

## بررسی کاربرد رگرسیون ریج در علوم پزشکی

نازنین فکری<sup>۱</sup> - حبیب‌الله اسماعیلی<sup>۲</sup> - حسن دوستی<sup>۳</sup> - رقیه پارسایی<sup>۴</sup>

<sup>۱</sup> گروه بهداشت عمومی، دانشکده علوم پزشکی نیشابور

<sup>۲</sup> گروه آمار زیستی و عضو مرکز تحقیقات علوم بهداشتی دانشکده بهداشت، دانشگاه علوم

پزشکی مشهد

<sup>۳</sup> گروه آمار زیستی، دانشکده بهداشت، دانشگاه علوم پزشکی مشهد

<sup>۴</sup> گروه آمار زیستی و اپیدمیولوژی دانشکده بهداشت، دانشگاه علوم پزشکی اصفهان

### چکیده

روش‌های رگرسیونی یکی از پرکاربردترین روش‌های آنالیز آماری در علوم پزشکی است. وجود همبستگی بین متغیرهای پیشگو باعث بوجود آمدن همخطی می‌شود. بنابراین برآورد پارامترها بر مبنای کمترین مربعات باقیمانده نادرست است. در این مقاله سعی بر این است تا رگرسیون ریج، به عنوان یکی از کاراترین و موثرترین روش‌های مقابله با مشکل هم خطی معرفی گردد. در انتها روش فوق را با یک مثال کاربردی که روی داده‌های مربوط به عوامل خطر بیماری‌های قلبی انجام شده است، توضیح می‌دهیم.

واژه‌های کلیدی: هم خطی، رگرسیون ریج، کمترین مربعات باقیمانده.

مدل‌های رگرسیون سعی بر این است که از طریق یک یا

### ۱ مقدمه

چند متغیر پیشگو، متغیر پاسخ پیش‌بینی شود [۱]. مدل

استاندارد برای رگرسیون خطی چندگانه به شکل

$$\underline{Y} = X\underline{\beta} + e \quad (1)$$

است که در آن بردار  $n \times 1$  از متغیر پاسخ،  $X$  یک

ماتریس  $n \times p$  از متغیرهای پیشگو با رتبه  $p$ ،  $\underline{\beta}$  بردار

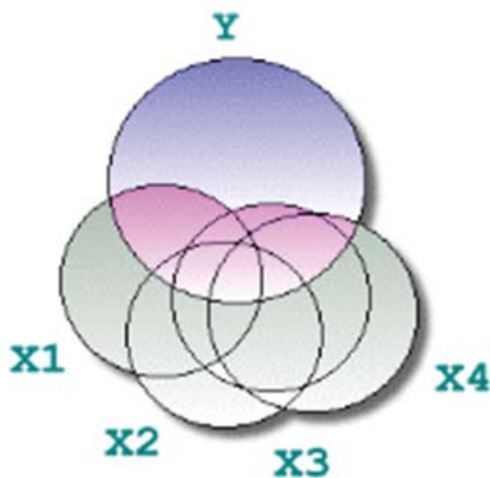
رگرسیون یکی از روش‌های متداول و قدیمی برای بررسی

رابطه‌ی بین متغیرهای قابل اندازه‌گیری است. به ویژه

رگرسیون خطی که تقریباً در همه‌ی شاخه‌های علوم مانند:

علوم اجتماعی اقتصادی، علوم فیزیکی و زیست‌شناسی،

پزشکی، تجارت، تکنولوژی و غیره کاربرد فراوان دارد. در



شکل ۱: وجود هم خطی بین متغیرهای توضیحی  $X_1$  تا  $X_4$

این بردارها مقادیر مشاهده شده برای هر کدام از متغیرهای پیشگو هستند. اگر به ازای ثابت‌های  $a_1, \dots, a_p$  یک رابطه‌ی خطی کامل، یعنی  $\sum_{i=1}^p a_i x_i = 0$  وجود داشته باشد، آن‌گاه می‌گوییم «هم خطی کامل» وجود دارد و اگر این رابطه‌ی خطی تقریبی باشد، یعنی  $\sum_{i=1}^p a_i x_i \approx 0$  آن‌گاه می‌گوییم «هم خطی ناقص» وجود دارد. در صورت وجود هم خطی کامل، معکوس ماتریس  $X^T X$  وجود نخواهد داشت و نیز کوچک‌ترین مقدار ویژه‌ی ماتریس  $X^T X$  برابر صفر است. هم خطی ناقص نیز بدین معنی است که  $X^T X$  تقریباً منفرد است. و کوچک‌ترین مقدار ویژه‌ی ماتریس  $X^T X$  نزدیک صفر می‌باشد.

### ۳ آثار هم خطی

برخی از مهم‌ترین نشانه‌های حضور هم خطی بین متغیرهای پیشگو به صورت زیر است:

- با حذف یا اضافه کردن یک متغیر پیشگو، یا با

$p \times 1$  ضرایب و مجهول هستند.  $e$  نیز بردار  $n \times 1$  باقی‌مانده‌های مدل با میانگین صفر و واریانس ثابت  $\sigma^2$  است [۳].

با کمک داده‌های جمع آوری شده و به کارگیری روش کمترین مربعات باقی‌مانده<sup>۱</sup> (روش گاوس مارکوف)، ضرایب مجهول  $\beta_1, \dots, \beta_p$  را می‌توان برآورد کرد و بدین ترتیب رابطه‌ی (۱) مشخص می‌شود. به عبارتی اگر تعداد مشاهدات  $n$  باشد آنگاه داریم [۱]:

$$y_i = b_1 x_{i1} + b_2 x_{i2} + \dots + b_p x_{ip} + e_i,$$

که  $i = 1, 2, \dots, n$  و یا به صورت ماتریس:

$$\underline{Y} = X\underline{b} + \underline{e}.$$

ضرایب  $b$  در این حالت برآورد ضرایب  $\beta$  و نارایب با کمترین واریانس هستند. در نظر نگرفتن پیش فرض‌ها در رگرسیون خطی می‌تواند روی نتیجه‌گیری‌ها اثر نامطلوب گذاشته و نتایجی نادرست بدست آید. یکی از این پذیره‌ها وجود هم خطی در ماتریس داده‌ها است.

### ۲ هم خطی چیست؟

$p$  متغیر را زمانی هم خط گویند که حداقل یکی از بردارهای ماتریس  $X$  یک ترکیب خطی از بردارهای دیگر باشد [۴]. به عبارتی دیگر، هرگاه دو متغیر با یکدیگر هم خطی داشته باشند، یعنی یکی از آنها، در معادله‌ی رگرسیون، حرف جدیدی برای گفتن ندارد و حضورش در معادله توجیه نمی‌شود. شکل ۱، نمونه‌ای از وجود هم خطی بین متغیرهای توضیحی را نشان می‌دهد. ستون‌های ماتریس  $X$  به صورت  $x_1, \dots, x_p$  هستند به طوری که هر یک از

<sup>۱</sup>Least square residuals

که در آن  $R_k^2$  مقدار ضریب تعیین چندگانه بین متغیرها است، زمانی که  $X_k$  را روی سایر متغیرها برگشت دهیم [۲]. با توجه به رابطه‌ی بالا چنانچه  $R_k^2 = 0$  باشد، آنگاه مقدار  $VIF_k = 1$  است. بنابراین اگر  $X_k$  رابطه‌ی خطی با متغیرهای پیشگوی دیگر نداشته باشد، مقدار  $VIF_k = 1$  خواهد بود.

## ۵ روش‌های مقابله با هم خطی

برای حذف یا کم کردن هم خطی بین متغیرهای پیشگو روش‌های مختلفی ارائه شده است که به طور کلی این روش‌ها شامل حذف متغیر پیشگویی که باعث ایجاد هم خطی شده، جمع آوری بیشتر و مجدد داده و یا استفاده از روش‌هایی غیر از کمترین مربعات باقیمانده است. استفاده از روش‌های غیر از کمترین مربعات باقیمانده از نظر کاربردی و دقت، دارای کارایی بیشتری هستند و از کاهش توان آزمون جلوگیری می‌کنند. یکی از این روش‌ها برای مقابله با اثرات نامطلوب هم خطی استفاده از برآوردهای اریب است که رگرسیون ریج، یکی از مهم‌ترین و کاراترین این روش‌هاست.

## ۶ رگرسیون ریج

هارل و کنارد [۳] اولین بار در مقاله‌ای در مجله‌ی مهندسی شیمی، بیان کردند زمانی که در مدل هم خطی وجود داشته باشد، برای برآورد  $\beta$ ، می‌توان به جای برآوردهای کمترین مربعات باقی‌مانده از برآوردهای زیر استفاده کرد:

$$\widehat{\beta}^* = [X'X + KI]^{-1} X'Y,$$

ورود و خروج یک مشاهده به مدل، تغییرات زیادی در ضرایب برآورد شده‌ی رگرسیون بوجود می‌آید.

- نتایج غیر معمول در آزمون معنی داری ضرایب رگرسیون، دیده می‌شود.
- ضرایب رگرسیون برآورد شده، علامتی مخالف آنچه مورد انتظار است دارند.

- فواصل اطمینان پرعرض برای ضرایب رگرسیونی مربوط به متغیرهای توضیحی مهم موجود در مدل، برآورد می‌شود.

- انحراف استاندارد برآورد شده‌ی ضرایب رگرسیون، زیاد می‌شود.

موارد ذکر شده، حضور هم خطی را به صورت کمی اندازه‌گیری نمی‌کنند [۴].

## ۴ روش‌های تشخیص هم خطی

از جمله روش‌های تشخیص هم خطی می‌توان به موارد زیر اشاره کرد:

۱. رسم نمودار leverage یا نمودار پراکنندگی

۲. استفاده از آزمون‌های نسبت  $F^2$  بالا و یا نسبت  $T^3$  کوچک

۳. استفاده از شاخص عامل تورم واریانس  $VIF^4$ . این شاخص برای  $k = 1, \dots, p$  برابر است با:

$$VIF_k = (1 - R_k^2)^{-1},$$

<sup>۲</sup>F-ratio

<sup>۳</sup>T-ratio

<sup>۴</sup>Variance Inflation Factor

بنابراین برای انتخاب  $k$  مقداری را برمی‌گزینیم که کاهش در واریانس برآوردگر اریب، بیش از افزایش مربع اریبی باشد. در نتیجه MSE آن کمتر از واریانس برآوردگر نااریب خواهد بود [۱].

## ۷ کاربرد

یک مطالعه مقطعی، برای بررسی ارتباط بین شاخص‌های چاقی و تیترا آنتی بادی شوک حرارتی (HSP-27) در دانشگاه علوم پزشکی مشهد روی سه گروه متشکل از ۵۰ نفر با وزن طبیعی، ۵۰ نفر دارای اضافه وزن و ۱۰۰ نفر که خیلی چاق هستند، انجام شده است، متغیرهای اندازه‌گیری شده در این مطالعه فشار خون، سن، جنس، چربی خون، قند خون و تعدادی از عوامل خطر بیماری‌های قلبی-عروقی بوده‌اند [۴]. در مطالعه فوق وجود همخطی بین متغیرها در نظر گرفته نشده است. بین متغیرها در گروه اول (افراد با وزن طبیعی)، VIF بررسی و همخطی شدید مشاهده شد. با استاندارد کردن متغیرهای پیشگویی سعی در از بین بردن هم خطی (در همان گروه اول) نمودیم که نتایج در جدول ۱ آمده است. با توجه به جدول ۱ و انجام رگرسیون روی داده‌های استاندارد شده مشخص شد که هنوز همخطی وجود داشته و در رگرسیون معمولی تنها متغیر فشار خون بر HSP-27، معنی‌دار شد.

بنابراین جهت رفع هم خطی، از رگرسیون ریب استفاده نمودیم. با توجه به مطالبی که در قسمت رگرسیون ریب آمده، ضریب  $K$  محاسبه و با استفاده از فرمول  $\tilde{b}(k) = (x'x + kI)^{-1}x'y$  پارامترهای مدل مجدداً برآورد شد.

جدول ۲ به‌ازای مقادیر مختلف  $k$  برآورد ضرایب

برآوردهایی که با  $K \geq 0$  بدست می‌آیند، شباهت‌های ریاضی زیادی با تعریف توابع پاسخ درجه دوم دارند. به این دلیل، برآوردها و تحلیل‌هایی که با  $\hat{\beta}^*$  انجام می‌شود ” رگرسیون ریب ” نامیده می‌شوند [۳]. رابطه‌ی برآورد ریب با برآورد کمترین مربعات باقی‌مانده، به شکل زیر است:

$$\tilde{b} = \hat{\beta}^* = Z\hat{\beta}.$$

به عبارتی برآوردگر ریب یک ترکیب خطی از برآوردگر کمترین مربعات باقی‌مانده است، به طوری که

$$Z = [X'X + KI]^{-1}X'X.$$

در روش کمترین مربعات باقی‌مانده، فرض کردیم  $\hat{b}$  برآوردی نااریب برای  $\beta$  است و خاصیت گاوس مارکوف، مینیمم بودن واریانس این برآوردگر در کلاس برآوردگرهای خطی نااریب را بیان می‌کند. در محاسبه‌ی برآوردگر ریب از فرض نااریب بودن چشم‌پوشی شده و برآوردگر اریبی برای  $\beta$  بدست می‌آید که دارای واریانس کمتری نسبت به برآوردگر کمترین مربعات باقی‌مانده است. بنابراین مقدار برآورد در این حالت پایدارتر خواهد بود [۴]، یعنی با حذف یا افزودن یک متغیر پیشگویی جدید به مدل و یا... ضرایب برآوردشده‌ی مدل تغییر چندانی نخواهند کرد.

در محاسبه‌ی برآوردگر ریب ( $\tilde{b}$ )، با فرض این که ستون‌های ماتریس  $X$  استاندارد شده باشند، از حل معادلات نرمال داریم:

$$\tilde{b}(k) = (x'x + kI)^{-1}x'y,$$

که در آن  $k$  مقداری مثبت است که بایستی توسط تحلیل‌گر با توجه به یک سری معیارها تعیین شود.

با توجه به این نکته که  $MSE(\tilde{b}) = \text{variance} + (\text{bias})^2$ .

رگرسیون و مقدار VIF را نشان می‌دهد. همانطور که ملاحظه می‌شود، مقدار VIF برای  $k = 0/16$  به میزان یک و یا کمتر رسیده که بیانگر عدم وجود همخطی است. از طرفی، ضرایب برآورد شده، در حالت  $k = 0/16$  به مقدار پایدارتری رسیده و تغییرات آن‌ها کمتر می‌باشد.

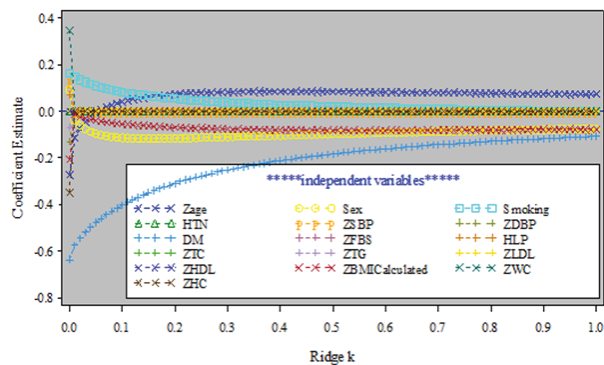
جدول ۱: برآورد ضرایب مدل رگرسیون معمولی در تاثیر متغیرهای استاندارد شده بر HSP-27

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t	Variance Inflation
intercept	۱	۰/۰۹۸۳۴	۱/۲۸۱۵۱	۰/۰۸	۰۹۳۹۳	۰
Zage	۱	-۰/۲۷۲۱۳	۰/۱۸۹۶۴	-۱/۴۳	۰/۱۶۰۷	۳/۷۶۵۲۱
Sex	۱	۰/۰۹۲۲۳	۰/۲۷۲۴۵	۰/۳۴	۰/۷۳۷۰	۱/۷۷۲۴۷
Smoking	۱	۰/۱۶۰۸۳	۰/۳۴۵۴۶	۰/۴۷	۰/۶۴۴۶	۱/۵۳۴۴۵
HTN	۱	۰/۰۰۲۶۸	۰/۰۰۱۲۲	۲/۲۰	۰/۰۳۴۷	۱۶/۹۵۹۵۶
ZSBP	۱	۰/۱۲۸۲۳	۰/۱۷۹۷۸	۰/۷۱	۰/۴۸۰۷	۴۵۵۳۶۰
ZDBP	۱	-۰/۱۳۰۱۹	۰/۱۷۹۹۸	-۰/۷۲	۰/۴۷۴۵	۴۵۶۶۰۴
DM	۱	-۰/۶۳۶۴۸	۰/۴۶۵۴۷	-۱/۳۷	۰/۱۸۰۸	۱/۵۲۵۵۷
ZFBS	۱	۰/۰۰۱۰۸	۰/۰۰۰۸۰۶۱۲	۱/۳۴	۰/۱۸۸۲	۶/۵۵۱۷۲
HLP	۱	۰/۰۶۹۳۹	۰/۱۲۶۶۹	۰/۵۵	۰/۵۸۷۶	۱۳۷۴۷۵
ZTC	۱	-۰/۰۰۰۴۵۴۷۷	۰/۰۰۰۳۰۴۲۸	-۱/۴۹	۰/۱۴۴۵	۲/۰۸۱۰۵
ZTG	۱	-۰/۰۷۰۵۴	۰/۱۲۶۲۸	-۰/۵۶	۰/۵۸۰۲	۱۳۷۰۹۰
ZLDL	۱	۰/۰۰۰۳۹۴۶۶	۰/۰۰۰۴۰۳۰۹	۰/۹۸	۰/۳۳۴۷	۲/۴۷۹۶۷
ZHDL	۱	-۰/۰۰۰۷۶۶۷۳	۰/۰۰۰۷۰۰۱۶	-۱/۰۹	۰/۲۸۱۶	۵/۶۳۰۶۲
Calculated ZBMI	۱	-۰/۲۰۴۵۷	۰/۳۷۰۹۹	-۰/۵۵	۰/۵۸۵۱	۱/۸۶۷۹۵
ZWC	۱	۰/۳۴۶۵۳	۰/۲۴۲۴۹	۱/۴۳	۰/۱۶۲۴	۷۵۵۹۳۳
ZHC	۱	-۰/۳۴۷۰۱	۰/۲۴۲۴۹	-۱/۴۳	۰/۱۶۱۸	۷۵۵۹۹۸

جدول ۲: برآورد ضرایب رگرسیونی به ازای مقادیر مختلف k.

obs	type	k	b <sub>o</sub>	Zage	sex	smoking	HTN	ZSBP	ZDBP	DM	ZFBS	HLP	ZTC	ZTG
۲۲	VIF	۰/۱۰		۱/۰۰۲۴۸	۱/۰۰۵۰۵	۰/۹۷۷۵۹	۱/۵۰۷۷	۰/۴۲	۰/۴۱	۰/۹۹۸۴۲	۵۰۶۲۸/۱	۰/۶۹	۱/۱۷۳۶۱	۰/۶۹
۲۳	b(k)	۰/۱۰	۰/۳۱۴۲۳	۰/۰۴۱۶۱	-۰/۱۱۳۶۹	۰/۰۸۲۲۳	۰/۰۰۰۹	-۰/۰۰	-۰/۰۰	-۰/۴۰۲۴۷	۰/۰۰۰۶۷	-۰/۰۰	-۰/۰۰۰۲۱	-۰/۰۰
۲۴	VIF	۰/۱۱		۰۹۵۸۰۰	۰/۹۷۶۳۹	۰/۹۵۰۰۲	۱/۳۳۵۸	۰/۳۸	۰/۳۸	۰/۹۶۹۲۷	۱/۶۸۳۰۳	۰/۶۳	۱/۱۳۳۳۸	۰/۶۳
۲۵	b(k)	۰/۱۱	۰/۲۹۴۹۷	۰/۰۴۷۱۰	-۰/۰۱۱۴۹۸	۰/۰۷۸۳۷	۰/۰۰۰۸	-۰/۰۰	-۰/۰۰	-۰/۳۹۰۵۶	۰/۰۰۰۶۴	-۰/۰۰	-۰/۰۰۰۲۰	-۰/۰۰
۲۶	VIF	۰/۱۲		۰۹۱۸۸۱	۰/۹۴۹۳۸	۰/۹۲۴۰۳	۱/۱۹۴۲	۰/۳۶	۰/۳۵	۰/۹۴۱۷۸	۱/۵۵۷۶۶	۰/۵۸	۱/۰۹۵۶۷	۰/۵۸
۲۷	b(k)	۰/۱۲	۰/۲۷۶۵۹	۰/۰۵۱۸۲	-۰/۱۱۵۸۸	۰/۰۷۴۸۳	۰/۰۰۰۸	-۰/۰۰	-۰/۰۰	-۰/۳۷۹۴۰	۰/۰۰۰۶۱	-۰/۰۰	-۰/۰۰۰۲۰	-۰/۰۰
۲۸	VIF	۰/۱۳		۰/۸۸۳۸۷	۰/۹۲۳۸۳	۰/۸۹۹۴۶	۱/۰۷۶۰	۰/۳۳	۰/۳۳	۰/۹۱۵۷۹	۱/۴۴۷۰۵	۰/۵۴	۱/۰۶۰۲۱	۰/۵۴
۲۹	b(k)	۰/۱۳	۰/۲۵۹۰۲	۰/۰۵۵۹۱	-۰/۱۱۶۴۷	۰/۰۷۱۵۶	۰/۰۰۰۷	-۰/۰۰	-۰/۰۰	-۰/۳۶۸۹۱	۰/۰۰۰۵۸	-۰/۰۰	-۰/۰۰۰۲۰	-۰/۰۰
۳۰	VIF	۰/۱۴		۰/۸۵۲۳۸	۰/۸۹۹۵۹	۰/۸۷۶۱۵	۰/۹۷۶۳	۰/۳۱	۰/۳۱	۰/۸۹۱۱۷	۱/۳۴۸۸۸	۰/۵۰	۱/۰۲۶۷۸	۰/۵۰
۳۱	b(k)	۰/۱۴	۰/۲۴۲۲۱	۰/۰۵۹۴۶	-۰/۱۱۶۸۱	۰/۰۶۸۵۳	۰/۰۰۰۷	-۰/۰۰	-۰/۰۰	-۰/۳۵۹۰۲	۰/۰۰۰۵۵	-۰/۰۰	-۰/۰۰۰۲۰	-۰/۰۰
۳۲	VIF	۰/۱۵		۰/۸۲۳۷۷	۰/۸۷۶۵۳	۰/۸۵۴۰۱	۰/۸۹۱۳	۰/۲۹	۰/۲۹	۰/۸۶۷۸۰	۱/۲۶۱۲۸	۰/۴۶	۰/۹۹۵۱۹	۰/۴۷
۳۳	b(k)	۰/۱۵	۰/۲۲۶۱۳	۰/۰۶۲۵۶	-۰/۱۱۶۹۵	۰/۰۶۵۷۱	۰/۰۰۰۷	-۰/۰۰	-۰/۰۰	-۰/۳۴۹۶۸	۰/۰۰۰۵۳	-۰/۰۰	-۰/۰۰۰۱۹	-۰/۰۰
۳۴	VIF	۰/۱۶		۰/۷۹۷۵۷	۰/۸۵۴۵۶	۰/۸۳۲۹۳	۰/۸۱۸۴	۰/۲۸	۰/۲۸	۰/۸۴۵۵۷	۱/۱۸۲۷۳	۰/۴۳	۰/۹۶۵۲۷	۰/۴۴
۳۵	b(k)	۰/۱۶	۰/۲۱۰۷۳	۰/۰۶۵۲۷	-۰/۱۱۶۹۳	۰/۰۶۳۰۷	۰/۰۰۰۶	-۰/۰۰	-۰/۰۰	-۰/۳۴۰۸۳	۰/۰۰۰۵۰	-۰/۰۰	-۰/۰۰۰۱۸	-۰/۰۰

همچنین نمودار اثر ریدج<sup>۵</sup> نیز برآورد ضرایب هر یک از متغیرها را به ازای مقادیر مختلف  $k$  نشان می‌دهد. این نمودار شکل دو بعدی است که با توجه به جدول ۲ رسم می‌شود و در آن، با افزایش مقدار  $k$  برآورد ضرایب پایدارتر می‌شود [۲].



شکل ۲: نمودار اثر ریدج

## ۸ بحث و نتیجه‌گیری

در صورت وقوع هم خطی، بهتر است از روش رگرسیون ریدج برای برآورد ضرایب مدل استفاده شود. که در این حالت، MSE، حتی کمتر از ضرایب در روش کمترین مربعات باقی‌مانده است و دقت مدل بیشتر می‌شود.

## تقدیر و تشکر

نویسندگان از تمام کارکنان کتابخانه‌ی دانشکده‌ی بهداشت، دانشگاه علوم پزشکی مشهد و نیز کارکنان کتابخانه‌ی استاد فاطمی، دانشکده‌ی علوم ریاضی، دانشگاه فردوسی مشهد تقدیر و تشکر می‌نمایند.

<sup>۵</sup>Ridge trace

## مراجع

- [۱] آذرنوش، ح. ع. و دوج، ی. (۱۳۷۹)، روش‌های جایگزین در رگرسیون. ترجمه. دانشگاه فردوسی مشهد.
- [2] Hoerl, A. E and Kennard, R. W. (1970a). Ridge Regression: Applications to Nonorthogonal Problems. *Technometrics*: 69-82.
- [3] Hoerl, A. E. and Kennard, R. W. (1970b). Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics*: 55-67.
- [4] Kutner, M., Nachtsheim, C., Neter, J. and Li, W. (2005). *Applied Linear Statistical Models*. 5<sup>th</sup>ed. New York: McGraw-Hill.
- [5] Tavallaie, S., Rahsepar, A. A. et al. (2012). Association between indices of body mass and antibody titers to heat-shock protein-27 in healthy subjects. *Clinical Biochemistry*, Vol. 45, 1-2, 144-147.