



دانشکده علوم ریاضی و کامپیوتر

گروه آمار

پایان نامه کارشناسی ارشد آمار

گرایش محض

عنوان:

تأثیر مشاهدات بر آمارهای آزمون با تمرکز روی مدل‌های رگرسیون خطی

نگارش:

فاطمه قنداق‌ساز

اساتید راهنما:

دکتر عبدالرحمن راسخ

دکتر قاسم قارمست

استاد مشاور:

دکتر صادق رضایی

بسمه تعالی

«چکیده پایان نامه»

نام خانوادگی دانشجو: قنداق‌ساز	نام: فاطمه
عنوان پایان نامه: تأثیر مشاهدات بر آماره های آزمون با تمرکز روی مدل های رگرسیون خطی	
اساتید راهنمای: دکتر عبدالرحمن راسخ و دکتر قاسم تارمیست	
درجه تحصیلی: کارشناسی ارشد	گرایش: محض رشته: آمار
محل تحصیل (دانشگاه): شهید چمران اهواز	
دانشکده: علوم ریاضی و علوم کامپیوتر	
تاریخ فارغ التحصیلی: ۱۳۹۰/۰۴/۰۶	
تعداد صفحات: ۱۳۲	
کلیدواژه‌ها: آماره‌ی نسبت درستنمایی، آنالیز تشخیص حذف موردی، باقیمانده‌های استیودنت شده، تابع تأثیر، ماتریس تصویر.	
چکیده: از جمله مسائلی که در روش‌های رگرسیونی از اهمیت ویژه‌ای برخوردار است، بررسی میزان تأثیری است که هر مشاهده بر هر یک از جنبه‌های مدل‌بندی دارد. اما از آنجا که تأثیر هر مشاهده بر هر یک از جنبه‌های رگرسیونی یکسان نیست، شناسایی مشاهدات مؤثر در هر زمینه مستلزم مطالعه‌ی تأثیر در همان زمینه خاص می‌باشد. هدف اساسی این پایان‌نامه بررسی آنالیز تشخیص حذف موردی برای ارزیابی تأثیر هر مشاهده بر آماره‌های آزمون است. در این راستا ابتدا مفاهیم مورد استفاده در پایان‌نامه را مرور می‌کنیم، سپس به مطالعه‌ی تابع تأثیر آزمون‌ها در مدل‌های رگرسیونی خطی پرداخته و برای این مدل‌ها، دو فرم تابع تأثیر نسبتی و تفاضلی آماره‌ی آزمون را بررسی می‌کنیم. سپس این توابع تأثیر را برای مدل‌های رگرسیونی وزنی به دست آورده و با استفاده از آن‌ها مشاهدات را براساس میزان تأثیری که در آزمون‌ها دارند، رتبه‌بندی و مؤثر بودن هر مشاهده در آماره‌ی آزمون را مورد سنجش و ارزیابی قرار می‌دهیم. سرانجام برای تأیید مباحث نظری ارائه شده، مثال‌هایی از داده‌ای واقعی را مطرح می‌کنیم.	

فهرست مطالب

صفحه	عنوان
------	-------

فصل اول: مقدمه و کلیات

۱	مقدمه.....	۱-۱
۲	مدل رگرسیونی خطی.....	۲-۱
۲	روش‌های آنالیز تشخیص.....	۳-۱
۴	مطالعه تأثیر.....	۴-۱
۵	تاریخچه اولیه.....	۵-۱

فصل دوم: معرفی مدل‌های رگرسیونی

۹	۱-۲ مقدمه.....
۱۰	۲-۲ فرم کلی مدل‌های رگرسیونی خطی.....
۱۲	۳-۲ مدل‌های رگرسیونی خطی.....
۱۲	۱-۳-۲ برآورد پارامترها و خواص آنها.....
۱۵	۲-۳-۲ برآورد پارامترها تحت محدودیت خطی.....
۱۷	۳-۳-۲ آزمون فرض‌های محدودیت خطی.....
۲۰	۴-۲ مدل‌های رگرسیونی خطی وزنی.....
۲۱	۱-۴-۲ برآورد پارامترها و خواص آنها.....
۲۴	۲-۴-۲ برآورد وزنها.....
۲۵	۳-۴-۲ برآورد پارامترهای مدل‌های رگرسیون وزنی تحت محدودیت‌های خطی.....
۲۷	۴-۴-۲ آزمون فرض‌های محدودیت‌های خطی وزنی برای مدل‌های رگرسیونی وزنی.....

فصل سوم: ماتریس تصویر و انواع باقیمانده‌ها

۲۹.....	۱-۳ مقدمه
۳۰	۲-۳ فرض‌های مدل‌بندی و انحرافات موجود در آن‌ها
۳۳.....	۳-۳ ماتریس تصویر
۳۵.....	۱-۳-۳ نقش ماتریس تصویر در تجزیه و تحلیل داده‌ها
۳۶.....	۴-۳ تعریف کلی باقیمانده‌ها
۳۷.....	۳-۵ باقیمانده‌های مدل‌های رگرسیونی خطی
۳۷.....	۱-۵-۳ باقیمانده‌های استیودنت شده
۳۹.....	۲-۵-۳ باقیمانده‌های استیودنت شده تحت وجود محدودیت خطی
۴۰.....	۳-۵-۳ باقیمانده‌های پیش‌بینی
۴۱.....	۴-۵-۳ باقیمانده‌های انحراف
۴۲.....	۳-۶ باقیمانده‌های رگرسیونی وزنی
۴۲.....	۱-۶-۳ باقیمانده‌های استیودنت شده وزنی
۴۳.....	۲-۶-۳ باقیمانده‌های استیودنت شده وزنی تحت وجود محدودیت خطی
۴۴.....	۳-۶-۳ باقیمانده‌های شرطی

فصل چهارم: تابع تأثیر و آزمون محدودیتهای خطی

۴۶.....	۱-۴ مقدمه
۴۷.....	۲-۴ ایده‌ی اصلی آنالیز تشخیص
۴۸.....	۴-۳ تابع تأثیر
۴۹.....	۴-۳-۴ تابع تأثیر در مدل‌های رگرسیونی خطی
۵۰.....	۴-۳-۴ تابع تأثیر تجربی بر اساس n مشاهده
۵۱.....	۴-۳-۴ تابع تأثیر تجربی براساس حذف یک مشاهده

۵۱	تابع تأثیر نمونه‌ای.....۴-۳-۴
۵۲	۴-۴ تابع تأثیر برای براورد پارامترهای مدل‌های رگرسیون خطی.....۴
۵۳	۱-۴-۴ فاصله کوک.....
۵۶	۵-۴ تابع تأثیر نسبتی آزمون فرض محدودیت‌های خطی.....
۵۶	۱-۵-۴ محاسبه تابع تأثیر نسبتی.....

فصل پنجم: تابع تأثیر آماره آزمون‌های نسبت درستنمایی در مدل‌های رگرسیونی

۶۹	۱-۵ مقدمه.....
۶۹	۲-۵ تابع تأثیرتفاضلی برای آماره آزمون.....
۷۱	۳-۵ ایده‌ی اصلی محاسبه‌ی توزیع تابع تأثیر تفاضلی آزمونها.....
۷۶	۴-۵ توزیع تقریبی d_A^B مدل‌های رگرسیونی خطی.....
۸۱	۱-۴-۵ توزیع حاصل ضرب دو متغیر تصادفی نرمال.....
۸۲	۲-۴-۵ اصلاحیه‌ای بر روشن ریتز و اسکوگارد.....

فصل ششم: تابع تأثیر آزمون‌های مدل‌های رگرسیونی وزنی

۹۰	۱-۶ مقدمه.....
۹۱	۲-۶ براورد β پس از حذف مشاهده‌ی نام در مدل‌های رگرسیونی وزنی.....
۹۲	۳-۶ تابع تأثیر نسبتی آماره‌ی آزمون برای مدل‌های رگرسیونی وزنی.....
۱۰۸	۴-۶ تابع تأثیر تفاضلی آماره آزمون برای مدل‌های رگرسیونی وزنی.....
۱۰۸	۱-۴-۶ توزیع مجانبی تابع تأثیر آماره‌ی آزمون در مدل‌های رگرسیونی وزنی.....

نتیجه‌گیری و پیشنهادات

۱۱۵	نتیجه نهایی.....
-----	------------------

۱۱۷	پیشنهادات
۱۱۹	پیوست
۱۲۴	فهرست واژه‌ها
۱۲۸	فهرست منابع و مأخذ
۱۳۲	چکیده پایان‌نامه به زبان انگلیسی

فصل اول:

مقدمه و کلیات

۱-۱ مقدمه

یکی از مسائلی که از دیرباز ذهن بشر را به خود مشغول کرده است، شناسایی پدیده‌های طبیعی و فیزیکی و تلاش برای معرفی آن‌ها می‌باشد. زمانی که عوامل مؤثر در تشکیل یک پدیده قابل اندازه‌گیری باشند، یکی از ساده‌ترین راه‌ها برای معرفی رفتار آن پدیده، استفاده از یک مدل ریاضی است. گاهی برای سادگی بیشتر ممکن است آگاهانه یا ناآگاهانه، با چشم‌پوشی از بررسی برخی فاکتورهای مهم، زمینه ایجاد نوعی خطا در مدل فراهم گردد و یا حتی زمانی که با اطمینان کامل تمامی فاکتورهای مهم در مدل گنجانده شوند، باز هم به دلیل خطاهای موجود در اندازه-گیری‌ها و عدم دقیقی، مدل‌بندی با نوعی خطای غیر قابل اندازه‌گیری همراه خواهد بود. این مسئله در عمل ما را به استفاده از مدل‌های آماری رهنمون می‌سازد که این عدم اطمینان با اضافه

کردن عبارت خطا در مدل گنجانده می‌شود. ضمن این‌که با استفاده از این مدل‌ها امکان بررسی و تجزیه و تحلیل‌های آماری، میسر می‌شود (رایان،^۱ ۱۹۹۷).

۲-۱ مدل رگرسیونی خطی

ساده‌ترین مدلی که به کمک آن می‌توان رابطه‌ی یک متغیر مشاهده شده را با m متغیر دیگر توصیف کرد، مدل رگرسیونی خطی^۲ است. این مدل که بر حسب پارامترها خطی است، به صورت زیر تعریف می‌شود:

$$Y = X\beta + \varepsilon \quad (1-1)$$

که در آن Y ، بردار $n \times 1$ از متغیر پاسخ، X ماتریس $n \times p$ از متغیرهای توضیحی، β بردار p بعدی از پارامترهای براوردپذیر و ε بردار n بعدی عبارت خطأ است. در عمل مدل‌های خطی کاربرد بیشتری نسبت به سایر مدل‌ها دارند؛ زیرا علاوه بر سادگی محاسبات، به راحتی نیز قابل تفسیرند. همچنین بسیاری از مدل‌های غیرخطی را می‌توان با کمک تبدیلهای ساده ریاضی به مدل‌های خطی تبدیل کرد.

۳-۱ روش‌های آنالیز تشخیص

هنگام استفاده از مدل‌های به‌فرم (۱-۱)، مدل‌بندی بر مبنای یک سری فرضیات صورت می‌گیرد و اعتبار استنباطهای ناشی از این مدل‌ها، مستلزم صحت این فرضیات است. در دهه‌های اخیر تأکید فراوانی بر ارزیابی فرض‌های آماری صورت گرفته و مطالعات گسترهای در این زمینه انجام شده است. همچنین به‌منظور بررسی تأثیر مشاهدات بر براوردهای پارامترها، مقادیر پیش‌بینی و نیز جنبه‌های دیگر مدل روش‌هایی تحت عنوان آنالیز تشخیص پیشنهاد شده است. اهمیت استفاده از آنالیز تشخیص در آن است که علاوه بر شناسایی جنبه‌هایی از مدل که با فرضیات مدل‌بندی هماهنگ نیستند، برای اصلاح آنالیز مدل نیز ایده‌های مطلوبی ارائه می‌دهد و به کمک آن

^۱ Ryan

^۲ Linear Regression Model

می‌توان به این پرسش پاسخ داد که؛ چگونه می‌توان به بررسی ثبات یا مطالعه انحراف در آنالیز نتایج پس از تعیین یک فرمول‌بندی پیچیده پرداخت؟ (کوک و ویزبرگ^۱، ۱۹۸۲)

کوک و ویزبرگ (۱۹۸۲) مطالعه وابستگی نتایج و استنباط‌ها به جنبه‌های مختلف فرمول‌بندی یک مسئله را مطالعه‌ی تأثیر^۲ نامیدند. ایده اصلی آنالیز تأثیر بسیار ساده و به این صورت است که ابتدا یک اغتشاش^۳(پرشیدگی) جزئی در فرمول‌بندی‌های پیچیده اعمال می‌شود و سپس بررسی می‌شود که این اغتشاش چگونه نتایج آنالیز را تغییر می‌دهد؟ یکی از انواع اغتشاشات که بیشترین مقبولیت را نسبت به سایر روش‌ها دارا می‌باشد، روش حذف موردی^۴ است که در آن پس از حذف هر مشاهده، تأثیر ایجاد شده توسط آن مشاهده بررسی می‌شود (کوک و ویزبرگ، ۱۹۸۲).

یکی از فرضیاتی که مدل‌بندی بر مبنای آن صورت می‌گیرد، عدم وجود داده‌ی مؤثر است. اما این فرض به سختی برقرار است و معمولاً نقش همهی مشاهدات در تعیین براوردها، آماره‌های آزمون و دیگر آماره‌ها یکسان نیست، به طوری که در برخی مسائل، ویژگی‌های رگرسیونی تنها تابع وجود مشاهداتی خاص بوده و عملاً از بیشتر مشاهدات چشم‌پوشی می‌شود. به همین خاطر شناسایی مشاهدات مؤثر و بررسی اثر آن‌ها بر روی جنبه‌های مختلف یک آنالیز، برای تحلیلگر از اهمیت ویژه‌ای برخوردار است. ضمن مطالعه تأثیر، علاوه بر شناسایی کامل این گونه مشاهدات، اطلاعاتی در رابطه با میزان اطمینان به نتایج حاصل و وابستگی آن‌ها به مدل مفروض در اختیار آماردان قرار می‌گیرد. به عنوان نمونه، اگر حذف یک مشاهده‌ی مؤثر از مجموعه‌ی داده‌ها علامت یک پارامتر براورده شده را تغییر دهد، به کار بردن استنباط‌های مربوط به آن پارامتر با تردید همراه خواهد بود (کوک و ویزبرگ، ۱۹۸۲).

علاوه بر این، با مطالعه‌ی تأثیر، بخش‌هایی که در آن‌ها داده‌ی گمشده وجود دارد و در آن قسمت از فضای مشاهدات پوشش کافی برای براورد و پیش‌بینی وجود ندارد، مشخص می‌شود. اگر مشاهدات مؤثر و یا مشاهدات دورافتاده ناشی از خطای اندازه‌گیری یا شرایط نامناسب

^۱ Cook & Weisberg

^۲ Study of influence

^۳ Perturbation

^۴ Case Deletion

آزمایشگاهی باشند، باید حذف شوند در غیر این صورت جمع آوری داده‌های بیشتر اقدامی مناسب خواهد بود (کوک و ویزبرگ، ۱۹۸۲).

۱-۴ مطالعه تأثیر

مطالعه‌ی تأثیر براساس هر یک از پنج معیار زیر می‌تواند صورت گیرد:

۱- باقیمانده‌ها

۲- ماتریس تصویر

۳- اندازه‌ی بیضی‌های اطمینان

۴- تابع تأثیر

۵- تأثیر موضعی

که از میان این معیارها، بررسی باقیمانده‌های مدل برآشش شده، از مهم‌ترین ابزارهای آنالیز تشخیص است که به کمک آن می‌توان درستی فرضیاتی چون وجود داده‌ی دورافتاده و یا همگنی واریانس^۱ را تحقیق کرد. یک داده‌ی دورافتاده در واقع مشاهده‌ای است که به طرز آشکاری از مجموعه‌ی اکثریت داده‌ها به دور است و ممکن است نتایج برآشش مدل را به مقدار زیادی تحت تأثیر قرار دهد به گونه‌ای که با حذف آن مشاهده از مجموعه داده‌ها نتایج کاملاً متفاوتی به دست می‌آید. یک ملاک برای شناسایی مشاهدات دور افتاده بررسی باقیمانده‌های مدل است. ضمن این‌که مواردی نیز وجود دارند که داده دورافتاده نبوده، اما نتایج کاملاً به وجود آن مشاهدات وابسته است که آن‌ها را مشاهدات مؤثر^۲ می‌نامند. به این ترتیب نه تنها مشاهدات دور افتاده بلکه هر داده‌ای ممکن است مؤثر تلقی شده و نتایج آنالیزی را تا حد زیادی، تحت تأثیر خود قرار دهد (چترجی و هادی^۳، ۱۹۸۶).

^۱ Homoskedasticity of variance

^۲ Influential observations

^۳ Chatterjee&Hadi

با توجه به این که هدف اصلی این مطالعه، شناسایی این گونه مشاهدات غیر معمول است، آنالیز باقی‌ماندها به تنها‌ی برای شناسایی این موارد کافی نیست. بدین جهت در این پایان‌نامه، علاوه بر مطالعه‌ی باقی‌ماندها به مطالعه‌ی تابع تأثیر نیز می‌پردازیم.

تابع تأثیر از جمله معیارهای مطالعه‌ی تأثیر است که تحلیلگر به کمک آن می‌تواند مشاهدات را براساس میزان تأثیرشان رتبه‌بندی کند. ایده‌ی اصلی تابع تأثیر اولین بار توسط همپل^۱ (۱۹۷۴) مطرح شد و طرح کلی آن به این صورت است که رفتار جزئی توابع آماری مختلف و نیز سهم تک تک مشاهدات بر روی آن‌ها مورد بررسی قرار گرفته و با توجه به آن، به هر مشاهده رتبه‌ای اختصاص می‌یابد، سپس مشاهده‌ای که بیشترین تأثیر را داشته باشد، به عنوان مؤثرترین مشاهده معرفی می‌شود (همپل، ۱۹۷۴).

۱-۵ تاریخچه‌ی اولیه

هر مشاهده از مدل مورد نظر می‌تواند تأثیر زیادی در برآورد پارامترها و به دنبال آن در استنباطهایی که براساس این برآوردها صورت می‌گیرد، داشته باشد. در زمینه‌ی شناسایی مشاهدات مؤثر در برآوردها، مقالات زیادی ارائه شده است که سرآغاز همه‌ی آن‌ها مقاله‌ای بود که توسط کوک (۱۹۷۷) ارائه شد و در آن آنالیز تشخیص حذف موردنی برای مدل‌های خطی-براساس تفاوت بین برآوردهای پارامتری حاصل از مجموعه‌ی داده‌های کامل و مجموعه‌ی داده‌هایی که یک مشاهده از آن‌ها حذف شده، تحت عنوان فاصله‌ی کوک^۲ معرفی شد (Ritze & Skogard, ۲۰۰۷).

به دنبال آن روش‌های متعددی برای مدل‌های مختلف به دست آمد. به عنوان مثال کوک (۱۹۷۹) به بررسی تأثیر هر مشاهده در مدل‌های رگرسیونی خطی پرداخت. پرگیبان^۳ (۱۹۸۱) ابزارهای آنالیز تشخیص را برای رگرسیون لجستیک مطرح کرد. ویلیامز (۱۹۸۷) مباحث آنالیز

^۱ Hample

^۲ Cook's Distance

^۳ Ritz & Skogard

^۴ Pregibon

تشخیص حذفی را برای مدل‌های خطی تعمیم یافته^۱ گسترش داد. کریستن سن^۲ و همکاران (۱۹۹۲) روشی برای مدل‌های آمیخته خطی^۳ ارائه دادند. راسن^۴ (۱۹۹۵) به بررسی مشاهدات مؤثر در مدل‌های رگرسیونی چندمتغیره پرداخت. شی و چن^۵ (۲۰۰۸) آنالیز تشخیص حذفی در مدل‌های چندسطحی را بررسی کردند.

آنچه مسلم است، این است که تأثیر یک مشاهده روی همهٔ کمیت‌های رگرسیونی یکسان نیست و شناسایی نقاط مؤثر در هر زمینه، مستلزم مطالعهٔ تأثیر، در آن زمینه است و این‌گونه نیست که اگر یک مشاهده مثلاً در برآورد پارامترها مؤثر باشد، بر آمارهٔ آزمون نیز مؤثر باشد. آنالیز تشخیص برای بررسی تأثیر یک مشاهده روی آزمون فرض‌ها چندان مورد توجه قرار نگرفته و مطالعات اندکی در این زمینه صورت گرفته است. به عنوان نمونه تغییرات آمارهٔ F جزئی ابتدا در کوک (۱۹۷۹) و سپس به صورت کلی تر، در کوک و ویزبرگ (۱۹۸۲) مورد بررسی قرار گرفت. لاس‌بدر و مول‌گاکر^۶ (۱۹۸۵) به شناسایی نقاط مؤثر در آزمون‌های امتیاز با استفاده از آنالیز تشخیص حذفی پرداختند. ویلیامز^۷ (۱۹۸۷) تقریب یک مرحله‌ای را برای تغییر در آمارهٔ نسبت لگاریتم درست‌نمایی محاسبه کرد. جالیف و لاکادا^۸ (۱۹۹۳) تأثیر تک مشاهدات را بر روی آزمون‌های تک مشاهده‌ای ساده برای میانگین و واریانس مورد مطالعه قرار دادند. گرمیت و ریدن-آور^۹ (۱۹۹۶) اثر تک مشاهدات را بر روی آزمون‌های معمولی میانگین بررسی کردند. لی و ژائو^{۱۰} ژائو^{۱۱} (۱۹۹۶ و ۱۹۹۷) تأثیر مشاهدات را بر آمارهٔ نیکویی برازش پیرسن و نیکویی تابع پیوندی در مدل‌های خطی تعمیم یافته مطالعه کردند. کیم^{۱۲} (۱۹۹۸) روش تأثیر موضعی را برای مطالعهٔ تأثیر مشاهدات روی آزمون فرض‌های خطی به کار برد و کیم (۲۰۰۲) به مطالعهٔ تأثیر مشاهدات

^۱ Generalized linear model

^۲ Christensen

^۳ Linear mixed models

^۴ Rosen

^۵ Shi & Chen

^۶ Lusbader & Moolgavkar

^۷ Williams

^۸ Jolliffe & Lukuda

^۹ Grimmett & Ridenhour

^{۱۰} Lee & Zhao

^{۱۱} Kim

بر روی آزمون فرض‌های خطی با به‌کار بردن آنالیز حذفی موردی یا گروهی پرداخت. ریتز و اسکوگارد (۲۰۰۷) با به‌کار بردن آنالیز تشخیص حذف موردی که در داویسون و اسنل^۱ (۱۹۹۱) تعریف شد، آنالیز حذفی را برای ارزیابی تأثیر یک مشاهده بر روی آزمون‌های نسبت درست‌نمایی بررسی کردند.

در این پایان‌نامه بر روی مدل‌های رگرسیونی وزنی تمرکز می‌کنیم و آنالیز تشخیص حذفی را برای آماره آزمون این مدل‌ها به‌دست می‌آوردم و با به‌کار بردن آن‌ها مشاهدات را براساس میزان تأثیری که در آماره آزمون دارند، رتبه‌بندی کنیم. اما براساس این معیار قادر به تشخیص مشاهدات مؤثر در آزمون نمی‌باشیم زیرا هیچ مدرکی دال بر این‌که مشاهدات با تأثیر زیاد قطعاً دارای تأثیر غیرمعمول بر آزمون می‌باشد، وجود ندارد! به همین دلیل به این معیار کفايت نکرده و برای حل این مشکل سعی می‌شود تا با تعریف معیارهای تأثیر برای آماره لگاریتم نسبت درست‌نمایی و محاسبه توزیع مجانبی مربوط به آن در مورد مؤثر بودن هر مشاهده در آزمون تصمیم‌گیری شود.

پیکره کلی این پایان‌نامه به این صورت است که در فصل بعد به معرفی مدل‌های رگرسیونی و مدل‌های رگرسیونی وزنی می‌پردازم و علاوه بر برآورد پارامترها، آزمون‌های فرضیه را نیز برای این مدل‌ها مورد مطالعه قرار خواهیم داد. در فصل سوم به تشریح مفاهیم و اصطلاحاتی که در فصل‌های بعد به طور مکرر از آن‌ها استفاده خواهیم کرد، می‌پردازم. در همین راستا نخست ماتریس تصویر و چگونگی تفسیر آن را بررسی خواهیم کرد. سپس به معرفی انواع باقی‌مانده‌ها برای مدل‌های رگرسیونی خطی و مدل‌های رگرسیونی وزنی می‌پردازم. در فصل چهارم ضمن تعریف کلی تابع تأثیر، مرور کوتاهی بر بررسی تأثیر هر مشاهده در برآورد پارامترها خواهیم داشت و سپس به مطالعه و بررسی تابع تأثیر نسبتی^۲ آماره نسبت درست‌نمایی برای مدل‌های رگرسیونی خواهیم پرداخت. فصل پنجم را با تعریف تابع تأثیر تفاضلی^۳ برای آماره لگاریتم درست‌نمایی برای مدل‌های رگرسیونی خطی آغاز می‌کنیم، پس از آن توجه خود را بر روی تابع

^۱ Davison & Snell

^۲ Ratio influence function

^۳ Difference influence function

توزیع مجانی این آماره معطوف کرده و از آن برای به دست آوردن معیاری برای تصمیم‌گیری در مورد مؤثر بودن هر مشاهده در آزمون، استفاده می‌کنیم. در فصل ششم تابع تأثیر آماره‌های آزمون را در مدل‌های رگرسیونی خطی وزنی به دست می‌آوریم. به این ترتیب که ابتدا تابع تأثیر نسبتی را برای این مدل‌ها به دست آورده، پس از آن ضمن تعریف تابع تأثیر تفاضلی برای آماره‌ی آزمون این مدل‌ها سعی می‌کنیم توزیع مجانی آن را نیز به منظور شناسایی مشاهدات مؤثر در آزمون محاسبه کنیم.

فصل دوم:

معرفی مدل‌های رگرسیونی

۱-۲ مقدمه

در این فصل از پایان‌نامه به معرفی مدل‌های رگرسیونی خطی کلی می‌پردازیم و برخی نتایج اساسی در زمینه براوردگرهای کمترین توان‌های دوم تعمیم یافته^۱ (*GLSEs*) را بیان خواهیم کرد. سپس برخی حالات خاص آن شامل مدل‌های رگرسیونی خطی و مدل‌های رگرسیونی وزنی را معرفی کرده و برخی خواص استنباطی در زمینه براورد پارامترها و آزمون‌های مربوط به این مدل‌ها را بیان خواهیم نمود.

در بخش بعد، فرم کلی مدل‌های رگرسیونی معرفی می‌شوند. پس از آن در بخش سوم به معرفی مدل‌های رگرسیونی خطی با واریانس همگن می‌پردازیم و برای این مدل‌ها، براورد

^۱ Generalized least square estimators

پارامترها، برآوردهای پارامترها تحت محدودیت‌های خطی و آزمون فرض‌های خطی را بیان خواهیم نمود. در بخش آخر بر روی مدل‌های رگرسیونی وزنی تمرکز می‌کنیم. در این مدل‌ها مشاهدات، دارای ساختارهای ناهمگن برای واریانس هستند. هدف این بخش چگونگی مدل‌بندی این داده‌ها و یافتن برآوردهای این مدل‌ها در حالت کلی و در حالت محدود شده و انجام آزمون فرض‌های خطی است.

۲-۲ فرم کلی مدل‌های رگرسیونی خطی

مدل رگرسیونی خطی کلی به فرم

$$Y = X\beta + \varepsilon \quad (1-2)$$

که در آن Y یک بردار $n \times 1$, X یک ماتریس $n \times p$ معلوم پرتبه ستونی و ε بردار تصادفی خطای غیر قابل مشاهده با

$$E(\varepsilon) = 0, \text{cov}(\varepsilon) = E(\varepsilon\varepsilon^T) = \Omega \in \mathcal{S}(n)$$

می‌باشد که $\Omega(n)$ نشان‌دهنده مجموعه‌ی تمام ماتریس‌های همیشه مثبت $n \times n$ است. اگر $\mathcal{P}_n(0, \Omega)$ کلاسی گسترده از توزیع‌هایی با میانگین صفر و ماتریس کوواریانس Ω تعریف شود، این کلاس به طور خاص توزیع نرمال $(0, \Omega)$ را در بر می‌گیرد. در این پایان‌نامه، مطالعه‌ی خود را بر روی این توزیع متمرکز می‌کنیم. اما ماتریس کوواریانس^۱ Ω معمولاً مجھول است و به عنوان تابعی از یک پارامتر مجھول اما برآوردهای θ به صورت زیر در نظر گرفته می‌شود:

$$\Omega = \Omega(\theta).$$

ساختار کلی ماتریس کوواریانس برای عبارت خطای ε ، به شکل $\Sigma = \sigma^2 \Omega$ است که به وسیله تابعی از یک پارامتر مجھول و برآوردهای θ مشخص می‌شود؛ یعنی

$$\Sigma = \Sigma(\theta).$$

^۱ Covariance matrix

حالات‌های خاص این مدل، با توجه به ساختارهای خاص ماتریس X و ماتریس کواریانس $\Sigma(\theta)$ مدل‌های مختلف را نتیجه می‌دهد که هر کدام در شرایط خاص کاربرد دارند. این مدل‌ها به مدل‌های تک متغیره محدود نمی‌شوند، بلکه مدل‌های چند متغیره مانند مدل‌های آنالیز واریانس چند متغیره^۱ (*MANOVA*)، مدل‌های رگرسیونی به‌ظاهر نامرتبه^۲ و مدل منحنی رشد^۳ را نیز در بر می‌گیرند (کاریا و کوراتا، ۲۰۰۴).

در این مدل‌ها، برای براورد پارامترها می‌توان از قضیه زیر که به قضیه گوس-مارکوف^۴ معروف است، استفاده کرد:

قضیه (۱-۲). در مدل‌هایی که به فرم (۱-۱) هستند، اگر ماتریس Σ معلوم باشد، براورد پارامترها طبق رابطه‌ی زیر به‌دست می‌آیند:

$$b(\Sigma) = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y. \quad (۲-۲)$$

براوردگرهایی که با استفاده از این قضیه به‌دست می‌آیند، بهترین براوردگرهای نالریب خطی^۵ (*BLUE*) می‌باشند و به براوردهای گوس-مارکوف (*GME*) معروف‌اند، اما در بیشتر موارد ماتریس Σ مجھول است و بنابراین نمی‌توان از این براوردگرهای استفاده کرد. در چنین مواردی از براوردگر کمترین توان‌های دوم تعمیم یافته^۶ (*GLSE*) استفاده می‌شود که همان *GME* بوده، با این تفاوت که براوردگر $\widehat{\Sigma}$ به جای Σ جایگذاری شده است

$$b(\widehat{\Sigma}) = (X^T \widehat{\Sigma}^{-1} X)^{-1} X^T \widehat{\Sigma}^{-1} Y. \quad (۳-۲)$$

از این براوردگر به طور گسترده برای داده‌های واقعی استفاده می‌شود (کاریا و کوراتا، ۲۰۰۴).

همان‌گونه که بیان شد، مدل (۱-۲) بیشتر مدل‌ها از جمله مدل‌های رگرسیون خطی ساده و مدل‌های رگرسیونی وزنی را نیز در بر می‌گیرد. در ادامه به معرفی این مدل‌ها می‌پردازیم.

^۱ Multivariate analysis of variance

^۲ Seemingly unrelated regression

^۳ Growth curve model

^۴ Kariya & Kurata

^۵ Gauss-Markov

^۶ Best linear unbias estimator

^۷ General least square estimator

۳-۲ مدل‌های رگرسیون خطی

در مدل‌های رگرسیون خطی، هدف مورد نظر یافتن رابطه‌ی خطی است که به کمک آن می‌توان یک متغیر پاسخ Y را بر مبنای چند متغیر توضیحی x_1, \dots, x_p توصیف کرده و از آن برای پیش‌بینی مقادیر دیگر متغیر Y استفاده کرد. یکی از حالت‌های خاص مدل (۱-۲) مدل رگرسیون خطی ساده به فرم

$$Y = X\beta + \varepsilon \quad (4-2)$$

است که در آن Y بردار $n \times 1$ متغیر پاسخ، X ماتریس $p \times n$ از متغیرهای توضیحی، β بردار p بعدی پارامترهای براوردپذیر و ε بردار n بعدی عبارت خطأ است. براورد بردار پارامترهای β بر مبنای یک نمونه‌ی n تایی از مقادیر مشاهده شده‌ی $y_i, i = 1, \dots, n$ و متغیرهای توضیحی مربوط به آن‌ها انجام می‌شود. برای هر مشاهده‌ی خاص، مثلاً مشاهده‌ی نام، اگر x_i^T ردیف نام ماتریس X باشد، می‌توان مدل را به صورت $y_i = x_i^T \beta + \varepsilon_i$ نوشت که در آن

$$\text{var}(\varepsilon_i) = \sigma^2, E(\varepsilon_i) = 0$$

و برای هر دو مشاهده‌ی i و j ($i \neq j$) $\text{cov}(\varepsilon_i, \varepsilon_j) = 0$. در این حالت امید ریاضی مشاهدات به فرم حاصل‌ضرب ماتریس X در بردار β به شکل $E(Y) = X\beta$ و ماتریس واریانس مشاهدات را به فرم حاصل‌ضرب $\sigma^2 I_n$ در ماتریس همانی می‌باشد؛ پس $\text{var}(Y) = \sigma^2 I_n$.

۱-۳-۲ براورد پارامترها و خواص آن‌ها

قضیه (۲-۲). اگر مدل، به صورت مدل (۴-۲) در نظر گرفته شود، آن‌گاه براورد کمترین توان‌های دوم بردار پارامترهای β عبارت است از

$$\hat{\beta} = (X^T X)^{-1} X^T Y. \quad (5-2)$$

اثبات: رنچر و اسکالچ (۲۰۰۷) صفحه ۱۴۲.

این براورده‌گرها در مقایسه با دیگر براوردها از خواص مطلوب‌تری برخوردارند از جمله این که دارای توزیع مجانبی

$$\sqrt{n}(\hat{\beta} - \beta) \rightarrow N_p(0, \sigma^2(X^T X)^{-1}) \quad (6-2)$$

می‌باشند. در ادامه به برخی دیگر از خواص این براورده‌گرها اشاره خواهیم کرد (رنچر و اسکالچ، ۲۰۰۷).

قضیه‌ی (۳-۲). اگر $var(Y) = \sigma^2 I_n$ و $E(Y) = X\beta$ باشد، آنگاه براورده‌گرهای کمترین مربعات $\hat{\beta}_j$ ، $j = 1, \dots, p$ ، در میان تمام براورده‌گرهای نازلیب خطی، دارای مینیمم واریانس هستند.

اثبات: رنچر و اسکالچ (۲۰۰۷) صفحه ۱۴۲.

نتیجه (۱-۲). اگر $var(Y) = \sigma^2 I_n$ و $E(Y) = X\beta$ باشد، بهترین براورده‌گر خطی نازلیب عبارت است از $a^T \hat{\beta}$ که در آن $\hat{\beta}$ در رابطه‌ی (۵-۵) تعریف شده است.

اثبات: رنچر و اسکالچ (۲۰۰۷) صفحه ۱۴۸.

قضیه (۴-۲). اگر $var(Y) = \sigma^2 I_n$ باشد، کوواریانس ماتریس $\hat{\beta}$ برابر

$$var(\hat{\beta}) = \sigma^2 (X^T X)^{-1} \quad (7-2)$$

خواهد بود.

اثبات: رنچر و اسکالچ (۲۰۰۷) صفحه ۱۴۵.

بردار مقادیر برازش شده عبارت است از:

$$\hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y = HY. \quad (8-2)$$

در تعریف فوق، H ، ماتریس تصویر^۱ و برابر

$$H = X(X^T X)^{-1} X^T \quad (9-2)$$

^۱ Projection matrix

می‌باشد که با استفاده از آن بردار باقی‌ماندها را به شکل

$$e = Y - \hat{Y} = Y - X\hat{\beta} = (I - H)Y \quad (10-2)$$

تعریف می‌شود و می‌توان از آن به عنوان جانشینی برای ε استفاده کرد. با استفاده از رابطه‌ی (۲-۱۰) مقدار باقی‌مانده‌ی آم برابر $e_i = y_i - \hat{y}_i$ خواهد بود. اگر طبق رابطه‌ی (۴-۲) به جای Y از مقدار معادل آن استفاده شود، رابطه‌ی میان e و ε به دست می‌آید:

$$\begin{aligned} e &= (I - H)(X\beta + \varepsilon) \\ &= (I - H)\varepsilon. \end{aligned}$$

این رابطه صراحتاً بیان می‌کند که رابطه‌ی میان e و ε تنها به ماتریس H بستگی دارد و با استفاده از آن، می‌توان e را به فرم اسکالر $e_i = \sum_{j=1}^n h_{ij}\varepsilon_j$ نوشت که

$$h_{ij} = x_i^T (X^T X)^{-1} x_j \quad (11-2)$$

عضو (i,j) ام ماتریس H می‌باشد و در آن x_i^T سطر i ام ماتریس X است. اگر h_{ij} ‌ها به اندازه‌ی کافی کوچک باشند، e جانشین مناسبی برای ε خواهد بود (رنچر و اسکالچ^۱، ۲۰۰۷).

در ک بهتر رابطه‌ی میان e و ε مستلزم مطالعه‌ی دقیق‌تر پیرامون ماتریس H می‌باشد. به دلیل اهمیت بردار باقی‌مانده‌ها و ماتریس تصویر، در فصل آینده به تفصیل درباره‌ی آن‌ها بحث خواهیم کرد.

از سوی دیگر آماره‌ی مجموع مربعات باقی‌مانده‌ها به صورت

$$s^2 = e^T e \quad (12-2)$$

تعریف شود. می‌توان از

$$\hat{\sigma}^2 = \frac{s^2}{n - p} \quad (13-2)$$

^۱ Rencher & Schaalje