

# کاربرد رگرسیون ستیغی کمترین توان‌های دوم پیراسته محدودشده تصادفی در مدل‌سازی مصرف آب

مهدی روزبه<sup>۱</sup>، ملیحه سادات ملک‌جعفریان<sup>۲</sup> و منیره معنوی<sup>۳</sup>

تاریخ دریافت: ۱۴۰۰/۰۲/۳۱

تاریخ پذیرش: ۱۴۰۰/۱۲/۲۶

## چکیده:

مهم‌ترین هدف علم آمار تجزیه و تحلیل داده‌های واقعی دنیای پیرامون بشر است. اگر این اطلاعات دقیق و درست تحلیل شوند، نتایج حاصل در بسیاری از تصمیمات مهم یاریگر ما خواهد بود. از جمله داده‌های واقعی پیرامون ما که تحلیل آن بسیار مهم است، داده‌های مربوط به مصرف آب می‌باشد. با توجه به اینکه کشور ایران در ناحیه نیمه‌خشک آب و هوایی از کره زمین قرار دارد، لازم است برای پیش‌بینی و برگزیدن بهترین و مناسب‌ترین مدل‌های دقیق مصرف آب گام‌های ژرفی برداشت که لازمه تصمیمات کلان‌کشوری می‌باشد. در تجزیه و تحلیل داده‌های واقعی ممکن است محقق با مشکل هم خطی و نقاط دورافتاده مواجه شود. روش‌های مقاوم (استوار) برای تحلیل مجموعه داده‌های دارای نقاط دورافتاده و رویکرد ستیغی روشی است که برای تحلیل مجموعه داده‌های دارای هم خطی استفاده می‌شوند. محدودیت روی مدل‌ها نیز ناشی از به‌کارگیری اطلاعات غیرنمونه‌ای در برآورد ضرایب رگرسیونی است. در این مقاله به مدل‌سازی داده‌های مصرف آب، با استفاده از رویکرد ستیغی محدودشده تصادفی استوار پرداخته می‌شود.

واژه‌های کلیدی: روش کمترین توان‌های دوم پیراسته ستیغی، محدودیت خطی تصادفی، مصرف آب، نقاط دورافتاده، هم خطی.

## ۱ مقدمه

نزدیک بین متغیرهای رگرسیونی است. وجود وابستگی خطی نزدیک، توانایی برآورد ضرایب مدل رگرسیونی با استفاده از روش کمترین توان‌های دوم معمولی را با مشکل مواجه می‌کند. در بیشتر کاربردهای رگرسیون، بین متغیرهای رگرسیونی مسئله هم خطی چندگانه وجود دارد و استنباط‌هایی نظیر مشخص کردن اثرات نسبی متغیرهای رگرسیونی، پیش‌گویی و برآورد و انتخاب یک مجموعه مناسب از متغیرها بر اساس این الگوی رگرسیون ممکن است تا حد زیادی گمراه‌کننده و پراشتباه باشد.

منابع اصلی هم خطی به‌صورت زیر است [۱۲]:

- شیوه جمع‌آوری داده‌ها.
- گذاشتن قیدهایی روی مدل.
- در نظر گرفتن یک مدل با متغیرهای رگرسیونی بیش‌ازحد نیاز.

تجزیه و تحلیل درست این داده‌های واقعی دنیای پیرامون بشر اطلاعات بسیار مهمی در مورد جامعه در اختیار محقق قرار می‌دهد. این اطلاعات و پیش‌بینی‌های حاصل از آن در تصمیم‌گیری‌های مهم و حیاتی در دنیای واقعی به کمک مسئولین خواهند آمد. اگر کوچک‌ترین اشتباهی در این تحلیل‌ها رخ دهد نتایج نامطلوب آن تا مدت‌ها گریبان‌گیر جامعه خواهد بود. بنابراین می‌بایست تحلیل‌ها با بیشترین دقت ممکن انجام شوند. اما حقیقت این است که داده‌های واقعی موجود در جامعه دارای پیچیدگی خاصی هستند و به‌سادگی و با روش‌های ساده موجود قابل تجزیه و تحلیل نیستند. یکی از رایج‌ترین مشکلات موجود در داده‌های واقعی هم خطی<sup>۴</sup> است. یک مسئله جدی که می‌تواند استفاده از مدل رگرسیونی را با اشکال مواجه کند، هم خطی چندگانه یا وابستگی خطی

<sup>۱</sup> هیئت‌علمی گروه آمار، دانشگاه سمنان، سمنان، ایران (نویسنده مسئول: mahdi.roozbeh@semnan.ac.ir)

<sup>۲</sup> شرکت آب و فاضلاب شهری استان سمنان، سمنان، ایران

<sup>۳</sup> دانش‌آموخته کارشناسی ارشد آمار، دانشگاه سمنان، سمنان، ایران.

در شکل ۱ ملاحظه می‌گردد وجود داده دورافتاده، برازش خط رگرسیون را تحت تأثیر قرار داده به طوری که با حذف این نقاط خط رگرسیونی به شدت تغییر می‌کند.

روش‌های متنوع و بسیاری برای شناسایی نقاط دورافتاده، نظیر تحلیل باقی‌مانده‌ها<sup>۱۰</sup> و ماتریس کلاهدار<sup>۱۱</sup>، پتانسیل‌ها<sup>۱۲</sup>، آماره نسبت کوواریانس<sup>۱۳</sup> و... وجود دارد. به طور کلی می‌توان روش‌های شناسایی نقاط را به دو گروه روش‌های عددی و روش‌های شهودی نیز تقسیم نمود که این روش‌ها به طور کامل در [۲] و [۶] قابل مشاهده است.

یکی دیگر از نکات بسیار مهم در برآورد ضرایب رگرسیونی یافتن محدودیت‌های موجود روی فضای پارامتر و توجه به آن‌ها در برآورد است. زیرا این محدودیت‌ها روی مقادیر به دست آمده و برآورد ضرایب تأثیر زیادی خواهند داشت. کاملاً واضح است که محدودیت فضای پارامتر فضای برآورد را نیز محدود خواهد کرد و می‌بایست برآوردهای به دست آمده نیز در این فضای محدود شده صدق کنند و نمی‌توان برآوردی یافت که در فضای محدود شده نباشند. بنابراین نادیده گرفتن این محدودیت‌ها، درستی برآوردهای به دست آمده را تحت الشعاع قرار می‌دهد. بدین جهت می‌بایست ضرایب رگرسیونی، تحت این محدودها روی مدل اعمال شوند [۱۵].

مدل‌های محدود شده به طور گسترده‌ای در مسئله آزمون فرضیات آماری، به ویژه آزمون نسبت درست نمایی تعمیم یافته در مدل‌های رگرسیون کاربرد دارند. محدودیت روی پارامترهای مدل می‌تواند به یکی از دلایل زیر رخ دهد:

۱. حقیقت ناشی از ملاحظات نظری و یا تجربی (آزمایشگاهی)،
۲. فرضیه‌ای که باید آزمون شود،
۳. نظر محقق یا متخصص،
۴. شرط غیرمعمول به منظور کاهش یا حذف فزونگی روی مدل مورد نظر [۲۰].

در دو بخش بعدی به معرفی دو روش کمترین توان‌های دوم معمولی و روش ستیغی می‌پردازیم.

وجود هم خطی ممکن است منجر به فواصل اطمینان عریض برای پارامترها، ناپایداری برآورد پارامترها<sup>۵</sup> یا تولید برآوردهایی با علامت اشتباه شود [۱]. روش‌های مختلفی برای غلبه برای هم خطی وجود دارد از جمله رگرسیون مؤلفه‌ی اصلی<sup>۶</sup>، رگرسیون ستیغی<sup>۷</sup> و ... . مفهوم رگرسیون ستیغی توسط هورل و کنارد [۹] برای مقابله با هم خطی در مسائل رگرسیون پیشنهاد شده است. از آن زمان تاکنون بسیاری از نویسندگان و محققان از این روش برای تحقیقات و پژوهش‌های خود استفاده کرده‌اند.

اولین و مهم‌ترین گام جهت انجام تحقیق، جمع‌آوری داده‌ها است. اغلب در بین مجموعه داده‌ها چند مشاهده وجود دارد که با سایر مشاهدات تفاوت دارند که تحت نام نقاط دورافتاده<sup>۸</sup> یاد می‌شوند. داده‌های دورافتاده مقادیر بسیاری از آماره‌ها و برآوردها از جمله میانگین و واریانس را تحت تأثیر قرار داده و موجب نتایج اشتباه می‌گردند. به عنوان مثال این داده‌ها با اثرگذاری بر برآورد ضرایب مدل‌های رگرسیون، مجموع توان‌های دوم خطا و غیره سبب ارائه نتایج غیرواقعی می‌گردند. حذف این‌گونه داده‌ها سبب از بین رفتن اطلاعات می‌شود. بنابراین بهترین راه برخورد با داده دورافتاده، استفاده از روش‌های استوار<sup>۹</sup> است به طوری که آماره‌ها و برآوردهای این روش تحت تأثیر داده دورافتاده قرار نگیرد.

زمانی که یک مشاهده تفاوت اساسی با دیگر مشاهدات داشته باشد، می‌تواند یک انحراف اساسی در نتایج تحلیل رگرسیونی ایجاد کند [۱۱]. محققان در مواجهه با این نقاط علاقه‌مند هستند که این مجموعه از داده‌ها را، که به آن‌ها نقاط دورافتاده نیز گفته می‌شود، تعیین و تأثیر آن‌ها را بر جنبه‌های مختلف تحلیل رگرسیونی ارزیابی کنند. نقاط دورافتاده ممکن است نتایج برازش مدل را تحت تأثیر قرار دهند، به گونه‌ای که حذف آن‌ها از مجموعه داده‌ها نتایج کاملاً متفاوتی به بار آورد. روش‌ها و معیارهای مختلفی به منظور شناسایی چنین مشاهداتی پیشنهاد شده است. البته برخی روش‌های استوار، به مشاهدات دورافتاده وزن صفر را اختصاص می‌دهند و عملاً مشاهده دورافتاده از مجموعه داده‌ها حذف می‌شود اما روش‌های دیگری نیز وجود دارد که به مشاهدات دورافتاده وزن غیرصفر ولی ناچیز اختصاص می‌دهند ([۱۴]، [۱۶] و [۱۷]).

<sup>5</sup>Instability of parameters estimation

<sup>6</sup>Principal component regression (PCR)

<sup>7</sup>Ridge regression

<sup>8</sup>Outlier

<sup>9</sup>Robust methods

<sup>10</sup>Analysis of residuals

<sup>11</sup>Hat matrix

<sup>12</sup>Potentials

<sup>13</sup>Covariance ratio statistics

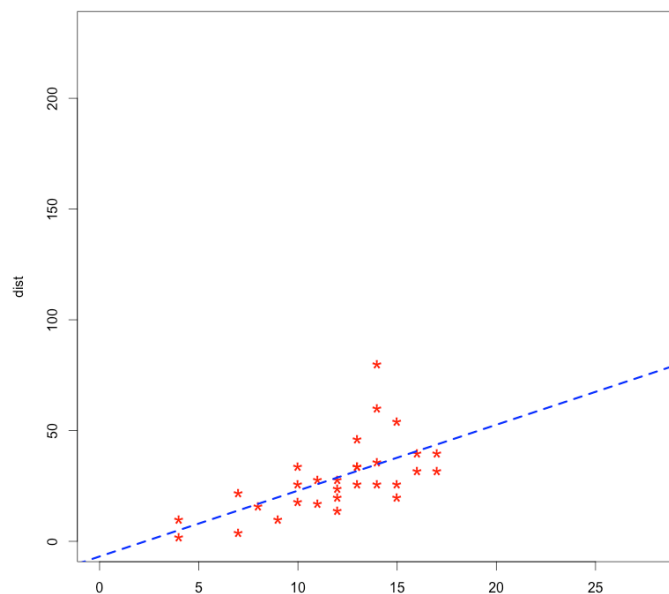
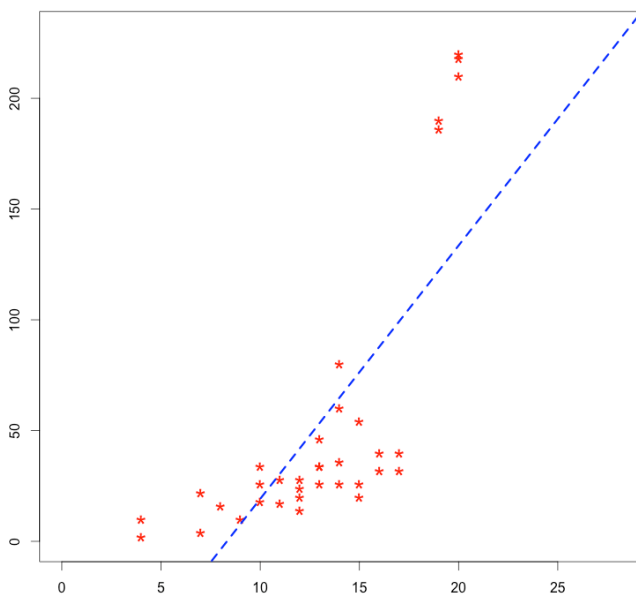
توضیحی  $j$  ام،  $\beta = (\beta_1, \dots, \beta_p)^T$  بردار ضرایب رگرسیونی و  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$  بردار خطای تصادفی با شرایط  $E(\varepsilon) = 0$  و  $E(\varepsilon\varepsilon^T) = \sigma^2 I_n$  است. لازم به ذکر است که منظور از  $0$  بردار صفر با  $n$  مؤلفه صفر و  $I_n$  ماتریس همانی  $n \times n$  است.

## ۱.۱ روش کمترین توان‌های دوم معمولی

مدل رگرسیونی خطی چندگانه را به صورت

$$y = X\beta + \varepsilon \quad (1)$$

تعریف می‌شود که در آن  $y = (y_1, \dots, y_n)^T$  متغیر پاسخ،  $x_j = (x_{1j}, \dots, x_{nj})^T, j = 1, \dots, p$  ستون  $j$  ام ماتریس  $X$  یعنی متغیر



شکل ۱: نمودار سمت چپ نمایش تأثیر داده‌های دورافتاده، نمودار سمت راست حذف داده‌های دورافتاده و برازش مجدد خط رگرسیون

## ۲ مدل‌های محدود شده

مدل‌های محدود شده، به دو گروه تصادفی و غیر تصادفی تقسیم می‌شوند که در ادامه معرفی می‌شوند.

تعریف ۱.۰۲. محدودیت خطی غیر تصادفی<sup>۱۴</sup>: اگر محدودیت به صورت

$$\hat{r} = \hat{R}\beta, \quad (2)$$

بیان شود، محدودیت خطی غیر تصادفی نام دارد به طوری که در آن  $\hat{r}$  بردار معلوم  $1 \times q$  و  $\hat{R}$  ماتریس  $q \times p$  اطلاعات پیشین معلوم روی بردار پارامترها است. ماتریس  $\hat{R}$  رتبه کامل سطری است.

لازم به ذکر است که منظور از اطلاعات پیشین، اطلاعات غیر نمونه‌ای در مورد پارامتر مجهول است به بیان ساده‌تر خلاصه‌ای از اطلاعات و دانسته‌های کاربر آمار در مورد پارامتر مجهول است.

برآورد ضرایب رگرسیونی با استفاده از روش کمترین با حل مسئله

رگرسیونی

$$\min_{\beta} (y - X\beta)^T (y - X\beta),$$

محاسبه می‌شود. طبق قضیه گاوس مارکوف برآوردگر به دست آمده با این روش، یعنی  $\hat{\beta}_{OLS} = (X^T X)^{-1} X^T y$  بهترین برآوردگر خطی نااریب با کمترین واریانس است.

## ۲.۱ روش ستیغی

در صورت وجود هم خطی از روش ستیغی برای برآورد ضرایب رگرسیونی استفاده می‌شود که این برآوردگر با حل مسئله

$$\min_{\beta} (y - X\beta)^T (y - X\beta) + k\beta^T \beta,$$

حاصل می‌شود. برآوردگر روش ستیغی، یعنی  $\hat{\beta}_{Ridge} = (X^T X + kI_p)^{-1} X^T y$  یک برآوردگر اریب است.

<sup>14</sup>Linear restricted

اطلاعات پیشین، دقیق نبوده و می‌بایست محدودیت به صورت تصادفی بیان می‌شود [۵]. محدودیت خطی تصادفی در مواردی نظیر روابط اقتصادی، ساختارهای صنعتی، برنامه‌ریزی‌های تولیدی و ... کاربرد دارند [۴]. اگر عدم اطمینان در مورد اطلاعات پیشین وجود داشته باشد، یک جایگزین مناسب برای محدودیت‌های خطی، محدودیت‌های خطی تصادفی است.

تعریف ۲۰۲. محدودیت خطی تصادفی<sup>۱۷</sup>: اگر محدودیت به صورت

$$\mathbf{r} = \mathbf{R}\beta + \nu, \quad (4)$$

بیان شود، محدودیت خطی تصادفی نام دارد به طوری که در آن  $\mathbf{r}$  بردار معلوم  $1 \times q$ ،  $\mathbf{R}$  ماتریس  $q \times p$  اطلاعات پیشین معلوم روی بردار پارامترها و  $\nu$  خطای تصادفی (اطلاعات پیشین مجهول) است. ماتریس  $\mathbf{R}$  سطری رتبه کامل است.

تیل و گلدبرگر [۲۱] با اعمال محدودیت خطی تصادفی (۴) در مدل (۱)، مدل جدیدی با فرض استقلال  $\varepsilon$  و  $\nu$  به صورت

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{r} \end{bmatrix} = \begin{bmatrix} \mathbf{X} \\ \mathbf{R} \end{bmatrix} \beta + \begin{bmatrix} \varepsilon \\ \nu \end{bmatrix},$$

به دست آوردند که می‌توان آن را به صورت

$$\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta + \tilde{\varepsilon}, \quad \tilde{\mathbf{y}} = \begin{bmatrix} \mathbf{y} \\ \mathbf{r} \end{bmatrix}, \quad \tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \mathbf{R} \end{bmatrix}, \quad \tilde{\varepsilon} = \begin{bmatrix} \varepsilon \\ \nu \end{bmatrix}, \quad (5)$$

بازنویسی کرد که در واقع، همان مدل رگرسیونی خطی ساده با شرایط  $E(\tilde{\varepsilon}\tilde{\varepsilon}^T) = \sigma^2 \mathbf{I}_{(n+q) \times p}$  و  $E(\tilde{\varepsilon}) = 0$  است.

برآورد ضرایب رگرسیونی مدل محدود شده تصادفی با فرض وجود هم خطی در مجموعه داده‌ها با روش ستیغی نیز به صورت

$$\min_{\beta} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta)^T (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta) + k\beta^T \beta,$$

قابل محاسبه است که در آن  $k$  پارامتر ستیغی نام دارد. برای محاسبه این پارامتر از روش اعتبارسنجی مقابل استفاده می‌شود. برآوردگر به روش ستیغی با محدودیت تصادفی برابر  $\hat{\beta}_{Ridge}^{SR} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + k\mathbf{I})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}$  است.

طبق [۱۵] با ترکیب محدودیت خطی غیر تصادفی (۲) و مدل (۱) برآوردگر محدود شده غیر تصادفی به صورت زیر قابل محاسبه است:

$$\hat{\beta}^R = \hat{\beta}_{ols} + (\mathbf{X}^T \mathbf{X})^{-1} \dot{\mathbf{R}} (\dot{\mathbf{R}} (\mathbf{X}^T \mathbf{X})^{-1} \dot{\mathbf{R}}^T)^{-1} (\dot{\mathbf{r}} - \dot{\mathbf{R}} \hat{\beta}_{ols}),$$

که در آن  $\hat{\beta}_{ols} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  است.

سرکار [۱۹] برآوردگر ستیغی محدود شده غیر تصادفی را برای مواردی که هم خطی در مجموعه داده‌ها وجود داشته باشد، به صورت  $\hat{\beta}_{Ridge}^{*R} = (\mathbf{I}_p + k(\mathbf{X}^T \mathbf{X})^{-1})^{-1} \hat{\beta}^R$  معرفی کرد که در آن  $k$  پارامتر ستیغی نام دارد. متأسفانه عبارت  $\dot{\mathbf{R}} \hat{\beta}_{Ridge}^{*R} = \dot{\mathbf{r}}$  برای این برآوردگر برای هر  $\beta$  برقرار نخواهد بود [۷]. بنابراین برآوردگر  $\hat{\beta}_{Ridge}^{*R}$  برآوردگر مناسبی نبوده و لذا او برآوردگر

$$\hat{\beta}_{Ridge}^R = \hat{\beta}_{Ridge} + (\mathbf{X}^T \mathbf{X} + k\mathbf{I})^{-1} \times (\dot{\mathbf{R}}^T (\mathbf{X}^T \mathbf{X} + k\mathbf{I})^{-1} \dot{\mathbf{R}})^{-1} (\dot{\mathbf{r}} - \dot{\mathbf{R}} \hat{\beta}_{Ridge}), \quad (3)$$

را پیشنهاد کرد به طوری که  $\hat{\beta}_{Ridge}^R = (\mathbf{X}^T \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$  است.

## ۱۰۲ تعیین پارامتر ستیغی

تعیین مقدار بهینه پارامتر ستیغی نقش بسیار مهمی در نتایج نهایی دارد. مقدار بهینه این پارامتر توازن بین اریبی و واریانس را برقرار می‌کند. در این شرایط اریبی و واریانس به طور هم‌زمان به کمترین حالت ممکن خود می‌رسند. برای محاسبه این پارامتر از روش‌های مختلفی می‌توان استفاده کرد که در این مقاله از روش اعتبارسنجی متقابل<sup>۱۵</sup> استفاده می‌شود. اعتبارسنجی متقابل انواع مختلفی دارد که در اینجا روش حذف تکی<sup>۱۶</sup> مورد استفاده قرار می‌گیرد. این روش برای مدل رگرسیونی (۱) به صورت

$$LOCV = \|\mathbf{y} - \hat{\mathbf{y}}_{(-i)}\|^2,$$

تعریف می‌شود که در آن  $\hat{\mathbf{y}}_{(-i)} = \mathbf{X}_{(-i)} \hat{\beta}_{(-i)}$  به طوری که  $\hat{\beta}_{(-i)}$  برآوردگر دلخواه است که با حذف  $i$  امین سطر ماتریس  $\mathbf{X}$  ( $\mathbf{X}_{(-i)}$ ) و  $i$  امین مشاهده بردار  $\mathbf{y}$  ( $\mathbf{y}_{(-i)}$ ) به دست می‌آید. به بیان ساده‌تر یعنی مدلی با حذف مشاهده  $i$  ام برازش داده می‌شود و بعد با استفاده از مدل برازش داده شده مشاهده  $y_i$  پیش‌بینی می‌شود. به عبارت دیگر  $\hat{y}_i$  مقدار پیش‌بینی شده مشاهده  $i$  ام از مدل حاصل با حذف مشاهده  $i$  ام است. در عمل، موقعیت‌های بسیاری وجود دارد که در آن فرض می‌شود که اطلاعات پیشین مناسب نیستند و یا به بیان ساده‌تر، در حقیقت

<sup>15</sup>Cross-validation (CV)

<sup>16</sup>Leave-one-out cross-validation (LOCV)

<sup>17</sup>Stochastic linear restricted

را حل کنیم که در آن  $Z$  ماتریس قطری با عناصر قطری  $z = (z_1, \dots, z_n)^T, z_i \in \{0, 1\}$  است. بدین طریق برآورد ضرایب رگرسیونی مدل محدودشده تصادفی با فرض وجود هم خطی و نقاط دورافتاده در مجموعه داده‌ها از رابطه

$$\hat{\beta}_{LTSRidge}^{SR} = (\tilde{X}^T Z \tilde{X} + kI)^{-1} \tilde{X}^T Z \tilde{y}$$

محاسبه می‌شود که به‌طور مشابه فرم غیرتصادفی نیز قابل بازنویسی است.

## ۴ نتایج عددی

در این بخش مطالعه عددی روی داده‌های واقعی انجام می‌شود. در هر دو مجموعه داده ضرایب رگرسیونی به چهار روش کمترین توان‌های دوم معمولی (OLS)، روش ستیغی (Ridge)، روش کمترین توان‌های دوم پیراسته ستیغی (LTS) و روش کمترین توان‌های دوم پیراسته ستیغی غیر تصادفی (LTSridge) در سه حالت بدون محدودیت، محدودشده خطی غیر تصادفی و محدودشده خطی تصادفی برآورد می‌شوند. نتایج نهایی با معیار مجموع توان‌های دوم خطا تحت ارزیابی قرار می‌گیرند.

### ۱.۴ مجموعه داده‌های واقعی

در این بخش مجموعه داده واقعی مربوط به مصرف آب موردبررسی قرار می‌گیرد.

#### ۱.۱.۴ مجموعه داده آب

شرکت‌های آب و فاضلاب کشور حاوی داده‌های بسیار با اهمیت از حیث و نگاه مشترکین آب می‌باشد. طبق آیین‌نامه عملیاتی وزارت نیرو برای تعریف اشتراک آب، هر شخص حقیقی یا حقوقی که انشعاب آب یا انشعاب‌های (آب و یا فاضلاب) مورد تقاضای وی، طبق مقررات برقرارشده باشد و در تعریف دیگر هر شخص حقیقی یا حقوقی که برای مصرف آب شرب و سایر خدمات دریافتی از شرکت برای آن قبض صادر می‌گردد، در لیست مشترکین آب در شرکت‌های آب و فاضلاب ثبت می‌گردد و از خدمات این شرکت‌ها بهره‌مند خواهد گردید. کنتورهای نصب‌شده در محل اشتراک آب به‌صورت دوره‌ای کنتورخوانی می‌شوند (حدود شش الی هفت بار در سال) و مصارف هر یک از مشترکین جهت محاسبه مبلغ آب‌بها در سامانه‌های آیفنا ثبت می‌گردد. داده‌های مورد استفاده در این پژوهش برحسب استانی بوده و به تفکیک سال‌های

<sup>18</sup>Least trimmed squares(LTS)

<sup>19</sup>Trimmed parameter

## ۳ روش کمترین توان‌های دوم پیراسته

یکی از مشهورترین روش‌های موجود برای برآورد ضرایب رگرسیونی در داده‌هایی که دارای نقاط دورافتاده هستند روش کمترین توان‌های دوم پیراسته<sup>۱۸</sup> [۱۸] است که در ادامه تعریف می‌شود. برای محاسبه برآوردگر روش کمترین توان‌های دوم پیراسته کافی است مسئله

$$\begin{aligned} \min_{\beta, Z} \varphi(\beta, Z) &= (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{Z}(\mathbf{y} - \mathbf{X}\beta) \\ \text{s.t.} \quad \mathbf{e}^T \mathbf{z} &= h, \quad z_i \in \{0, 1\}, \quad i = 1, \dots, n, \end{aligned}$$

حل شود. که در آن  $\mathbf{z} = (z_1, \dots, z_n)$  و  $\mathbf{e} = (1, \dots, 1)_{n \times 1}^T$  است. همچنین پارامتر  $h$ ، پارامتر پیراسته<sup>۱۹</sup> است. لازم به ذکر است که ماتریس قطری  $Z$  با عناصر قطری  $z_i$  به‌منظور تعیین مشاهده دورافتاده و از بین بردن تأثیر آن روی مدل، در مسئله بهینه‌سازی حضور دارد، به‌طوری‌که در آن به مشاهدات دورافتاده وزن صفر اختصاص می‌یابد. یعنی،

$$z_i = \begin{cases} 1 & \text{مشاهده } i \text{ ام دورافتاده نباشد} \\ 0 & \text{مشاهده } i \text{ ام دورافتاده باشد} \end{cases}$$

تعریف می‌شود. لازم به ذکر است که در صد شکست (میزان استواری و مقاومت روش ذکرشده در برابر نقاط دورافتاده) این روش ۵۰ درصد است. یعنی حتی اگر نیمی از داده‌ها دورافتاده باشند بازهم این روش در برابر نقاط دورافتاده، استوار است [۱۳].

### ۱.۳ مدل‌های رگرسیونی کمترین توان‌های دوم

#### پیراسته محدودشده

در شرایطی که به‌طور هم‌زمان با هم خطی و نقاط دورافتاده در داده‌ها مواجه باشیم ناچار به ترکیب روش ستیغی با روش کمترین توان‌های دوم پیراسته خواهیم بود. در این شرایط برآورد ضرایب رگرسیونی با حل

$$\min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{Z}(\mathbf{y} - \mathbf{X}\beta) + k\beta^T \beta,$$

به دست می‌آید که در آن  $Z$  ماتریس قطری با عناصر قطری  $z = (z_1, \dots, z_n)^T, z_i \in \{0, 1\}$  است. در این روش  $z_i = 0$  یعنی مشاهده  $i$  ام یک نقطه دورافتاده است.

به همین ترتیب برای برآورد ضرایب رگرسیونی مدل محدودشده تصادفی (۵) با فرض وجود هم خطی و نقطه دورافتاده در مجموعه داده کافی است مسئله

$$\min_{\beta} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta)^T \mathbf{Z}(\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta) + k\beta^T \beta.$$

پراکنش نقاط دورافتاده نمایش داده شده است. نقاط دورافتاده با علامت “+” مشخص شده اند. نمودار انواع نقاط در مجموعه داده‌ها نیز در شکل ۳ قابل مشاهده است. طبق این نمودار مشاهدات ۱۲، ۱۴، ۱۶، ۱۸ و ۱۹ نقاط دورافتاده (پرت و اهرمی بد) هستند. مشاهدات ۱۹ و ۱۸ که به عنوان نقاط پرت معرفی شده‌اند، مربوط به داده‌های سال‌های ۱۳۹۷ و ۱۳۹۸ می‌باشد که با توجه به ثبات تقریبی مصرف آب، افزایش چند درصدی تعرفه آب، میزان آب‌بها رشد به سزایی داشته است. همچنین داده‌های ۱۲، ۱۴، ۱۶ مربوط به سال‌های قبل از هدفمندی یارانه‌ها است که مبلغ آب‌بها با تعرفه‌های قدیمی است.

معیار عامل تورم واریانس<sup>۲۱</sup> برای متغیرهای توضیحی در جدول ۱ نشان داده شده است. عامل تورم واریانس تمام متغیرهای توضیحی به جز متغیر *waterproduc* بزرگتر از ۱۰ است. همچنین مقدار آماره کاپا<sup>۲۲</sup> برابر  $10^{12} \times 10^{17124}$  می‌باشد. نمودار ۴ همبستگی میان متغیرهای مجموعه داده را نمایش می‌دهد. این موارد نشان‌دهنده وجود هم خطی بالا در مجموعه داده‌ها است.

بر اساس برآوردگر کمترین توان‌های دوم معمولی و اطلاعات و مشاهدات موجود، ماتریس محدودیت به صورت

$$R = \begin{bmatrix} 0 & 1 & 3 & 50 & 1 \\ 0 & 4 & 1 & -60 & 2 \\ 1 & 0 & -1 & -0.5 & -1 \end{bmatrix}, \quad (7)$$

در نظر گرفته شد و بردار  $r$  به صورت تقریبی از رابطه

$$r \simeq R\hat{\beta}_{ols} = \begin{bmatrix} 1809800 \\ 3945460 \\ -1820759 \end{bmatrix},$$

به دست آمد. بنابراین در مسائل کاربردی و عملی در نظر گرفتن محدودیت تصادفی خطی به جای محدودیت دقیق واقع‌بینانه‌تر به نظر می‌آید. اکنون لازم است فرض  $H_0: R\beta \simeq r$  برای مدل (۶) آزمون شود. آماره آزمون این فرضیه به صورت

$$\chi^2 \simeq (R\hat{\beta} - r)^T (R\hat{\Sigma}R^T)^{-1} (R\hat{\beta} - r) = 0.004352015,$$

است که در آن

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y, \\ \hat{\sigma}^2 &= \frac{1}{n-p} (y - X\hat{\beta})^T (y - X\hat{\beta}) \\ \hat{\Sigma} &= \hat{\sigma}^2 (X^T X)^{-1}. \end{aligned}$$

۱۳۸۱ الی ۱۳۹۸ یعنی بازه ۱۹ ساله است، بدین معنی که هر یک از متغیرها برحسب سال موردنظر جمع‌آوری گردیده است. قابل ذکر است که این داده‌ها مربوط به شهرهای استان سمنان می‌باشد و مشترکین روستایی شامل این داده‌ها نیست. در ادامه این بخش به تفسیر هر یک از متغیرها پرداخته می‌شود. متغیر مصرف آب نشان‌دهنده میزان مصرف کلیه مشترکین آب در کلیه کاربری‌ها طی یک سال است که برحسب واحد مترمکعب در سال ثبت می‌شود. این مقدار عددی میزان دقیق مصرف آب هر مشترک می‌باشد که توسط کنترل‌رسان صورت گرفته و در سیستم جهت صدور قبوض ثبت می‌گردد. متغیر تعداد مشترکین نشان‌دهنده تعداد مشترکین موجود در استان در کلیه کاربری‌ها در پایان سال مذکور است که برحسب واحد فقره می‌باشد. لازم به ذکر است که فقط مشترک‌های آب تنها موردنظر می‌باشد. متغیر جمعیت تعداد نفرات ساکن در شهرهای استان می‌باشد که برحسب پیش‌بینی حاصله از سرشماری‌های مرکز آمار برحسب نفر حاصل شده است. از آنجاکه در سال‌های ۱۳۸۵، ۱۳۹۰ و ۱۳۹۵ در کشور سرشماری صورت گرفته است، لذا در این سال‌ها اعداد جمعیتی اعداد سرشماری است و مابقی سال‌ها برآورد جمعیت می‌باشد. متغیر تولید آب به معنای آب تولیدی از منابع تأمین آب برای شهرهای استانی است که برحسب مترمکعب در سال می‌باشد. به‌طورکلی منابع تأمین آب به دو دسته زیرزمینی و سطحی تقسیم می‌شود که در استان سمنان اکثر منابع تأمین آب زیرزمینی است. متغیر آب‌بها، بهای خالص دریافتی از مشترکین برحسب آب مصرفی مشترکین می‌باشد که برحسب پلکان‌های تعرفه آب مصرفی لحاظ شده است. لازم به ذکر است این متغیر شامل هیچ‌یک از موارد دیگر ذکر شده در قبوض نمی‌باشد. به‌عنوان مثال قبوض آب شامل مواردی نظیر آب‌نمان، تبصره‌ها، ماده ۱۱ فاضلاب‌بها و ... است. مدل این مجموعه داده‌ها به صورت

$$\begin{aligned} (Wateruse)_i &= \beta_1 (Waterproduce)_i + \beta_2 (Money)_i \\ &+ \beta_3 (Year)_i + \beta_4 (Crowd)_i + \beta_5 (Branch)_i + \epsilon_i, \\ i &= 1, \dots, 19, \end{aligned} \quad (8)$$

است که در آن متغیر وابسته میزان مصرف آب (*Wateruse*) است و متغیرهای توضیحی عبارت‌اند از سال آبی (*Year*)، مبلغ آب‌بها (*Money*)، جمعیت (*Crowd*)، میزان تولید آب (*Waterproduce*) و تعداد انشعاب‌ها (*Branch*)<sup>۲۰</sup>. در ادامه به بررسی وجود هم خطی و نقاط دورافتاده در مجموعه داده‌ها پرداخته می‌شود. در شکل ۲ نحوه<sup>۲۰</sup> منبع این داده‌ها شرکت آب و فاضلاب شهری استان سمنان می‌باشد.

<sup>21</sup> Variance Inflation Factor

<sup>22</sup> Kappa

بنابراین فرض صفر پذیرفته می‌شود. لازم به ذکر است که محدودیت‌های در نظر گرفته شده در مثال واقعی صرفاً بر مبنای اطلاعات پیشین و با استفاده از برآوردگر کمترین توان‌های دوم معمولی به دست آمده و به دلیل مطمئن شدن در مورد آن، از آزمون فرض محدودیت خطی که در [۱۲] آمده نیز استفاده شده است، تا علاوه بر پیاده‌سازی برآوردگرهای پیشنهادی جنبه آموزشی نیز داشته باشد.

نتایج در جداول ۲ و ۳ گزارش شده‌اند. طبق جدول ۲ بهترین مدل برازش شده بر اساس معیار میانگین توان‌های دوم خطا، برآورد کمترین توان‌های دوم پیراسته ستیغی است. به عبارت دیگر برآورد رگرسیونی کمترین توان‌های دوم پیراسته ستیغی نقش پررنگ‌تری در

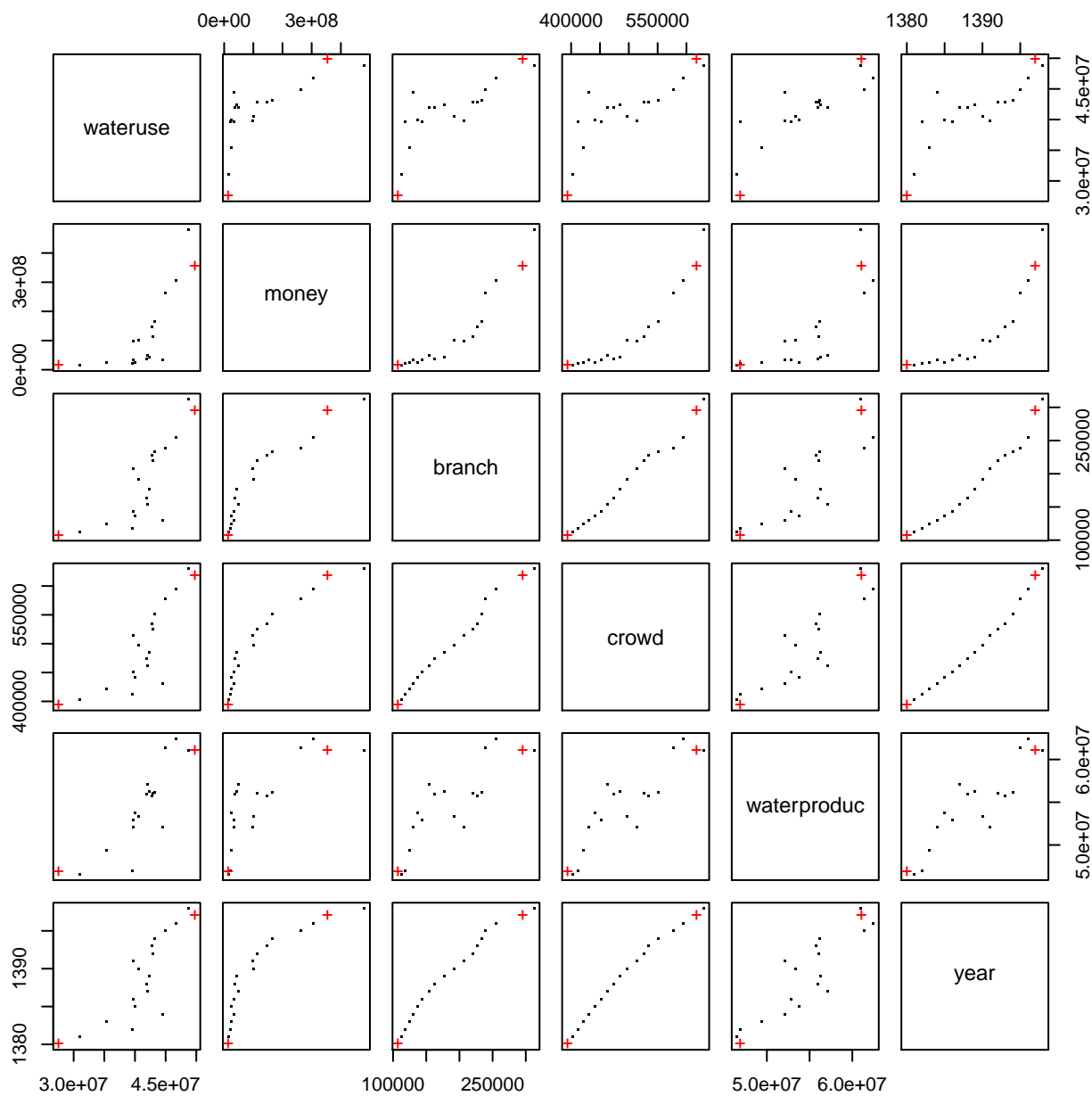
توجیه واریانس کل دارد و میزان خطا در این مدل کمتر است که این امر کاملاً مورد انتظار بود چراکه در این مجموعه از داده‌ها به طور هم‌زمان نقاط دورافتاده و هم خطی وجود دارد. همچنین طبق جدول ۳ بهترین برآورد، برآورد کمترین توان‌های دوم پیراسته ستیغی محدود شده است که طبق معیار میانگین توان‌های دوم خطا نسبت به حالت غیر محدود شده، عملکرد بهتری دارد. بنابراین، بهترین برازش داده شده به متغیرهای آبی تعریف شده در این پژوهش که می‌تواند به عنوان برآورد پیش‌بینی مصرف آب به کار رود، برازش حاصل از روش محدود شده تصادفی کمترین توان‌های دوم پیراسته ستیغی است.

جدول ۱: عامل تورم واریانس برای متغیرها در مجموعه داده آب

$VIF_{money}$	$VIF_{branch}$	$VIF_{crowd}$	$VIF_{waterproduc}$	$VIF_{year}$
۳۷۸۱۲۶	۱۳۱,۲۷۲۴	۴۹۰,۱۸۵۲	۹,۶۰۶۰	۲۷۳,۴۱۹۶

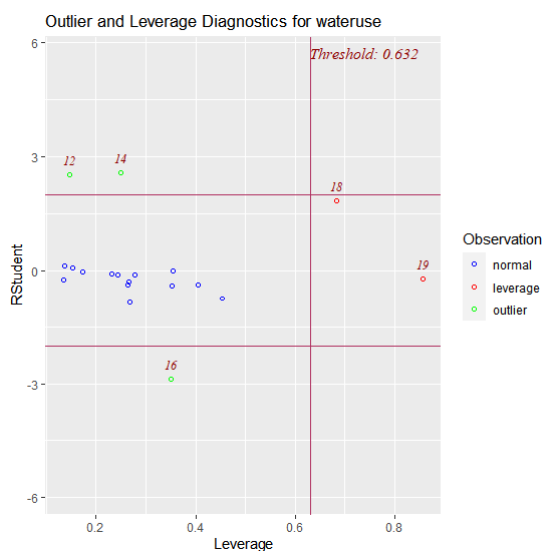
جدول ۲: برآوردگر ضرایب رگرسیونی برای مدل محدود نشده در مجموعه داده آب

LTSRidge	LTS	Ridge	OLS	روش برآورد
-۰,۴۵۳۷	-۰,۰۰۶۴	-۰,۳۸۵۸	۰,۳۹۲۷	$\hat{\beta}_1$
-۰,۶۶۳۶	۱,۱۳۸۴	-۰,۵۱۲۹	۰,۵۰۱۶	$\hat{\beta}_2$
-۰,۰۶۳۴	-۱,۹۳۳۰۸	۰,۱۱۵۵	-۲,۱۶۰۸	$\hat{\beta}_3$
۰,۷۱۹۲	۰,۸۸۴۴	۰,۲۱۳۳	۰,۷۰۸۸	$\hat{\beta}_4$
۰,۶۸۹۴	۰,۹۹۳۰	۰,۶۸۸۷	۱,۴۹۷۲	$\hat{\beta}_5$
۴,۳۰۴	۱۰,۹۹۷۷۳	۶,۰۵۱	۱۵,۵۳۰۱۲	MSE
۰,۰۰۳	-	۰,۰۰۴	-	k

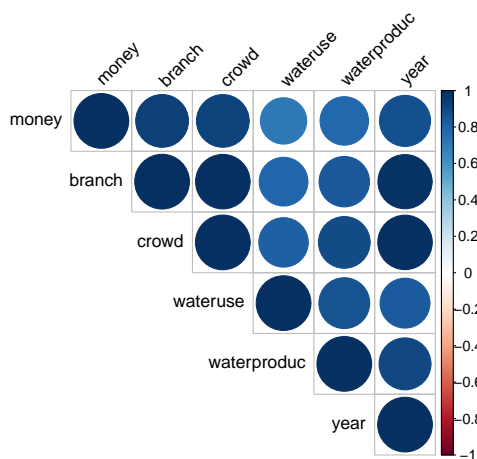


شکل ۲: نمودار پراکنش متغیرها در مجموعه داده مصرف آب





شکل ۳: نمودار انواع نقاط مجموعه داده مصرف آب



شکل ۴: نمودار همبستگی متغیرهای مجموعه داده آب

جدول ۳: برآوردگر ضرایب رگرسیونی برای مدل محدودشده خطی تصادفی در مجموعه داده آب

LTSRidge	LTS	Ridge	OLS	روش برآورد
-۰٫۴۶۸۳	۰٫۰۰۳۵	-۰٫۴۷۳۵	-۰٫۰۳۵۹	$\hat{\beta}_1$
-۰٫۲۳۰۹	۰٫۵۴۴۸	-۰٫۲۱۷۵	۰٫۳۸۵۳	$\hat{\beta}_2$
-۰٫۰۲۷۹	-۰٫۸۹۶۷	۰٫۰۵۶۰	-۰٫۴۹۸۳	$\hat{\beta}_3$
-۰٫۵۱۳۵	۱٫۰۰۱۲	-۰٫۶۴۷۳	۱٫۱۰۱۸	$\hat{\beta}_4$
۱٫۳۱۱۸	۰٫۳۹۳۳	۰٫۱۸۲۱	۰٫۰۵۹۲	$\hat{\beta}_5$
۰٫۲۸۸	۰٫۳۱۶۱۸۱	۰٫۷۴	۲٫۹۵۱۰۷۹	MSE
۰٫۰۰۱	-	۰٫۰۰۲	-	k

## ۵ بحث و نتیجه‌گیری

داده‌های واقعی اعمال محدودیت روی مدل و به‌کارگیری اطلاعات غیر نمونه‌ای باعث بهبود نتایج و کاهش میانگین مجموع توان‌های دوم خطا شده است و این امر بدان معنی است که اعمال محدودیت باعث افزایش دقت مدل و نتایج خواهد شد.

همان‌طور که ذکر شد در مواردی که هم خطی و نقاط دورافتاده در مجموعه داده واقعی وجود داشته باشد، روش کمترین توان‌های دوم معمولی عملکرد مناسبی نخواهد داشت. روش ستیغی برای رفع مشکل هم خطی مناسب است اما زمانی که به همراه هم خطی نقاط دورافتاده در مدل وجود داشته باشند نیز روش مناسبی نخواهد بود و در این حالت لازم است که از روش کمترین توان‌های دوم پیراسته ستیغی که به‌طور هم‌زمان بر این دو مشکل غلبه می‌کند، استفاده کرد. طبق نتایج مجموعه داده‌های واقعی روش کمترین توان‌های دوم پیراسته ستیغی بهترین عملکرد را نسبت به سه روش دیگر داشته است. طبق نتایج

## تقدیر و تشکر

نویسندگان مقاله ضمن تشکر از اعضای محترم هیئت تحریریه مجله، از پیشنهادها و نظرات ارزشمند داوران و ویراستار محترم مقاله که موجب ارتقاء سطح آن گردید کمال تشکر و قدردانی را دارند.

## مراجع

- [۱] امینی، م.، روزبه، م.، زمانی، ح. (۱۳۹۷)، تحلیل رگرسیون پیشرفته با R، انتشارات علمی پارسیان، تهران.
- [۲] معنوی، م.، (۱۳۹۸)، روش‌های تحلیل رگرسیونی در داده‌های با بعد بالا، پایان‌نامه کارشناسی ارشد، دانشگاه سمنان، سمنان.
- [3] Alheety, M. and Kibria, G. (2019). A new version of unbiased ridge regression estimator under the stochastic restricted linear regression model, *Communications in statistics: Simulation and computation*.
- [4] Arashi, M., Kibria, B.M.G and Valizadeh T., (2017). On Ridge Parameter Estimators under Stochastic Subspace Hypothesis, *Journal Statistical computation and simulation*, **5**, 966-983.
- [5] Belsley, D.A., Kuh, E., Welsch, R.E., (2004). *Regression Diagnostics Identifying Influential Data and Sources of Collinearity*, Wiley and Sons, New Jersey.

- [6] Chatterjee, S. and Hadi, A. S., (1986). "Influential observation, High leverage points, and outliers in linear regression," *Journal Statistical Science*. **1** no.3, 379-416.
- [7] Grob, J. (2003) Restricted ridge estimation, *Statistics & Probability Letters*, **65**, 57-64.
- [8] Hald, A. (1952). *Statistical Theory with Engineering Applications*, NewYork: John Wiley.
- [9] Hoerl A. E. and Kennard R. W. (1970), Ridge regression: biased estimation for non-orthogonal problems, *Thechnometrics*, **12**, 55-67.
- [10] Kibria, B. M. G., (2005). Applications of some improved estimators in linear regression, *Journal Modern Applied Statistical Methods*, **5(2)**, 367-380.
- [11] Liu H., Shah S. and Jiang W. (2004). On-line outlier detection and data cleaning, *Computers and Chemical Engineering*, **28 (9)**, 1635-1647.
- [12] Montgomery, D.C., Peck, E. A. and Vining G. G. (2012) Introduction to Linear Regression Analysis, 5th Edition, John Wiley & Sons, New Jersey.
- [13] Roozbeh M. (2016) Robust ridge estimator in restricted semi parametric regression models, *Journal of Multivariate Analysis*, **147**, 127-144.
- [14] Roozbeh M. and Aishah N. A. (2017). Feasible robust estimator in restricted semiparametric regression models based on the LTS approach, *Communications in Statistics – Simulation & Computation*, **46**, 7332-7350.
- [15] Roozbeh M. and Aishah N. A. (2020). Uncertain stochastic ridge estimation in partially linear regression models with elliptically distributed errors, *Statistics: A Journal of Theoretical and Applied Statistics*, **54**, 494-523.
- [16] Roozbeh, M. and Arashi, M. (2017). Least-trimmed squares: asymptotic normality of robust estimator in semiparametric regression models, *Journal of Statistical Computation & Simulation*. **147**, 1130-1147.
- [17] Roozbeh, M. and Babaie-Kafaki, S. (2016). Extended least trimmed squares estimator in semiparametric regression models with correlated errors, *Journal of Statistical Computation and Simulation*. **86(2)**, 357-372.
- [18] Rousseeuw, P.J. and Leroy, A.M. (1987), *Robust Regression and Outlier Detection*, John Wiley, New York.
- [19] Sarkar, N. (1992). A new estimator combining the ridge regression and the restricted least squares methods of estimation, *Communications in Statistics Theory Methods*, **21**, 1987-2000.
- [20] Sengupta, D. and Jammalamadaka, S.R.(2003). *Linear Models: An Integrated Approach*. World Scientific Publishing Company.
- [21] Theil, H., Goldberger, A.S., (1961). On pure and mixed statistical estimation in economics, *International Economic Review*, **2**, 65-78.

## Application of stochastic restricted least trimmed squares ridge regression in water consumption modeling

Mahdi Roozbeh<sup>1</sup>, Mlihe Malekjafarian<sup>2</sup> and Monireh Manavi<sup>3</sup>

### Abstract:

The most important goal of statistical science is the analysis of the real data of the world around us. If this information is analyzed accurately and correctly, the obtained results will help us in many important decisions. Among the real data around us which its analysis is very important, is the water consumption data. Considering that Iran is located in a semi-arid climate area of the earth, it is necessary to take big steps for predicting and selecting the best and the most appropriate accurate models of water consumption, which is necessary for the macro-national decisions. In the analysis of the real data set, we usually encounter with the problems of multicollinearity and outliers points. Robust methods are used for analyzing the data sets with outliers and ridge method is used for analyzing the data sets with multicollinearity. Also, the restriction on the models is resulted from using non-sample information in estimation of regression coefficients. In this paper, it is proceeded to model the water consumption data using robust stochastic restricted ridge approach.

**Keywords:** Multicollinearity, Outliers, Ridge least trimmed squares method, Stochastic linear restriction, Water consumption.

---

<sup>1</sup>Faculty of mathematics, Semnan university, Semnan, Iran.

<sup>2</sup>Facing the Governor of Semnan Province City Water and Wastewater Company, Semnan, Iran.

<sup>3</sup>Master's degree graduate, statistics and Computer science, Semnan university, Semnan, Iran.