

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه شهید چمران اهواز

دانشگاه شهید چمران اهواز

دانشکده علوم ریاضی و کامپیوتر

شماره پایان نامه :

۹۲۱۳۶۱۰

پایان نامه کارشناسی ارشد

گرایش آمار ریاضی

عنوان :

مطالعه‌ی برآورد ریج و برخی مباحث تشخیصی در مدل‌های رگرسیونی تحت محدودیت‌های

خطی تصادفی

استاد راهنما:

دکتر عبدالرحمن راسخ

استاد مشاور:

دکتر بهزاد منصوری

نگارنده :

نرگس هدایت‌پور

دی ماه سال ۱۳۹۲

چکیده

نام خانوادگی: هدایت پور	نام: نرگس	شماره دانشجویی: ۹۰۱۳۶۱۱
عنوان پایان نامه: مطالعه‌ی برآورد ريج و برخی مباحث تشخیصی در مدل‌های رگرسیونی تحت محدودیت- های خطی تصادفی		
استاد راهنما: دکتر عبدالرحمن راسخ		
استاد مشاور: دکتر بهزاد منصوری		
درجه تحصیلی: کارشناسی ارشد	رشته: آمار	گرایش: آمار ریاضی
دانشگاه: شهید چمران اهواز	دانشکده: علوم ریاضی و کامپیوتر	گروه: آمار
تاریخ فارغ التحصیلی: ۱۳۹۲/۱۰/۳۰		تعداد صفحه: ۱۳۵
کلید واژه ها: هم خطی، برآورد ريج، برآورد ريج تحت محدودیت‌های خطی، مباحث تشخیصی، مشاهدات موثر، مشاهدات پرت.		
<p>مطالعه‌ی مباحث تشخیصی در مدل‌های مختلف رگرسیونی دارای سابقه‌ی طولانی است. گاهی مجموعه‌ی کوچکی از داده‌ها اثرات نامتناسبی را بر روی نتایج حاصل از آنالیز رگرسیونی اعمال می‌کنند. به طوری که برآورد پارامترها یا مقادیر پیش‌بینی شده بیشتر تحت نفوذ این مشاهدات قرار می‌گیرند. شناسایی چنین مشاهداتی از طریق روش‌ها و معیارهایی که مبتنی بر آنالیز تأثیر است امکان‌پذیر بوده که در مقوله‌ی مباحث تشخیصی تجلی پیدا می‌کند. باهدف تعیین مشاهدات مؤثر و پرت معیارهایی همچون ماتریس پیش‌بینی، باقیمانده‌های استاندارد، $DFBETAS$، $DFFITSS$، فاصله‌ی کوک و روش انتقال میانگین پیشنهاد شده است. از سوی دیگر وجود هم‌خطی در میان متغیرهای مستقل پیامدهای نامطلوبی به همراه دارد. این پدیده روش کمترین توان‌های دوم را به چالش کشیده و موجب عدم کارایی برآورد حاصل می‌گردد. راه‌بردهایی همچون به‌کارگیری برآورد اریب ريج، برآورد آمیخته و ترکیبی از این دو تحت عنوان برآورد ريج آمیخته باهدف اصلاح مشکل هم‌خطی پیشنهاد شده است. حضور توأم هم‌خطی و مشاهدات مؤثر در مجموعه‌ی داده‌ها امر نامعمول و غیرمنطقی نیست بلکه یک موضوع پیچیده است. در این رساله ضمن مطالعه‌ی رگرسیون ريج و رگرسیون ريج تحت محدودیت‌های خطی تصادفی اقدام به بررسی روش‌ها و معیارهای تشخیصی مذکور جهت تعیین مشاهدات تأثیرگذار و پرت در روش رگرسیون ريج می‌کنیم. در پی آن این معیارها را به روش رگرسیونی ريج تحت محدودیت‌های خطی تصادفی تعمیم داده و از این طریق مشاهدات تأثیرگذار و پرت را تشخیص می‌دهیم. در نهایت به منظور شفاف سازی و توضیح نتایج حاصل شده یک مثال عددی که طی آن مجموعه داده‌های واقعی سیمان پرتلند را مورد بررسی قرار داده‌ایم ذکر می‌کنیم.</p>		

فهرست

عنوان	صفحه
فصل ۱	مقدمه و تاریخچه‌ای از پژوهش‌های پیشین..... ۱
۱-۱	مقدمه..... ۱
۲-۱	تاریخچه‌ی پژوهش‌های پیشین..... ۴
۱-۲-۱	برآوردگرهای اریب..... ۵
۲-۲-۱	محدودیت‌های خطی..... ۶
۳-۲-۱	ترکیب برآوردهای اریب و محدودیت‌های خطی..... ۷
۴-۲-۱	آنالیز تشخیصی..... ۹
۳-۱	ساختار پایان‌نامه..... ۱۰
فصل ۲	هم‌خطی و برآورد ریج..... ۱۲
۱-۲	مقدمه..... ۱۲
۲-۲	هم‌خطی..... ۱۳
۱-۲-۲	آثار هم‌خطی..... ۱۴
۲-۲-۲	روش‌های تشخیص هم‌خطی..... ۱۶
۱-۲-۲-۲	ماتریس همبستگی (R)..... ۱۶
۲-۲-۲-۲	عامل تورم واریانس (VIF)..... ۱۷
۳-۲-۲-۲	تحلیل سیستم مقادیر ویژه $X'X$ ۱۷
۴-۲-۲-۲	تجزیه‌ی مقادیر منفرد..... ۱۸

۱۹تجزیه‌ی واریانس ضرایب رگرسیونی
۲۱۳-۲ رگرسیون ریج
۲۴۱-۳-۲ ویژگی‌های برآورد ریج
۲۶۲-۳-۲ روش‌هایی برای انتخاب k
۲۷۳-۳-۲ مقایسه‌ی برآورد کمترین توان‌های دوم با برآورد ریج
۲۹فصل ۳ برآورد ریج تحت محدودیت‌های خطی
۲۹۱-۳ مقدمه
۳۰۲-۳ مفروضات و اطلاعات تکمیل کننده
۳۱۳-۳ برآورد کمترین توان‌های دوم تحت محدودیت‌های خطی دقیق
۳۳۱-۳-۳ آزمون سازگاری محدودیت‌ها در مدل رگرسیونی
۳۵۴-۳ برآورد کمترین توان‌های دوم تحت محدودیت‌های خطی تصادفی
۳۸۵-۳ هم‌خطی و اطلاعات تکمیل کننده
۳۹۶-۳ برآورد ریج تحت محدودیت‌های خطی دقیق
۴۲۷-۳ برآورد ریج تحت محدودیت‌های خطی تصادفی
۴۴۸-۳ مقایسه برآورد ریج آمیخته با برآورد ریج و برآورد آمیخته
۴۴۱-۸-۳ مقایسه برآورد ریج آمیخته با برآورد ریج
۴۶۲-۸-۳ مقایسه برآورد ریج آمیخته با برآورد آمیخته
۵۰فصل ۴ مباحث تشخیصی در رگرسیون ریج
۵۰۱-۴ مقدمه

- ۵۱ ۲-۴ مشاهدات تأثیرگذار
- ۵۲ ۳-۴ روش‌های تشخیص مشاهدات تأثیرگذار و پرت در رگرسیون کمترین توان‌های دوم
- ۵۲ ۱-۳-۴ ماتریس پیش‌بینی در رگرسیون کمترین توان‌های دوم
- ۵۴ ۲-۳-۴ معرفی انواع باقیمانده‌ها
- ۵۶ ۳-۳-۴ معیار $DFBETA$
- ۵۹ ۴-۳-۴ معیار $DFFITS$
- ۶۰ ۵-۳-۴ معیار فاصله‌ی کوک
- ۶۱ ۶-۳-۴ روش انتقال میانگین
- ۶۳ ۴-۴ هم‌خطی و مشاهدات مؤثر
- ۶۳ ۵-۴ روش‌های تشخیص مشاهدات تأثیرگذار و پرت در رگرسیون ریج
- ۶۴ ۱-۵-۴ ماتریس پیش‌بینی در رگرسیون ریج
- ۶۵ ۲-۵-۴ بردار باقیمانده در رگرسیون ریج
- ۶۷ ۳-۵-۴ معیار $DFBETA^*$
- ۶۸ ۴-۵-۴ معیار $DFFITS^*$
- ۶۹ ۵-۵-۴ فاصله‌ی کوک
- ۷۰ ۶-۵-۴ روش انتقال میانگین رگرسیون ریج
- ۸۲ فصل ۵ مباحث تشخیصی در رگرسیون ریج با محدودیت‌های خطی تصادفی
- ۸۲ ۱-۵ مقدمه
- ۸۳ ۲-۵ روش‌های تشخیص مشاهدات تأثیرگذار و پرت تحت برآورد ریج آمیخته

۸۳	۱-۲-۵ بردار باقیمانده و ماتریس پیش‌بینی در برآورد رنج آمیخته.....
۸۵	۲-۲-۵ معیار $DFBETAS^{**}$
۸۶	۳-۲-۵ معیار $DFFITSS^{**}$
۸۷	۴-۲-۵ معیار فاصله‌ی کوک.....
۸۸	۵-۲-۵ روش انتقال میانگین رگرسیون رنج تحت محدودیت‌های خطی تصادفی.....
۹۸	۳-۵ مثال کاربردی.....
۱۰۰	۱-۳-۵ محاسبه و مقایسه‌ی برآوردها.....
۱۰۴	۲-۳-۵ محاسبه و مقایسه‌ی معیارهای حذف موردی.....
۱۰۹	۳-۳-۵ محاسبه و مقایسه‌ی روش انتقال میانگین.....
۱۱۱	۴-۵ نتیجه‌گیری و پیشنهادات.....
۱۱۶	پیوست الف.....
۱۲۱	پیوست ب.....
۱۲۵	منابع.....
۱۳۳	واژه‌نامه.....

فصل اول

مقدمه و تاریخچه‌ای از پژوهش‌های پیشین

۱-۱ مقدمه

آنالیز رگرسیونی یکی از روش‌های آماری برای تحلیل داده‌ها است. این روش برای بررسی و به مدل در آوردن ارتباط بین متغیرها به کار می‌رود. کاربردهای رگرسیون متعدد است و تقریباً در هر زمینه‌ای از جمله علوم مهندسی، فیزیک، اقتصاد و... کاربرد دارد. تقریباً در همه‌ی کاربردهای رگرسیون معادله رگرسیون فقط یک تقریب از رابطه‌ی درست بین متغیرها است. هدف از تحلیل رگرسیونی و به‌کارگیری مدل رگرسیونی استنباط‌هایی چون مشخص کردن اثرات نسبی متغیرهای مستقل، پیش‌بینی و یا برآورد و انتخاب یک مجموعه مناسب از متغیرها برای مدل است.

زمانی که بین متغیرهای مستقل وابستگی خطی وجود نداشته باشد به سادگی می‌توان استنباط‌های ذکر شده را نتیجه گرفت. ولی متأسفانه اغلب در کاربرد مدل‌های رگرسیونی بین متغیرهای مستقل وابستگی خطی وجود دارد. وجود وابستگی خطی بین متغیرهای مستقل در مدل رگرسیونی تحت عنوان پدیده‌ی هم‌خطی^۱ مطرح می‌شود. در این حالت نتایجی که براساس مدل به دست خواهیم آورد ممکن است گمراه‌کننده باشد. وجود هم‌خطی در مدل مشکلاتی چون ناپایداری و عدم کیفیت در برآوردهای آماری را به همراه دارد و همچنین دلیلی بر عدم اعتبار برآورد کمترین توان-

^۱ Collinearity

های دوم است. به این ترتیب برای برخورد با اثرات نامطلوب هم‌خطی راه‌های مختلفی چون جمع‌آوری داده اضافی، تخصیص مدل مجدد و استفاده از روش‌هایی متفاوت با روش کمترین توان‌های دوم برای برآورد ضرایب رگرسیونی ارائه شده است. با توجه به اینکه پدیده‌ی هم‌خطی روش کمترین توان‌های دوم را به چالش می‌کشد، بنابراین برای برآورد ضرایب رگرسیونی روش‌های دیگری پیشنهاد شده است. برآورد حاصل از این روش‌ها بر خلاف روش کمترین توان‌های دوم اریب است. پس با نادیده گرفتن ضرورت نااریبی برآوردی به دست خواهیم آورد که واریانس آن کمتر از واریانس برآورد کمترین توان‌های دوم است. به کارگیری کلاس برآوردگرهای اریب یکی از با اهمیت‌ترین روش‌ها جهت برآورد ضرایب تحت شرایط حضور هم‌خطی و با هدف رفع این مشکل محسوب می‌گردد. از جمله این برآوردها می‌توان به برآورد ریج^۱، برآورد استاین^۲ و برآورد لیو^۳ اشاره کرد.

راه‌کار دیگری که برای اصلاح اثرات ایجاد شده‌ی ناشی از پدیده‌ی هم‌خطی ارائه شده است، استفاده از اطلاعات کمکی^۴ و یا پیشین^۵ علاوه بر اطلاعات نمونه در قالب محدودیت‌های خطی^۶ بر روی ضرایب رگرسیونی است. اعمال محدودیت‌های خطی یعنی اضافه کردن این اطلاعات به مجموعه‌ی داده‌ها باعث بهبود برآورد و به دنبال آن تقلیل اثر هم‌خطی می‌گردد. برآورد ضرایب در این شرایط تحت عنوان برآورد آمیخته^۷ معرفی می‌شود. قابل ذکر است که برآورد استاین، ریج و لیو حالت خاصی از برآورد آمیخته می‌باشند.

¹ Ridge Estimator

² Stein Estimator

³ Liu Estimator

⁴ Auxiliary information

⁵ Prior

⁶ Linear Restriction

⁷ Mixed Estimator

باید توجه داشت که می‌توان از ترکیب دو روش مذکور یعنی کلاس برآوردهای اریب همراه با اعمال محدودیت‌های خطی برای رفع مشکل هم‌خطی بهره جست. از مقایسه‌ی معیار ماتریس میانگین توان دوم خطای^۱ برآورد حاصل از ادغام دو روش با دیگر برآوردها در می‌یابیم که این برآورد عمل‌کرد بهتری دارد.

از سوی دیگر مطالعه مباحث تشخیصی^۲ در مدل‌های مختلف رگرسیونی سابقه‌ی طولانی دارد. بررسی روش‌هایی چون یافتن مشاهدات پرت^۳ و مشاهدات تأثیرگذار^۴ از جمله این روش‌ها می‌باشند. با توجه به تأثیری که مشاهدات تأثیرگذار بر نتایج حاصل از تحلیل رگرسیونی دارند؛ بنابراین باید به دقت مشاهدات را مورد بررسی قرار داده و میزان تأثیرگذاری آن‌ها را تعیین کرد. آماردانان و محققان مطالعات گسترده‌ای در این زمینه انجام داده‌اند و برای تشخیص این مشاهدات معیارها و روش‌هایی پیشنهاد کرده‌اند که می‌توان به درایه‌های قطری ماتریس پیش‌بینی^۵ (آماره‌ی نفوذ)، باقیمانده‌های استاندارد، معیار $DFBETAS$ ، $DFFITSS$ ، فاصله‌ی کوک^۶ و روش انتقال میانگین^۷ اشاره کرد.

حضور هم‌زمان هم‌خطی و مشاهدات با اثرات ناروا یکی از مسائل مهم ولی پیچیده است و توجه بسیاری از آماردانان به این موضوع معطوف شده است. همچنین بیان شده که وجود هم‌خطی می‌تواند اثر مشاهدات ناروا را مستور کند. با توجه به جایگاه ویژه‌ی مباحث تشخیصی پژوهش‌گران تحقیقات وسیعی جهت بسط این مباحث در شرایط وجود هم‌خطی در مدل انجام دادند.

¹ Mean Square Error Matrix (MSEM)

² Diagnostic Methods

³ Outliers

⁴ Influential observations

⁵ Prediction matrix

⁶ Cook's Distance

⁷ Mean shift

۲-۱ تاريخچه‌ي پژوهش‌هاي پيشين

آناليز رگرسيوني در مباحث آماري پيشينه‌اي طولاني دارد. در آغاز مفهوم رگرسيون در نظريه‌ي کمترین توان‌هاي دوم که توسط برخي محققين طراحي شده بود، ظهور پيدا کرد. اصطلاح رگرسيون به وسيله‌ي گالتون^۱ (۱۸۷۷؛ ۱۹۸۹) براي توصيف يک مسأله‌ي زيست‌شناختي ابداع گرديد. براي گالتون رگرسيون مفهومي زيست‌شناختي داشت، اما يول^۲ (۱۸۹۷) و پيرسون^۳ (۱۹۰۳) کارهاي او را براي مفاهيم آماري توسعه دادند. بعدها فيشر^۴ (۱۹۲۲) با اشاره به ضعيف بودن فرض‌هاي آنها اين مفهوم را به طور گسترده‌تري مطرح ساخت. امروزه رگرسيون با داشتن نتايج مطلوب يکي از روش‌هاي پرکاربرد جهت آناليز داده‌ها محسوب مي‌گردد.

با توجه به اين موضوع که وجود هم‌خطي در ميان متغيرهاي مستقل نتايج حاصل از تحليل رگرسيوني را تحت تأثير قرار مي‌دهد بايد جهت اصلاح اين مشکل برآييم. مفهوم هم‌خطي، روش‌هاي تشخيص و برخورد با اين مشکل و همچنين پيامدهاي ناشي از آن به‌طور وسيعي توسط مونتگمري و پک^۵ (۱۹۸۲) بحث شده است. راه‌حل‌هاي زيادي براي مقابله با اين مسأله ارائه شده است. توجه به کلاس برآوردگرهاي اريب و بهره‌مندی از آنها يکي از راه‌حل‌هاي متداول براي مقابله با اين مشکل است.

¹ Galton

² Yule

³ Pearson

⁴ Fisher

⁵ Montgomery and Peck

۱-۲-۱ برآوردگرهاي اريب

استاين^۱ (۱۹۵۶) برآورد اريبي جهت رفع مشكل هم‌خطي معرفي كرد. برآورد اريب ديگري تحت عنوان برآورد ريچ توسط هورل و كنارد^۲ (۱۹۷۰ b,a) مطرح شد. اين برآوردگر يكي از روش‌هاي رايچ به منظور تقليل اثرات هم‌خطي است. به همين جهت جايگاه ويژه‌اي در ميان شيوه‌هاي رفع اين مسأله داشته و توجه بسياري از آماردانان به اين موضوع جلب شده است. دمپستر^۳ (۱۹۷۷) و وينود^۴ (۱۹۷۸) نيز برآورد ريچ را مورد بررسي قرار داده و مطلبي را در اين زمينه ارائه كردند. روش رگرسيون ريچ توسط برخي محققان بسط داده شد و گونه‌هاي مختلف از اين روش ارائه گرديد. هورل و كنارد (۱۹۷۰) گسترشي را در روش رگرسيون ريچ پيشنهاد كردند به اين صورت كه براي هر متغير رگرسيوني پارامتر اريب^۵ جداگانه‌اي را اجازه مي‌دهد. اين روش به عنوان رگرسيون ريچ تعميم يافته^۶ شناخته مي‌شود. رگرسيون ريچ وزني^۷ نيز توسط هولند^۸ (۱۹۷۳) مطرح گرديد. ترنكلر^۹ (۱۹۸۴) برآورد ريچ را براي مدلي با خطاهاي وابسته و غير همگن به دست آورد. اسويندل^{۱۰} (۱۹۷۶) برآورد ريچ اصلاح شده^{۱۱} را با استفاده از يك بردار اطلاع پيشين معرفي كرد و پليسكين^{۱۲} (۱۹۸۷) به مقايسه‌ي اين برآورد با برآورد ريچ پرداخت. برآورد ريچ ناريب^{۱۳} توسط سينگ^{۱۴} (۱۹۸۶) ارائه شد. نسخه‌ي ديگري از برآورد ريچ ناريب توسط كراس و

¹ Stein

² Horel and Kennard

³ Dempster

⁴ Vinod

⁵ Bias parameter

⁶ Generalized Ridge regression

⁷ Weighted Ridge Regression

⁸ Holland

⁹ Trenkler

¹⁰ Swindel

¹¹ Modified Ridge Estimator

¹² Pliskin

¹³ Unbiased Ridge Estimator

¹⁴ Singh

همکاران^۱ (۱۹۹۵) مطرح گردید. به دنبال آن فلاح و سلام^۲ (۲۰۱۱) برآورد ريج ناريب اصلاح شده^۳ را معرفی نموده و در مورد ویژگی‌های آن بحث کردند. مسی^۴ (۱۹۶۵) برآورد اریب مؤلفه-های اصلی^۵ را ارائه کرد. لیو^۶ (۱۹۹۳) برآورد جدیدی با ادغام برآورد ريج و استاین و به منظور ترقی این دو برآورد پیشنهاد داد. آکدنیز و کسیرنلر^۷ (۱۹۹۵) این برآورد را با لفظ برآورد لیو معرفی کردند و ساکالی‌الگو و همکاران^۸ (۲۰۰۱) نیز ماتریس میانگین توان دوم خطای برآورد ريج و لیو را باهم قیاس کرد. همچنين آکدنیز و ارول^۹ (۲۰۰۳) پس از بررسی برخی برآوردگرهای اریب ذکر شده به مقایسه‌ی آن‌ها با یکدیگر اقدام کرد.

۱-۲-۲ محدودیت‌های خطی

گاهی اطلاعاتی راجع به برخی ضرایب رگرسیونی در دسترس است. برخی آماردانان چون بلسلی و همکاران^{۱۰} (۱۹۸۰) و راتو و توتنبرگ^{۱۱} (۱۹۹۵) راتو و همکاران (۲۰۰۸) بیان کردند که اضافه کردن این اطلاعات به نمونه موجب بهبود فرایند برآورد ضرایب می‌شود. همچنین آن‌ها فرض کردند که در طی فرایند، اطلاعات پیشین و مشاهدات، وزن یکسانی را دریافت می‌کنند. در عمل موقعیت‌هایی ممکن است اتفاق بیفتد که چنین فرضی امکان‌پذیر نباشد. سیبر و لی^{۱۲} (۲۰۰۲) و تیل^{۱۳} (۱۹۶۳) آزمون آماری جهت بررسی سازگاری و معنی‌دار بودن این اطلاعات به فرم

¹ Crouse et al

² Fallah and salam

³ Modified Unbiased Ridge Estimator

⁴ Massy

⁵ Principal Components Regression

⁶ Liu

⁷ Akdeniz and Kaciranlar

⁸ Sakallioğlu et al

⁹ Akdeniz and Erol

¹⁰ Belsley et al

¹¹ Rao and Toutenburg

¹² Seber and Lee

¹³ Theil

محدودیت‌های دقیق و یا تصادفی ارائه دادند. آن‌ها بیان کردند که در صورت سازگاری، ترکیبی از اطلاعات نمونه با این اطلاعات به فرم محدودیت‌های خطی برای برآورد ضرایب به‌کار برده می‌شود. گاهی معلومات موجود دقیق نبوده و با اندکی خطا همراه می‌باشند، از طرفی نادیده گرفتن این اطلاعات عملکرد ناصحیحی است. بنابراین تیل و گولدرگر^۱ (۱۹۶۱) و تیل (۱۹۶۳) به‌کارگیری این معلومات به فرم محدودیت‌های خطی تصادفی را پیشنهاد کردند.

۱-۲-۳ ترکیب برآوردهای اریب و محدودیت‌های خطی

راه دیگر پاسخ‌گویی به مسأله‌ی هم‌خطی استفاده از اطلاعات کمکی موجود افزون بر نمونه است. این روش توسط بلسلی و همکاران (۱۹۸۰) مطرح گردید. آن‌ها بیان کردند که استفاده از این اطلاعات به دو شیوه‌ی بیزی و برآورد آمیخته امکان‌پذیر است. استفاده از روش بیزی برای بهبود شرایط حاکم توسط زلنر^۲ (۱۹۷۱) و لیمیر^۳ (۱۹۷۳، ۱۹۷۸) بررسی شده است. مشکل اصلی این روش این است که نیاز به شرح دقیقی از توزیع پیشین دارد. روش برآورد آمیخته مشابه با روش بیزی است. اولین بار دوربین^۴ (۱۹۵۳) به‌طور هم‌زمان اطلاعات پیشین و نمونه را برای بهبود برآورد به‌کار گرفت. سپس برآورد آمیخته به وسیله‌ی تیل و گولدرگر (۱۹۶۱) و تیل (۱۹۶۳) معرفی شد. در این روش اطلاعات کمکی به‌طور مستقیم در قالب محدودیت‌های خطی تصادفی به مجموعه‌ی داده‌ها اضافه می‌شود و این شیوه نیاز به توزیع پیشین ندارد. سارکار^۵ (۱۹۹۲) ایده‌ی ترکیب برآورد ریح و برآورد کمترین توان‌های دوم محدود شده^۶ را ارائه نموده و در پی آن برآورد

^۱ Theil and Goldberger

^۲ Zellner

^۳ Leamer

^۴ Durbin

^۵ Sarkar

^۶ Restricted least-Squares Estimator

ريج محدود شده^۱ را معرفي کرد. وجه تمايز و برتري برآورد ريج اصلاح شده نسبت به برآورد ريج و برآورد ريج محدود شده توسط کسيرنلر و همکاران (۱۹۹۸) و ويجکون^۲ (۱۹۹۸) نشان داده شد. کسيرنلر و همکاران (۱۹۹۹) ديدگاه برآورد ليو اصلاح شده را از طريق متحد کردن برآورد ليو و کمترین توان‌هاي دوم محدود شده ارائه دادند و برآورد ليو محدود شده^۳ را معرفي کردند. گروب^۴ (۲۰۰۳) بيان کرد که برآورد پيشنهادي سارکار (۱۹۹۲) براي کليه محدوديت‌هاي هاي خطي در نظر گرفته شده برآوردي رضايت بخش نيست. به همين منظور پس از مطالعه و بررسي‌هاي لازم برآورد ريج محدود شده‌ي تازه‌اي را که از ادغام روش کمترین توان‌هاي دوم محدود شده با برآورد پيشنهادي اسويندل (۱۹۷۶) حاصل شده بود ارائه داد. هابرت^۵ و جيکون (۲۰۰۶) برآورد ليو تحت محدوديت تصادفي^۶ را با ايجاد اصلاحاتي در برآورد ليو محدود شده پيشنهاده دادند اين برآورد تعميمي از برآورد ليو و آميخته است. ريوان اوزکال^۷ (۲۰۰۹) نظريه‌ي بهره‌مندی هم‌زمان از رگرسيون ريج و محدوديت‌هاي خطي تصادفي را ارائه کرد و همچنين برآورد ريج آميخته^۸ را براي موقعيتي که در آن بردار خطاي مدل همبسته و غير همگن باشد به دست آورد. نسخه‌ي ديگري از برآورد ريج آميخته در سال‌هاي اخير توسط لي و يانگ^۹ (۲۰۱۰) پيشنهاده و مورد مطالعه قرار گرفته است.

¹ Restricted Ridge Estimator

² Wijekoon

³ Restricted Liu Estimator

⁴ Grob

⁵ Hubert

⁶ Stochastic Restricted Liu Estimator

⁷ Revan Ozkale

⁸ Stochastic Restricted Ridge Estimator

⁹ Li and Yang

۱-۲-۴ آناليز تشخيصي

در تحليل رگرسيون مفروضات عمده‌اي در نظر گرفته مي‌شود. براي آزمون مناسب بودن مدل بايد به اعتبار اين مفروضات توجه داشت. انحرافات شديد از مفروضات مي‌تواند به مدل ناپايدار منجر شود. براي تشخيص و برخورد با انحراف از اين مفروضات، روش‌هاي مختلفي پيشنهاد شده است که مجموعه‌ي اين روش‌ها را به عنوان مباحث تشخيصي بيان مي‌کنيم (مونتگمري و پک، ۲۰۰۱). مطالعه‌ي مباحث تشخيصي به ويژه يافتن مشاهدات پرت و تأثيرگذار در مدل‌هاي رگرسيوني از اهميت وافري برخوردار است. نويسندگان و محققان زيادي در اين زمينه به تحقيق و بررسي پرداختند. به‌طور مشخص تعيين مشاهدات مؤثر و پرت توجه بسياري از آماردانان از جمله کوک^۱ (۱۹۷۷)، بلسلي و همکاران (۱۹۸۰)، کوک و ويسبرگ^۲ (۱۹۸۲)، چترجي و هادي^۳ (۱۹۸۶) و (۱۹۸۸) و آتکينسون^۴ (۱۹۸۸) را جلب کرده است. بلسلي و همکاران (۱۹۸۰) بيان کردند که وجود هم‌خطي در مدل مي‌تواند حالت هم‌پوشاني براي اين مشاهدات ايجاد کند. همچنين بيان شد که در صورت پديدار شدن هم‌خطي ممکن است مشاهدات تأثيرگذار در برآورد ريج متمايز با برآورد کمترین توان‌هاي دوم باشد (والکر و بيرچ^۵ ۱۹۸۸). روش‌هاي تشخيص اين مشاهدات براي رگرسيون ريج توسط والکر و بيرچ (۱۹۸۸) ارائه شد. در مقاله‌ي ارائه شده توسط استیک^۶ (۱۹۸۶) ذکر شد که اثر مشاهدات بانفوذ^۷ بر برآورد ريج کمتر است از اثری که همان مشاهدات بر برآورد کمترین توان‌هاي دوم دارند. رويکرد تأثير مکاني^۸ تعميمي از روش‌هاي تشخيصي است که توسط

^۱ Cook

^۲ Weisberg

^۳ Chatterjee and Hadi

^۴ Atkinson

^۵ Walker and Birch

^۶ Steece

^۷ Leverage Observations

^۸ Local Influence

توسط کوک (۱۹۸۶) ارائه شد. شی^۱ (۱۹۹۷) و شی و وانگ^۲ (۱۹۹۹) به ترتیب به بررسی آنالیز تأثیر مکانی در روش مؤلفه‌های اصلی و روش رگرسیون ریح پرداختند. روش‌های تشخیصی موردی^۳ برای رگرسیون ریح اصلاح شده توسط جاهوفر و چن^۴ (۲۰۰۹) بسط داده شد. جاهوفر و چن (۲۰۱۱) به مطالعه‌ی آنالیز تأثیر مکانی در رگرسیون ریح اصلاح شده پرداختند.

۳-۱ ساختار پایان‌نامه

تشخیص مشاهدات تأثیرگذار به دلیل اثرات نامناسب و غیرمعقولی که بر نتایج حاصل از تحلیل رگرسیونی دارند در مباحث آماری از اهمیت ویژه‌ای برخوردار است. بنابراین در این رساله بر مطالعه‌ی مباحث تشخیصی چون یافتن مشاهدات تأثیرگذار و داده‌های پرت در شرایط حضور هم‌خطی و استفاده از رگرسیون ریح تمرکز می‌شود و در ادامه نتایج به دست آمده را برای رگرسیون ریح تحت محدودیت‌های خطی تصادفی تعمیم می‌دهیم. با این هدف در فصل اول مقدمه و تاریخچه‌ای در این زمینه ارائه گردیده است. فصل دوم این رساله، به بیان پدیده‌ی هم‌خطی و مشکلات حاصل از بروز آن می‌پردازد. همچنین مفصلاً روش‌های شناسایی هم‌خطی شرح داده می‌شود. در ادامه فصل، به بیان روش رگرسیون ریح و برآورد حاصل از این روش خواهیم پرداخت، همچنین ویژگی‌های برآورد ریح را برمی‌شماریم. در فصل سوم استفاده از اطلاعات اضافی به فرم محدودیت‌های خطی با جزئیات بیشتر شرح داده می‌شود. به دنبال آن برآورد محدود شده و برآورد آمیخته معرفی می‌گردد. در پی آن نیز برآورد ریح تحت محدودیت دقیق و برآورد ریح تحت محدودیت تصادفی و خواص آن شرح داده می‌شود. فصل چهارم مفهوم مباحث تشخیصی را تشریح کرده و برخی معیارهای تشخیصی را جهت شناسایی نقاط تأثیرگذار

¹ Shi

² Wang

³ Case deletion

⁴ Jahufer and Chen

معرفي مي‌کنيم. در ادامه بررسي‌هاي انجام شده در جهت تخصيص اين معيارها براي رگرسيون ريج را بيان کرده و معيارهاي تشخيصي تعميم يافته را شرح خواهيم داد. فصل پنجم، روش‌هاي تشخيص مشاهدات تأثيرگذار و پرت تحت شرايط استفاده هم‌زمان از رگرسيون ريج و اعمال محدوديت‌هاي خطي تصادفي ارائه شده و همچنين نتايج تئوري بيان شده در اين فصل و فصل -هاي گذشته را با ارائه يک مثال عددي تشریح مي‌کنيم و در پايان نتايج و پيشنهادات را مطرح خواهيم کرد.

هم خطی و برآورد ريج

۱-۲ مقدمه

مدل رگرسیونی خطی بیان کننده‌ی یک رابطه‌ی خطی بین متغیرهای مستقل و متغیر پاسخ است. چه بسا رابطه‌ی قوی تجربی بین دو یا چند متغیر وجود داشته باشد؛ اما دلیلی بر وجود رابطه‌ی علت و معلولی بین متغیرها نیست. تحلیل رگرسیونی می‌تواند در تصدیق و تأیید این رابطه کمک کننده باشد. در تحلیل و به‌کارگیری مدل رگرسیونی برآورد ضرایب یکی از اهداف و استنباط‌های رایج است؛ که محققان به منظور نائل آمدن به این هدف بسیار کوشیده‌اند. یکی از چالش‌های پیش‌رو در این زمینه مشکل هم‌خطی است. بدین منظور در اینجا به طرح این معضل پرداخته و بیان می‌کنیم که وجود هم‌خطی در میان متغیرهای مستقل چه مشکلاتی را به همراه خواهد داشت. همچنین شاخص‌هایی را که برای تشخیص هم‌خطی پیشنهاد شده ارائه می‌دهیم. از طرفی برای رفع مشکل هم‌خطی راه‌حل‌هایی مطرح شده است که در این فصل به بیان یکی از این روش‌ها یعنی به‌کارگیری کلاس برآوردگرهای اریب و به‌طور ویژه برآورد ريج می‌پردازیم. در ادامه ویژگی‌های برآورد حاصل را ذکر کرده، در نهایت به مقایسه‌ی برآورد ريج و برآورد کمترین توان-های دوم می‌پردازیم.

۲-۲ هم خطی

در آغاز پیش از طرح موضوع بنا به ضرورت و پیش برد رساله، یک مدل رگرسیون خطی به عنوان پیش فرض قرار می دهیم؛ بنابراین مدل رگرسیونی زیر را در نظر می گیریم.

$$Y = X\beta + \varepsilon \quad , \quad \varepsilon \sim N(0, \sigma^2 I) \quad (1.2)$$

که در آن Y یک بردار $n \times 1$ از مشاهدات، X یک ماتریس $n \times p$ از متغیرهای مستقل با رتبه p β یک بردار $p \times 1$ از ضرایب رگرسیونی نامعلوم و ε یک بردار $n \times 1$ از خطاهاست. برآورد ضرایب را با استفاده از روش متداول کمترین توان های دوم بر مبنای کمینه کردن مجموع توان دوم خطا به صورت زیر محاسبه می کنیم.

$$\hat{\beta} = (X'X)^{-1}X'Y.$$

در مدل های رگرسیونی با ماتریس متغیرهای مستقل پر رتبه، برآورد ضرایب رگرسیونی به سادگی به دست می آید؛ اما در مواردی علی رغم مستقل بودن این متغیرها به دلیل وابستگی زیاد آنها برآوردهای حاصل قابل اعتماد نخواهند بود و در مواردی می توانند به استنتاج نادرستی منجر شوند.

موقعیتی را که در آن بین متغیرهای مستقل وابستگی خطی وجود داشته باشد تحت عنوان مسأله هم خطی مطرح می کنند. اغلب بین هم خطی کامل و ناقص تمایز وجود دارد. وقتی که برخی از ستون های ماتریس متغیرهای مستقل ارتباط خطی دقیق و کامل داشته باشند آنگاه هم خطی کامل خواهیم داشت. به بیان دیگر می توان هم خطی کامل را به این صورت تعریف کرد که اگر برای زیر مجموعه ای از ستون های X رابطه ی زیر برقرار باشد که در آن t_j مقادیر ثابت غیر صفرند و X_j نشان دهنده ی j امین ستون X است، آنگاه در مدل هم خطی کامل وجود دارد.