مدل‌های رگرسیون خطی معمولی با توزیع لاپلاس - چوله معرفی خواهند شد و سپس آن‌ها را به مدل‌های آمیخته با اثرهای تصادفی، که در تحلیل داده‌های همبسته به وفور به کار می‌روند، تعمیم می‌دهیم. با توجه به آنکه استنباط پارامترهای مدل بر اساس روش حداکثر درست‌نمایی حاشیه‌ای منجر به محاسبات جبری پیچیده‌ای می‌شود ما از روش‌های عددی پیشرفته مانند الگوریتم EM از دیدگاه فراوانی‌گرا و نیز رهیافت مونت کارلوی زنجیر مارکوفی از دیدگاه بیزی بهره می‌گیریم. در این راستا از نمایش تصادفی سلسله‌مراتبی توزیع لاپلاس - چوله‌ی چندمتغیره که آمیخته‌ای از نرمال و گاماست استفاده خواهیم کرد. اهمیت نظری نتایج حاصل را با ارائه‌ی تحلیل داده‌های واقعی نشان می‌دهیم.

واژه‌های کلیدی: استنباط بیزی، الگوریتم EM، نمونه‌گیری گیبز، حداکثر درست‌نمایی، مدل رگرسیونی با اثرهای آمیخته، نمایش تصادفی سلسله‌مراتبی.

فهرست مطالب

صفحه

عنوان

فصل اول: مقدمه

- ۱-۱ موضوع و پیشینه‌ی تحقیق ۱
- ۲-۱ اهداف تحقیق ۳
- ۳-۱ اهمیت و کاربرد نتیجه‌های تحقیق ۳
- ۴-۱ جنبه‌های محاسباتی ۴
- ۵-۱ ساختار پایان‌نامه ۴

فصل دوم: کلیات

- ۱-۲ مقدمه ۵
- ۲-۲ آمار بیز ۶
- ۱-۲-۲ توزیع پیشین ۶
- ۳-۲ روش‌های مونت کارلوی زنجیر مارکوفی ۷
- ۱-۳-۲ روش مونت کارلو ۸
- ۲-۳-۲ زنجیر مارکوف ۹
- ۴-۲ الگوریتم متروپلیس - هستینگس ۱۰
- ۵-۲ نمونه‌گیری گیبز ۱۱
- ۶-۲ معیارهای تشخیص همگرایی در روش‌های مونت کارلوی زنجیر مارکوفی ۱۲
- ۷-۲ مقایسه‌ی مدل‌ها با استفاده از معیارهای اطلاع آکائیک و اطلاع انحراف ۱۴
- ۸-۲ پیش‌بینی در آمار بیز ۱۵
- ۱-۸-۲ برآورد پیش‌بینی پسین ترتیبی و پیش‌بینی شرطی ترتیبی ۱۷
- ۹-۲ مقایسه‌ی مدل‌ها با استفاده از تابع چگالی پیش‌بینی پسین ۱۸
- ۱۰-۲ روش داده‌افزایی ۱۸
- ۱۱-۲ مدل‌های سلسله‌مراتبی ۱۸

۱۲-۲ الگوریتم EM	۱۹
۱-۱۲-۲ کاربردهای دیگر الگوریتم EM	۲۱
۱۳-۲ معرفی چند توزیع	۲۱

فصل سوم: توزیع لاپلاس - چوله و کاربرد آن در مدل‌های رگرسیونی

۱-۳ مقدمه	۲۴
۲-۳ توزیع لاپلاس استاندارد	۲۵
۳-۳ توزیع لاپلاس - چوله‌ی کاتز و همکاران (۲۰۰۱)	۲۵
۱-۳-۳ تابع چگالی توزیع لاپلاس - چوله	۲۶
۲-۳-۳ ویژگی‌های توزیع	۲۷
۳-۳-۳ نمایش‌های تصادفی توزیع لاپلاس - چوله	۲۹
۴-۳ توزیع لاپلاس - چوله‌ی یو و زانگ (۲۰۰۵)	۳۱
۱-۴-۳ ویژگی‌های توزیع	۳۲
۵-۳ توزیع لاپلاس - چوله‌ی ارسلان (۲۰۰۹)	۳۳
۱-۵-۳ ویژگی‌های توزیع	۳۴
۶-۳ کاربرد توزیع لاپلاس - چوله در مدل‌های رگرسیونی	۳۵
۱-۶-۳ برآوردیابی مدل رگرسیونی لاپلاس - چوله با استفاده از الگوریتم EM	۳۶
۲-۶-۳ برآوردیابی مدل رگرسیونی لاپلاس - چوله از دیدگاه آمار بیز	۳۸
۳-۶-۳ الگوریتم نمونه‌گیر گیبز	۴۲
۷-۳ مثال کاربردی	۴۳
۱-۷-۳ برازش مدل لاپلاس - چوله	۴۶
۲-۷-۳ مقایسه‌ی مدل‌های لاپلاس - چوله و نرمال	۴۹
۸-۳ نتیجه‌گیری	۴۹

فصل چهارم: توزیع لاپلاس - چوله‌ی چندمتغیره و کاربرد آن

۱-۴ مقدمه	۵۰
۲-۴ معرفی توزیع لاپلاس - چوله‌ی چندمتغیره	۵۱
۳-۴ نمایش تصادفی توزیع لاپلاس - چوله‌ی چندمتغیره	۵۲
۴-۴ ویژگی‌های توزیع لاپلاس - چوله‌ی چندمتغیره	۵۳
۵-۴ برآورد حداکثر درست‌نمایی	۵۵
۶-۴ مثال کاربردی	۵۷

فصل پنجم: مدل‌های چندسطحی

۱-۵ مقدمه	۵۹
۲-۵ مدل‌های چندسطحی	۶۰
۳-۵ مدل عرض از مبدأ تصادفی	۶۱
۱-۳-۵ برآورد حداکثر درست‌نمایی پارامترهای مدل عرض از مبدأ تصادفی با فرض نرمال	۶۳
۴-۵ مدل دوسطحی با ضرایب تصادفی	۶۶
۵-۵ تحلیل بیزی مدل رگرسیونی عرض از مبدأ تصادفی با توزیع لاپلاس - چوله	۶۸
۱-۵-۵ توزیع پیشین پارامترهای مدل	۶۹
۲-۵-۵ توزیع پسین توام	۶۹
۳-۵-۵ توزیع‌های پسین شرطی کامل	۷۰
۶-۵ برازش مدل رگرسیونی عرض از مبدأ تصادفی با توزیع لاپلاس - چوله‌ی چندمتغیره	۷۴
۱-۶-۵ برآورد حداکثر درست‌نمایی پارامترهای مدل	۷۵
۷-۵ تحلیل بیزی مدل رگرسیون خطی با اثرهای آمیخته با توزیع لاپلاس - چوله‌ی چندمتغیره	۷۸
۱-۷-۵ توزیع پیشین پارامترهای مدل	۷۹
۲-۷-۵ توزیع پسین توام	۷۹
۳-۷-۵ توزیع‌های پسین شرطی کامل	۸۰
۸-۵ مثال کاربردی (بدهی‌های مالیاتی)	۸۳

۸۴.....	۱-۸-۵ برازش مدل رگرسیونی عرض از مبدأ تصادفی نرمال
۸۵.....	۲-۸-۵ بررسی درستی فرض‌های مدل
۸۷.....	۳-۸-۵ برازش مدل رگرسیونی عرض از مبدأ تصادفی لاپلاس - چوله
۸۹.....	۹-۵ بحث و نتیجه‌گیری
۹۰.....	منابع و مأخذ

فهرست شکل‌ها

صفحه	عنوان
۲۷	شکل ۳-۱ تابع چگالی لاپلاس - چوله ی کاتز
۳۱	شکل ۳-۲ تابع چگالی لاپلاس - چوله ی یو و زانگ
۳۳	شکل ۳-۳ تابع چگالی لاپلاس - چوله ی ارسلان
۴۴	شکل ۳-۴ بافت‌نگار مانده‌ها
۴۵	شکل ۳-۵ نمودار مانده‌ها در برابر مقادیر برازش شده
۴۶	شکل ۳-۶ نمودار احتمال نرمال مانده‌های استاندارد شده
۴۷	شکل ۳-۷ نمودار سری زمانی تاریخچه پارامتر γ
۴۸	شکل ۳-۸ نمودار خودهمبستگی نگار پارامترهای σ^2 و γ
۴۸	شکل ۳-۹ نمودار هسته‌ی چگالی پارامترهای β_2 و β_5
۴۹	شکل ۳-۱۰ نمودار گل‌من - رایین پارامترهای σ^2 و γ
۵۲	شکل ۴-۱ تابع چگالی توزیع لاپلاس - چوله‌ی دومتغیره با $\gamma = (0.5, 0.5)'$
۵۲	شکل ۴-۲ تابع چگالی توزیع لاپلاس - چوله‌ی دومتغیره با $\gamma = (-0.5, -0.5)'$
۸۶	شکل ۵-۱ بافت‌نگار مانده‌های سطح اول
۸۶	شکل ۵-۲ بافت‌نگار مانده‌های سطح دوم
۸۶	شکل ۵-۳ نمودار مانده‌های سطح اول در برابر مقادیر برازش شده
۸۶	شکل ۵-۴ نمودار مانده‌های سطح دوم در برابر مقادیر برازش شده

فهرست جدول‌ها

صفحه	عنوان
۴۴	جدول ۱-۳ برآورد پارامترهای مدل نرمال
۴۵	جدول ۲-۳ آزمون شاپیرو-ویلک برای نرمال بودن مانده‌های مدل
۴۷	جدول ۳-۳ برآورد بیز پارامترهای مدل لاپلاس - چوله
۴۹	جدول ۴-۳ مقادیر معیار اطلاع انحراف دو مدل نرمال و لاپلاس - چوله
۵۸	جدول ۱-۴ برآورد پارامترهای توزیع لاپلاس - چوله و نرمال
۵۸	جدول ۲-۴ مقادیر معیار اطلاع انحراف دو مدل نرمال و لاپلاس - چوله
۸۵	جدول ۱-۵ برآورد پارامترهای مدل عرض از مبدأ تصادفی با توزیع نرمال
۸۶	جدول ۲-۵ آزمون شاپیرو-ویلک برای نرمال بودن خطاهای مدل
۸۷	جدول ۳-۵ مقادیر $p, D(\bar{\theta})$ و AIC چهار مدل مختلف
۸۸	جدول ۴-۵ برآورد بیز پارامترهای مدل با خطاهای لاپلاس - چوله و اثرهای تصادفی نرمال

فصل اول

مقدمه

۱-۱ موضوع و پیشینه‌ی تحقیق

در تحلیل مدل‌های رگرسیونی فرض معمول این است که توزیع متغیر پاسخ به شرط متغیرهای توضیحی نرمال است. اما از آنجا که در بسیاری از مثال‌های تجربی، توزیع داده‌ها از توزیع نرمال انحراف زیادی دارد، لذا استفاده از توزیع‌های جایگزین دیگر، از زمان‌های پیش مورد توجه آماردانان قرار گرفته است. بعضی از این توزیع‌های جایگزین به جهت داشتن ویژگی‌های خاصی چون کشیدگی بیشتر و دم‌های سنگین‌تر و این که قابلیت تحت پوشش قرار دادن مشاهدات دورافتاده را دارند، به توزیع‌های مقاوم^۱ معروف هستند. یکی از این توزیع‌ها، لاپلاس است که آمیخته‌ی نامتناهی از توزیع‌های نرمال و نمایی است (کاتز و همکاران^۲ ۲۰۰۱). این توزیع که علاوه بر داشتن دم‌های سنگین‌تر، نسبت به نرمال از کشیدگی بیشتری نیز برخوردار است، در استنباط‌های مقاوم اهمیت فراوانی دارد. کاربردهایی از این توزیع در مباحث امور مالی، علوم اقتصادی و زیستی یافت می‌شود. نظر به این که در تحلیل داده‌های واقعی، مثال‌هایی نیز وجود دارند که توزیع مشاهدات متقارن نیست، لذا مطالعه بر روی تعمیم‌هایی از توزیع لاپلاس به حالت چوله‌ی آن نیز مورد توجه محققان قرار گرفته است. از بین آن‌ها، کاتز و همکاران (۲۰۰۱) توزیع لاپلاس - چوله‌ی تعمیم‌یافته را به صورت زیر تعریف می‌کنند

^۱ Robust

^۲ Kotz

تعریف: فرض کنید متغیر تصادفی Z دارای توزیع نرمال p متغیره با میانگین صفر و واریانس I_p و متغیر تصادفی نامنفی V که مستقل از Z است، دارای توزیع گاما با پارامترهای θ و 1 باشد. در این صورت متغیر تصادفی Y با نمایش تصادفی

$$Y = \mu + V\gamma + \sqrt{V}\Sigma^{1/2}Z$$

که در آن $\mu \in \mathbb{R}^p$ ، $\gamma \in \mathbb{R}^p$ ، Σ ماتریس معین مثبت و $\Sigma^{1/2}$ جذر Σ است، دارای توزیع لاپلاس - چوله‌ی تعمیم‌یافته است. در تعریف فوق، به متغیر تصادفی V که موجب تصادفی کردن میانگین و واریانس توزیع نرمال گردیده است، متغیر آمیختگی^۱ گفته می‌شود. همچنین به پارامتر γ که نقش چوله کردن توزیع لاپلاس را بر عهده دارد، پارامتر چولگی^۲ می‌گویند (اگر $\gamma = 0$ آن‌گاه توزیع لاپلاس تعمیم‌یافته حاصل می‌شود). کاتز و همکاران (۲۰۰۱) نشان دادند که تابع چگالی متناظر با توزیع لاپلاس - چوله‌ی تعمیم‌یافته دارای یک نمایش پیچیده بر حسب تابع بسل تعدیل‌شده‌ی نوع سوم^۳ می‌باشد. در حالت خاص، اگر $\theta = 1$ آن‌گاه توزیع لاپلاس - چوله‌ی چندمتغیره به دست می‌آید که ویژگی‌های آن از سوی محققان زیادی مطالعه شده است. از جمله کاتز و همکاران (۲۰۰۱) ضمن بررسی ویژگی‌های این توزیع، به محاسبه‌ی برآورد حداکثر درستمایی پارامترهای آن در حالت یک‌متغیره پرداخته‌اند. ویسک^۴ (۲۰۰۹) موضوع برآوردیابی پارامترهای توزیع فوق را با استفاده از روش گشتاوری مورد بررسی قرار داده است. کاربردهایی از توزیع لاپلاس - چوله‌ی چندمتغیره در امور مالی و علوم زیستی توسط کزوبسکی و پدگرسکی^۵ (۲۰۰۱)، لیندسی و لیندسی^۶ (۲۰۰۶) مطالعه شده است.

ارسلان^۷ (۲۰۰۹) نشان داده است که اگر در تعریف فوق $\theta = \frac{p+1}{2}$ باشد (p بعد متغیر تصادفی است) آن‌گاه تابع چگالی لاپلاس - چوله‌ی تعمیم‌یافته شکل ساده‌تری پیدا می‌کند که می‌تواند در مباحث مدل‌سازی مفید واقع شود. وی همچنین با استفاده از الگوریتم EM، برآورد حداکثر درستمایی پارامترها را به دست آورده و الگوریتم شبیه‌سازی ساده‌ای برای این توزیع ارائه داده است.

توزیع لاپلاس - چوله به لحاظ انعطاف‌پذیری زیادی که نسبت به نرمال دارد، می‌تواند در مدل‌های رگرسیونی به‌عنوان جایگزینی مقاوم مورد استفاده قرار گیرد. با توجه به پیچیدگی محاسبات در ارتباط با مسئله‌ی برآوردیابی

¹ Mixing variable

² Skewness parameter

³ Modified Bessel function of the third kind

⁴ Visk

⁵ Kozubowski and Podgorski

⁶ Lindsey and Lindsey

⁷ Arslan

پارامترها در مدل‌های فوق، روش‌های جایگزین همچون نمونه‌گیری گیز که در تحلیل‌های بیزی مورد استفاده قرار می‌گیرند، در این پایان‌نامه در نظر گرفته می‌شود. در این راستا با توجه به از نمایش سلسله‌مراتبی^۱ مدل‌های فوق، از روش داده‌افزایی^۲ که توسط تنر و ونگ^۳ (۱۹۸۷) معرفی گردید و موجب سهولت در انجام محاسبات آماری می‌شود، استفاده می‌کنیم.

یکی از مدل‌های معروف آماری که برای تحلیل داده‌های وابسته و آشیانه‌ای به کار می‌رود، مدل چندسطحی^۴ است (فریز^۵ ۲۰۰۴). از آنجا که این مدل‌ها به دلیل وجود اثرهای تصادفی دارای ساختار پیچیده‌تری هستند، لذا استنباط آن‌ها بر مبنای توزیع‌های مقاوم، توسط روش‌های معمول کلاسیک به آسانی امکان‌پذیر نمی‌باشد. از این‌رو استفاده از روش‌های بیزی در مطالعه‌ی مدل‌های چندسطحی، به موضوع اصلی بسیاری از تحقیقات کاربردی تبدیل شده است. از جمله آرلانو- واله^۶ و همکاران (۲۰۰۷) به مطالعه‌ی مدل رگرسیون خطی با اثرهای آمیخته^۷ با استفاده از توزیع نرمال - چوله پرداختند. در این پایان‌نامه نیز مدل رگرسیون خطی با اثرهای آمیخته با توزیع لاپلاس - چوله‌ی چندمتغیره معرفی می‌شود.

۲-۱ اهداف تحقیق

هدف از انجام این پایان‌نامه، ضمن معرفی توزیع لاپلاس - چوله و مروری بر شکل‌های متفاوتی از این توزیع که در دهه‌ی اخیر معرفی شده است، بررسی کاربرد این توزیع در مدل‌های رگرسیونی می‌باشد. همچنین با توجه به کاربرد گسترده‌ی مدل‌های چندسطحی در تحلیل داده‌ها با ساختار سلسله‌مراتبی، در این پایان‌نامه به جای توزیع نرمال که به‌طور متداول به کار می‌رود، توزیع لاپلاس - چوله‌ی چندمتغیره را منظور کرده و استنباط پارامترهای این مدل‌ها را با استفاده از دو دیدگاه کلاسیک و بیزی انجام می‌دهیم.

۳-۱ اهمیت و کاربرد نتیجه‌های تحقیق

با توجه به اهمیت مطالعه‌ی توزیع‌های مقاوم در مباحث مدل‌سازی، نتایج حاصل از این پایان‌نامه که شامل معرفی مدل‌های رگرسیونی با توزیع لاپلاس - چوله‌ی چندمتغیره می‌باشد، می‌تواند در زمینه‌ی تحلیل داده‌های تجربی

¹ Hierarchical representation

² Data augmentation

³ Tanner and Wong

⁴ Multilevel model

⁵ Frees

⁶ Arellano-Valle

⁷ Linear mixed- effect model

مفید واقع شود. از سوی دیگر با توجه به فراگیر شدن کاربرد روش‌های بیزی مبتنی بر روش‌های مونت کارلوی زنجیر مارکوف^۱ (MCMC) در دهه‌های اخیر، که موجب سهولت انجام محاسبات آماری گردیده است، در این پایان‌نامه نیز از روش‌های فوق بهره برده‌ایم. همچنین دو مثال کاربردی در این پایان‌نامه ارائه می‌شوند تا اهمیت استفاده از مدل‌های رگرسیونی با توزیع مذکور نشان داده شود.

۴-۱ جنبه‌های محاسباتی

در این پایان‌نامه برای رسم توابع چگالی مختلف، نرم‌افزار Maple به کار رفته است. همچنین جهت برازش مدل رگرسیونی معمولی از نرم‌افزار SPSS و جهت برازش مدل رگرسیون خطی با اثرهای آمیخته از نرم‌افزار Stata استفاده شده است. برآوردیابی حداکثر درستنمایی پارامترهای توزیع لاپلاس - چوله‌ی چندمتغیره توسط نرم‌افزار SAS (فصل چهارم را ببینید) و محاسبات مربوط به تحلیل بیزی مثال‌های ارائه شده در پایان‌نامه توسط نرم‌افزار وین‌باگز^۲ انجام شده است.

۵-۱ ساختار پایان‌نامه

در فصل دوم به معرفی آمار بیز و مباحث مرتبط با آن می‌پردازیم. همچنین الگوریتم EM را که یک روش عددی در محاسبه‌ی برآورد حداکثر درستنمایی محسوب می‌شود، در این فصل شرح می‌دهیم. در فصل سوم توزیع لاپلاس - چوله را معرفی می‌کنیم و از این توزیع در برازش مدل‌های رگرسیونی استفاده می‌کنیم. در فصل چهارم به معرفی توزیع لاپلاس - چوله‌ی چندمتغیره پرداخته و موضوع برآوردیابی پارامترهای آن را توسط روش حداکثر درستنمایی بررسی می‌کنیم. در فصل پنجم مدل‌های چندسطحی را معرفی می‌کنیم و از توزیع لاپلاس - چوله‌ی چندمتغیره جهت تحلیل این مدل‌ها استفاده می‌کنیم. در این راستا مدل رگرسیونی عرض از مبدأ تصادفی^۳ تصادفی^۳ را با فرض آن که اثرها و مؤلفه‌های خطا دارای توزیع لاپلاس - چوله‌ی چندمتغیره باشند، معرفی کرده و استنباط بیزی پارامترهای این مدل را با استفاده از رهیافت مونت کارلوی زنجیر مارکوفی انجام می‌دهیم. علاوه بر آن، این موضوع را از دیدگاه فراوانی‌گرا، توسط روش EM، بررسی می‌کنیم. در پایان، مطالعه‌ای بر برازش مدل رگرسیونی خطی با اثرهای آمیخته، توسط توزیع لاپلاس - چوله‌ی چندمتغیره، انجام خواهیم داد. قابل توجه آن که در کل پایان‌نامه، نماد * به منظور نشان دادن نتایج مستخرج از پایان‌نامه به کار رفته است.

^۱ Markov Chain Monte Carlo

^۲ WinBUGS

^۳ Random intercept model

فصل دوم

کلیات

۲-۱ مقدمه

در این فصل ابتدا به معرفی آمار بیز و روش‌های مونت کارلوی زنجیر مارکوف که از معروف‌ترین روش‌های محاسباتی در آمار بیز محسوب می‌شوند، می‌پردازیم. سپس الگوریتم متروپلیس - هستینگس^۱ و نمونه‌گیری گیبز^۲ را به‌عنوان دو روش مهم از روش‌های مونت کارلوی زنجیر مارکوف معرفی می‌کنیم. همچنین مسئله‌ی همگرایی و روش‌های تشخیص آن را که از مباحث مهم در استفاده از این روش‌ها هستند، بررسی می‌کنیم. در ادامه به موضوع پیش‌بینی در مدل‌های بیزی و همچنین معرفی معیارهایی جهت مقایسه‌ی مدل‌ها می‌پردازیم. در پایان فصل نیز الگوریتم EM را به‌عنوان یک روش برآوردیابی عددی معرفی می‌کنیم که در فصل‌های بعد از آن استفاده خواهیم کرد.

¹ Metropolis-Hastings

² Gibbs Sampling

۲-۲ آمار بیز

تفاوت اصلی آمار بیز با آمار کلاسیک این است که در آمار بیز، پارامترهای مجهول نیز به عنوان کمیت‌های تصادفی در نظر گرفته می‌شوند و توزیع آماری مناسبی به آن‌ها تخصیص داده می‌شود که به آن توزیع پیشین^۱ می‌گویند. این توزیع، بیانگر اعتقاد محقق در مورد پارامتر قبل از جمع‌آوری داده‌هاست. در این حالت استنباط آماری بر مبنای قضیه‌ی بیز انجام می‌گیرد. بر این اساس، اطلاعات موجود در توزیع پیشین و اطلاعات حاصل از جمع‌آوری داده‌ها با هم ترکیب شده و توزیع جدیدی حاصل می‌شود که به آن توزیع پسین^۲ می‌گویند. به بیان دیگر، ادغام اطلاعات پیشین و کنونی موجب به‌روز شدن اطلاعات محقق از پارامتر می‌گردد و در نتیجه توزیع پسین حاصل می‌شود که استنباط بر مبنای آن انجام می‌گیرد.

فرض کنید متغیر تصادفی Y دارای تابع چگالی $f(y|\theta)$ باشد که در آن متغیر تصادفی θ دارای توزیع پیشین $\pi(\theta)$ است. اگر Y_1, \dots, Y_n را یک نمونه‌ی تصادفی از $f(y|\theta)$ در نظر بگیریم و تابع درستنمایی را با $L(\theta|y)$ نشان دهیم داریم

$$L(\theta|y) = f(y_1, \dots, y_n) = \prod_{i=1}^n f(y_i|\theta). \quad (۱-۲)$$

در این صورت با توجه به قضیه‌ی بیز، توزیع پسین به صورت

$$f(\theta|y) = \frac{L(\theta|y)\pi(\theta)}{f(y)} \quad (۲-۲)$$

به دست می‌آید که می‌توان آن را متناسب با حاصل ضرب تابع درستنمایی در توزیع پیشین به صورت

$$f(\theta|y) \propto L(\theta|y)\pi(\theta) \quad (۳-۲)$$

در نظر گرفت. با استفاده از توزیع پسین می‌توان بر آورد نقطه‌ای پارامتر θ را محاسبه کرد که متداول‌ترین آن، میانگین این توزیع است.

۱-۲-۲ توزیع پیشین

انتخاب توزیع پیشین به دلیل آن‌که توزیع پسین را تحت تأثیر قرار می‌دهد، اهمیت زیادی دارد. در این قسمت به معرفی انواع توزیع‌های پیشین می‌پردازیم.

^۱ Prior distribution

^۲ Posterior distribution

۱- توزیع پیشین ناآگاهی بخش^۱

در اغلب مثال‌های کاربردی که هیچ‌گونه اطلاعاتی از پارامتر در دست نیست، توزیع پیشین را طوری در نظر می‌گیرند که توزیع پسین را تحت تأثیر قرار ندهد. چنین توزیع‌هایی را توزیع پیشین مبهم^۲ یا ناآگاهی بخش می‌گویند (نتزوفراس^۳ ۲۰۰۹).

یک توزیع پیشین ناآگاهی بخش می‌تواند به صورت

$$\pi(\theta) \propto 1$$

روی فضای پارامتر در نظر گرفته شود که در این صورت به آن توزیع یکنواخت یا مسطح^۴ می‌گویند. در این حالت توزیع پسین با تابع درستنمایی معادل می‌شود و در نتیجه برآورد مد توزیع پسین با برآورد حداکثر درستنمایی در حالت کلاسیک یکسان می‌شود. همچنین اگر توزیع پیشین طوری در نظر گرفته شود که

$$\int \pi(\theta) d\theta = \infty$$

آن‌گاه به آن، توزیع ناسره^۵ می‌گویند. توزیع‌های ناسره در صورتی قابل استفاده هستند که منجر به یک توزیع پسین سره شوند.

جزئیات بیشتر در مورد توزیع‌های پیشین ناآگاهی بخش در یانگ و برگر^۶ (۱۹۹۶) موجود است.

۲- توزیع پیشین آگاهی بخش

به توزیع‌های پیشینی که حاوی اطلاعات پیشین هستند و توزیع پسین را تحت تأثیر قرار می‌دهند پیشین آگاهی بخش می‌گویند. به‌عنوان یک مثال از این توزیع‌ها می‌توان توزیع پیشین مزدوج^۷ را در نظر گرفت. هرگاه توزیع‌های پیشین و پسین هر دو از یک خانواده باشند، توزیع پیشین را مزدوج می‌گویند (نتزوفراس ۲۰۰۹).

۲-۳ روش‌های مونت کارلوی زنجیر مارکوفی

نظر به این که محاسبه‌ی توزیع پسین و برآورد بیز همیشه به سادگی امکان‌پذیر نمی‌باشد و در اغلب موارد مستلزم حل انتگرال‌های پیچیده است، لذا در این مواقع باید از روش‌های عددی استفاده شود. یکی از روش‌های

¹ Noninformative

² Vague

³ Ntzoufras

⁴ Flat

⁵ Improper

⁶ Yang and Berger

⁷ Conjugate