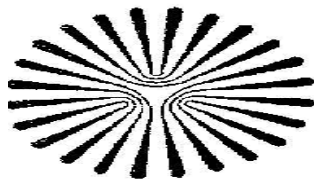


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه پیام نور

دانشکده علوم پایه و کشاورزی

مرکز تهران

پایان نامه

برای دریافت درجه کارشناسی ارشد

در رشته آمار ریاضی

گروه علمی آمار

عنوان پایان نامه:

مقایسه خطای نوع اول برخی آزمونهای ناپارامتری روی ضرائب مدلهای

رگرسیون چندگانه

طاهره محمدی

استاد راهنما: جناب آقای دکتر علی شادرخ

استاد مشاور: جناب آقای دکتر مسعود یارمحمدی

شهریور ۱۳۹۰

اینجانب طاهره محمدی دانشجوی ورودی سال ۸۶ مقطع کارشناسی ارشد رشته آمار ریاضی گواهی می‌نمایم چنانچه در پایان نامه خود از فکر، ایده 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 آن نباشد بر عهده خویش می‌دانم و جوابگوی آن خواهم بود.

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

نام و نام خانوادگی: طاهره محمدی

تاریخ و امضاء: ۱۳۹۰/۶/۱۳

اینجانب طاهره محمدی دانشجوی ورودی سال ۸۶ مقطع کارشناسی ارشد رشته آمار ریاضی گواهی می‌نمایم چنانچه براساس مطالب پایان نامه خود اقدام به انتشار مقاله، کتاب، و ... نمایم ضمن مطلع نمودن استاد راهنما، با نظر ایشان نسبت به نشر مقاله، کتاب، و ... و به صورت مشترک و با ذکر نام استاد راهنما مبادرت نمایم.

نام و نام خانوادگی : طاهره محمدی

تاریخ و امضاء: ۱۳۹۰/۶/۱۳

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

شهریور ۱۳۹۰

تقدیم به

امید چشمان منتظر

و

همسر عزیزم

تشکر و قدر دانی:

با تشکر از استاد گرانقدر جناب آقای دکتر شادرخ که بنده حقیر را با راهنماییهای ارزنده خود در تحریر این رساله یاری فرمودند و همچنین جناب آقای دکتر یار محمدی که نکات مفیدی را به بنده یادآوری نمودند. جا دارد از جناب آقای دکتر احتشامی همسرم که با صبر و حوصله فراوان در این امر، مرا یاری نمودند نیز تشکر بنمایم. امید است این رساله چراغ راهی برای محققان شود.

چکیده

روش های ناپارامتری متفاوتی برای انجام آزمون فرض روی ضرائب رگرسیونی ارائه شده است. در این پایان نامه ابتدا مهمترین آنها که از نتیجه مقایسات سایر آماردانان بدست آمده است را معرفی می کنیم. سپس در جستجوی روشی هستیم که از بقیه بهتر عمل کند. در این راستا هو و جون (۲۰۰۱) روش کندی را که در سال ۱۹۹۵ ارائه شده بود اصلاح نموده و با شبیه سازی نشان دادند که روش کندی اصلاح شده که ما آن را به نام روش هو - جون می نامیم، دارای خطای نوع اول کمتری است. از طرف دیگر شادرخ (۲۰۰۵) بطوری تئوری و شبیه سازی نشان داد که روش فریدمن - لان که یک روش نسبتاً بهتری نسبت به روش کندی است، دارای عملکرد بهتری می باشد. ما در این پایان نامه اصلاحاتی را که هو-جون روی روش کندی انجام دادند روی روش فریدمن-لان انجام دادیم و این روش اصلاح شده را با خود روش فریدمن-لان و هو-جون مقایسه کردیم و نتیجه گرفتیم که با افزایش حجم نمونه روش فریدمن-لان اصلاح شده دارای خطای نوع اول کمتری توان بالاتری نسبت به روش فریدمن-لان و هو-جون است.

واژگان کلیدی: آزمونهای جایگشتی، خطای نوع اول جایگشتی، تعویض پذیری، ضریب رگرسیونی

صفحه	فهرست
۱	مقدمه
۲	فصل اول: رگرسیون خطی
۳	مقدمه
۴	۱-۱) ضریب همبستگی
۵	۲-۱) معرفی رگرسیون خطی ساده
۹	۳-۱) رگرسیون خطی چندگانه
۱۳	۴-۱) بررسی مانده ها
۱۴	۵-۱) روشهای ناپارامتری
۱۶	فصل دوم: معرفی آزمون های جایگشتی
۱۷	مقدمه
۱۷	۱-۲) معرفی روش بوت استراپ و جک نایف
۱۹	۲-۲) روشهای انجام آزمون های جایگشتی
۲۱	۳-۲) آزمون های جایگشتی روی ضریب رگرسیون در مدل خطی ساده
۲۲	۴-۲) آزمون روی یک ضریب رگرسیونی در مدل خطی چندگانه
۲۷	۵-۲) روش فریدمن - لان در حالت چندگانه
۳۴	فصل سوم: شبیه سازی
۳۵	مقدمه

۳۵	۳-۱) روش شبیه سازی
۱۰۸	۳-۲) نتایج شبیه سازی
۱۱۲	نتیجه گیری کلی
۱۱۳	پیوست ۱) برنامه های شبیه سازی
۱۲۰	پیوست ۲) دستور گرام-اشمیت
۱۲۰	پیوست ۳) تجزیه طیفی
۱۲۲	منابع
۱۲۳	واژه نامه

مقدمه

درحالتی که فرضهای زیربنائی برای انجام آزمون فرض روی ضرائب مدلهای رگرسیونی خطی برقرار نباشد، انجام آزمونهای کلاسیک روی ضرائب رگرسیونی پاسخی مطمئن به ما نخواهد داد و لذا می‌بایست از روشهای ناپارامتری استفاده کرد. یک گروه از این آزمونها، آزمونهای دوباره نمونه‌گیری هستند که آزمونهای جایگشتی جزئی از این گروه می‌باشند که هدف این پایان نامه مقایسه برخی روشهای آزمونهای جایگشتی با روشهای اصلاح شده آنها می‌باشد.

آزمونهای جایگشتی در حدود سال ۱۹۳۰ توسط فیشر معرفی شد. اما به علت نیاز به محاسبات زیاد چندان مورد توجه قرار نگرفت. تا اینکه کامپیوترهای با سرعت بالا در دسترس قرار گرفت، با توجه به اینکه در اکثر مواقع فرضهای زیربنائی روی داده‌های حقیقی برقرار نبودند لذا آزمونهای جایگشتی جایگاه ویژه‌ای یافته و مورد توجه محققین رشته‌های اقتصاد، روانشناسی، علوم آزمایشگاهی و... قرار گرفت.

در فصل اول این پایان نامه مدل رگرسیون خطی معرفی شده و در رابطه با خطا و آزمونهای مربوط به آن می‌پردازد، در فصل دوم آزمونهای متفاوت جایگشتی و نیز روشهای اصلاح شده آنها را معرفی می‌کنیم، در فصل سوم به شیوه شبیه‌سازی روشهای معرفی شده را از دیدگاه خطای نوع اول مورد مقایسه قرار می‌دهیم.

فصل اول

رگرسیون خطی

مقدمه

واژه رگرسیون^۱ را نخستین بار فرانسیس گالتن^۲ (۱۸۲۲-۱۹۱۱) دانشمند انگلیسی در کتاب وراثت خود بکار برد (بهبودیان، جواد، ۱۳۸۳، تهران، پیام نور). در این فصل بنا داریم به طور مختصر رگرسیون خطی ساده و چندگانه^۳ را بیان و به بررسی مانده‌های^۴ حاصل از برازش خط رگرسیون پردازیم و توضیح بسیار کوتاهی در مورد آزمونهای ناپارامتری دهیم. آنهایی که در زمینه علوم زیستی یا تجربی یا اجتماعی مطالعه و پژوهش می‌کنند معمولاً با دو نوع متغیر برخورد دارند یک نوع آن اغلب قابل کنترل و نوع دیگر تصادفی می‌باشد مثلاً سن نوزاد (برحسب روز) یک نوع متغیر قابل کنترل ولی وزن او (برحسب گرم) متغیر تصادفی می‌باشد. کشف ارتباط و مدل سازی میان این گونه متغیرها، به منظور پیش بینی مقدار متغیر تصادفی وابسته و مطالعه‌ی عوامل مؤثر در مقدار آن، برای پژوهشهای علمی ضرورت دارد. البته ممکن است حالت‌های پیش بیاید که هر دو متغیر مستقل^۵ و وابسته^۶ متغیرهای تصادفی باشند که ما در این پایان نامه بدان نمی‌پردازیم، همچنین ممکن است متغیر وابسته متغیر کیفیت^۷ باشد که منجر به رگرسیون لوژستیک^۸ می‌گردد که پرداختن به آن خارج از بحث مربوط به این رساله می‌باشد.

^۱ Regression

^۲ Francis Galton

^۳ Multi linear regression

^۴ Residual

^۵ Independent

^۶ Dependent

^۷ Qualitative variable

^۸ Logistic regression

۱-۱ ضریب همبستگی

کواریانس^۱ بین X و Y معیاری عددی است که برای اندازه گیری تغییرات توأم دو متغیر تصادفی بکار می‌رود و به صورت زیر تعریف می‌شود:

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] = E(XY) - \mu_X \mu_Y \quad (1-1)$$

بطور شهودی می‌توان گفت که اگر احتمال رخ دادن مقادیر بزرگ X همراه با مقادیر بزرگ Y و مقادیر کوچک X همراه با مقادیر کوچک Y زیاد باشد، X و Y در یک جهت تغییر می‌کنند، در چنین وضعیتی هر دو مقدار انحراف $(X - \mu_X)$ و $(Y - \mu_Y)$ با احتمال زیادی مثبت یا منفی خواهند بود، به طوری که حاصلضرب $(X - \mu_X)(Y - \mu_Y)$ به احتمال قوی مثبت است. در نتیجه، امید ریاضی حاصلضرب مثبت است. از سوی دیگر، اگر X و Y گرایش به تغییر در جهت عکس هم داشته باشند، مقادیر مثبت $(X - \mu_X)$ اغلب متناظر با مقادیر منفی $(Y - \mu_Y)$ می‌شوند و بعکس. در نتیجه حاصلضرب به احتمال قوی منفی است و امید ریاضی منفی خواهد بود. با توجه به مطالب مذکور علامت و مقدار کواریانس نشان دهنده جهت و میزان بستگی بین X و Y است. این معیار به واحدهای اندازه گیری X و Y بستگی دارد برای داشتن معیاری که رابطه بین دو متغیر را اندازه گیری می‌کند که بستگی به واحدهای متغیرها نداشته باشد، کواریانس را بر انحراف معیار دو متغیر تقسیم می‌کنیم که این معیار را ضریب همبستگی بین X و Y می‌نامیم و بصورت ذیل محاسبه می‌نمائیم:

$$\rho = \rho(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} \quad (2-1)$$

که برآورد کننده گشتاوری آن به صورت ذیل است:

$$R = S_{XY} / S_X S_Y \quad (3-1)$$

$$S_{XY} = 1/n \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \overline{XY} - \bar{X} \cdot \bar{Y}$$

$$S_X^2 = 1/n \sum_{i=1}^n (X_i - \bar{X})^2 = \overline{X^2} - \bar{X}^2$$

$$S_Y^2 = 1/n \sum_{i=1}^n (Y_i - \bar{Y})^2 = \overline{Y^2} - \bar{Y}^2$$

^۱ Covariance

شایان ذکر است که محاسبه ضریب همبستگی بین دو متغیر با توجه به نوع متغیر متفاوت است انواع دیگری از ضریب همبستگی هم موجود است که با توجه به خلاصه بودن این مبحث از آوردن آنها خودداری می‌نمائیم.

۲-۱) معرفی رگرسیون خطی ساده

هر گاه بتوانیم به کمک یک تابع، از روی مقادیر یک متغیر، مقدار یک متغیر تصادفی دیگر را پیش بینی کنیم می‌گوییم یک تابع رگرسیون یا یک تابع برگشت داریم. به کمک نمودار پراکنش^۱ وجود رابطه (خطی) و پراکندگی داده‌ها بطور قابل تشخیص است. در صورتی که وجود رابطه خطی بین دو متغیر از روی نمودار پراکنش مشاهده شود می‌توان خط رگرسیون زیر را به داده‌ها برازش داد.

$$Y|x_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad (۱-۴)$$

که در فرمول فوق متغیر قابل کنترل را با مقدار ثابت x نشان می‌دهیم و آنرا متغیر مستقل^۲ می‌نامیم. متغیر دوم را که به متغیر x بستگی دارد و مقدار آن تصادفی می‌باشد با متغیر تصادفی $(Y|x)$ نشان می‌دهیم و آنرا متغیر وابسته^۳ می‌گوئیم. برای اینکه کلمه مستقل با مفهوم احتمالی آن اشتباه نشود، گاهی x را متغیر کنترل شده یا پیش بینی کننده یا بازگشت دهنده^۴ می‌گویند، متغیر دوم $(Y|x)$ را هم اغلب متغیر پاسخ^۵ می‌نامند.

یک متغیر تصادفی دیگر که با ε نشان داده می‌شود و آنرا خطای تصادفی یا به سادگی خطا^۶ می‌نامیم. فرضهای زیر بنایی مدل رگرسیونی ذکر شده به صورت ذیل می‌باشد:

(۱) ε_i ها ناهمبسته باشند.

(۲) ε_i ها دارای میانگین صفر و واریانس ثابت σ^2 هستند.

در صورتی که ε_i ها دارای توزیع نرمال با میانگین صفر و واریانس ثابت σ^2 باشند مدل (۱-۴) را یک مدل رگرسیونی نرمال (گاوس^۱) می‌نامند.

^۱ Scatter plot

^۲ Independent

^۳ Dependent

^۴ Controlled or predictor or regressor

^۵ Response

^۶ Error

۱-۲-۱) برآورد پارامترهای رگرسیون خطی ساده

برای برآورد کردن خط رگرسیونی لازم است پارامترهای آنرا با استفاده از داده‌های نمونه تصادفی $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ برآورد کنیم. برای انجام این کار دو روش وجود دارد. یکی روش حداقل مربعات^۲ (LS) و دیگری روش ماکسیمم درستنمایی^۳ است. روش حداقل مربعات را نخستین بار ریاضیدان برجسته آلمانی گاوس در محاسبات خود استفاده کرده است. به روایتی هم لژنارد^۴ ریاضیدان فرانسوی درباره این روش سه سال قبل از گاوس مقاله‌ای منتشر کرده است. به طور کلی در روش LS سعی می‌کنیم با مینیم کردن مجموع توانهای دوم خطا، پارامترهای مجهول را برآورد کنیم. دلیل مینیم کردن مجموع یاد شده این است که می‌خواهیم یک خط مستقیم مناسب بر داده‌های دوتائی برازش دهیم. در اینجا از ذکر محاسبات در بدست آوردن پارامترها با این روش صرف نظر کرده و مستقیماً برآوردها را بصورت ذیل بیان می‌کنیم.

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x} \quad \hat{\beta}_1 = \frac{S_{XY}}{S_X^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (5-1)$$

در نتیجه برآورد خط رگرسیون ساده به صورت زیر می‌شود:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad (6-1)$$

می‌دانیم که برآوردهای فوق برآوردهای نارایی برای پارامترهای خط رگرسیونی هستند. به منظور برآورد واریانس‌ها از مشاهده مانده‌ها استفاده می‌کنیم که تفاضل بین دو مقدار مشاهده شده y_i و پیش‌بینی شده \hat{y}_i بر اساس خط رگرسیونی می‌باشد این برآورد به صورت زیر است:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\epsilon}_i^2 = \frac{n}{n-2} (S_Y^2 - \hat{\beta}_1^2 S_X^2) \quad (7-1)$$

$$S_Y^2 = \sum (y_i - \bar{y})^2 \quad \text{که در آن}$$

درواقع تعبیر هندسی روش LS بدین صورت است که خط برازش داده شده به داده‌ها به گونه‌ای از میان آنها عبور می‌کند که خطاها دارای کمترین مقدار شود.

^۱ Gauss

^۲ Least squares Method

^۳ Maxim Likelihood Estimator

^۴ Legendr

۲-۲-۱) توزیع برآوردها و آماره های آزمون

با توجه به فرض های که درباره مدل ε_i ها داشتیم و Y_i ها هم مستقل و دارای توزیع نرمال هستند داریم:

$$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2) \quad (۸-۱)$$

اگر تابع درستنمایی^۱، یعنی چگالی توأم را برای Y_1, Y_2, \dots, Y_n را تشکیل دهیم و ماکزیم کنیم، برآوردهای بدست آمده در رابطه مذکور برآوردهای درستنمایی ماکزیم هستند که هر سه با برآوردهای که به روش حداقل مربعات خطا بدست آمده یکی هستند و نیز این برآوردها دو به دو مستقل اند و با توجه به اینکه $\hat{\beta}_1$ و $\hat{\beta}_0$ ترکیب خطی از Y_i ها که نرمال هستند پس آنها دارای توزیع نرمال می باشند و داریم:

$$\left\{ \begin{array}{l} \hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{n s_x^2}\right) \\ \hat{\beta}_0 \sim N\left(\beta_0, \frac{\sigma^2}{n} \left(1 + \frac{\bar{x}^2}{s_x^2}\right)\right) \\ \frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{(n-2)} \end{array} \right. \quad (۹-۱)$$

با توجه به توزیع های فوق و تعریف توزیع t-استودنت^۲ آماره های ذیل دارای توزیع t-استودنت با n-2 درجه آزادی^۳ می باشند.

$$T_1 = \frac{\sqrt{n}(\hat{\beta}_1 - \beta_1)}{\frac{\hat{\sigma}}{s_x}} = T_{(n-2)} \quad (۱۰-۱)$$

$$T_2 = \frac{\sqrt{n}(\hat{\beta}_0 - \beta_0)}{\hat{\sigma} \sqrt{1 + \frac{\bar{x}^2}{s_x^2}}} = T_{(n-2)} \quad (۱۱-۱)$$

بنابراین با توجه به نکات فوق می توان به راحتی در مورد پارامترهای مدل خطی ساده آزمونهای آماری ذیل را انجام دهیم و یا برای آنها فواصل اطمینان بدست آوریم.

^۱ Likelihood function

^۲ Students t distribution

^۳ Degrees of freedom

$$\begin{cases} H_0: \beta_0=0 \\ H_1: \beta_0 \neq 0 \end{cases} \quad \begin{cases} H_0: \beta_1 =0 \\ H_1: \beta_1 \neq 0 \end{cases}$$

۳-۲-۱) تحلیل واریانس

در این بخش به اختصار تحلیل واریانس برای رگرسیون خطی ساده را بیان می‌کنیم. برای اینکه مدل رگرسیون خطی ساده را با پارامترهای کامل و با دو مدل که هر کدام از حذف یکی از پارامترهای β_0, β_1 بدست آمده‌اند، مقایسه کنیم از تحلیل واریانس استفاده می‌نمائیم، با فرض اینکه مدل کامل است جمله زیر را به دو قسمت تفکیک می‌نمائیم.

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n [(Y_i - \hat{Y}_i) + (\hat{Y}_i - \bar{Y})]^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 \quad (12-1)$$

که جملات را با نمادهای زیر نشان می‌دهیم:

$$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2 = nS_{\bar{Y}}^2 \quad (13-1)$$

تغییر پذیری نسبت به معدل را می‌سنجد درجه آزادی آن (n-1) در مدل شامل هر دو پارامتر تغییر پذیری خطا با (n-2) درجه آزادی است.

$$SSE = \sum_{i=1}^n \hat{E}_i^2 = (n-2) \hat{\sigma}^2 \quad (14-1)$$

و جمله دیگر را تغییر پذیری رگرسیون با ۱ درجه آزادی می‌نامیم هر چقدر نسبت تغییر پذیری رگرسیون نسبت به خطا بیشتر باشد برازش داده‌ها با این خط رگرسیونی موفقیت آمیزتر بوده است.

$$SSR = n\hat{\beta}_1^2 s_x^2 \quad (15-1)$$

نسبت SSE به واریانس و نسبت SSR به واریانس دارای توزیع کی دو است و با توجه به رابطه توزیع F-فیشر^۱ و کی دو^۲، به منظور آزمون اینکه $\beta_1=0$ از آماره زیر استفاده می‌نمائیم:

$$MSE = \frac{SSE}{df_E} \quad MSR = \frac{SSR}{df_R} \quad F = \frac{MSR}{MSE} \quad (16-1)$$

^۱ F-Fisher distribution

^۲ Chi-square distribution

ناحیه رد(بحرانی)^۱ H_0 به صورت ذیل می‌باشد.

$$C. R_{H_0}: F > F_{(1-\alpha)}(1, n - 2) \quad (17-1)$$

محاسبه نسبت SSR به SST نیز برای درک این مطلب که Y را می‌توان به کمک خط رگرسیون خطی پیش بینی کرد لازم است، این نسبت را ضریب تعیین^۲ می‌نامند.

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \quad (18-1)$$

در واقع استراتژی کلی به صورت ذیل است که:

۱. اجرای مدل کامل و بدست آوردن SSE

۲. اجرای مدل کاهش یافته بر اساس H_0 و بدست آوردن SSE

۳. محاسبه آماره آزمون

۴. تصمیم گیری

(بزرگ نیا، آذرنوش، ۱۳۸۲، ۳۹۳ الی ۴۳۲)

۳-۱) رگرسیون خطی چندگانه

در دو بخش قبل در مورد اینکه یک متغیر وابسته تنها به یک متغیر کنترل شده به صورت خطی ارتباط دارد مطالبی را ارائه کرده ایم، اینک می‌خواهیم متغیر وابسته‌ای را مورد بررسی قرار دهیم که به صورت خطی به چندین متغیر کنترل شده ارتباط دارد. به عبارتی دیگر می‌خواهیم در مورد رگرسیون خطی چند متغیره یا رگرسیون چند گانه به اختصار صحبت کنیم.

تعریف مدل خطی: فرض کنید x_1, x_2, \dots, x_k متغیرهای کنترل شده و متغیر تصادفی شرطی $(Y|x_1, x_2, \dots, x_k)$ متغیر وابسته باشد همچنین متغیرهای کنترل شده بدون خطای اندازه گیری و متغیر وابسته دارای خطای اندازه گیری و خطای تصادفی که روی هم با ε نشان می‌دهیم باشند رابطه

^۱ Critical level

^۲ Coefficient of determination

^۳ Multiple linear regression

زیر یک مدل خطی k گانه با پارامترهای $\beta_1, \beta_2, \dots, \beta_k$ را نشان می‌دهد. ε دارای میانگین صفر و واریانس ثابت σ^2 می‌باشد که اگر نرمال باشد یک مدل نرمال داریم.

$$\{Y|X_1, X_2, \dots, X_n\} = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon \quad (19-1)$$

اگر n بار آزمایش را انجام دهیم داده‌ها را می‌توان صورت ماتریسی ذیل نمایش داد با این فرض که آزمایشها مستقلا انجام می‌گیرد ε_i ها مستقل و هم توزیع اند.

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_k \end{pmatrix} \quad X = \begin{pmatrix} x_{11} & \dots & x_{1k} \\ \vdots & \dots & \vdots \\ x_{n1} & \dots & x_{nk} \end{pmatrix}$$

اینک می‌توان n رابطه ی (۱۹-۱) را که بصورت ماتریسی بیان شد، بطور خلاصه ذیل نشان دهیم:

$$Y = X\beta + \varepsilon \quad E(\varepsilon) = 0, \quad \sum \varepsilon = \sigma^2 I \quad (21-1)$$

جهت برآورد پارامترهای مدل با روش مجموع توانهای دوم خطا و استفاده از معادلات نرمال^۱ به شرط اینکه ماتریس XX' وارون پذیر باشد (عدم وجود هم خطی) داریم:

$$(XX')^{-1} XY = \hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} \quad (22-1)$$

با توجه به فرضیات مدل، ماتریس کوواریانس $\hat{\beta}$ به صورت زیر است:

$$\sum \hat{\beta} = \sigma^2 (XX')^{-1} \quad (23-1)$$

اما در رابطه (۲۳-۱) مقدار σ^2 مجهول است یک برآورد نااریب برای آن به صورت زیر است:

$$\hat{\sigma}^2 = \frac{\|Y - X\hat{\beta}\|^2}{n-k} \quad (24-1)$$

و بردار میانگین و برآورد آنرا بصورت زیر نشان می‌دهیم:

$$\mu = E(Y) = X\beta \quad \hat{\mu} = X\hat{\beta} = X(XX')^{-1}XY \quad (25-1)$$

اگر رابطه (۲۱-۱) را، مدل (۱) در نظر بگیریم. مدل زیر را با عنوان (۲) معرفی می‌کنیم.

^۱ Normal equations

$$Y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I) \quad (26-1)$$

برآورد پارامترها در مدل (۲) همانند مدل (۱) می باشد با ذکر این نکته که برآوردها دارای توزیع های به صورت زیر می باشند:

$$\left[\begin{array}{l} \hat{\beta} \sim N(\beta, \sigma^2 (X'X)^{-1}) \\ \hat{\mu} \sim N(\mu, \sigma^2 X(X'X)^{-1}X') \\ \frac{(n-k)\hat{\sigma}^2}{\sigma^2} \Rightarrow \chi^2_{(n-k)} \end{array} \right. \quad (27-1)$$

(بزرگ نیا، آذرنوش، ۱۳۸۲، ۱۶۹، ۱۷۴)

۱-۳-۱) استنباط آماری در مدل های خطی چند متغیره

مدل خطی (۲) و ماتریس M راکه $h \times k$ و پررتبه سطری^۱ (سطرها مستقل باشند) می باشد در نظر می گیریم تحت فرض آماری $H_0: M\beta = 0$ که آنرا مدل (۳) در نظر می گیریم. برآورد پارامترهای مدل با استفاده از روش نسبت درستنمایی تعمیم یافته^۲ به صورت ذیل است:

$$\hat{\sigma}^2 = \frac{\|Y - \hat{\mu}\|^2}{n-k-h}, \quad \hat{\mu} = \text{proj}_W^Y \quad (28-1)$$

که در آن برآورد میانگین تصویر بردار Y روی فضای ایجاد شده توسط سطرهاى ماتریس M است و طبق قضیه تصویری بردار میانگین یکتاست.

۱-۳-۲) تحلیل واریانس

تحلیل واریانس در مدل چندمتغیره به همان صورت مدل خطی ساده صورت می گیرد، فقط به جای اسکالرها از معادل برداری یا ماتریسی آنها استفاده می شود که بصورت زیر بیان می نمائیم.

الف) مجموع توانهای دوم خطا در مدل (۲)

$$SSE = \|Y - \hat{\mu}\|^2 = \|Y\|^2 - \|\hat{\mu}\|^2 \quad (29-1)$$

^۱ Full row rank

^۲ Generalized likelihood ratio

ب) مجموع توانهای دوم خطا در مدل (۳) تحت فرض H_0

$$SST = \|Y - \hat{\mu}\|^2 = \|Y\|^2 - \|\hat{\mu}\|^2 \quad (30-1)$$

ج) مجموع توانهای دوم رگرسیونی

$$SSR = \|\hat{\mu} - \hat{\mu}\|^2 = \|\hat{\mu}\|^2 - \|\hat{\mu}\|^2 \quad (31-1)$$

با استفاده از فرم کانونی Y^1 ، تحت فرض صفر داریم:

$$SSE = \sum_{i=k+1}^n Z_i^2 \stackrel{d}{\Rightarrow} \sigma^2 x_{(n-k)}^2 \quad (32-1)$$

$$SST = \sum_{i=k-h+1}^n Z_i^2 \stackrel{d}{\Rightarrow} \sigma^2 x_{(n-k+h)}^2 \quad (33-1)$$

$$SSR = \sum_{i=k-h+1}^k Z_i^2 \stackrel{d}{\Rightarrow} \sigma^2 x_{(h)}^2 \quad (34-1)$$

ملاحظه می‌شود با تعاریف فوق SSR و SSE همواره مستقل هستند. بدین ترتیب برآورد نارایب σ^2 در دو مدل (2) و (3) تحت فرض صفر به صورت زیر است:

I) در مدل (2): $\hat{\sigma}^2 = \frac{SSE}{n-k} = MSE$ دارای توزیع χ^2 با $n-k$ درجه آزادی

II) در مدل (3): $\hat{\sigma}^2 = \frac{SST}{n-k+h} = MST$ دارای توزیع χ^2 با $n-k+h$ درجه آزادی

باز هم مانند بخش (۳-۲-۱) جدول تحلیل واریانس را تشکیل داده و با توجه به آماره F تعریف شده در رابطه (۱۶-۱) در مورد فرض صفر تصمیم می‌گیریم. در مورد زمانی که بردار X پرتبه ستونی^۲ نباشد معادلات نرمال باز هم دارای پاسخ هستند، زیرا به هر حال تصویر y روی فضای ستونی^۳ X همواره وجود دارد. اما پاسخ آن یکتا نیست. این موضوع منجر به بحث وارون گسترده^۴ می‌گردد که از حوصله این بخش خارج است.

(بهبودیان، جواد، ۱۳۸۳، رگرسیون، تهران، پیام نور)

^۱ Canonical form
^۲ Full column rank
^۳ column space
^۴ Generalized inverse