

الله
كريم



دانشگاه اصفهان

دانشکده علوم

گروه آمار

پایان نامه‌ی کارشناسی ارشد رشته‌ی آمار گرایش آمار ریاضی

توزیع‌های آمیخته – مقیاس نرمال – چوله و کاربرد آن‌ها در مدل‌های رگرسیونی غیرخطی

استاد راهنما:

دکتر ایرج کاظمی

پژوهشگر:

مهرسا عابدینی

آبان ماه ۱۳۹۱

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

پایان نامه کارشناسی ارشد رشته آمار گرایش آمار ریاضی
برخاسته شده است



دانشگاه امیرکبیر
دانشکده علوم
گروه آمار

پایان نامه‌ی کارشناسی ارشد رشته‌ی آمار گرایش آمار ریاضی

خانم مهسا عابدینی

تحت عنوان

کاربرد توزیع‌های آمیخته - مقیاس فرمال - چوله در مدل‌های رگرسیونی ثیروخطی

در تاریخ ۹۱/۸/۲۸ توسط هیأت داوران زیر بررسی با درجه عالی به تصویب نهایی رسید.

امضاء
امضاء
امضاء

۱- استاد راهنمای پایان نامه دکتر ایرج کاظمی با مرتبه‌ی علمی استادیار

۲- استاد داخل پایان نامه دکتر محمد بهرامی با مرتبه‌ی علمی استادیار

۳- استاد داور خارج از گروه دکتر سروش علیمرادی با مرتبه‌ی علمی استادیار

امضای مدیر گروه

خدای را بسی شاکرم که از روی کرم، پدر و مادری فداکار به من عطا فرمود تا در سایه درخت پربار وجودشان بیاسایم و از ریشه آنها شاخ و برگ گیرم و از سایه‌ی وجودشان در راه کسب علم و دانش قدم بردارم، والدینی که وجودشان تاج افتخاری است بر سرم و نامشان دلیلی است بر بودنم چرا که این دو وجود، پس از پروردگار مایه‌ی هستی‌ام بوده‌اند، دستم را گرفته‌اند و راه رفتن را در این وادی زندگی پر فراز و

نشیب به من آموخته‌اند؛ اساتیدی که برایم زندگی کردن، بودن و انسان بودن را معنا

کردند و حال این برگ سبزی است تحفه‌ی درویش که تقدیم می‌کنم به

پدر و مادر عزیزم و خانواده‌ی مهربانم

به نام یگانه خالق هستی بخش جان آفرین که در ذات احديتش شکی نیست، به عون و قوهی الهی اين پروژه را با راهنمایی های جناب آقای دکتر کاظمی به پایان رساندم و بنا به فرمودهی مولای متقيان علی(ع): «من علّمنی حرفًا فقد يسّرني عبداً» همواره افتخار شاگردی ايشان را با خود خواهم داشت و کمال تشکر را از ايشان دارم.

چکیده

یکی از مباحث مهم در تحلیل رگرسیون غیرخطی، فرض درباره توزیع خطای مدل است که در چند دهه اخیر موضوع تحقیق بسیاری از محققان بوده است. پیش از این، فرض بر این بوده که خطا از توزیع نرمال پیروی می‌کند. بسیاری از محققان توزیع‌های جدیدی را با خواص منحصر به فردی برای توزیع خطا در نظر گرفته‌اند، از جمله می‌توان به خانواده توزیع آمیخته-مقیاس چوله-نرمال اشاره نمود. مبحث مهم دیگر در تحلیل رگرسیون غیرخطی، تعیین شکل تابع انتظار است که با توجه به ساختار داده‌های واقعی مشخص می‌شود.

هدف از انجام این پایان‌نامه، معرفی خانواده توزیع آمیخته مقیاس چوله-نرمال و توزیع‌های چوله-نرمال، چوله-تی و چوله-اسلش متعلق به آن و بررسی خواص مطلوب آن‌ها در حالت کلی و در رگرسیون غیرخطی است. از خواص مطلوب این خانواده از توزیع‌ها می‌توان به نمایش سلسله مراتبی آن که امکان محاسبه ساده‌تر برآورد پارامترها به روش الگوریتم EM ، تعمیم‌های آن و رهیافت نمونه‌گیر گیز را فراهم می‌کند، اشاره کرد. با برآش دادن مدل رگرسیون غیرخطی مناسب با توزیع خطای مختلف متعلق به خانواده توزیع آمیخته-مقیاس چوله-نرمال بر داده‌های واقعی به این نتیجه رسیدیم که توزیع‌های چوله در مواردی بهتر از توزیع نرمال به داده‌ها برآش می‌شود. با انجام مطالعات شبیه‌سازی از یک مدل رگرسیون غیرخطی با خطای چوله-تی و مقایسه با سایر توزیع‌های معرفی شده برای خطا، قابلیت توزیع آمیخته مقیاس چوله-نرمال در برآش مدل‌های رگرسیون غیرخطی مشخص می‌شود.

کلید واژه: شکل سلسله مراتبی، نمونه‌گیر گیز، الگوریتم EM ، توزیع‌های نامتقارن، توزیع‌های دم-سنگین.

فهرست مطالب

عنوان	صفحه
فصل اول: مقدمه	
۱- موضوع و پیشینه تحقیق	۱
۲- اهمیت و اهداف تحقیق	۳
۳- روش‌های محاسباتی	۳
۴- ساختار پایان‌نامه	۳
فصل دوم: مفاهیم و تعاریف مقدماتی	
۱-۲ مقدمه	۵
۲-۲ تابع مشخصه	۶
۳-۲ خانواده توزیع چوله-متقارن	۷
۴-۲ توزیع SN	۸
۴-۲ توزیع SMN	۱۲
۴-۲ توزیع تی-استیودنت	۱۴
۴-۲ توزیع اسلش	۱۶
۵-۲ الگوریتم EM	۱۷
۶-۲ استنباط بیزی	۲۰
۶-۲ توزیع پیشین	۲۱
۶-۲ توزیع پسین	۲۳
۶-۲ برآورد بیزی	۲۳
۷-۲ رهیافت $McMC$	۲۵
۷-۲ زنجیر مارکف	۲۶
۷-۲ انتگرال گیری مونت کارلو	۲۶
۷-۲ نمونه‌گیری مهم	۲۶
۷-۲ الگوریتم متropolیس-هستینگر	۲۷
۷-۲ نمونه‌گیری گیبز	۲۸

۲۹.....	۶-۷-۲ تشخیص همگرایی الگوریتم
۳۰.....	۸-۲ معیار انتخاب مدل
فصل سوم: خانواده توزیع‌های SMSN	
۳۲.....	۱-۳ مقدمه
۳۳.....	۲-۳ خانواده توزیع <i>SMSN</i>
۳۷.....	۱-۲-۳ تابع مشخصه
۳۹.....	۲-۲-۳ برآورد پارامترها توسط الگوریتم <i>ECME</i>
۴۲.....	۳-۳ توزیع چوله-تی (<i>ST</i>)
۵۰.....	۱-۳-۳ تابع مشخصه
۵۲.....	۲-۳-۳ برآورد پارامترها توسط الگوریتم <i>ECM</i>
۵۵.....	۴-۳ : توزیع چوله-اسلش (<i>SSL</i>)
۶۰.....	۱-۴-۳ تابع مشخصه
۶۲.....	۲-۴-۳ برآورد پارامترها توسط الگوریتم <i>ECM</i>
فصل چهارم: رگرسیون غیرخطی	
۶۵.....	۱-۴ مقدمه
۶۶.....	۲-۴ مدل رگرسیون خطی
۶۷.....	۳-۴ مدل رگرسیون غیرخطی
۶۹.....	۴-۴ تابع درستنمایی
۶۹.....	۵-۴ برآورد پارامترها
۷۱.....	۱-۵-۴ روش تکرار گاوس-نیوتون
۷۵.....	۴-۶ آزمون فرض و فاصله اطمینان
۷۹.....	۷-۴ محاسبات بیزی
۷۹.....	۱-۷-۴ چگالی پیشین
۸۰.....	۲-۷-۴ چگالی پسین
فصل پنجم: کاربرد توزیع‌های SMSN در رگرسیون غیرخطی	
۸۲.....	۱-۵ مقدمه
۸۳.....	۲-۵ مدل رگرسیون غیرخطی با خطای <i>SMSN</i>

۱-۲-۵ تابع درستنمایی	۸۴
۲-۲-۵ رهیافت بیز	۸۴
۳-۲-۵ الگوریتم <i>ECM</i>	۸۶
۳-۵ مدل رگرسیون غیرخطی با خطای <i>ST</i>	۸۹
۱-۳-۵ تابع درستنمایی	۸۹
۲-۳-۵ رهیافت بیز	۹۰
۳-۳-۵ الگوریتم <i>ECM</i>	۹۱
۴-۵ مدل رگرسیون غیرخطی با خطای <i>SSL</i>	۹۴
۱-۴-۵ تابع درستنمایی	۹۵
۲-۴-۵ رهیافت بیز	۹۶
۳-۴-۵ الگوریتم <i>ECM</i>	۹۷
۵-۵ تحلیل داده‌های فارماکوکینتیکز	۱۰۰
۶-۵ مطالعه شبیه‌سازی	۱۰۶
۷-۵ بحث و نتیجه‌گیری	۱۱۰
پیوست ۱	۱۱۱
پیوست ۲	۱۱۳
واژه‌نامه فارسی-انگلیسی	۱۱۵
واژه‌نامه انگلیسی-فارسی	۱۲۰
منابع	۱۲۵

فهرست جداول‌ها

صفحه	عنوان
۷۳	جدول ۴-۱: نتایج محاسبه عددی روش تکرار گاوس-نیوتون
۷۴	جدول ۴-۲: نتایج روش تکرار گاوس-نیوتون
۷۴	جدول ۴-۳: برآورد پارامتر به روش گاوس-نیوتون
۷۶	جدول ۴-۴: نتایج تحلیل مدل ۶-۴
۷۷	جدول ۴-۵: برآورد پارامترهای مدل ۶-۴
۷۷	جدول ۴-۶: نتایج تحلیل مدل ۷-۴
۷۸	جدول ۴-۷: نتایج تحلیل مدل ۸-۴
۷۸	جدول ۴-۸: نتایج تحلیل مدل ۹-۴
۷۹	جدول ۴-۹: برآورد پارامترهای مدل ۹-۴
۱۰۱	جدول ۵-۱: نتایج تحلیل مدل ۴۲-۵
۱۰۱	جدول ۵-۲: برآورد پارامترهای مدل ۴۲-۵
۱۰۲	جدول ۵-۳: نتایج تحلیل مدل ۴۳-۵
۱۰۴	جدول ۵-۴: برآوردهای بیز پارامترها با توزیع خطای نرمال و چوله-نرمال
۱۰۴	جدول ۵-۵: برآوردهای بیز پارامترها با توزیع خطای تی و چوله-تی
۱۰۵	جدول ۵-۶: برآوردهای بیز پارامترها با توزیع خطای اسلش و چوله-اسلش
۱۰۵	جدول ۵-۷: معیارهای انتخاب مدل دادهای فارماکوکینتیکز
۱۰۷	جدول ۵-۸: برآوردهای بیز و MSE پارامترها با توزیع خطای نرمال و چوله-نرمال در مطالعه شبیه‌سازی
۱۰۸	جدول ۵-۹: برآوردهای بیز و MSE پارامترها با توزیع خطای تی و چوله-تی در مطالعه شبیه‌سازی
۱۰۹	جدول ۵-۱۰: برآوردهای بیز و MSE پارامترها با توزیع خطای اسلش و چوله-اسلش در مطالعه شبیه‌سازی

فهرست شکل‌ها

صفحه	عنوان
۱۰	شکل ۲-۱: توزیع‌های نرمال و چوله-نرمال
۱۵	شکل ۲-۲: توزیع‌های نرمال و تی-استیودنت
۱۷	شکل ۲-۳: توزیع‌های نرمال و اسلش
۴۹	شکل ۳-۱: توزیع‌های چوله-تی و تی بریده شده
۵۲	شکل ۳-۲: توزیع چوله-تی
۶۰	شکل ۳-۳: توزیع‌های چوله-اسلش و اسلش بریده شده
۶۲	شکل ۳-۴: توزیع چوله-اسلش
۱۰۳	شکل ۳-۵: هیستوگرام باقیمانده‌ها

مخفف یا کوتاه نوشت

E	امیدریاضی
$ease$	برآورد خطای استاندارد مجانبی
\approx	تقریب
\propto	تناسب
\sim	توزیع
ase	خطای استاندارد مجانبی
$diag$	ماتریس قطری
obs	مشاهده
MSE	میانگین توان دوم خطای
Var	واریانس
$\underline{\underline{d}}$	هم توزیع

فصل اول

مقدمه

۱- موضوع و پیشینه تحقیق

در بسیاری از آزمایش‌های بالینی در علوم پزشکی و یا تحلیل داده‌های وابسته و پانلی که عموماً از مطالعه‌های اقتصادی گردآوری شده‌اند مشاهده می‌شود که مدل رگرسیونی مناسب برای داده‌ها مدلی غیرخطی است. بنابراین در این موارد اهمیت مدل رگرسیونی غیرخطی آشکار می‌شود. گاهی نیز در تحلیل رگرسیونی غیرخطی، محقق با این مسئله مواجه می‌شود که ساختار داده‌ها متقارن نیست و یا دارای مشاهده‌های دورافتاده است. بنابراین برای این مدل‌های غیرخطی در نظر گرفتن توزیع خطای نرمال مناسب نیست و توزیع‌های متعلق به خانواده توزیع *SMSN* انتخابی می‌تواند مناسب باشد. با توجه به خواص مطلوب این خانواده از توزیع‌ها، پژوهش‌های جدید بسیاری بر مبنای خواص این خانواده بنا شده‌اند. در این پایان‌نامه خانواده توزیع آمیخته-مقیاس چوله-نرمال^۱ (*SMSN*) که تعمیمی از توزیع چوله-نرمال^۲ (*SN*) (آزالینی^۳، ۱۹۸۵) است و دو روش برآوردهایی

¹ Scale mixture of skew normal

² Skew normal

³ Azzalini

تکراری عددی تحت عنوان‌های الگوریتم ماکسیمم سازی امید^۱ (EM) (دempster و همکاران^۲، ۱۹۷۷) و رهیافت مونت کارلوی زنجیر مارکفی^۳ ($MCMC$) برای پارامترهای آن معرفی می‌شود. همچنین پس از یادآوری مختصر مدل رگرسیون خطی، به بررسی مدل رگرسیون غیرخطی پرداخته می‌شود. در نهایت نیز با در نظر گرفتن توزیع خطای $SMSN$ در مدل رگرسیون غیرخطی بخشی از کاربرد این خانواده توزیع نشان داده می‌شود.

در سال ۲۰۰۱ برانکو و دی^۴ استفاده از خانواده توزیع $SMSN$ را برای تحلیل مشاهده‌های نامتقارن^۵ پیشنهاد دادند. حالت خاصی از توزیع‌های فوق، برای مدل‌های غیرخطی توسط لین و همکاران^۶ (۲۰۰۹) ارائه شد. لاجس و همکاران^۷ (۲۰۱۰) با شبیه‌سازی داده‌های چوله از توزیع‌های دم-سنگین^۸، کاربرد این توزیع‌ها را نشان داده‌اند.

با سو و همکاران^۹ (۲۰۱۰) نیز توزیع‌های $SMSN$ را در مدل‌های آمیخته پیشنهاد دادند. مدل پیشنهادی آن‌ها به طور همزمان عدم تقارن توزیع داده‌ها و تأثیر مشاهده‌های دورافتاده^{۱۰} را در نظر می‌گیرد. گری و همکاران^{۱۱} (۲۰۱۱) برآش مدل‌های رگرسیونی غیرخطی را توسط توزیع‌های $SMSN$ و با به کارگیری الگوریتم EM انجام دادند. آنان همچنین با مطالعه‌های شبیه‌سازی این مدل‌ها برای داده‌هایی با ساختار پخش نامتقارن نتیجه گرفتند که استنباط آماری پارامترها توسط این توزیع‌های آمیخته-مقیاس در برابر مشاهدات دورافتاده نسبت به چوله-نرمال استوارتر است. همچنین استفاده از روش‌های بیزی مانند $MCMC$ را در برآش مدل‌های فوق پیشنهاد دادند. کانچو و همکاران^{۱۲} (۲۰۱۱) با استفاده از رهیافت $MCMC$ ، کاربردی از خانواده توزیع $SMSN$ را در مدل‌های غیر خطی ارائه دادند. این مدل‌ها با لحاظ کردن چولگی و دم-سنگین بودن در تحلیل داده‌ها مدل‌های بسیاری، از جمله چوله-تی^{۱۳} (ST) و چوله-اسلش^{۱۴} (SSL)، را پوشش می‌دهند. آن‌ها با استفاده از روش‌های بیزی بر اساس معیار کولبک-لیبر به تشخیص داده‌های پرنفوذ پرداختند. با تحلیل داده‌های واقعی نشان دادند که مدل ST متعلق به خانواده توزیع $SMSN$ برآزنده‌تر از مدل SN است.

^۱ Expectation maximization algorithm

^۲ Dempster et al.

^۳ Markov chain

^۴ Branco and Dey

^۵ Asymmetric

^۶ Lin et al.

^۷ Lachos et al.

^۸ Heavy-tailed

^۹ Basso et al.

^{۱۰} Outlier

^{۱۱} Garay et al.

^{۱۲} Cancho et al.

^{۱۳} Skew t

^{۱۴} Skew slash

۱-۲ اهمیت و اهداف تحقیق

در تحلیل داده‌هایی که مدل‌های رگرسیون با فرض نرمال بودن مؤلفه‌های خطا به آنها خوب برازش نمی‌شوند اهمیت استفاده از توزیع‌های دیگر، از جمله توزیع‌های $SMSN$ ، روشن است. اهمیت ویژه خانواده توزیع‌های معرفی شده در کاربرد آنها جهت مدل‌سازی داده‌هایی با ساختار نامتقارن و دارای مشاهده‌های دورافتاده یا پرنفوذ است. نمایش سلسله مراتبی مطلوب توزیع‌های $SMSN$ امکان برآوردهایی توسط روش‌های EM و $MCMC$ برای پارامترهای این توزیع‌ها را فراهم می‌کند.

هدف از انجام این پایان‌نامه، بررسی توزیع‌های $SMSN$ و به دست آوردن خواص احتمالی و استباطی آنها در مدل‌سازی آماری است. با توجه به خواص مطلوب این توزیع‌ها توجه بیشتری به کاربرد آنها در برازش مدل‌های خطی و غیر خطی جهت برطرف کردن مشکلات مطرح شده خواهد شد. در این رابطه مدل‌های فوق بر مجموعه داده‌های واقعی برازش داده می‌شوند و با استفاده از معیارهای انتخاب مدل مناسب، درستنمایی توزیع‌های متعلق به خانواده توزیع $SMSN$ جهت حصول نتایج معتبر مقایسه می‌شوند. در این زمینه جهت برآورد پارامترها روش‌های تحلیل بیزی، به ویژه رهیافت نمونه‌گیرگیز^۱، استفاده خواهد شد. به علاوه برای شناخت بهتر خواص این توزیع‌ها در مدل‌های غیرخطی، شبیه‌سازی از آنها صورت می‌گیرد و نتایج به دست آمده تحلیل می‌شوند.

۱-۳ روش‌های محاسباتی

برازش مدل‌های رگرسیونی و تحلیل آنها با استفاده از نرم‌افزارهای نسخه sas 9.2 و اپنباگز^2 نسخه 3.2.1 انجام شده است. با استفاده از نرم‌افزار $maple$ نمودار توزیع‌ها رسم شده‌اند. برای برآورد بیزی پارامترهای مدل از نرم‌افزار اپنباگز و برای سایر تحلیل‌ها از دستورهای $proc nlin$ و یا $proc nlmixed$ از نرم‌افزار sas استفاده شده است.

۱-۴ ساختار پایان‌نامه

ساختار این پایان‌نامه به این صورت است که در فصل دوم مباحثی که در فصل‌های بعد مورد استفاده قرار می‌گیرند مطرح می‌شوند. در فصل سوم خانواده توزیع $SMSN$ و توزیع‌های SSL و ST متعلق به آن معرفی و

¹ Gibbs sampler

² Open Bugs

پارامترهای آنها از طریق روش‌های گیز مبتنی بر $McMC$ و تعمیم‌های الگوریتم EM برآورد می‌شوند. فصل چهارم در بردارنده مطالبی در خصوص مدل رگرسیون غیرخطی است. در نهایت در فصل پنجم کاربرد توزیع‌های متعلق به خانواده توزیع $SMSN$ در مدل غیرخطی بررسی می‌شود. در انتهای این فصل نیز داده‌های واقعی را تحلیل خواهیم کرد.

فصل دوم

مفاهیم و تعاریف مقدماتی

۱-۲ مقدمه

در این فصل تعاریف و مفاهیم به کار گرفته شده در این پایان نامه را شرح می‌دهیم. خانواده توزیع $SMSN$ (برانکو و دی، ۲۰۰۱) که در این پایان نامه مورد بررسی قرار می‌گیرد آمیخته-مقیاسی از چوله-نرمال است که برای معرفی این خانواده توزیع ابتدا آشنایی مختصری با خانواده توزیع‌های چوله-متقارن (ونگ و همکاران^۱، ۲۰۰۴) و آمیخته-مقیاس نرمال (SMN) (آندره و مالوز^۲، ۱۹۷۴) ارائه می‌شود. از طرف دیگر، برای محاسبه برآورد پارامترها نیاز به معرفی روش‌های تکراری-عددی مانند الگوریتم EM (دمپستر و همکاران، ۱۹۷۷) مبتنی بر رهیافت فراوانی‌گرا و روش‌های $MCMC$ از جمله متropolis-هستینگز^۳ (هستینگز، ۱۹۷۰) و گیز (گمن و گمن^۴، ۱۹۸۰) بر اساس رهیافت بیز است.

¹ Wang et al.

² Andrew and Mallows

³ Metropolis-Hastings

⁴ Geman and Geman