

دانشگاه پیام نور
دانشکده علوم پایه

پایان نامه
برای دریافت مدرک کارشناسی ارشد
رشته آمار ریاضی
گروه آمار

عنوان :

مقایسه کاربرد رگرسیون لجستیک با روش پروبیت در تجزیه
و تحلیل متغیرهای دوتایی

جهانچهر جواهری

استاد راهنما :

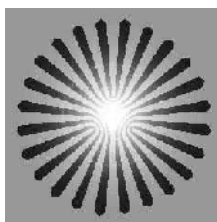
دکتر پرویز نصیری

استاد مشاور :

دکتر مسعود یارمحمدی

دی ماه ۱۳۸۹

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه پیام نور
دانشکده علوم پایه
مرکز تهران

پایان نامه
برای دریافت مدرک کارشناسی ارشد
رشته آمار ریاضی
گروه آمار

عنوان:

مقایسه کاربرد رگرسیون لجستیک با روش پروبیت در تجزیه
و تحلیل متغیرهای دوتایی

جهانچهر جواهری

استاد راهنما:

دکتر پرویز نصیری

استاد مشاور:

دکتر مسعود یارمحمدی

دی ماه ۱۳۸۹



شماره
تاریخ
پیوست

صور تجلسه دفاع از پایان نامه کارشناسی ارشد

جلسه دفاع از پایان نامه کارشناسی ارشد خانم جهانچهر جواهری
دانشجوی رشته آمار ریاضی به شماره دانشجویی: ۸۷ ۸۶۷۱۰۴۶
تحت عنوان:

"مقایسه کاربرد رگرسیون لجستیک با روش پروبیت در تجزیه و تحلیل داده های دوتایی"

جلسه دفاع با حضور داوران نامبرده ذیل در روز چهارشنبه مورخ: ۸۹/۱۰/۸ ساعت: ۹/۳۰-۸/۳۰ در محل
مجمع علوم پایه و کشاورزی برگزار شد. و پس از بررسی پایان نامه مذکور با نمره به عدد ۷۰...
به حروف هفتاد و با درجه ارزشیابی خوب... مورد قبول واقع شد نشد

ردیف	نام و نام خانوادگی	هیات داوران	مرتبه دانشگاهی	دانشگاه/ موسسه	امضاء
۱	دکتر پرویز نصیری	استاد راهنما	استاد	پیام نور	
۲	دکتر مسعود بیار محمدی	استاد مشاور	استاد	پیام نور	
۳	دکتر صادق رضایی	استاد داور	استاد	دانشگاه شهید عمرالاحقاز	
۴	دکتر مسعود بیار محمدی	نماینده علمی گروه	استاد	پیام نور	

۵۰۱۳

اینجانب **جهانچهر جواهری** دانشجوی ورودی سال ۱۳۸۶ مقطع کارشناسی ارشد رشته آمار ریاضی گواهی می‌نمایم چنانچه در پایان نامه خود از فکر، ایده و اندیشه و نوشته دیگری بهره گرفته‌ام با نقل قول مستقیم یا غیر مستقیم منبع و ماخذ آن را نیز در جای مناسب ذکر کرده‌ام. بدیهی است مسئولیت تمامی مطالبی که نقل قول دیگران نباشد بر عهده خویش می‌دانم و جوابگوی آن خواهم بود.

دانشجو تأیید می‌نماید که مطالب مندرج در این پایان نامه (رساله) نتیجه تحقیقات خودش می‌باشد و در صورت استفاده از نتایج دیگران مرجع آن را ذکر نموده‌است.

نام و نام خانوادگی دانشجو: **جهانچهر جواهری**

تاریخ و امضاء

اینجانب **جهانچهر جواهری** دانشجوی ورودی سال ۱۳۸۶ مقطع کارشناسی ارشد رشته آمار ریاضی گواهی می‌نمایم چنانچه بر اساس مطالب پایان نامه خود اقدام به انتشار مقاله، کتاب ونمایم ضمن مطلع نمودن استاد راهنما، با ایشان نسبت به نشر مقاله، کتاب وو به صورت مشترک و با نام ذکر نام استاد راهنما مبادرت نمایم.

نام و نام خانوادگی دانشجو: **جهانچهر جواهری**

تاریخ و امضاء

کلیه حقوق مادی مترتب از نتایج مطالعات، آزمایشات و نوآوری ناشی از تحقیق موضوع این پایان نامه متعلق به دانشگاه پیام نور می‌باشد.

دی ماه ۱۳۸۹

چکیده

آنالیز رگرسیون ابزاری بسیار مفید برای بررسی رابطه و همبستگی بین دو و یا چندمتغیر است و به علت قابلیت های خوب کاربردی، امروزه وسیله مفیدی برای رشته های مختلف علوم است. در آنالیز رگرسیون در حالت خاص با متغیر های وابسته ای که فقط ۲ مقدار دارند سروکار داریم. در این حالت می توان از رگرسیون لجستیک و رگرسیون پروبیت برای آنالیز روابط بین متغیر وابسته و چند متغیر مستقل استفاده کرد. هدف این آنالیز تنظیم عوامل موثر، بخش بندی کواریت های مهم مربوط به متغیر وابسته و پیشگویی مقدار متغیر وابسته است. این دو رگرسیون از نظر کیفی نتایج مشابهی دارند، ولی برآورد پارامترهای دو مدل را نباید به طور مستقیم مورد مقایسه قرار داد. لجیت و پروبیت هر دو شکلی از تبدیل داد های دو حالتی (دوتایی) هستند. هر دوی این آنالیزها از مدل های احتمال خطی بمنظور برآورد پارامترهای مدل استفاده می کنند و حالت خاصی از مدل های عمومی خطی هستند. برای محاسبه پارامترهای مدل، تخمین حداکثر درستنمایی از پارامترهای رگرسیون را محاسبه می کنند و برای پیشینه کردن و رسیدن به حداکثر درستنمایی روش Fisher's scoring مورد استفاده قرار می گیرد. در نهایت در هر دو روش در صورت رسیدن به حداکثر درستنمایی پارامترهای مدل برآورد می شود. هر دو روش برای رسیدن به حداکثر درستنمایی از رویه مکرر استفاده می کنند.

در حالتی که توزیع تجمعی پاسخها در برابر دادهها از توزیع نرمال تبعیت نکند پیشنهاد می شود بجای تبدیل پروبیت از تبدیل لجیت استفاده کنیم و در صورتی که توزیع تجمعی دادهها از توزیع نرمال تبعیت کنند، آنالیز پروبیت ترجیح داده می شود. در حضور تعداد زیادی سطوح متغیر مستقل تبدیل لجیت برازش بهتری را فراهم می کند و بر عکس پروبیت مدل های با اثرات تصادفی و مجموعه داده های در حد متوسط را بهتر برازش می کند. در صورتی که طرح مدونی بکار برده نشده باشد خطای برآورد احتمالات در آنالیز پروبیت نسبت به لجستیک بیشتر است و لذا در این حالت آنالیز لجستیک توصیه می شود.

از مشکلات مهم همراه با این نوع رگرسیونها می توان، نرمال نبودن متغیرها، عدم پوشش کامل نتایج و احتمال واقع شدن \hat{Y} خارج از دامنه (۰ و ۱)، ناهمگنی واریانس توزیع ها و مشکوک بودن مقدار R^2 به عنوان معیار نکویی برازش را نام برد. مدل های احتمال خطی از نظر منطقی خیلی جالب نیستند زیرا این مدلها فرض می کنند احتمال همزمان با افزایش X بطور خطی افزایش می یابد و اثر افزایشی X بطور ثابت باقی می ماند. این موضوع گاهی اوقات غیر واقعی است.

فهرست مطالب

عنوان مطلب..... صفحه

I.....چکیده

II.....پیش گفتار

فصل اول: کلیات و مفاهیم

- ۱-۱ ارتباط آماری متغیرها..... ۱
- ۲-۱ تحلیل رگرسیون..... ۱
- ۱-۲-۱ مدل‌های خطی..... ۲
- ۲-۲-۱ برآورد پارامترها به روش کمترین مربعات..... ۳
- ۳-۲-۱ برآورد پارامترها بروش کمترین مربعات..... ۴
- ۳-۱ رگرسیون خطی ساده..... ۵
- ۱-۳-۱ تاریخچه..... ۵
- ۴-۱ رگرسیون خطی چندگانه..... ۱۰
- ۵-۱ رگرسیون غیرخطی..... ۱۳
- ۱-۵-۱ نمونه‌هائی از مدل‌های غیر خطی..... ۱۵
- ۲-۵-۱ روش نیوتن رافسن..... ۱۶
- ۶-۱ خانواده مدل‌های خطی تعمیم یافته..... ۱۸
- ۱-۶-۱ مولفه‌های مدل‌های خطی تعمیم یافته..... ۱۹
- ۱-۶-۲ مولفه تصادفی..... ۱۹
- ۱-۶-۳ مولفه سیستماتیک .. ۱۹
- ۱-۶-۴ مولفه ربط..... ۲۰
- ۱-۶-۵ ویژگی‌های مدل خطی تعمیم یافته..... ۲۱
- ۱-۶-۶ برآورد برای الگوهای خطی تعمیم یافته..... ۲۲

۲۳	۷-۶-۱ خانواده نمائی توزیع ها.....
۲۶	۸-۶-۱ معادلات درستنمائی برای مدل خطی تعمیم یافته
۲۸	۹-۶-۱ برخی یافته های مهم در باره GLM
۲۸	۱۰-۶-۱ روش برآورد MLE
۲۹	۱۱-۶-۱ روش برآورد MME
۲۹	۷-۱ رگرسیون خاص
۲۹	۱-۷-۱ رگرسیون لجستیک
۳۲	۱-۱-۷-۱ مثال کاربردی از رگرسیون لجستیک
۳۳	۲-۱-۷-۱ مدل رگرسیون لجستیک
۳۴	۳-۱-۷-۱ مدل رگرسیون لجستیک دوتائی
۳۴	۴-۱-۷-۱ مدل رگرسیون لجستیک ترتیبی
۳۴	۵-۱-۷-۱ کاربرد رگرسیون لجستیک
۳۴	۶-۱-۷-۱ مشخص کردن عوامل ریسک مرتبط
۳۵	۷-۱-۷-۱ اصلاح و تعدیل عوامل
۳۵	۸-۱-۷-۱ پیشگوئی و تشخیص
۳۶	۲-۷-۱ آنالیز پروبیت
۳۷	۱-۲-۷-۱ کاربرد آنالیز پروبیت
۳۸	۲-۲-۷-۱ چگونگی انجام پروبیت
۳۸	۳-۲-۷-۱ مراحل انجام آنالیز پروبیت
۴۰	۴-۲-۷-۱ نکات قابل توجه در آنالیز پروبیت
۴۱	۸-۱ آزمون کای مربع
۴۲	۱-۸-۱ مفروضات
۴۲	۲-۸-۱ محدودیت ها
۴۲	۳-۸-۱ آزمون فیشر

فصل دوم: رگرسیون لجستیک و رگرسیون پروبیت

- ۴۴ ۱-۲ رگرسیون لجستیک و رگرسیون پروبیت
- ۴۵ ۱-۱-۲ رگرسیون لجستیک
- ۴۵ ۱-۱-۱-۲ مدل آماری
- ۴۵ ۲-۱-۱-۲ بررسی مدل
- ۴۶ ۳-۱-۱-۲ معیارهای نکوئی برازش
- ۴۶ ۴-۱-۱-۲ نکوئی برازش با آماره پیرسون
- ۴۷ ۵-۱-۱-۲ آماره نسبت درستنمائی
- ۴۹ ۶-۱-۱-۲ آماره والد
- ۵۰ ۷-۱-۱-۲ آماره اسکور
- ۵۱ ۸-۱-۱-۲ odds ratio برآوردهای
- ۵۱ ۹-۱-۱-۲ روش گام به گام
- ۵۲ ۱۰-۱-۱-۲ حجم نمونه برای پیشگوئی کمی
- ۵۲ ۲-۱-۲ رگرسیون پروبیت
- ۵۳ ۱-۱-۲-۲ مدل‌های پروبیت
- ۵۴ ۳-۱-۲ مثال مقایسه نتایج رگرسیون لجستیک و رگرسیون پروبیت

فصل سوم: نتایج رگرسیون لجستیک و رگرسیون پروبیت

- ۵۶ ۱-۳ مقدمه
- ۵۶ ۲-۳ مواد و روش‌ها
- ۵۶ ۱-۲-۳ جامعه آماری
- ۵۶ ۲-۲-۳ نمونه مورد مطالعه
- ۵۶ ۳-۲-۳ پرسشنامه
- ۵۷ ۴-۲-۳ آزمون‌های روائی و پویائی
- ۵۷ ۵-۲-۳ روش‌های آماری مورد استفاده

۵۷.....	۳-۳ نتایج
۵۷.....	۱-۳-۳ ترسیم جداول فراوانی.....
۶۵.....	۲-۳-۳ استفاده از جداول توافق آزمون کای مربع.....
۶۷.....	۳-۳-۳ تعیین پایائی و سنجش.....
۶۷.....	۴-۳-۳ نتایج آنالیز لجستیک.....
۷۶.....	۴-۳ نتایج رگرسیون گام به گام.....
۸۱.....	۳-۵ نتایج آنالیز پروبیت

فصل چهارم: مقایسه رگرسیون لجستیک و پروبیت

۸۸.....	۴-۱ اختلاف رگرسیون لجستیک و پروبیت.....
۸۸.....	۴-۱-۱ از نظر نوع تبدیل بکار برده شده.....
۸۹.....	۴-۱-۲ از نظر توزیع متغیر پاسخ.....
۸۹.....	۴-۱-۳ از نظر تعداد سطوح متغیر مستقل.....
۸۹.....	۴-۱-۴ از نظر نوع کاربرد.....
۸۹.....	۴-۲ تشابه رگرسیون لجستیک و پروبیت.....
۹۰.....	۴-۳ مشکلات مرتبط با نتایج و برآوردهای متغیر وابسته دو حالتی.....
۹۰.....	۴-۴ نتیجه گیری کلی

منابع

۹۲.....	منابع.....
---------	------------

فصل ۱

کلیات و مفاهیم

۱-۱ ارتباط آماری متغیرها

در اغلب آزمایش های آماری در هر یک از واحدهای آزمایشی، متغیرهای مختلفی اندازه گیری می شوند. ولی در اغلب تکنیک های تجزیه و تحلیلی فقط یک متغیر به صورت جداگانه مورد بررسی قرار می گیرد. ارتباط آماری به سنجش ارتباط بین دو یا چند متغیر کمی مربوط می شود. روش های مختلفی نیز جود دارد که می توان ارتباط بالقوه بین متغیرها را بررسی کرد. بوسیله روش های یک و چند متغیره می توان وابستگی بین متغیرها را بررسی کرد ولی روش ها چند متغیره خیلی پیچیده هستند (اوریت و دان ۱۹۹۱- موریسون ۱۹۹۰). یکی از روشهای ساده بمنظور بررسی ارتباط آماری بین متغیرها استفاده از ضرایب همبستگی است. ضریب همبستگی شاخصی است ریاضی که جهت و مقدار رابطه ی بین دو متغیر را توصیف می کند (پلاتا ۲۰۰۶). ضریب همبستگی درمورد توزیع های دو یا چند متغیره به کار می رود. اندازه های مختلفی از همبستگی وجود دارد ولی همیشه مقدار آن بین ۱- تا ۱+ قرار دارد. اگر مقادیر دو متغیر شبیه هم تغییر کند یعنی با کم یا زیاد شدن یکی، دیگری هم کم یا زیاد شود، به گونه ای که بتوان رابطه آنها را به صورت یک معادله بیان کرد، گوئیم بین این دو متغیر همبستگی وجود دارد. برای سنجش همبستگی بین متغیرها، ضرایب گوناگونی به کار می رود که مهمترین آنها ضریب همبستگی ساده پیرسون، ضریب همبستگی اسپیرمن و ضریب همبستگی کندال است (هاک و کساوسکی ۲۰۰۷).

روش دیگری که برای بررسی ارتباط آماری بین متغیرها مورد استفاده قرار می گیرد معادلات رگرسیون است. در اینجا انواع رگرسیون بمنظور سنجش ارتباط آماری بین متغیرها مورد بررسی قرار خواهد گرفت (چان ۲۰۰۳ و لیوایز ۱۹۹۳).

۲-۱ تحلیل رگرسیون

تحلیل رگرسیون، مجموعه ای از تکنیک های آماری است که برای مدل سازی و بررسی رابطه بین یک متغیر پاسخ Y و مجموعه ای از متغیرهای پیشگو (مستقل) X_1, X_2, \dots, X_k می باشد. کاربردهای رگرسیون گسترده بوده و تقریباً در هر زمینه کاربردی چون پزشکی، مدیریت، اقتصاد و استفاده می شود. مدل های رگرسیون خطی بطور گسترده بعنوان مدل های تجربی برای تقریب یک رابطه پیچیده و معمولاً نامعلوم بین متغیرهای پاسخ و پیشگو بکار می رود.

نلدر و دربرن در سال ۱۹۷۲ مدل‌های خطی عمومی را به مدل‌های خطی تعمیم یافته توسعه دادند. در حقیقت مدل‌های خطی تعمیم یافته، مدل‌های خطی عمومی هستند که متعلق به خانواده نمایی می باشند و خط‌هایی با توزیع نرمال دارند. GLM ها به فرض مشخص بودن توزیع کامل نیاز دارد، اما در برخی موقعیت‌ها مخصوصاً در بررسی داده‌های گسسته، گاهی اوقات فرض توزیع کامل محدود می شود. برای رفع نیاز به وجود توزیع کامل و دربرن در سال ۱۹۷۴ تابعی به نام تابع شبه درستنمایی معرفی نمود که فقط بر اساس فرض وجود گشتاور مرتبه دوم متناهی متغیر پاسخ برقرار می باشد. (استفورد، ۱۹۹۶). از روش تابع شبه درستنمایی فقط در برآورد کردن میانگین یا ضرایب رگرسیونی استفاده می شود.

۱-۲-۱ مدل‌های خطی

مدل‌های خطی بصورت زیر بیان می شود:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \varepsilon \quad (1-1)$$

که در آن Y متغیر پاسخ، X ها متغیر پیشگو، $\beta_0, \beta_1, \dots, \beta_k$ مجموعه پارامترهای نامعلوم و ε جمله خطای تصادفی است. مدل بالا را اغلب مدل رگرسیون خطی می نامند. جمله خطا دارای میانگین صفر است. پس میانگین پاسخ در مدل رگرسیون خطی عبارت است از:

$$E(Y) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k \quad (2-1)$$

در این مدل میانگین پاسخ یک تابع خطی از پارامترهای مجهول β_0 و β_1 و \dots و β_k است. مدل‌های رگرسیون خطی به دلایل گوناگونی کاربرد وسیعی دارند. دلیل اول اینکه این مدلها تقریب طبیعی چندجمله‌ای برای روابط تابعی پیچیده ترند. یعنی اگر $E(Y)=f(X)$ رابطه دقیق یا قطعی بین متغیر پاسخ و متغیر پیشگو باشد، آنگاه تقریب اول سری تیلور این رابطه در نقطه‌ای مانند x_0 با رابطه زیر بدست می آید.

$$E(Y) \cong f(x_0) + \left. \frac{df(x)}{dx} \right|_{x=x_0} (x-x_0) + R \quad (3-1)$$

$$\cong \beta_0 + \beta_1 (x-x_0)$$

که صرف نظر از باقیمانده R (بجز جمله خطا) یک الگوی رگرسیون خطی یک متغیری است. وقتی K متغیر پیشگو داشته باشیم، تقریب مرتبه اول تیلور مستقیماً به یک مدل رگرسیونی خطی K متغیری منجر می شود. چون اغلب از الگوهای رگرسیون خطی (به طور موفقیت آمیزی) بعنوان تقریب چند جمله ای ها استفاده می شود، لذا بعضی اوقات این مدلها را مدل تجربی می نامند. دومین دلیل استفاده از مدلهای رگرسیون خطی اینست که از طریق آنها پارامترهای مجهول مدل یعنی $\beta_0, \beta_1, \dots, \beta_k$ مستقیماً برآورد شوند.

در ادامه به دو روش مهم برآورد پارامترهای مدل رگرسیون خطی که روشهای کمترین مربعات معمولی و درستمائی ماکزیمم هستند، اشاره می نمائیم.

۱-۲-۲-۲ برآورد پارامترها

۱-۲-۲-۱ برآورد پارامترها به روش کمترین مربعات^۱

این روش نوعاً برای برآورد ضرایب رگرسیون در یک مدل رگرسیون خطی چندگانه بکار می رود. فرض می کنیم $n > k$ مشاهده از متغیر پاسخ Y بصورت y_1, y_2, \dots, y_n داریم. برای هر پاسخ مشاهده شده y_i یک مشاهده برای هر متغیر پیشگو داریم و فرض می کنیم $x_{i1}, x_{i2}, \dots, x_{ik}$ مشاهده i ام متغیر x_i باشد. جمله خطای ε در مدل را دارای میانگین صفر و واریانس ثابت σ^2 و $\{\varepsilon_i\}$ را متغیرهای تصادفی ناهمبسته فرض می کنیم. معادله ۳ را بر حسب مشاهدات بصورت معادله زیر می نویسیم:

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \varepsilon_i \\ &= \beta_0 + \sum_{t=1}^k \beta_t x_{it} + \varepsilon_i \end{aligned} \quad (4-1)$$

روش کمترین مربعات، ها را در معادله بالا طوری انتخاب می کند که مجموع مربعات خطاهای ε_i مینیمم شود. برآورد کمترین مربعات β ها بصورت زیر است:

$$b = (X'X)^{-1} X'Y$$

۱-۲-۲ برآورد پارامترها به روش درستنمایی ماکزیمم^۱

از روش کمترین مربعات برای برآورد پارامترهای یک مدل رگرسیون خطی بدون در نظر گرفتن شکل توزیع متغیر پاسخ Y می توان استفاده کرد. اگر شکل توزیع اسخ معلوم باشد، آنگاه از روش برآورد دیگری به نام درستنمایی ماکزیمم می توانیم استفاده کنیم.

مدل رگرسیون خطی $Y = X\beta + \varepsilon$ را در نظر می گیریم. فرض می کنیم خطاها در این مدل دارای توزیع نرمال مستقل با میانگین صفر و واریانس ثابت σ^2 باشند. در این صورت مشاهده در نمونه

$$(y_i, x_{i1}, \dots, x_{ik}) \text{ دارای توزیع نرمال مستقل با میانگین } \beta_0 + \sum_{t=1}^k \beta_t x_{it} \text{ و واریانس } \sigma^2$$

خواهد بود. تابع درستنمایی را از توزیع احتمال توام مشاهدات پیدا می کنیم. اگر این توزیع توأم را با مشاهدات داده شده و پارامترهای مجهول در نظر بگیریم، آنگاه تابع درستنمایی را به صورت زیر داریم:

$$L(y, \beta, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2} [(y - x\beta)'(y - x\beta)]\right\} \quad (5-1)$$

برآوردگرهای درستنمایی ماکزیمم مقادیری از پارامترهای β و σ^2 هستند که تابع درستنمایی را ماکزیمم می کنند. ماکزیمم کردن تابع درستنمایی L معادل ماکزیمم کردن تابع لگاریتم درستنمایی یعنی $\ln L$ می باشد.

تابع لگاریتم درستنمایی بصورت زیر است :

$$\ln[L(y, \beta, \sigma^2)] = \frac{-n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} [(y - x\beta)'(y - x\beta)] \quad (6-1)$$

مشتق تابع لگاریتم درستنمایی را تابع امتیاز می نامند. اگر از تابع لگاریتم درستنمایی نسبت به پارامترهای β و σ^2 مشتقات جزئی گرفته و مساوی صفر قرار دهیم خواهیم داشت:

$$\frac{1}{\sigma^2} X'(Y - X\beta) = 0 \quad (7-1)$$

¹-Maximom Liklihood

دستگاه معادلات $k \times k$ بالا را معادلات امتیاز درستنمایی ماکزیمم^۱ می نامند. برآورد درستنمایی ماکزیمم، پاسخ دستگاه معادلات امتیاز است.

$$\beta = (XX)^{-1}XY$$

اگر نسبت به σ^2 مشتق بگیریم و مساوی صفر قرار دهیم، برآوردگر σ^2 بصورت زیر بدست می آید:

$$\hat{\sigma}^2 = \frac{1}{n}(y - X\beta)'(y - X\beta)$$

بطور کلی برآوردگرهای درستنمایی نسبت به برآوردگرهای کمترین مربعات دارای خواص آماری بهتری هستند که بخاطر فرض های اضافی می باشد که برای آنها در نظر گرفته می شود. مثلاً برآوردگرهای درستنمایی ماکزیمم نیاز به نرمال بودن توزیع مشاهدات دارند ولی در روش برآورد کمترین مربعات

برقراری این فرض لازم نیست. برآوردگرهای درستنمایی ماکزیمم نارایب بطور مجانبی نارایب هستند و با افزایش حجم نمونه اریب می شوند و نیز مجموعه ای از آماره های بسنده را تشکیل می دهند و سازگارند.

۳-۱ رگرسیون خطی ساده

۱-۳-۱ تاریخچه رگرسیون خطی

رگرسیون شاخه ای از علم آمار است که به بررسی رابطه بین متغیرها می پردازد (کاتنر ۲۰۰۵). موضوع رگرسیون ابتدا در قرن هیجدهم با استفاده از نجوم با هدف هدایت کشتی ها مورد توجه قرار گرفت. لژاندر^۲ در سال ۱۸۰۵ روش کمترین توان های دوم را بسط و توسعه داد. گوس^۳ ادعا می کند که او این روش را چند سال قبل توسعه داده و در سال ۱۸۰۹ نشان داد که وقتی خطها دارای

1-Score Equations Maximum Likelihood
1-catner
2- Legendre
3-Guss-۱

توزیع نرمال اند، روش کمترین توانهای دوم به جواب بهینه منجر می شود. تا قرن نوزدهم این روش شناسی تقریباً در علوم فیزیک مورد استفاده قرار می گرفت. فرانسیس گالتن^۱ در سال ۱۸۷۵ اصطلاح برگشت به مقدار متوسط را در مورد یک معادله رگرسیونی ساده در نظر گرفت. گالتن از این اصطلاح برای بیان این پدیده که پسران پدران قد بلند، بلند قد بوده ولی نه به بلندی قد پدرشان، در حالی که پسران پدران کوتاه قد میل به کوتاه قد بودن دارند ولی این تمایل به اندازه کوتاه قد بودن پدرانشان نیست، استفاده کرد. برای اطلاع بیشتر در مورد تاریخچه استفاده از رگرسیون می توانید به (استیگر^۲ ۱۹۸۶) مراجعه کنید. کاربرد های رگرسیون متعدد است و همانند سایر شاخه های علم آمار، تقریباً در هر زمینه ای از جمله مهندسی، فیزیک، اقتصاد، مدیریت، علوم زیستی و علوم اجتماعی صورت می پذیرد. تحلیل رگرسیونی فن و تکنیکی آماری برای بررسی و به مدل در آوردن ارتباط بین متغیرهاست. در حقیقت این تحلیل ممکن است فن آماری با بیشترین و وسیعترین کاربرد بین فنون آماری باشد. (مونت گمری و پک^۱، ۱۹۹۲). هدف در این تحلیل یافتن یک مدل مناسب برای نمایش ارتباط بین متغیرهای وابسته و مستقل می باشد.

ضریب کورولاسیون فقط یک معیار از ارتباط خطی را می دهد و هیچ نشانی از این که کدام خط راست را ارائه نمی دهد (اگر یک خط راست مناسب باشد) و نشان دهنده بهترین ارتباط است. آنالیز رگرسیون ساده خطی به معنی تعیین این ارتباط است. ساده در اینجا اشاره به این حقیقت دارد که یک متغیر مستقل در آنالیز قرار داده شده است و ربطی به پیچیده یا ساده بودن آنالیز ندارد. بوسیله نمودار داده ها نیز می توان نشان داد که مدل خط راست مناسب است یا نه.

دو متغیر موجود در رگرسیون ساده یکسان نیستند. تغییر در یکی (متغیر مستقل) باعث تغییر در دیگری (متغیر وابسته) می شود و رگرسیون y برحسب x همانند رگرسیون x برحسب y نیست و باید متوجه این اختلاف بود. به عنوان مثال چون تغییر در سطوح بیماری می تواند روی

4-Galton-۲
5-Stigler-²
6-gomeri

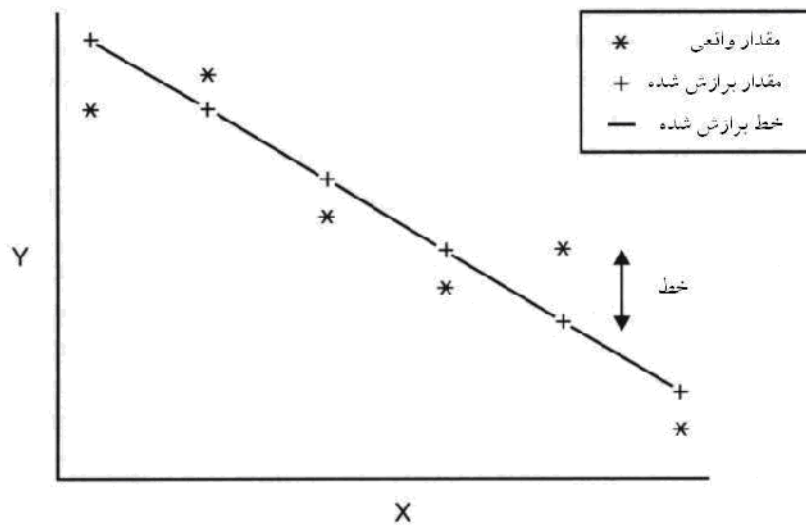
عملکرد تأثیر بگذارد، امکان پیدا کردن رگرسیون عملکرد محصول (به عنوان متغیر وابسته y) بر حسب سطوح بیماری در آن محصول (به عنوان متغیر مستقل x) وجود دارد. ولی این احتمال وجود ندارد که عملکرد روی سطح بیماری تأثیر بگذارد. لذا رگرسیون بیماری بر حسب عملکرد غیر ضروری خواهد بود. اگر x و y همبستگی کامل داشته باشند، در این صورت مدل رگرسیون y بر حسب x و x بر حسب y یکسان خواهد بود ولی تفسیر آنها هنوز متفاوت خواهد بود.

اگر متغیر مستقل به وسیله x و متغیر وابسته را به وسیله y نشان دهیم، آنالیز رگرسیون ساده یک مدل رگرسیون خطی ساده را به شکل زیر برآزش خواهد نمود.

$$y = a + bx$$

که فقط یک x برای پیش بینی y وجود دارد. و خطی یعنی اینکه مدل بر حسب β_0 و β_1 خطی است. این مدل کاملاً مناسب داده ها نخواهد بود ولی همان طوری که در شکل نشان داده شده است، مشاهده ها به صورت تصادفی از خط راست انحراف دارند. مقادیر y محاسبه شده به وسیله مدل در واقع مقادیر برآزش شده هستند و تشکیل یک خط راست را می دهند. اختلاف بین مقادیر برآزش شده و مقادیر y مشاهده شده خطاها هستند. آنالیز رگرسیون مقادیر پارامترهای b, a را برای مدل تخمین می زند و بنابراین مدل خط راست را برآزش می کند. روش های مختلفی برای تخمین پارامترهای b, a وجود دارد. متداول ترین روش «روش حداقل مربعات» است که در آن، مقادیر b, a را به گونه ای پیدا می کنند که حداقل مجموع مربعات خطا را داشته باشیم. مفروضات مورد نیاز برای آنالیز رگرسیون خطی همان

مفروضات آنالیز واریانس، یعنی ثابت بودن واریانس، توزیع تصادفی و استقلال خطاها است. همانند آنالیز واریانس، یک روش خوب برای ارزیابی اعتبار مدل برآزش شده، ترسیم انحراف های مدل در برابر مقادیر برآزش شده است. مجدداً توزیع تصادفی، به صورت یک نوار افقی در اطراف خطای صفر نشان دهنده مناسب بودن مدل مورد استفاده است. نباید الگو یا تمایل خاصی در نمودار وجود داشته باشد.



یک مثال از مفهوم مقادیر برازش شده، مقادیر واقعی و مدل برازش شده
(در این حالت یک مدل خط راست) برگرفته از دراپر و اسمیت ۱۹۸۱

اگ

ر نمودار داده های خام اولیه یک عکس العمل منحنی خط را به متغیر مستقل نشان دهد، ممکن است هنوز رگرسیون خطی مناسب باشد. به وسیله یک تابع ریاضی مناسب یا تبدیل کردن، می توان با استفاده از رگرسیون خطی یک مدل برای داده ها برازش نمود (کولین ۱۹۹۹). اصطلاح رگرسیون خطی را می توان برای همه مدل هایی که از نظر ریاضی پارامترهای خطی دارند به کار برد ولی منعکس کننده شکل تابعیت نیست. هر دو متغیر وابسته و مستقل را می توان به منظور برازش مدل تبدیل کرد. به عنوان مثال در مدل های زیر می توان متغیر مستقل x را تبدیل نمود:

$$y = a + bx + cx^2$$

$$y = a + bx + cx^2 + dx^3$$

$$y = a + b \ln x$$

معادلات خط راست، درجه دو و درجه سه که در بالا ذکر شد، اولین سه معادله چند جمله ای مناسبی هستند که وابسته به مقایسه های چند جمله ای می باشند.

اگر یک مجموعه داده را برای آنالیز رگرسیون مورد استفاده قرار می گیرد، باید یک دامنه مناسب و کافی از مقادیر x موجود باشد. برای دامنه مقادیر x موجود فقط می توان یک خط یا منحنی خط برازش نمود و هر نتیجه بدست آمده درباره این داده ها، فقط برای محدوده همین دامنه است. مدلی که برای مقادیر از یک دامنه داده شده برازش می شود را نباید برای پیش بینی یا تخمین مقادیری که خارج از این دامنه واقع شده اند به کار برد چون به دلیل موجود نبودن هیچ اطلاعاتی در زمینه این مقادیر نمی دانیم که این مدل برای خارج از این دامنه نیز مناسب است یا خیر. خط برازش شده با اطمینان بالاتری در نواحی که داده های جفتی بیشتری وجود دارد، برآورد می گردد.

بنابراین اگر این امکان وجود دارد که مقادیر y را برای مقادیر x انتخاب شده اندازه گیری کنیم، انتخاب مقادیر x ، بسته به هدف تحقیق تغییر خواهد کرد. اگر همه نقاط از اهمیت یکسانی برخوردار هستند که این امر اغلب در مورد خط راست صدق می کند، یک توزیع یکنواخت مقادیر x در کل دامنه x مطلوب خواهد بود. اگر هدف آزمایش برازش یک مدل ریاضی برای یک مجموعه داده است، و دانش اولیه در باره شکل مدل وجود ندارد، یک توزیع یکنواخت از مقادیر مورد نیاز است. اگر مقادیر خاصی به عنوان مثال یک مقدار مطلوب مورد نظر است و تا حدودی می دانیم که ممکن است این مقدار کجای محور x ها واقع شود. اگر بیشتر مقادیر x در اطراف محل مورد انتظار باشند و تعداد کمی از مقادیر در دو منتهی الیه واقع شوند، محل واقعی این مقدار با قطعیت بیشتری یافت خواهد شد.

برای حصول اطمینان از برازش مدل صحیح، تعداد کافی از داده ها در دامنه مقادیر x مورد نیاز است. تعداد مقادیر مختلف موجود برای x پیچیدگی مدلی را که می توان برای داده ها برازش نمود تعیین می کند (دراپر و اسمیت ۱۹۸۱). بین دو نقطه داده همیشه فقط یک خط راست عبور می کند ولی می توان تعداد نامحدودی منحنی های پیچیده تر، بین همان دو نقطه کشید. برای برازش