



دانشگاه صنعتی شاهرود

دانشکده علوم ریاضی

پایان نامه جهت اخذ درجه کارشناسی ارشد

آمار ریاضی

بررسی رفتار برآوردگرهای رگرسیونی ریج در مدل های خطی منفرد

نگارش

سمانه نجاریان

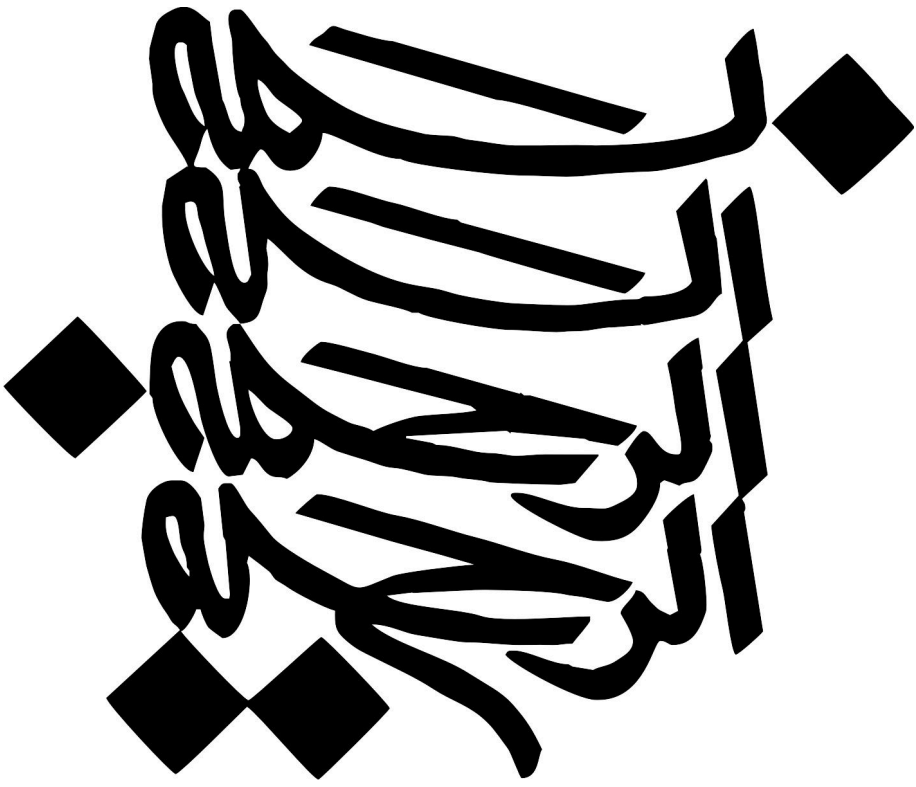
استاد راهنما

دکتر محمد آرشی

استاد مشاور

دکتر جعفر فتحعلی

شهریور ۱۳۹۰



تقدیم به:

همه آنانی که

به من اندیشیدن آموختند

تا امروز

کورسوی ظلمت را از روشنایی معرفت

بازشناسم.

فهرست مقالات مستخرج از متن پایان نامه

[۱] نجاریان. س. و آرشی. م. (۱۳۸۹)، شبیه سازی در مدل رگرسیون خطی ریج منفرد، ششمین همایش ملی آمار - دانشگاه پیام نور اهواز.

[2] Najarian. S. and Arashi. M. (2010), Ridge Estimation In Singular Linear Model: New Approach, The 10th Iranian Statistical Conference.

[3] Najarian. S., Arashi. M. and Golam Kibria. B. M. (2011), A Simulation Study On Some Ridge Regression Estimators, Communication in statistics- simulation Restricted and Computation.

پیشگفتار

در دهه‌های اخیر، تلاش‌های زیادی در راستای یافتن برآوردهای بهبودیافته انجام گرفته است. یکی از مسائلی که زمینه‌ساز این تلاش بوده، همخطی چندگانه می‌باشد. وقتی که ستون‌های ماتریس طرح وابسته خطی باشند، برآوردهای کمترین توان دوم بدست آمده به روش گوس-مارکف به دلیل نزدیک به صفر بودن دترمینان ماتریس $X'X$ ، کارایی لازم را نخواهند داشت. این مطلب منجر به ارائه برآوردهای بهبودیافته از جمله برآوردهای نوع ريج شده است. هدف از این پایان‌نامه، ارائه برآوردهای ريج در مدل‌های خطی محدودشده و محدودنشده و تحلیل و تفسیر ویژگی‌های برآوردهای ريج شده، است. این مجموعه شامل ۵ فصل و ۴ ضمیمه می‌باشد. مطالب هر فصل بطور مختصر عبارتند از:

در فصل ۱، مقدمه‌ای بر تحلیل رگرسیون، ارائه و چگونگی برآورد ضرایب رگرسیونی و همچنین چگونگی بروز همخطی در مدل‌های رگرسیون چندگانه خطی و روش‌های تشخیص همخطی بررسی شده است.

در فصل ۲، برآوردهای رگرسیونی ريج در مدل خطی محدودشده مورد بررسی قرار گرفته و خطای برآوردهای محاسبه شده‌است و به کمک شبیه‌سازی و مثال عددی کارایی برآوردهای ريج و برتری آنها بر برآوردهای کمترین توان‌های دوم محدودشده نشان داده شده است.

در فصل ۳، برآوردهای رگرسیونی ريج محدودشده و محدودنشده در مدل خطی منفرد ارائه شده و ویژگی‌های برآوردهای ريج بررسی شده است.

در فصل ۴، برآوردهای رگرسیونی ريج در مدل خطی منفرد با در نظر گرفتن محدودیتی که دارای جمله خطای

تصادفی است، ارائه و ویژگی‌های برآوردگر بررسی و به کمک شبیه‌سازی کارایی برآوردگر تایید شده است. در فصل ۵، برآوردگر رگرسیونی ریج منفرد در مدل به ظاهر نامرتبط (SUR) ارائه و ویژگی‌های برآوردگر بطور خلاصه بیان شده است.

در ضمیمه A ، بحث و نتیجه‌گیری و پیشنهادها برای ادامه کار در آینده، مطرح شده‌اند. در ضمیمه B ، عملیات جبری ماتریس‌ها بمنظور سادگی درک اثبات قضایا و لم‌ها در متن پایان‌نامه بیان شده است.

در ضمیمه C ، توضیحاتی در خصوص عامل تورم واریانس به عنوان عاملی برای شناخت همخطی در مدل رگرسیونی خطی داده شده است.

در ضمیمه D ، توضیحات کاملتری درباره عدد شرطی که عاملی دیگر برای تشخیص مسئله همخطی است، ارائه شده است.

* در این مجموعه، قضایایی که برهان آن از نویسنده این پایان‌نامه می‌باشد با علامت ستاره مشخص شده است.

مدلهای ارائه شده در این مجموعه

$Y = X\beta + E$ $R\beta = r$ $E(E) = 0 \quad Cov(E) = \sigma^2 I$ $Y \sim n \times 1 \quad X \sim n \times p \quad \beta \sim p \times 1$ $E \sim n \times 1 \quad R \sim q \times p \quad r \sim q \times 1$	فصل دوم
$Y = X\beta + E$ $R\beta = r$ $E(E) = 0 \quad Cov(E) = \sigma^2 \Sigma$ $Y \sim n \times 1 \quad X \sim n \times p \quad \beta \sim p \times 1$ $E \sim n \times 1 \quad R \sim q \times p \quad r \sim q \times 1$	فصل سوم
$Y = X\beta + E$ $R\beta + v = r$ $E(E) = 0 \quad Cov(E) = \sigma^2 I$ $Y \sim n \times 1 \quad X \sim n \times p \quad \beta \sim p \times 1$ $E \sim n \times 1 \quad v \sim q \times 1$ $R \sim q \times p \quad r \sim q \times 1$	فصل چهارم
$Y = X\beta + E$ $R\beta = r$ $E(E) = 0 \quad Cov(E) = \Sigma \otimes I_q$ $Y \sim n \times 1 \quad X \sim n \times k \quad \beta \sim k \times 1$ $E \sim n \times 1 \quad R \sim q \times p \quad r \sim q \times 1$ $k = \sum_{i=1}^p k_i \quad n = pq$	فصل پنجم

قدردانی

منت خدای را عزوجل که طاعتش موجب قربت است و به شکر اندرش مزید نعمت

در پایان لازم می‌دانم از استاد گرامی، جناب آقای دکتر آرشی که راهنمایی و هدایت این پایان‌نامه را عهده‌دار بودند و تلاش زیادی در راستای جهت‌دهی به تفکرات و تلاش‌های اینجانب داشته‌اند، قدردانی کنم. همچنین

از جناب آقای دکتر فتحعلی جهت مشاوره و راهنمایی‌های بی‌دریغشان کمال تشکر را دارم.

از اساتید داور محترم، آقایان دکتر داوود شاهسونی (دانشگاه صنعتی شاهرود) و دکتر مهدی عمادی (دانشگاه فردوسی مشهد) که با حضور گرمشان تصحیح و داوری پایان‌نامه را بر عهده گرفتند، سپاسگزارم.

از جناب آقای پروفسور کیبیریا (دانشگاه فلوریدای آمریکا) بخاطر راهنمایی‌هایشان در استخراج مقاله‌ای از فصل یک این پایان‌نامه، تشکر می‌کنم.

از جناب آقای دکتر فکور و آقای دکتر روزبه بخاطر همکاری‌شان با بنده کمال تشکر را دارم.

همچنین از مسئولین کتابخانه دکتر فاطمی دانشگاه فردوسی مشهد بخصوص آقای مسئله‌گو بخاطر همکاری صمیمانه‌شان متشکرم.

سمانه نجاریان

تابستان ۱۳۹۰

علائم اختصاری

Mean Square Error میانگین توان دوم خطا : *MSE*

Mean Square Error Matrix ماتریس میانگین توان دوم خطا : *MSEM*

Least Squares Estimator برآوردگر کمترین توان‌های دوم : *LSE*

Restricted Least Squares Estimator برآوردگر کمترین توان‌های دوم محدودشده : *RLSE*

Generalized Least Squares Estimator برآوردگر کمترین توان‌های دوم محدودشده : *GLSE*

Ridge Estimator برآوردگر ریج : *RE*

Restricted Ridge Estimator برآوردگر ریج محدودشده : *RRE*

Restricted Linear Model مدل خطی محدودشده : *RLM*

Stein Estimator برآوردگر استاین : *SE*

Lui Estimator برآوردگر لویی : *LE*

Principal Components Estimator برآوردگر مولفه‌های اصلی : *PCE*

Partial Least Squares Estimator برآوردگر کمترین توان‌های دوم جزئی : *PLSE*

Almost Unbiased Estimator برآوردگر تقریباً نارایب : *AUE*

Best Linear Unbiased Estimator بهترین برآوردگر نارایب خطی : *BLUE*

Seemingly Unrelated Regression مدل به ظاهر نامرتبط : *SUR*

Asymptotic Normal *AN* : نرمال مجانبی

Variance Inflation Factor *VIF* : عامل تورم واریانس

Condition Number *CN* : عدد شرطی

فهرست مطالب

۱	مقدمه و دورنما	۱
۳	۱.۱ مدل رگرسیونی چندگانه	۳
۴	۲.۱ برآورد پارامترهای مدل	۴
۶	۳.۱ همخطی در رگرسیون خطی چندگانه	۶
۸	۴.۱ آثار همخطی چندگانه	۸
۹	۵.۱ شاخص‌های همخطی چندگانه	۹
۹	۱.۵.۱ محک ماتریس همبستگی	۹
۱۰	۲.۵.۱ عامل تورم واریانس	۱۰
۱۰	۳.۵.۱ عدد شرطی	۱۰
۱۲	۲ برآوردگر رگرسیونی رنج محدودشده	۱۲
۱۶	۱.۲ معرفی برآوردگر	۱۶
۲۰	۲.۲ محاسبه خطای برآوردگر	۲۰
۲۷	۳.۲ برآورد مقادیر k_i	۲۷
۲۹	۴.۲ شبیه‌سازی	۲۹
۳۵	۵.۲ کاربرد	۳۵

۳۶	مثال ۱	۱.۵.۲
۳۸	مثال ۲	۲.۵.۲
۳۹	مثال ۳	۳.۵.۲
۴۰	نتیجه‌گیری	۶.۲
۴۲	۳ برآوردگر رگرسیونی ریح منفرد	
۴۴	معرفی برآوردگر	۱.۳
۴۶	ویژگی‌های برآوردگر	۲.۳
۵۸	شبیه‌سازی	۳.۳
۶۰	نتیجه‌گیری	۴.۳
۶۱	۴ برآوردگر ریح منفرد در مدل با محدودیت تصادفی	
۶۲	معرفی برآوردگر	۱.۴
۶۶	ویژگی‌های برآوردگر	۲.۴
۶۹	شبیه‌سازی	۳.۴
۷۱	نتیجه‌گیری	۴.۴
۷۲	۵ برآوردگر ریح منفرد در مدل <i>SUR</i>	
۷۶	معرفی برآوردگر	۱.۵
۷۹	ویژگی‌های برآوردگر	۲.۵
۸۰	نتیجه‌گیری	۳.۵
۸۱	الف خلاصه و پیشنهادات برای تحقیقات آینده	
۸۱	الف.۱ خلاصه	

۸۲	الف. ۲. پیشنهادات برای پژوهش‌های آینده (ارائه زمینه‌های تحقیق)
۸۴	ب جبر ماتریس‌ها
۸۹	پ عامل تورم واریانس
۹۱	ت عدد شرطی
۹۴	مراجع

فصل ۱

مقدمه و دورنما

مقدمه

تحلیل رگرسیونی یکی از روش‌های آماری برای تحلیل داده‌های چندعاملی است که حوزه کاربرد آن بیشترین وسعت را دارد، نتایج پرجاذبه آن از نظر مفهومی، فرآیند ساده بکارگیری یک معادله است که ارتباط بین دسته‌ای از متغیرها را بیان می‌کند. تحلیل رگرسیونی همچنین از جهت نظری به لحاظ ظرافت و زیبایی ریاضیات دارای جذابیت است. موفقیت در تحلیل، نیاز به درک و تیزبینی در دو مقوله تئوری و مسائل عملی دارد و زمانی بروز می‌کند و خود را نشان می‌دهد که تکنیک در بکارگیری داده‌ها و اطلاعات جهان واقع قرار می‌گیرد.

تحلیل رگرسیونی فن و تکنیکی آماری برای بررسی و مدل‌بندی ارتباط بین متغیرهاست. کاربردهای رگرسیون متعدد است و تقریباً در هر زمینه‌ای از جمله مهندسی، فیزیک، اقتصاد، مدیریت، علوم زیستی و بیولوژی و علوم اجتماعی صورت می‌پذیرد.

کلمه "رگرسیون" در لغت به معنای "برگشت" است. دلیل این نامگذاری به مطالعات "گالتون" در مورد رابطه قد پسر و قد پدر مربوط می‌شود. گالتون با رسم نمودار قد پسر در مقابل قد پدر (نمودار پراکنش) براساس حجم بزرگی از داده‌ها به این نتیجه رسید که اغلب افراد بلندقد، پسرانی کوتاه‌قد و پدران کوتاه‌قد، پسرانی بلندقدتر از خود دارند و این پدیده را "برگشت به میانگین" نامید و از آن پس نام رگرسیون بر روش‌های بررسی روابط بین متغیرهای آماری باقی ماند.

به عنوان مثالی از یک مسئله که در آن تحلیل رگرسیونی می‌تواند مفید واقع شود، فرض می‌کنیم یک مهندس صنایع، توسط یک سازنده نوشابه استخدام شده که محصول تحویلی و عملیات سرویس و خدمت رسانی برای ماشین‌های فروش را تجزیه و تحلیل کند. او گمان می‌کند که زمان لازم برای اینکه یک فرد تحویل دهنده، این ماشین را سرویس دهد بستگی به تعداد موارد محصول تحویل شده دارد. مهندس ۲۵ فروشنده جزئی را که دارای ماشین فروش می‌باشند، به تصادف بازدید کرده است و زمان تحویل (برحسب دقیقه) و حجم محصول تحویل داده شده (بر حسب مورد) برای هر یک مشاهده می‌شود. بوسیله نمودار پراکنش، ارتباط بین

زمان تحویل و حجم تحویل مشخص می‌شود. نقاطی که نشان‌دهنده داده‌ها هستند در امتداد یک خط مستقیم قرار می‌گیرند، نه دقیقا روی خط. اگر Y نشان‌دهنده زمان تحویل و X را برای حجم تحویل شده قرار دهیم، در این صورت معادله یک خط مستقیم که این دو متغیر را به هم مربوط می‌سازد، چنین است

$$Y = \beta_0 + \beta_1 X \quad (1.1)$$

که β_0 عرض از مبدا و β_1 شیب خط است. از آنجایی که نقاط داده‌ها دقیقا روی یک خط مستقیم قرار نمی‌گیرند بنابراین (۱.۱) بایستی به عنوان برآوردکننده آن اصلاح و تعدیل شود. اگر اختلاف بین مقدار مشاهده‌شده Y و خط مستقیم $(\beta_0 + \beta_1 X)$ را خطای E قرار می‌دهیم، مناسب است که خطای E به عنوان یک خطای آماری تلقی شود. بدین معنی که متغیری تصادفی است که عدم برازش را بیان و اندازه ناتوانی مدل از برازش دقیق داده‌ها را برآورد می‌کند. این خطا ممکن است به لحاظ اثرات دیگر متغیرها، خطاهای اندازه‌گیری و غیره روی زمان تحویل صورت پذیرد. بنابراین مدلی با قابلیت پذیرش بیشتر برای داده‌های زمان تحویل عبارتست از

$$Y = \beta_0 + \beta_1 X + E \quad (2.1)$$

معادله (۲.۱) یک مدل رگرسیون خطی نامیده می‌شود. در این مدل X متغیر مستقل و Y متغیر وابسته نامیده می‌شود. در هر حال این نامگذاری باعث می‌شود که این مفهوم با مفهوم استقلال آماری اشتباه شود. چون در (۲.۱) فقط یک متغیر رگرسیونی وجود دارد، این مدل را مدل رگرسیونی ساده نامیده‌اند.

۱.۱ مدل رگرسیونی چندگانه

در حالت کلی ممکن است متغیر پاسخ به k متغیر رگرسیونی X_1, X_2, \dots, X_k مربوط باشد، بنابراین مدل رگرسیونی به صورت زیر نوشته می‌شود

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + E \quad (3.1)$$

چون در این مدل رگرسیونی بیش از یک متغیر درگیر است، رگرسیونی چندمتغیره نام دارد. پارامترهای β_j ($j = 0, 1, \dots, k$) ضرایب رگرسیون نامیده می‌شود. بکارگیری صفت خطی بودن در مدل فوق اشاره بر این دارد که مدل، برحسب پارامترهای $\beta_0, \beta_1, \dots, \beta_k$ خطی است و نه به دلیل اینکه Y تابعی خطی از X هاست. این مدل یک ابرصفحه در فضای k بعدی از متغیرهای رگرسیونی X_j است. پارامتر β_j نشان‌دهنده تغییرات مورد انتظار متغیر پاسخ به ازای یک واحد تغییر در X_j است و قتیکه همه متغیرهای رگرسیونی باقیمانده دیگر X_i ($i \neq j$) ثابت نگهداشته شوند. به همین جهت پارامترهای β_j ($j = 1, 2, \dots, k$) ضرایب جزئی رگرسیون نامیده می‌شوند. مدل‌های رگرسیون چندگانه، اغلب به عنوان تقریب توابع بکار می‌روند. بدین معنی که ارتباط تابعی واقعی بین X_1, X_2, \dots, X_k و Y شناخته شده نیست، اما مدل رگرسیون خطی روی دامنه تغییرات متغیرهای رگرسیونی، تقریب مناسبی از ارتباط تابعی مذکور می‌باشد.

۲.۱ برآورد پارامترهای مدل

برای برآورد ضرایب رگرسیون، روش کمترین توان‌های دوم بکار می‌رود. فرض می‌کنیم $n > k$ مشاهده در دسترس است و Y_i نمایش دهنده i امین پاسخ مشاهده شده و X_{ij} نمایش دهنده i امین مشاهده در j امین سطح متغیر رگرسیونی باشد. فرض می‌کنیم جمله خطای E در مدل دارای $E(E) = 0$ و $Var(E) = \sigma^2$ و خطاها ناهمبسته هستند. می‌توان مدل متناظر با (۳.۱) را بصورت زیر نوشت

$$\begin{cases} Y = X\beta + E \\ E(E) = 0, \quad Cov(E) = \sigma^2 I \end{cases} \quad (4.1)$$

تابع کمترین توان‌های دوم خطا بصورت زیر نوشته می‌شود

$$\begin{aligned} S(\beta_0, \beta_1, \dots, \beta_k) &= \sum_{i=1}^n E_i^2 \\ &= \sum_{i=1}^n \left(Y_i - \beta_0 - \sum_{j=1}^k \beta_j X_{ij} \right)^2 \\ &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

با مشتقگیری از S نسبت به $\beta_0, \beta_1, \dots, \beta_k$ و مساوی با صفر قرار دادن مشتق حاصل، در صورت نامنفرد بودن

ماتریس $\mathbf{X}'\mathbf{X}$ ، برآوردگرهای ضرایب مدل رگرسیونی را بصورت زیر بدست می‌آوریم

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \quad (5.1)$$

بردار مقادیر برازش شده \hat{Y}_i با مقادیر مشاهده شده Y_i چنین خواهد بود

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{H}\mathbf{Y} \quad (6.1)$$

ماتریس $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$ ، ماتریسی $p \times p$ ، متقارن و خودتوان است و بنابراین $\text{Rank}(\mathbf{H}) =$

$\text{tr}(\mathbf{H}) = p$ که ماتریس هت یا برازش نامیده می‌شود زیرا که بردار مقادیر مشاهدات را به بردار مقادیر

برازش شده تصویر می‌کند. ماتریس برازش و خواص آن نقشی مرکزی در تحلیل رگرسیونی بازی می‌کند.

ماتریس واریانس-کوواریانس $\hat{\boldsymbol{\beta}}$ عبارتست از

$$\text{Cov}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \quad (7.1)$$

همچنین برآوردگر ناریب σ^2 عبارتست از میانگین توان‌های دوم باقیمانده‌ها

$$\hat{\sigma}^2 = \frac{SSE}{n-p} \quad (8.1)$$

که در آن $SSE = (Y - X\beta)'(Y - X\beta)$ ، مجموع توان‌های دوم باقیمانده‌هاست. حال اگر در مدل (۴.۱)، $Cov(E) = \sigma^2 \Sigma$ باشد، که در آن Σ ماتریسی معیت مثبت است. برآوردگر ضرایب رگرسیونی و ماتریس واریانس-کوواریانس برآوردگرها و همچنین برآوردگر ناریب σ^2 به ترتیب به صورت زیر خواهند بود

$$\hat{\beta} = (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}Y \quad (9.1)$$

و

$$Cov(\hat{\beta}) = \sigma^2(X'\Sigma^{-1}X)^{-1} \quad (10.1)$$

و

$$\hat{\sigma}^2 = \frac{(Y - X\hat{\beta})'\Sigma^{-1}(Y - X\hat{\beta})}{n - p} \quad (11.1)$$

در این صورت مقادیر برازش شده به فرم زیر می‌باشند

$$\hat{Y} = X\hat{\beta} = H\Sigma^{-1/2}Y \quad (12.1)$$

که در آن $H = X\Sigma^{-1/2}(X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1/2}$

۳.۱ همخطی در رگرسیون خطی چندگانه

مدل‌های رگرسیونی در حوزه کاربردی وسیعی مورداستفاده قرار می‌گیرند. یک مسئله جدی که می‌تواند استفاده از مدل رگرسیونی را با اشکال مواجه کند، همخطی چندگانه یا وابستگی خطی نزدیک بین متغیرهای رگرسیونی است. وجود وابستگی خطی نزدیک، توانایی برآورد ضرایب مدل رگرسیون را با مشکل مواجه می‌کند.

بکارگیری و تعبیر یک مدل رگرسیون چندگانه بطور ضمنی یا بطور صریح به برآوردهای ضرایب مدل رگرسیون بستگی دارد. بعضی از استنباط‌هایی که بطور معمول در مدل رگرسیونی صورت می‌پذیرد، شامل

موارد زیر است:

۱- مشخص کردن اثرات نسبی متغیرهای رگرسیونی

۲- پیش‌بینی و یا برآورد

۳- انتخاب دسته‌ای مناسب از متغیرها برای مدل

اگر هیچ رابطه خطی بین متغیرهای رگرسیونی نباشد، گفته می‌شود که متغیرهای رگرسیونی مستقلند. اگر متغیرهای رگرسیونی مستقل باشند، استنباط‌هایی مانند مواردی که در بالا ذکر شد، نسبتاً بسادگی صورت می‌پذیرد. متأسفانه در بیشتر کاربردهای رگرسیون، متغیرهای رگرسیونی مستقل نیستند. هنگامی که ارتباط خطی نزدیکی بین متغیرهای رگرسیونی وجود دارد، گفته می‌شود مسئله همخطی چندگانه وجود دارد.

چهار منبع اولیه همخطی وجود دارد:

۱- شیوه جمع‌آوری داده‌ها

۲- قیدهای روی مدل یا درون جامعه

۳- نوع مدل

۴- مدلی با متغیرهای رگرسیونی بیش از حد

درک اختلاف بین این منابع همخطی چندگانه اهمیت دارد، به طوری که توصیه برای تحلیل داده‌ها و تعبیر مدل نتیجه‌شده تا اندازه‌ای به علت مسئله همخطی بستگی دارد. (برای اطلاعات بیشتر در مورد همخطی

چندگانه رجوع کنید به میسون، گونست و وبستر^۱ (۱۹۷۵).

^۱Mason, Gunst and Webster