

الله



دانشکده علوم

پایان نامه کارشناسی ارشد در رشته آمار ریاضی

بررسی کارایی برآوردهای تقریباً نااریب در مدل
رگرسیون خطی، در حضور همخطی چندگانه

توسط:

مرضیه توانگر

استاد راهنما:

دکتر مینا توحیدی

دکتر عبدالرسول برهانی حقیقی

مهر ماه ۱۳۹۱

به نام خدا

اظهارنامه

اینجانب مرضیه توانگر (۸۹۰۸۹۲) دانشجوی کارشناسی ارشد رشته‌ی آمار ریاضی اظهار می‌کنم که این پایان نامه حاصل پژوهش خودم بوده و در جاهایی که از منابع دیگران استفاده کرده‌ام، نشانی دقیق و مشخصات کامل آن را نوشتهم. همچنین اظهار می‌کنم که تحقیق و موضوع پایان نامه‌ام تکراری نیست و تعهد می‌نمایم که بدون مجوز دانشگاه دستاوردهای آن را منتشر ننموده و یا در اختیار غیر قرار ندهم. کلیه حقوق این اثر مطابق با آئین نامه مالکیت فکری و معنوی متعلق به دانشگاه شیراز است.

نام و نام خانوادگی : مرضیه توانگر

تاریخ و امضا : ۱۳۹۱/۰۷/۱۲



به نام خدا

بررسی کارایی برآوردهای تقریباً نااریب در مدل رگرسیون خطی، در حضور همخطی چندگانه

به وسیله :

مراضیه توانگر

پایان نامه

ارائه شده به تحصیلات تکمیلی دانشگاه به عنوان بخشی از فعالیت‌های
لازم برای اخذ درجه کارشناسی ارشد

در رشته :

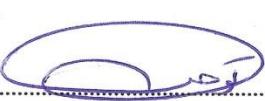
آمار ریاضی

از دانشگاه شیراز

شیراز

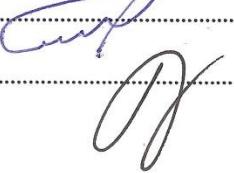
جمهوری اسلامی ایران

ارزیابی کمیته‌ی پایان نامه، با درجه‌ی : عالی

.....

دکتر مینا توحیدی، دانشیار بخش آمار (رئیس کمیته)

.....

دکتر عبدالرسول برهانی حقیقی، استادیار بخش آمار (رئیس کمیته)

.....

دکتر سلطان محمد صدوqi الوندی، استاد بخش آمار

.....

دکتر عبدالرضا بازرگان لاری، استادیار بخش آمار

مهر ماه ۱۳۹۱

تقدیم به

او که مرا آفرید

و آنکه بودن تا کشتم پدید

و ایشان که دادند از علم بر من نوید

و مایی که داریم عشق و امید

و شما اکنون

تقدیم به درود و مادر عزیزم که

وجودشان برایم به عشق

وجودم برایشان به رنج

تو انشان رفت تا به تو امایی رسم

و مویشان کرد سپیدی گرفت تارویم سپید ماند

سپاسگزاری

با سپاس بیکران به درگاه یگانه معبد هستی که این توفیق را نصیب من کرد تا در مرحله ای دیگر، گامی هرچند کوچک در تحقق اهداف خویش بردارم. اکنون که به یاری پروردگار، مراحل پژوهش و تدوین این پایان نامه به اتمام رسیده است، بر خود لازم می دانم که از تلاش و کمک تمام کسانی که مرا در این امر یاری نموده اند، صمیمانه سپاسگزاری نمایم. از مادر و پدر عزیزم به پاس محبت بی دریغشان که همواره مشوق من بوده اند و در سختی ها و دشواری های زندگی، صمیمانه و با خلوص قلب مرا یاری داده اند، خاضعانه تشکر می کنم .

مراتب سپاس و قدردانی بیکران خود را به محضر استاد بزرگوارم سرکار خانم دکتر مینا توحیدی که این پژوهش مرهون رهنماوهای ایشان می باشد و همواره با سعه‌ی صدر و ظرافت طبع در دوران تحصیل مرا یاری داده اند، تقدیم می دارم. همچنین از جناب آقای دکتر دکتر عبدالرسول برهانی حقیقی که در پایان نامه مرا یاری کردند نهایت سپاس را دارم. از اساتید محترم جناب آقای دکتر محمد صدوقی و جناب آقای دکتر عبدالرضا بازرگان لاری که زحمت مشاوره‌ی این تحقیق را تقبل نمودند و مرا مدیون محبت‌ها و راهنمایی‌های ارزنده‌ی خود کردن، سپاسگزارم. امید آن دارم که با یاری خداوند متعال در جهت پیشرفت ایران عزیز قدم بردارم.

چکیده

بررسی کارایی برآوردهای تقریبا ناریب در مدل رگرسیون

خطی، در حضور هم خطی چندگانه

به وسیله‌ی :

مرضیه توانگر

برآوردهای حداقل مربعات معمولی (OLS) $\hat{\beta} = (X'X)^{-1}X'Y$ برای برآورد ضرایب رگرسیونی در مدل رگرسیون خطی $Y = X\beta + \epsilon$ استفاده می‌شود. اما این برآوردهای به شدت به خصوصیات ماتریس $X'X$ بستگی دارد. هم خطی چندگانه بین متغیرهای توضیحی در مدل رگرسیون خطی، یک مسئله مهم در بکارگیری این مدل می‌باشد. در این حالت برآوردهای کمترین مربعات ($\hat{\beta}$) دارای واریانس بزرگی است. در این پایاننامه، ابتدا با الهام از کلاس برآوردهای ریج، $\hat{\beta}_R = (X'X + kI)^{-1}X'Y$ جدیدی از برآوردهای اریب، $\hat{\beta}_d = (X'X + dI)^{-1}(X'Y + d\hat{\beta})$ معرفی می‌شود. هر عضو این کلاس جدید یک برآوردهای لیو می‌باشد. برآوردهای بهینه، با استفاده از محک مینیمم توان دوم خطای MSE ، انتخاب می‌شود. با ادغام اعضای دو کلاس از برآوردهای مذکور، کلاسی دیگر از برآوردهای به نام کلاس برآوردهای دوپارامتری (TP) ساخته می‌شود. با استفاده از فرآیند جک-نایف، برآوردهای تقریبا ناریب لیو ($AULE$) و برآوردهای تقریبا ناریب دوپارامتری ($AUTP$) به دست آورده می‌شوند. برآوردهای ساخته شده در این پایاننامه، با برآوردهای LS با معیار MSE مقایسه می‌گردند و در آخر، یک مطالعه شبیه‌سازی، برای مشاهده کارایی این برآوردهای ساخته شده انجام گرفته است.

فهرست مطالب

عنوان	صفحه
فصل اول : مقدمه	
۱-۱- مقدمه	۲
۱-۲- برآوردهای انقباضی	۳
۱-۳- همخطی چندگانه	۴
۱-۳-۱- نتایج وجود همخطی چندگانه	۵
۱-۳-۲- علل وقوع همخطی چندگانه	۸
۱-۳-۳- تشخیص همخطی چندگانه	۱۱
۱-۴- روش‌هایی برای برخورد با همخطی چندگانه	۱۶
۱-۴-۱- آشنایی با روش جک نایف	۲۰
۱-۴-۲- روش جک نایف	۲۱
فصل دوم : کلاس جدیدی از برآوردهای اریب β در رگرسیون خطی	
۲-۱- مقدمه	۲۶
۲-۲- انتخاب بهینه پارامتر d در برآوردهای βd	۳۰
۲-۳- خواص آماری $\beta 1d$	۳۹
۲-۴- مثال عددی	۴۳
فصل سوم : برآوردهای تقریباً نااریب لیو	
۳-۱- مقدمه	۴۸

۲-۳- برآوردگر تعمیم یافته تقریبا ناریب لیو ۵۳	
۳-۳- برآورد ناریب برای اریبی برآوردگر تعمیم یافته لیو (<i>GLE</i>) ۵۶	
۴-۳- مقایسه بین <i>OLSE</i> , <i>GLE</i> و <i>AUGLE</i> ۵۹	
۳-۵- برآوردگر تقریبا ناریب لیو تحت مدل خطی تعمیم یافته ۶۴	
۳-۶- عملکرد برآوردگرهای تقریبا ناریب لیو تحت مدل خطی تعمیم یافته ۶۶	
۳-۷- مقایسه دو برآوردگر <i>AUL</i> و <i>GLS</i> ۶۶	
۳-۸- مقایسه دو برآوردگر <i>AUL</i> و <i>Liu</i> ۶۷	
۳-۹- مقایسه بین برآوردگرهای <i>AUL</i> و <i>ORR</i> ۶۹	
۳-۱۰- شبیه‌سازی مونت کارلو ۷۲	
۳-۱۱- نتایج پایانی ۷۵	
فصل چهارم : کارایی برآوردگر تقریبا ناریب دوپارامتری	
۴-۱- مقدمه ۷۷	
۴-۲- برآوردگر دوپارامتری (<i>TP</i>) ۷۹	
۴-۳- عملکرد برآوردگر <i>TP</i> بر حسب ملاک <i>MSE</i> ۸۱	
۴-۴- مقایسه بین برآوردگرهای <i>OLS</i> و <i>TP</i> بر حسب ملاک <i>MSE</i> ۸۲	
۴-۵- انتخاب پارامترهای اریب <i>k</i> و <i>d</i> در برآوردگر <i>TP</i> ۸۵	
۴-۶- برآوردگر تقریبا ناریب دوپارامتری (<i>AUTP</i>) ۸۹	
۴-۷- عملکرد برآوردگر <i>AUTP</i> بر حسب ملاک <i>MSE</i> ۹۱	
۴-۸- مقایسه برآوردگر <i>AUTP</i> و برآوردگر <i>LS</i> ۹۲	
۴-۹- مقایسه برآوردگر دوپارامتری و برآوردگر <i>AUTP</i> ۹۴	
۴-۱۰- انتخاب پارامترهای اریب <i>k</i> و <i>d</i> در برآوردگر <i>AUTP</i> ۹۶	

۱۰۰	-۴ مطالعه شبیه سازی.....
۱۰۶	فهرست منابع.....
۱۱۱	پیوست الف: (زبان برنامهنویسی)
۱۱۱	الف- برنامه R برای داده‌های فصل دوم
۱۱۳	ب- برنامه R برای داده‌های شبیه سازی شده
۱۱۹	واژه‌نامه(انگلیسی به فارسی)
۱۲۴	واژه‌نامه(فارسی به انگلیسی)

فهرست جداول

صفحه

عنوان و شماره

٤٥.....	جدول ١-٢
٤٥.....	جدول ٢-٢
٤٦.....	جدول ٣-٢
٤٦.....	جدول ٤-٢
٤٦.....	جدول ٥-٢
٧٥.....	جدول ١-٣
١٠١.....	جدول ١-٤
١٠١.....	جدول ٢-٤
١٠١.....	جدول ٣-٤
١٠٣.....	جدول ٤-٤
١٠٣.....	جدول ٥-٤
١٠٤.....	جدول ٦-٤

فهرست شکل‌ها

عنوان	صفحة
-------	------

شکل شماره ۱-۲	۳۲
---------------	----

فصل اول

مقدمه

۱-۱ مقدمه

مدل رگرسیون خطی زیر را در نظر بگیرید:

$$Y = X\beta + \varepsilon, \quad \varepsilon \sim (0, \sigma^2 I_n) \quad (1.1.1)$$

که Y بردار پاسخ $n \times 1$, X ماتریس $n \times p$ از متغیرهای توضیحی, β بردار $p \times 1$ از ضرایب رگرسیونی نامعلوم و ε بردار $n \times 1$ از خطاهای تصادفی مستقل و هم توزیع هستند که دارای میانگین صفر و واریانس $\sigma^2 I_n$ می‌باشند.

حال علاقمند هستیم تا پارامتر β را جهت استنباطهای پیش‌بینی کننده برآورد کنیم.

با توجه به روش حداقل مربعات بردار پارامتر β , به صورت زیر برآورد می‌شود:

$$S = (Y - X\beta)'(Y - X\beta) = \varepsilon' \varepsilon$$

$$S = Y'Y - Y'X\beta - \beta'X'Y + \beta'X'X\beta$$

$$\frac{\partial S}{\partial \beta} = -2X'Y + 2X'X\beta = 0$$

$$X'X\beta = X'Y$$

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

برآوردگر $\hat{\beta}$ را، برآوردگر حداقل مربعات معمولی (*OLS*)^۱ می‌نامند، که برآوردگر ناریب می‌باشد و

با توجه به قضیه گوس مارکوف^۲ در بین برآوردگرهای خطی ناریب، دارای کوچکترین واریانس می‌باشد و به این ویژگی BLUE ^۳, به معنای بهترین برآوردگر ناریب خطی، گویند.

¹ Ordinary least square

² Gauss Markov Theorem

³ Best Linear Unbiased Estimator

۱-۲- برآوردگر انقباضی^۱

حال می‌خواهیم برآوردگرهای انقباضی، که برآوردگرهای اریبی می‌باشند را معرفی کنیم. با توجه به اینکه برآوردگر (OLS) در بین برآوردگرهای خطی ناریب کوچکترین واریانس را دارد، اما نمی‌توان گفت که در بین تمام برآوردگرهای دارای کوچکترین MSE می‌باشد. بنابراین اگر فرض ناریبی را نادیده بگیریم، می‌توان برآوردگرهای اریبی را یافت که در بعضی موقعیت‌ها به دلیل MSE کوچکتر عملکرد بهتری داشته باشند.

در آمار برآوردگر انقباضی، برآوردگری است که با بهبود بخشیدن برآوردگر اولیه با استفاده از سایر اطلاعات درباره پارامتر، به دست می‌آید. در واقع برآوردگر بهبود یافته نسبت به برآوردگر اولیه به مقدار واقعی نزدیک‌تر می‌باشد و دارای ریسک کمتری می‌باشد.

فرض کنید برآوردگر اولیه صفر نباشد. با ضرب برآوردگر اولیه در یک پارامتر مشخص می‌توان برآوردگر دیگری به دست آورد و سپس مقدار پارامتر مشخص را طوری تعیین کرد که MSE برآوردگر جدید مینیمم شود. به ازای مقدار به دست آمده برای پارامتر، برآوردگر جدید MSE کوچکتری نسبت به برآوردگر اولیه خواهد داشت و در نتیجه برآوردگر جدید بهبود یافته است. یک مثال قابل ذکر در این رابطه، برآورد واریانس جامعه بر اساس یک نمونه ساده می‌باشد. برای نمونه‌ای به حجم n ، واریانس نمونه به صورت زیر می‌باشد:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

و می‌دانیم که s^2 ، برآوردگری ناریب برای واریانس جامعه می‌باشد اما اگر به جای کسر $\frac{1}{n-1}$ در s^2 ، کسر $\frac{1}{n+1}$ را جایگزین کنیم، برآورد اریبی با MSE مینیمم حاصل می‌شود. مبحث انقباض در استنباط بیزی ضمنی و استنباط درستنمایی توانیده^۲ و استنباط از نوع $James - Stein$ مطرح می‌شود. در مقابل فرآیند برآوردگر ماکسیمم درستنمایی و حداقل مربعات شامل تاثیرات انقباضی نمی‌باشند اگرچه در طرح‌های برآورد انقباضی مورد استفاده قرار می‌گیرند.

¹ Shrinkage Estimator

² Penalized

در استفاده از برآوردهای انقباضی در زمینه تحلیل رگرسیونی که ممکن است تعداد زیادی متغیرهای توضیحی وجود داشته باشند، توسط copas (۱۹۸۳) توصیف شده است. در این حالت مقدار ضرایب رگرسیونی برآورد شده به سمت صفر منقبض می‌شوند و تاثیرات آن در کاهش MSE مقادیر پیش‌بینی از مدل مشخص شده است.

یک نوع دیگر از برآوردهای اریب که برآوردهای انقباضی می‌باشند و مورد توجه بسیاری از پژوهشگران قرار گرفته است، برآوردهای ridge^۱ و برآوردهای Stein^۲ می‌باشند که به ترتیب به صورت زیر می‌باشند:

$$\hat{\beta}_R = (X'X + kI)^{-1}X'Y$$

$$\hat{\beta}_S = c\hat{\beta}$$

که $1 < c < 0$ و k پارامتر می‌باشد و در بخش‌های بعدی به طور کامل معرفی می‌شوند.

۱-۳-۱- هم خطی چندگانه^۳

بکارگیری و تعبیر یک مدل رگرسیون خطی چندگانه به طور ضمنی یا به طور صریح به برآوردهای ضرایب منفرد رگرسیون بستگی دارد. استنباطهایی که به طور معمول انجام می‌شود، شامل موارد زیر است:

۱. مشخص کردن اثرات نسبی متغیرهای رگرسیونی.

۲. پیش‌بینی و یا برآورد

۳. انتخاب یک مجموعه مناسب از متغیرها برای مدل

اگر هیچ رابطه خطی بین متغیرهای رگرسیونی نباشد، گفته می‌شود که متعامدند، اگر متغیرهای رگرسیونی بر هم عمود باشند، استنباطهایی که در بالا ذکر شد نسبتاً ساده صورت می‌پذیرد. متاسفانه در بیشتر کاربردهای رگرسیون متغیرهای رگرسیونی متعامد نیستند. گاهی اوقات متعامد نبودن چیز جدی نیست. در حالی که در بعضی وضعیت‌ها متغیرهای رگرسیونی تقریباً به طور کامل ارتباط خطی دارند و در چنین مواردی استنباطهای بر اساس این

¹ Ridge

² Stein

³ multicollinearity

رگرسیون می‌تواند گمراه کننده یا غلط باشد. هنگامی که ارتباط خطی نزدیکی بین متغیرهای رگرسیونی وجود دارد، گفته می‌شود مسئله هم خطی چندگانه وجود دارد (به کتاب مقدمه‌ای بر تحلیل رگرسیون خطی از C Montgomery, Douglas ترجمه سید ابراهیم رضوی پاریزی مراجعه شود).

در این بخش علل وقوع هم خطی چندگانه و بعضی از اثرات آن روی استنباط، بیان می‌شود و روش‌هایی برای تشخیص وجود هم خطی و بعضی فنون برخورد با این مسئله، مورد بحث و بررسی قرار خواهد گرفت.

۱-۳-۱- نتایج وجود هم خطی چندگانه

وجود هم خطی چندگانه بالقوه آثار جدی متعددی بر برآوردهای حداقل مربعات ضرایب رگرسیون دارد. بعضی از این اثرها را به سادگی می‌توان نمایش داد. فرض کنیم فقط دو متغیر رگرسیونی X_1 و X_2 وجود داشته باشد و با این فرض که X_1 و X_2 و Y به طول واحد مقیاس‌سازی شده‌اند. مدل را به صورت زیر در نظر می‌گیریم :

$$Y = \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

و معادلات حداقل مربعات نرمال عبارتند از :

$$(X'X)\hat{\beta} = X'Y$$

$$\begin{pmatrix} 1 & r_{12} \\ r_{12} & 1 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} r_{1y} \\ r_{2y} \end{pmatrix}$$

که در آن r_{12} همبستگی ساده بین X_1 و X_2 و r_{jy} همبستگی ساده بین X_j و Y می‌باشد و $j = 1, 2$. معکوس $(X'X)$ عبارتست از :

$$C = (X'X)^{-1} = \begin{pmatrix} \frac{1}{(1 - r_{12}^2)} & \frac{-r_{12}}{(1 - r_{12}^2)} \\ \frac{-r_{12}}{(1 - r_{12}^2)} & \frac{1}{(1 - r_{12}^2)} \end{pmatrix} \quad (1.3.1)$$

و برآوردهای ضرایب رگرسیون عبارتست از :

$$\hat{\beta}_1 = \frac{r_{1y} - r_{12}r_{2y}}{(1 - r_{12}^2)} \quad , \quad \hat{\beta}_2 = \frac{r_{2y} - r_{12}r_{1y}}{(1 - r_{12}^2)} \quad (1.3.2)$$

اگر بین X_1 و X_2 همخطی چندگانه شدید موجود باشد در این صورت ضریب همبستگی $V(\hat{\beta}_j) = C_{jj}\sigma^2 \rightarrow \infty$ و r_{12} بزرگ خواهد بود. با توجه به اینکه اگر $|r_{12}| > 1$ آنگاه $Cov(\hat{\beta}_1, \hat{\beta}_2) = C_{12}\sigma^2 \rightarrow \pm\infty$ بسته به اینکه $r_{12} \rightarrow +1$ یا $r_{12} \rightarrow -1$ یا r_{12} بنابراین همخطی زیاد بین X_1 و X_2 منجر به واریانس ها و کواریانس های بزرگ برای برآوردهای حداقل مربعات ضرایب رگرسیون خواهد شد.

در نتیجه اگر از سطوح یکسان X نمونه های مختلف گرفته شود، به برآوردهای با تفاوت زیاد از هر پارامتر منجر خواهد شد.

تعريف ضریب تعیین¹ چندگانه:

ضریب تعیین چندگانه یک محک مناسب برای ارزیابی مناسبت مدل می باشد و با R_p^2 نشان داده می شود. ضریب تعیین چندگانه برای یک مدل زیر مجموعه با p جمله می باشد و به صورت زیر تعریف می شود:

$$R_p^2 = \frac{SS_R(p)}{SS_{yy}} = 1 - \frac{SS_E(p)}{SS_{yy}}$$

که $(SS_R(p)$ و $SS_E(p)$ ، به ترتیب نشان دهنده مجموع مربعات رگرسیون و مجموع مربعات باقیمانده برای یک مدل زیر مجموعه با p جمله می باشد.

وقتی که بیش از دو متغیر رگرسیونی وجود دارد، همخطی اثرات مشابهی ایجاد می کند، می توان نشان داد که اعضای قطر ماتریس $(X'X)^{-1}$ عبارتند از

$$C_{jj} = \frac{1}{(1-R_j^2)} \quad j = 1, \dots, p \quad (1.3.3)$$

که در آن R_j^2 ضریب تعیین چندگانه از رگرسیون X_j نسبت به $p-1$ متغیر رگرسیونی باقیمانده است. اگر همخطی شدید بین X_j و هر زیرمجموعه ای از $p-1$ متغیر رگرسیونی دیگر وجود داشته باشد، در این صورت مقدار R_j^2 نزدیک به واحد خواهد بود. چون واریانس $j\hat{\beta}$ برابر است با $V(\hat{\beta}_j) = C_{jj}\sigma^2 = (1 - R_j^2)^{-1}\sigma^2$ برآورد حداقل مربعات ضریب رگرسیونی $j\hat{\beta}$ بسیار زیاد شود. بطور کلی اگر X_i و X_j در یک رابطه همخطی چندگانه درگیر باشند کواریانس $i\hat{\beta}_i$ و $j\hat{\beta}_j$ نیز بزرگ خواهد بود.

¹ Coefficient of Determination

هم خطی چندگانه همچنین به ایجاد برآوردهای حداقل مربعات $\hat{\beta}$ که از نظر قدرمطلق خیلی بزرگ می‌باشند منجر خواهد شد. برای ملاحظه این مطلب، مربع فاصله $\hat{\beta}$ تا بردار واقعی β را مورد توجه قرار می‌دهیم، یعنی فاصله β از $\hat{\beta}$ را به صورت زیر در نظر می‌گیریم:

$$L_1^2 = (\hat{\beta} - \beta)'(\hat{\beta} - \beta) \quad (1.3.4)$$

امید مربع فاصله، $E(L_1^2)$ ، برابر است با :

$$\begin{aligned} E(L_1^2) &= E(\hat{\beta} - \beta)'(\hat{\beta} - \beta) \\ &= \sum_{j=1}^p E(\hat{\beta}_j - \beta_j)^2 \\ &= \sum_{j=1}^p V(\hat{\beta}_j) \\ &= \sigma^2 \text{Tr}(X'X)^{-1} \end{aligned} \quad (1.3.5)$$

که اثر^۱ یک ماتریس (با Tr نمادگذاری می‌شود) با مجموع اعضای قطر اصلی برابر است. وقتی که هم خطی چندگانه وجود داشته باشد، بعضی از مقادیر ویژه $X'X$ کوچک خواهند بود. چون اثر یک ماتریس نیز با مجموع مقادیر ویژه آن برابر است (1.3.5) به صورت زیر در می‌آید:

$$E(L_1^2) = \sigma^2 \sum_{j=1}^p \frac{1}{\lambda_j} \quad (1.3.6)$$

که در آن $\lambda_j > 0, j = 1, 2, \dots, p$ مقادیر ویژه $X'X$ می‌باشند. بنابراین اگر ماتریس $X'X$ به لحاظ هم خطی شرایط بیمارگونه داشته باشد حداقل یکی از روابطها کوچک خواهد بود و (1.3.6) نتیجه می‌دهد که فاصله برآورد حداقل مربعات $\hat{\beta}$ تا مقدار واقعی پارامتر β می‌تواند بزرگ باشد. به طور معادل می‌توان نشان داد که

$$\begin{aligned} E(L_1^2) &= E(\hat{\beta} - \beta)'(\hat{\beta} - \beta) \\ &= E(\hat{\beta}'\hat{\beta} - 2\hat{\beta}'\beta + \beta'\beta) \end{aligned} \quad (1.3.7)$$

یا

$$E(\hat{\beta}'\hat{\beta}) = \beta'\beta + \sigma^2 \text{Tr}(X'X)^{-1} \quad (1.3.8)$$

یعنی به طور کلی طول بردار $\hat{\beta}$ بزرگتر از بردار β می‌باشد. این بدین معنی است که روش حداقل مربعات، ضرایب رگرسیونی را ایجاد می‌کند که از نظر قدرمطلق خیلی بزرگ می‌باشند.

^۱ Trace

به طور کلی هرگاه هم خطی چندگانه شدید وجود داشته باشد، روش حداقل مربعات برآوردهای ضعیفی از هر یک از پارامترهای منفرد مدل بدست می‌دهد. اما این لزوما نتیجه نمی‌دهد که مدل برازش شده یک پیش‌بینی ضعیفی می‌باشد. اگر پیش‌بینی‌ها محدود به ناحیه‌هایی از فضای X باشد که هم خطی چندگانه تقریباً برقرار است، مدل برازش شده اغلب پیش‌بینی قابل قبولی به دست می‌دهد. این بدین علت می‌تواند اتفاق بیفتد که ترکیب خطی $\sum_{j=1}^p \beta_j x_{ij}$ کاملاً خوب برآورد می‌شود، هرچند پارامترهای منفرد β_j به طور ضعیفی برآورد می‌شوند. بدین معنی که اگر داده‌های اصلی در امتداد ابر صفحه تعریف شده در رگرسیون قرار گیرند در این صورت مشاهدات آینده که نزدیک این ابر صفحه قرار می‌گیرند علیرغم برآوردهای نامناسب تک پارامترها می‌توانند بصورت دقیق پیش‌بینی شوند.

به طور خلاصه می‌توان پیامدهای وجود هم خطی را به صورت زیر فهرست کرد.

۱- چون در حالت هم خطی اطلاعات مستقل در مورد هر یک از متغیرهای مستقل وجود ندارد؛ لذا نمی‌توان اثرات جزئی متغیرهای مذکور روی متغیر وابسته را برآورد کرد.

۲- هنگامی که همبستگی شدید بین متغیرهای مستقل وجود داشته باشد کواریانس و واریانس ضرایب بزرگ‌تر برآورد خواهند شد.

۳- برآورد ضرایب رگرسیونی نالریب همچنان خاصیت *Blue* را دارند. گرچه کمیت به دست آمده برای آنها غیر قابل اعتماد است.

۴- در حالی که با هم خطی شدید در مدل مواجه هستیم، پیش‌بینی‌های صورت گرفته در آن غیر قابل اعتماد خواهد بود. در این حالت پیش‌بینی‌ها بر اساس مدلی که دارای زیر مجموعه‌ای از متغیرهای مستقل مدل اصلی است، بهتر صورت می‌گیرد.

۱-۳-۲- علل وقوع هم خطی چندگانه

مدل رگرسیون چندمتغیره (1.1.1) را در نظر بگیرید. فرض کنید متغیرهای رگرسیونی و پاسخ، مرکزی شده و به طول واحد مقیاس‌سازی شده‌اند. در نتیجه $X'X$ یک ماتریس $p \times p$ از همبستگی‌های بین متغیرهای رگرسیونی و y' یک بردار $1 \times p$ از همبستگی‌های بین متغیرهای رگرسیونی و پاسخ است. اگر z امین ستون ماتریس X با X_z نمایش داده شود،