

## مدل سازی سن تقویمی به روش رگرسیون ستیغی کمترین توان های دوم پیراسته

مهدی روزبه، منیره معنوی

گروه آمار، دانشکده علوم ریاضی، آمار و علوم کامپوتر، دانشگاه سمنان

تاریخ دریافت: ۱۳۹۸/۰۲/۲۶ تاریخ آخرین بازنگری: ۱۳۹۸/۱۰/۳۰

**چکیده:** متداول ترین روش برای برآورد پارامترهای مدل رگرسیون خطی، روش کمترین توان های دوم معمولی است که علی رغم سادگی محاسبه و دستیابی به بهترین برآورد خطی نااریب از پارامترها، گاهی منجر به جواب های گمراه کننده می شود. به عنوان مثال می توان به مشکلات ناشی از وجود همخطی و داده های دورافتاده در مجموعه داده ها اشاره کرد. روش کمترین توان های دوم پیراسته که یکی از معروف ترین روش های رگرسیون استوار است، تاثیر داده های دورافتاده را تا حد امکان کم می کند. هدف اصلی این مقاله ارائه یک برآورد ستیغی استوار در مدل سازی مربوط به داده های سن دندانی است. در بین روش هایی که برای تعیین سن استفاده می شود، رایج ترین روش در سراسر دنیا، روش نوین تعمیم یافته دمیرجیان است که بر اساس سخت شدگی دندان دائمی در رادیوگرافی پانورامیک بنا شده است. نشان داده شده است که استفاده از برآوردگر ستیغی استوار منجر به کاهش میانگین توان دوم خطای برآورد در مقایسه با برآوردگر کمترین توان های دوم معمولی می شود. البته برآوردگرهای پیشنهادی در داده های شبیه سازی شده نیز مورد ارزیابی قرار گرفتند.

**واژه های کلیدی:** برآوردگر ستیغی استوار، سن دندانی دمیرجیان، کمترین توان های دوم پیراسته.

## ۱ مقدمه

مدل رگرسیونی خطی به صورت

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n,$$

است، که در آن  $y_i$  متغیر پاسخ،  $\mathbf{x}_i^T = (x_{i1}, \dots, x_{ip})$  بردار متغیر توضیحی،  $\boldsymbol{\beta}$  بردار مجهول پارامترها و  $\varepsilon_i$  جمله خطا با شرایط

$$E(\varepsilon_i) = 0, \quad \text{Var}(\varepsilon_i) = \sigma^2, \quad \text{Cov}(\varepsilon_i, \varepsilon_j) = 0, \quad i \neq j,$$

است. برآوردگر بردار ضرایب رگرسیونی به روش کمترین توان‌های دوم به صورت  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  است، که در آن  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$  و  $\mathbf{y} = (y_1, \dots, y_n)^T$ . این برآورد، بهترین برآوردگر ناریب خطی با کمترین واریانس<sup>۱</sup> است. مانده‌ها بر اساس برآوردگر کمترین توان‌های دوم به صورت  $e_i = y_i - \hat{y}_i$  تعریف می‌شود که در آن  $\hat{y}_i = \mathbf{x}_i^T \hat{\boldsymbol{\beta}}$  مقدار پیش‌بینی شده برای مشاهده  $i$ ام است. این روش بسیار ساده و کاربردی است اما در مواردی نظیر وجود نقطه پرت در مجموعه داده و همخطی بین متغیرهای توضیحی بسیار ضعیف عمل می‌کند.

## ۲ روش‌های استوار

اولین برآورد استوار توسط **اجورث** (۱۸۸۷) معرفی شد. وی استدلال کرد که نقاط دورافتاده تاثیر زیادی روی روش کمترین توان‌های دوم دارند. زیرا این روش در برآورد پارامترها به هر مشاهده وزن یکسان می‌دهد و بنابراین در نظرگرفتن توان دوم مانده‌ها باعث افزایش تاثیر نقاط دورافتاده روی برآوردگرها می‌شود. به عبارت دیگر، برآوردگر کمترین توان‌های دوم معمولی از کمینه‌سازی عبارت  $\sum_{i=1}^n e_i^2$  حاصل می‌شود. **اجورث** (۱۸۸۷) پیشنهاد کرد، به جای استفاده از مسئله کمینه‌سازی کمترین توان‌های دوم عبارت  $\sum_{i=1}^n |e_i|$  کمینه شود. یعنی به جای استفاده از توان دوم مانده‌ها، از قدرمطلق آن‌ها استفاده شود. **مارونا و همکاران** (۲۰۰۶) معیار استواری یک برآوردگر را با استفاده از نقطه شکست<sup>۲</sup> تعریف کردند.

<sup>۱</sup>Best linear Unbiased Estimator (BLUE)

<sup>۲</sup>Breakdown point

قدیمی‌ترین تعریف نقطه شکست را هاگ (۱۹۷۴) ارائه داد. البته همپل (۱۹۷۱) رابطه کاربردی‌تری تعریف کرد. دونهو و هوبر (۱۹۸۳) نیز تعریف ساده‌ای ارائه کردند که کاربردی‌ترین تعریف برای نقطه شکست است که در ادامه بیان می‌شود. البته قبل از آن نیاز به ارائه مفهوم نرم مرتبه دوم است.

تعریف ۰.۱ تابع حقیقی  $\|\cdot\|_2$  نرم مرتبه دوم نامیده می‌شود، اگر به ازای هر بردار دلخواه  $\mathbf{a} = (a_1, \dots, a_p)$  رابطه  $\|\mathbf{a}\|_2 = (\sum_{i=1}^p |a_i|^2)^{\frac{1}{2}}$  برقرار باشد.

تعریف ۰.۲ (مارونا و همکاران، ۲۰۰۶) نقطه شکست برآوردگر  $T = T(\mathbf{W})$  برای نمونه  $\mathbf{W}$  به حجم  $n$  به صورت  $\text{BP}(T; \mathbf{W}) = \min_m \{m : \sup_{\mathbf{W}^*} \|T(\mathbf{W}^*)\|_2 = \infty\}$  تعریف می‌شود، که در آن  $\mathbf{W}^*$  نمونه آلوده شده<sup>۱</sup> با جایگزینی  $n$  با  $m \leq n$  نقطه از نمونه  $\mathbf{W}$  با مقادیر دلخواه است.

بنابراین درصد شکست برآوردگر  $T$  برای یک نمونه متناهی مانند  $\mathbf{W}$  برابر با  $\frac{m}{n}$  خواهد بود. برای رگرسیون کمترین توان‌های دوم معمولی، وجود یک داده غیر عادی برای تحت تاثیر قراردادن برآورد ضرایب کافی است. بنابراین نقطه شکست آن به صورت  $\frac{1}{n}$  است. هر چه  $n$  بزرگتر باشد،  $\frac{1}{n}$  به سمت صفر میل می‌کند، به این معنی که درصد نقطه شکست برآوردگر کمترین توان‌های دوم معمولی صفر درصد است. یعنی تنها حضور یک مشاهده دورافتاده می‌تواند تاثیر جدی و مخربی روی برآورد کمترین توان‌های دوم معمولی داشته باشد. روش‌های استوار مختلف دیگری برای غلبه بر آثار بد نقاط دورافتاده توسط محققان پیشنهاد شده است که از جمله آن‌ها می‌توان به روش‌های والد (۱۹۴۰)، نیر و سریواستاوا (۱۹۴۲)، بارتلت (۱۸۸۷)، براون ومود (۱۹۵۱)، توکی (۱۹۷۰)، جرکوا (۱۹۷۱)، جیگل (۱۹۷۲)، بیگل (۱۹۷۳)، اندروز (۱۹۷۴)، کونکر و باست (۱۹۷۸) و ولمن و باست (۱۹۸۱) اشاره نمود. متاسفانه درصد شکست این روش‌ها حداکثر ۳۰ درصد است و همچنین برخی از آن‌ها برای  $p > 2$  (تعداد متغیرهای توضیحی) تعریف نمی‌شوند. سیگل (۱۹۸۲) برآوردگر میانه مکرراً را با ضریب شکست ۵۰ درصد معرفی نمود، ولی ضعف این روش در عدم تشخیص بخش‌های «خوب» و «بد» نمونه بود. روسو (۱۹۸۴) برآوردگر کمترین میانه توان‌های دوم<sup>۲</sup> را که براساس طرحی از همپل (۱۹۷۵) بود، به صورت  $\min_{\beta} \{\text{median}_{1 \leq i \leq n} e_i^{\beta}\}$  معرفی کرد. این روش با نسبت به تبدیلات خطی در متغیرهای توضیحی هم‌وردا<sup>۴</sup> (اندازه تغییرات هماهنگ دو متغیر تصادفی) است. عدم سادگی مینیمم‌سازی تابع هدف

<sup>1</sup>Contaminated

<sup>2</sup>Repeated median estimator

<sup>3</sup>Least median squares

<sup>4</sup>Equivariant

در این مورد منجر به معرفی برآوردگر کمترین توان های دوم پیراسته<sup>۱</sup> توسط روسو و لروی (۱۹۸۷) شد که در ادامه مورد بررسی قرار می گیرد.

در بخش ۲ رگرسیون استوار توان های دوم پیراسته معرفی شده است. مروری بر رگرسیون ستیغی و معرفی آن در بخش ۳ آمده است. در بخش ۴ کاربردی از این روش برای داده های سن دندانی مورد مطالعه قرار گرفته است.

### ۳ رگرسیون استوار کمترین توان های دوم پیراسته

فرض کنید  $(e^x)_{1:n} \leq (e^x)_{2:n} \leq (e^x)_{i:n} \leq (e^x)_{n:n}$  توان دوم مانده های مرتب شده هستند. لازم به ذکر است که ابتدا مانده ها به توان دو می رسند، سپس مرتب می شوند. در روش رگرسیون استوار کمترین توان های دوم پیراسته، مجموع  $h$  کوچکترین توان های دوم مانده های مدل به صورت  $\sum_{i=1}^h (e^x)_{i:n}$  کمیته می شوند.  $h$  در بازه  $[\frac{n}{p}, n]$  تغییر می کند و  $\alpha = \frac{n-h}{h}$  درصد مشاهدات دور افتاده را نشان می دهد. روسو و لروی (۱۹۸۷) نشان داد که با انتخاب  $h = (\frac{n}{p}) + (\frac{p+1}{p})$  نقطه شکست روش کمترین توان های دوم پیراسته به ماکسیمم مقدار ممکن خواهد رسید. اما برای تعیین اینکه مشاهده  $i$ ام، یک مشاهده خوب است یا نه لازم است تابع نشانگر  $z_i$  را به صورت

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

تعریف می شود. در نهایت ماتریس  $\mathbf{Z}$  که یک ماتریس قطری است به صورت

$$\mathbf{Z} = \begin{pmatrix} z_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & z_n \end{pmatrix},$$

به دست می آید. شکل ماتریسی مدل رگرسیون خطی به صورت

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon, \quad E(\varepsilon) = 0, \quad E(\varepsilon^T \varepsilon) = \sigma^2 \mathbf{I}_p,$$

<sup>1</sup>Least trimmed squares

خواهد بود، به طوری که  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$  بردار  $n \times 1$  از خطاهای تصادفی است. بنابراین می توان مسئله کمینه سازی را به صورت

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

بازنویسی کرد، که در آن  $\mathbf{e} = (1, \dots, 1)_{n \times 1}^\top$  و  $\mathbf{z} = (z_1, \dots, z_n)$ . با حل مسئله بهینه سازی (۱)، برآوردگر کمترین توان های دوم پیراسته به صورت

$$\hat{\beta}(\mathbf{Z}) = \mathbf{C}_z^{-1} \mathbf{X}^\top \mathbf{Z} \mathbf{y}, \quad \mathbf{C}_z = \mathbf{X}^\top \mathbf{Z} \mathbf{X},$$

به دست می آید. در صد شکست این روش ۵۰ درصد است. یعنی حتی اگر نیمی از داده ها دورافتاده باشند باز هم این روش در برابر نقاط دورافتاده، استوار است (روزبه، ۲۰۱۶؛ امینی و روزبه، ۲۰۱۶؛ روزبه و آرشی، ۲۰۱۷).

#### ۴ رگرسیون ستیغی استوار

وجود ارتباط خطی بین متغیرهای توضیحی در مدل رگرسیون خطی چندگانه، همخطی نامیده می شود. دو نوع همخطی وجود دارد: همخطی کامل و همخطی ناقص. همخطی کامل، زمانی رخ می دهد که یکی از متغیرهای توضیحی، تابعی دقیق از یک یا چند متغیر توضیحی باشد و همخطی ناقص، زمانی رخ می دهد که یکی از متغیرهای توضیحی، تابعی تقریبی از یک یا چند متغیر توضیحی باشند. در عمل به دلیل وجود خطای اندازه گیری، همخطی کامل رخ نمی دهد. وجود همخطی ممکن است منجر به فاصله های اطمینان پهن برای پارامترها، ناپایداری برآوردها یا تولید برآوردهایی با علامت اشتباه شود. روش های مختلفی برای غلبه برای همخطی وجود دارد که از آن جمله رگرسیون مولفه اصلی و رگرسیون ستیغی اشاره کرد. **هول و کنارد (۱۹۷۰)** برای رفع مشکل همخطی برآوردگر ستیغی را با کمینه سازی تابع هدف به صورت

$$\varphi(\beta) = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta), \quad \beta^\top \beta \leq \phi^2, \quad (2)$$

معرفی کردند، که در آن  $\phi^2$  پارامتر منظم‌سازی<sup>۱</sup> بوده و برای کنترل ضرایب برآورد شده بکار می‌رود. از کمینه‌سازی تابع هدف (۲) برآوردگر ستیغی به صورت

$$\hat{\beta}(k) = \mathbf{C}_k^{-1} \mathbf{X}^T \mathbf{y}, \quad \mathbf{C} = \mathbf{X}^T \mathbf{X}, \quad \mathbf{C}_k = \mathbf{C} + k \mathbf{I}_p,$$

به‌دست می‌آید. این برآوردگر، یک برآوردگر اریب است و به  $k$  پارامتر ستیغی یا انقباضی گفته می‌شود. نکته مهم در روش ستیغی یافتن مقدار  $k$  است چرا که مقدار زیاد  $k$  باعث افزایش اریبی و کاهش واریانس و مقدار کم  $k$  باعث کاهش اریبی و افزایش واریانس می‌شود (امینی و روزبه، ۲۰۱۵) و برای حداقل نمودن تابع مخاطره لازم است  $k$  به گونه‌ای تعیین شود که موازنه بین اریبی و واریانس برقرار کند. روش‌های مختلفی مانند استفاده از معیارهای آکائیک، اطلاع بیزی،  $C_p$  مالوز، اعتبارسنجی متقابل، اعتبارسنجی متقابل تعمیم‌یافته و مانند آن‌ها برای انتخاب مقدار  $k$  وجود دارد. برای اطلاعات بیشتر در مورد برآوردگر رگرسیون ستیغی می‌توان به راسخ و همکاران (۱۳۹۸) و امامی (۱۳۹۷) مراجعه کرد. واسرمن (۲۰۰۶) از بین روش‌های بیان‌شده روش اعتبارسنجی و اعتبارسنجی متقابل را به دلیل یافتن همزمان مقدار بهینه  $k$  و آزمون میزان دقت برگزیدند.

مجموع مجذور خطا با روش اعتبارسنجی<sup>۲</sup> متقابل به صورت  $CV = \|\mathbf{y} - \hat{\mathbf{y}}_{(-i)}\|^2$  محاسبه می‌شود. که در آن  $\hat{\mathbf{y}}_{(-i)} = \mathbf{X}_{(-i)} \hat{\beta}_{(-i)}$ ،  $\hat{\beta}_{(-i)} = (\mathbf{X}_{(-i)}^T \mathbf{X}_{(-i)})^{-1} \mathbf{X}_{(-i)}^T \mathbf{y}_{(-i)}$  و  $\mathbf{X}_{(-i)}$  از حذف  $i$  امین سطر ماتریس  $\mathbf{X}$  و  $\mathbf{y}_{(-i)}$  نیز از حذف  $i$  امین مشاهده بردار  $\mathbf{y}$  بدست می‌آید. برای مشاهده جزئیات بیشتر به واسرمن (۲۰۰۶) مراجعه نمایید.

در روش اعتبارسنجی متقابل در هر مرحله باید یکی از مشاهدات حذف و مدل بدون آن مشاهده برازش شود یعنی در واقع باید  $n$  بار برازش انجام شود تا بتوان مقدار  $k$  بهینه را یافت و این امر باعث دشواری محاسبات و زمان‌گیر شدن آن خواهد شد. به همین دلیل روش اعتبارسنجی متقابل تعمیم‌یافته<sup>۳</sup> که محاسبات آن ساده‌تر از اعتبارسنجی متقابل بوده و به صورت

$$GCV_k = \frac{\frac{1}{n} \|(\mathbf{I}_n - \mathbf{H}_k) \mathbf{y}\|^2}{\left(1 - \frac{1}{n} \text{tr}(\mathbf{H}_k)\right)^2},$$

تعریف می‌شود، پیشنهاد شد (واسرمن، ۲۰۰۶)، که در آن  $\mathbf{H}_k = \mathbf{X}(\mathbf{X}^T \mathbf{X} + k \mathbf{I}_p)^{-1} \mathbf{X}^T$ . در بسیاری

<sup>۱</sup>Regularization parameter

<sup>۲</sup>Cross-validation (CV)

<sup>۳</sup>Generalized cross-validation (GCV)

از مسائل رگرسیونی ممکن است همزمان هم مشکل نقاط دورافتاده و هم مشکل همخطی وجود داشته باشد. ایده‌ای که در این‌گونه موارد به ذهن می‌رسد، استفاده همزمان از روش‌های ستیغی و استوار است، یعنی

$$\min_{\beta, \mathbf{Z}} \varphi(\beta, \mathbf{Z}) = (\mathbf{y} - \mathbf{X}\beta)^\top \mathbf{Z}(\mathbf{y} - \mathbf{X}\beta), \quad \beta^\top \beta \leq \phi^2, \quad \mathbf{e}^\top \mathbf{z} = h, \quad z_i \in \{0, 1\}. \quad (3)$$

با حل مسئله بهینه‌سازی (۳)، برآوردگر ستیغی استوار به صورت  $\hat{\beta}(k, \mathbf{Z}) = \mathbf{C}_{k, \mathbf{Z}}^{-1} \mathbf{X}^\top \mathbf{Z} \mathbf{y}$  می‌آید، که در آن  $\mathbf{C}_{k, \mathbf{Z}} = \mathbf{C}_Z + k \mathbf{I}_p$ . برای محاسبه برآوردگر ستیغی استوار در مسائل عددی، ابتدا برآوردگر کمترین توان‌های دوم پیراسته را با استفاده از رابطه (۱) محاسبه نمود و ماتریس  $\mathbf{Z}$  ساخته می‌شود. سپس، معیار اعتبارسنجی متقابل تعمیم‌یافته را برای برآوردگر ستیغی استوار به صورت

$$\text{GCV}_{k, \mathbf{Z}} = \frac{\frac{1}{n} \|(\mathbf{I}_n - \mathbf{H}_{k, \mathbf{Z}}) \mathbf{y}\|^2}{(1 - \frac{1}{n} \text{tr}(\mathbf{H}_{k, \mathbf{Z}}))^2}, \quad \mathbf{H}_{k, \mathbf{Z}} = \mathbf{X} \mathbf{C}_{k, \mathbf{Z}}^{-1} \mathbf{X}^\top,$$

بازنویسی می‌شود. در انتها با کمینه‌سازی  $\text{GCV}_{k, \mathbf{Z}}$  به کمک دستور optim در نرم‌افزار R مقدار بهینه  $k$  را به روش عددی L-BFGS-B بدست آورده و در رابطه (۳) جایگزین می‌شود.

## ۵ مطالعه شبیه‌سازی

در این بخش در یک مطالعه شبیه‌سازی به بررسی و مقایسه چهار برآوردگر کمترین توان‌های دوم، ستیغی، کمترین توان‌های دوم پیراسته و ستیغی استوار پرداخته می‌شود. در این مطالعه بردار پارامترها به صورت  $\beta = (-3, 4, 5.5, 2, -4, 7)^\top$  در نظر گرفته شده است. به منظور ایجاد همخطی بین متغیرهای توضیحی برای تولید هر یک از درایه‌های ماتریس طرح از رابطه

$$x_{ij} = (1 - \gamma^2)^{\frac{1}{2}} z_{ij} + \gamma z_{ip}, \quad i = 1, \dots, n, \quad j = 1, \dots, p,$$

استفاده شده است. به طوری که در آن اعداد شبه‌تصادفی نرمال استاندارد مستقل بوده و  $\gamma$  طوری تعیین می‌شود که بین هر دو متغیر توضیحی همبستگی ایجاد شود. در این جا  $\gamma = 0.75, 0.87$  فرض می‌شود. از طرفی برای ایجاد نقاط پرت در مجموعه داده‌ها، ۷۵ درصد خطاها از توزیع نرمال استاندارد و ۲۵ درصد

بقیه از توزیع  $t$  - استیودنت غیرمرکزی تولید می‌شود، یعنی

$$\varepsilon = (\varepsilon_1^T, \varepsilon_2^T)^T \quad \varepsilon_1 \sim N(0, \mathbf{I}), \quad \varepsilon_2 \sim t_2(\lambda),$$

بطوریکه  $t_{\nu}(\delta)$  نشان دهنده توزیع  $t$  - استیودنت غیر مرکزی با  $\nu$  درجه آزادی و پارامتر غیرمرکزی  $\delta$  است. در نهایت متغیر پاسخ با  $n = 100$  مشاهده و  $1000$  تکرار از مدل  $\mathbf{y} = \mathbf{X}\beta + \varepsilon$  تولید می‌شود. لازم به ذکر است که محاسبات در سیستمی با ویژگی‌های زیر انجام شده است:

Intel(R) Pentium(R) CPU G3250 @ 3.20 GHz 3.20 GHz, 8 GB RAM,  
System type 64-bit Operating System, x64-based processor.

مقادیر برآورد پارامترها، مقادیر اریبی، انحراف استاندارد و همچنین  $\| \mathbf{y} - \hat{\mathbf{y}} \|^2$  SSE برآوردگرهای پیشنهادی در جدول‌های ۱، ۲، ۳، ۴ و ۵ گزارش شده است. با توجه به این جدول‌ها کاملاً واضح است که روش ستیغی استوار بسیار بهتر از روش‌های دیگر عمل می‌کند، چرا که کمترین میزان SSE متعلق به این روش بوده است. همچنین میزان اریبی و انحراف استاندارد برآوردگر ستیغی استوار در مقایسه با روش‌های دیگر کمتر است.

جدول ۰.۱ ارزیابی برآوردگرهای کمترین توان‌های دوم معمولی و پیراسته برای  $\gamma = 0.75$

روش	معمولی		پیراسته	
	برآورد	اریبی	برآورد	انحراف استاندارد
$\beta_1$	-۲,۶۹۶۶	۰,۳۰۳۳	-۲,۸۸۲۲	۲,۰۷۸۳
$\beta_2$	۴,۱۷۷۶	۰,۱۷۷۶	۳,۹۵۱۴	۱,۴۸۰۸
$\beta_3$	۵,۵۶۵۸	۰,۰۶۵۸	۵,۵۴۷۶	۱,۲۵۸۸
$\beta_4$	۲,۳۷۷۳	۰,۳۷۷۳	۱,۸۶۰۲	۱,۷۶۶۷
$\beta_5$	-۳,۸۴۷۹	۰,۱۵۲۰	-۳,۹۴۲۲	۲,۵۳۴۸
$\beta_6$	۶,۴۰۱۳	-۰,۵۹۸۶	۶,۹۷۹۸	۲,۵۴۹۳

جدول ۰.۲ ارزیابی برآوردگرهای ستیغی و ستیغی استوار برای  $\gamma = 0.75$

روش	ستیغی		ستیغی استوار	
	برآورد	اریبی	برآورد	انحراف استاندارد
$\beta_1$	-۲,۱۸۰۸	۰,۸۱۹۱	-۲,۸۰۶۳	۰,۶۴۹۷
$\beta_2$	۴,۰۰۵۰	۰,۰۰۵۰	۳,۹۴۵۱	۰,۵۷۴۲
$\beta_3$	۵,۲۴۲۰	-۰,۲۵۷۹	۵,۵۲۱۹	۰,۵۷۸۹
$\beta_4$	۲,۴۰۴۹	۰,۴۰۴۹	۱,۸۸۰۹	۰,۷۴۰۸
$\beta_5$	-۳,۱۹۸۴	۰,۸۰۱۵	-۳,۸۵۲۲	۰,۶۳۱۰
$\beta_6$	۵,۹۰۴۲	-۱,۰۹۵۷	۶,۸۸۱۲	۰,۸۶۳۸



جدول ۳. ارزیابی برآوردگرهای کمترین توان‌های دوم معمولی و پیراسته برای  $\gamma = ۰.۸۷$ .

روش	کمترین توان دوم معمولی			کمترین توان دوم پیراسته		
	برآورد	اریبی	انحراف استاندارد	برآورد	اریبی	انحراف استاندارد
$\beta_1$	-۲۸۳۲۸	۰.۱۶۷۱	۲.۱۵۸۷	-۲۹۷۰۴	۰.۲۹۵	۰.۷۷۳۲
$\beta_2$	۳۹۸۱۰	-۰.۱۸۹	۲.۰۶۶۵	۳۹۱۴۹	-۰.۰۸۵۰	۰.۸۴۳۷
$\beta_3$	۵۶۵۳۹	۰.۱۵۳۹	۲.۴۲۲۶	۵۵۶۳۶	۰.۰۶۳۶	۰.۸۱۵۱
$\beta_4$	۱۹۷۳۲	-۰.۰۲۶۷	۳.۴۶۸۸	۱۹۹۴۵	-۰.۰۰۵۴	۰.۸۱۶۷
$\beta_5$	-۳۸۸۷۶	۰.۱۱۲۳	۲.۰۶۸۵	-۴۱۳۶۲	-۰.۱۳۶۲	۰.۸۵۶۸
$\beta_6$	۷۰۳۷۳	۰.۰۳۷۳	۳.۷۴۷۷	۷۰۷۶۶	۰.۰۷۶۶	۱.۱۹۳۷

جدول ۴. ارزیابی برآوردگرهای ستیغی و ستیغی استوار برای  $\gamma = ۰.۸۷$ .

روش	ستیغی		ستیغی استوار	
	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
$\beta_1$	-۱۸۵۵۹	۱.۱۴۴۰	-۲۸۳۸۲	۱.۷۰۷۱
$\beta_2$	۳۹۰۰۰	-۰.۰۹۹۹	۳۹۱۶۹	۱.۵۹۱۹
$\beta_3$	۵۲۷۸۱	-۰.۲۲۱۸	۵۵۳۳۵	۱.۹۲۰۳
$\beta_4$	۲۱۶۸۰	۰.۱۶۸۰	۲۰۳۲۴	۲.۴۸۴۸
$\beta_5$	-۲۶۹۱۰	۱.۳۰۸۹	-۳۹۷۶۵	۱.۸۰۶۱
$\beta_6$	۵۶۷۷۱	-۱.۳۲۲۸	۶۸۶۷۶	۲.۳۰۵۴

جدول ۵. خطا و سرعت محاسبه برآوردگرهای پیشنهادی.

روش برآورد کمترین توان‌های دوم	$\gamma = ۰.۸۷$		$\gamma = ۰.۸۵$	
	SSE	سرعت محاسبات	SSE	سرعت محاسبات
کمترین توان‌های دوم پیراسته	۴۴۷۷۲۷۴۰	۰۹ : ۰۰ : ۰۰	۳۶۷۵۵۸۰۲	۰۶ : ۰۰ : ۰۰
ستیغی	۴۷۸۹۳۰۴	۰۷ : ۰۶ : ۰۰	۲۸۵۷۴۸۸	۵۴ : ۰۵ : ۰۰
ستیغی استوار	۲۸۵۰۸۷۰۲	۰۱ : ۰۱ : ۰۰	۲۶۶۴۵۹۵۱	۵۱ : ۰۱ : ۰۰
	۴۶۷۲۹۵۵	۰۶ : ۰۶ : ۰۰	۲۸۳۲۸۰۰	۴۶ : ۰۶ : ۰۰

## ۶ مطالعه کاربردی: تحلیل داده‌های سن دندانی

برآورد سن، نقش مهمی در پزشکی قانونی، بیماری‌های غدد، دندان پزشکی، باستان‌شناسی و بسیاری علوم دیگر دارد. در کشورهای در حال توسعه که ثبت دقیق زمان تولد افراد وجود ندارد، این روش راه بسیار مناسبی برای تخمین سن تقویمی افراد است. اغلب بین بلوغ فیزیکی و سن فرد ارتباط وجود دارد و برخی از نشانه‌های بلوغ فیزیکی مانند سن اسکلتی، قاعدگی، اندازه و قد افراد و کلسیفیکاسیون (جمع شدن کلسیم در یکی از بافت‌های بدن) می‌تواند در نبود اطلاعاتی نظیر سن و جنسیت دقیق افراد، برای تعیین سن افراد مورد استفاده قرار گیرد. سن دندانی می‌تواند بر اساس رویش یا مراحل شکل‌گیری دندان در رادیوگرافی ارزیابی شود. در این بین مدل‌سازی برای برآورد سن تقویمی با استفاده از روش‌های نوین و غنی آماری دارای اهمیت است. در بین روش‌هایی که برای تعیین سن استفاده می‌شود، رایج‌ترین روش

مورد استفاده در سراسر دنیا، روش دمیجیان است. این روش بر اساس کلسیفیکاسیون دندان دائمی در رادیوگرافی پانورامیک است. در این نوع رادیوگرافی یک نوع اسکن پانورامیک فک بالا و پایین به وسیله اشعه ایکس دندان پزشکی است که دیدی دو بعدی از نیمکره فرضی تشکیل شده در حد فاصل بین گوش‌ها می‌دهد، است. مشکلی که در بررسی سن دندانی با روش دمیجیان وجود دارد زمانی است که تقریباً تمام ریشه‌های دندان‌های دائمی بیمار تکامل یافته و ریشه‌های دندان‌های عقل بیمار باید به عنوان ملاک مورد بررسی قرار گیرند. دمیجیان و همکارانش روش جدیدی ابداع کردند (دمیجیان و همکاران، ۱۹۷۶). پس از رویش دندان مولر دوم دائمی در حدود سنین ۱۲ تا ۱۳ سالگی تعیین سن با مشکل مواجه می‌شود. این مسئله، دندان مولر سوم را در این سنین در مرکز توجه قرار می‌دهد چرا که هنوز در حال شکل‌گیری است. این دندان در سنین سالگی در دهان ظاهر می‌شود و با توجه به فرآیند تشکیل ریشه می‌توان مراحل تکاملی آن ۱۷ تا ۲۱ را تعریف کرد که تقریباً از ۱۴ سالگی آغاز می‌شود. در کنار استفاده از روش دمیجیان در تعیین سن با کمک دندان مولر سوم، روش‌های دیگری نیز وجود دارد که برپایه تحقیقات صورت گرفته روش دمیجیان بر سایر روش‌ها عملکرد بهتری داشته است. اردکانی و همکاران (۲۰۰۷) مطالعه‌ای با هدف مطالعه مقطعی از نوع تشخیصی با همکاری سازمان پزشکی قانونی و دانشکده دندانپزشکی شهید صدوقی یزد در سال ۸۴ - ۱۳۸۳ درباره ۵۸ بیمار با محدوده سنی ۱۵ تا ۲۵ سال که به منظور جراحی دندان عقل به مطب خصوصی مراجعه کرده بودند، صورت گرفت. محاسبات آماری از طریق آزمون‌های تی زوجی، ضریب همبستگی، ویلکاکسون و کلموگروف اسمیرینف انجام پذیرفت. کمترین خطاهای به دست آمده در تخمین سن با سن تقویمی بیمار، مربوط به دندان عقل کشیده شده سمت راست فک پایین بود.

در این مقاله، بیماران مراجعه کننده در طی دوره ۱۲ ماهه سال ۹۴ - ۱۳۹۳ به درمانگاه خاتم‌الانبیا استان یزد، جامعه مورد بررسی تحقیق را تشکیل داده است. تعداد ۸۷ نفر از بیماران مراجعه کننده، که به صورت روش تصادفی ساده انتخاب شده اند. معیارهای ورود و خروج نمونه‌ها در مطالعه عبارتند از عدم وجود مشکلات رشدی- تکاملی، کیفیت قابل قبول رادیوگرافی پانورامیک، وجود حداقل یک دندان عقل و عدم وجود شکستگی و یا ضایعاتی پیرامون دندان عقل. تعداد ۸۷ رادیوگرافی پانورامیک دیجیتال مربوط به بیماران نوجوان ارتودنسی مراجعه‌کننده به کلینیک خاتم‌الانبیا از آرشیو کلینیک بررسی شد. جنسیت و سن کرونولوژیک دقیق افراد (دارای روز و ماه تولد علاوه برسال) از روی پرونده بیماران ثبت شد، در صورت عدم وجود در پرونده از طریق تماس با بیمار به طور مستقیم پرسیده شد. رادیوگرافی‌ها دارای دندان‌های عقل بوده و بین سنین ۱۴ تا ۱۹ قرار داشت.

در جدول‌های ۶ و ۷ به ترتیب درصد فراوانی بیماران به تفکیک جنسیت و سن تقویمی و برآورد شده

جدول ۰۶. فراوانی داده‌ها به تفکیک جنسیت.

جنسیت	زن	مرد
فراوانی نسبی	۶۷/۹	۳۲/۱

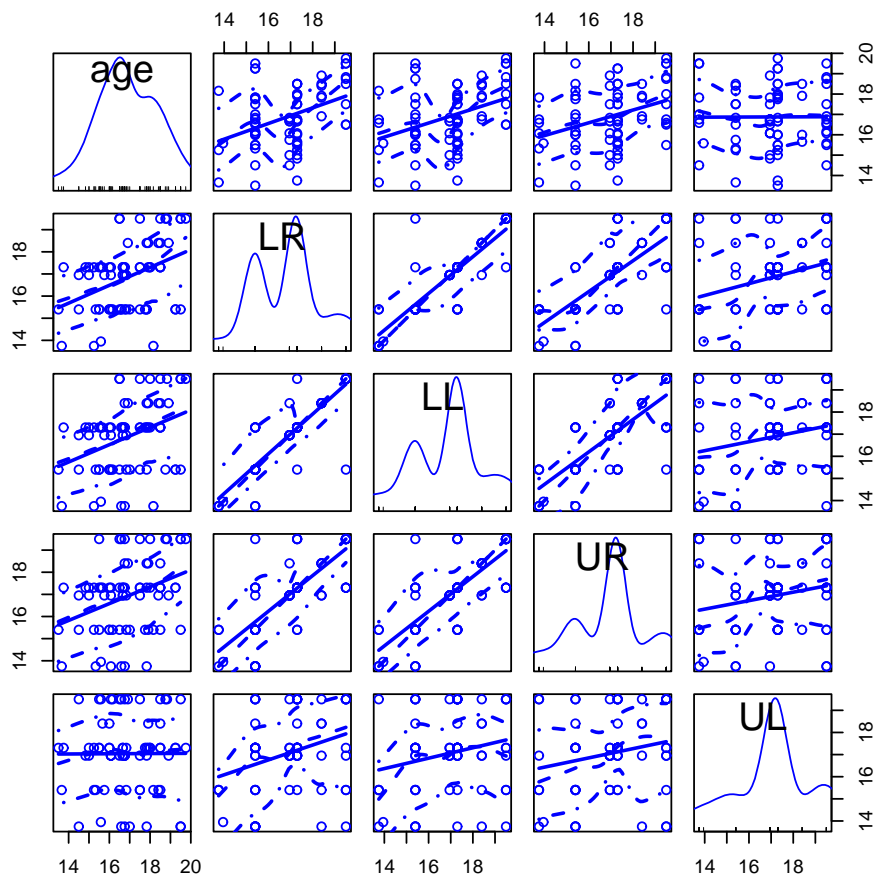
به نمایش گذاشته شده است. مدل رگرسیون خطی چندگانه به صورت

جدول ۰۷. آماره‌های توصیفی به تفکیک سن تقویمی و آزمون تخمین دمیرجیان.

-	مینیمم	ماکسیمم	میانگین	انحراف معیار
سن تقویمی	۱۳/۵	۲۰/۵۸	۱۷/۲۱	۱/۱۶۶
آزمون دمیرجیان	۱۳/۷۵	۱۹/۵	۱۶/۸۸	۱/۰۸

$$(age)_i = \beta_1(LR)_i + \beta_2(LL)_i + \beta_3(UR)_i + \beta_4(UL)_i + \varepsilon_i \quad i = 1, \dots, 70, \quad (4)$$

نوشته می‌شود، که در آن  $age$  متغیر پاسخ سن تقویمی،  $LR$  سن برآوردشده دمیرجیان دندان راست مولر فک پایین،  $LL$  سن برآوردشده دمیرجیان دندان چپ مولر فک پایین،  $UR$  سن برآوردشده دمیرجیان دندان راست مولر فک بالا و  $UL$  سن برآوردشده دمیرجیان دندان چپ مولر فک بالا هستند. همان‌طور که در شکل ۱ ملاحظه می‌شود، متغیرهای توضیحی دارای رابطه خطی معنی‌داری با متغیر وابسته نیستند. با انجام تبدیلات باکس - کاکس روی داده‌ها، تا حد ممکن رابطه خطی معنی‌داری بین متغیرهای توضیحی و وابسته ایجاد می‌شود، بطوریکه دست‌یابی به یک مدل رگرسیون قوی‌تر را امکان‌پذیر می‌کند. تبدیل باکس کاکس روی داده‌های اکیدا مثبت انجام می‌شود. در واقع این تبدیل توزیع متغیرها را به نرمال چند متغیره نزدیک می‌کند (کاکس و باکس، ۱۹۶۴). با انجام این تبدیل که به صورت نمودارهای باکس - کاکس در شکل ۲ به نمایش گذاشته شده‌است، وجود رابطه خطی معنی‌داری بین متغیرهای توضیحی و وابسته نسبت به حالت قبل مشهود است. برای انجام این تبدیلات از نرم‌افزار مینی‌ت‌ب استفاده شد. برای بررسی نحوه ارتباط خطی بین متغیرهای توضیحی از نمودار ۳ با استفاده از بسته "PerformanceAnalytics" در نرم‌افزار R 3.4.4 رسم شده است. همان‌طور که ملاحظه می‌شود، بیشترین مقدار همبستگی ۰/۸۷ است که مربوط به دو متغیر توضیحی  $LL_{new}$  و  $LR_{new}$  برای بررسی وجود همخطی در مدل از آماره کاپا

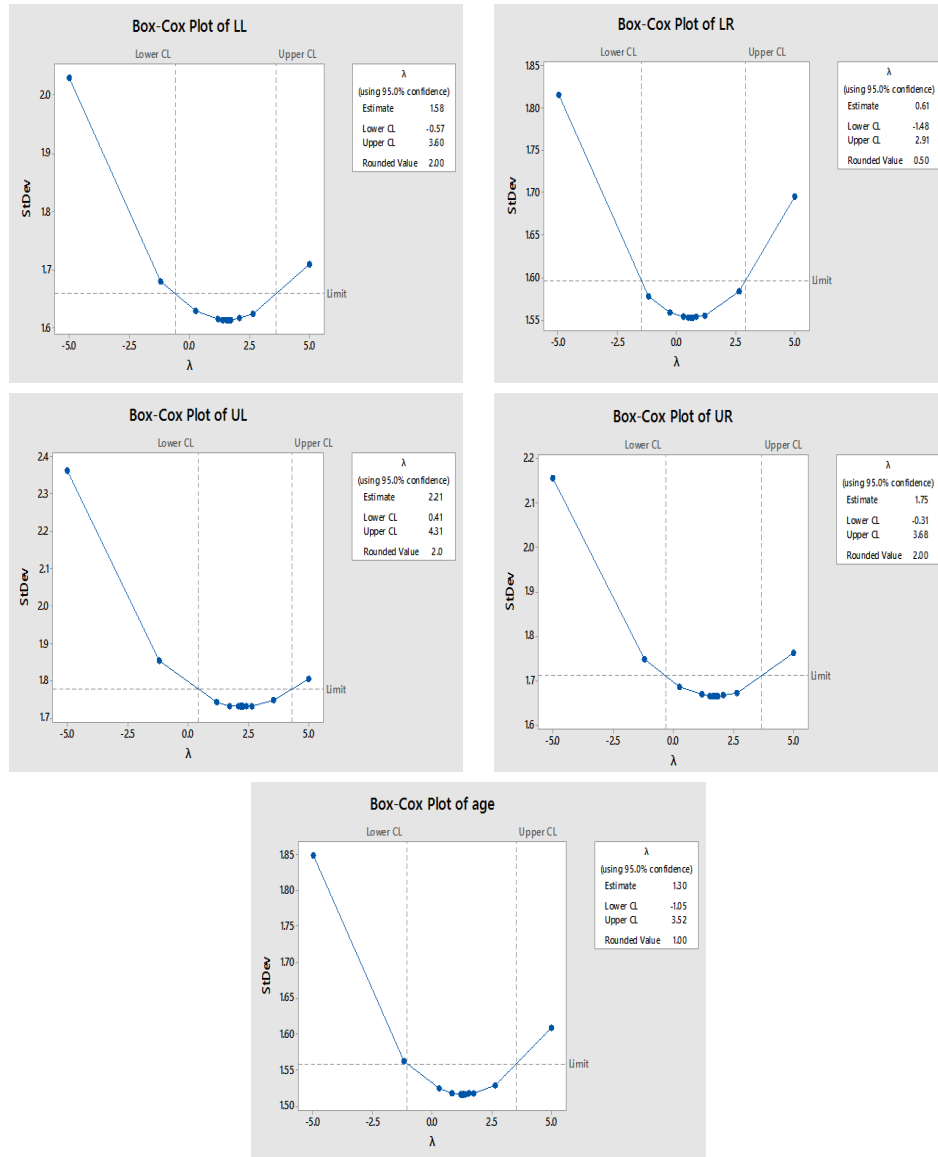


شکل ۰۱. نمودار پراکنش متغیرهای پاسخ و توضیحی.

به صورت

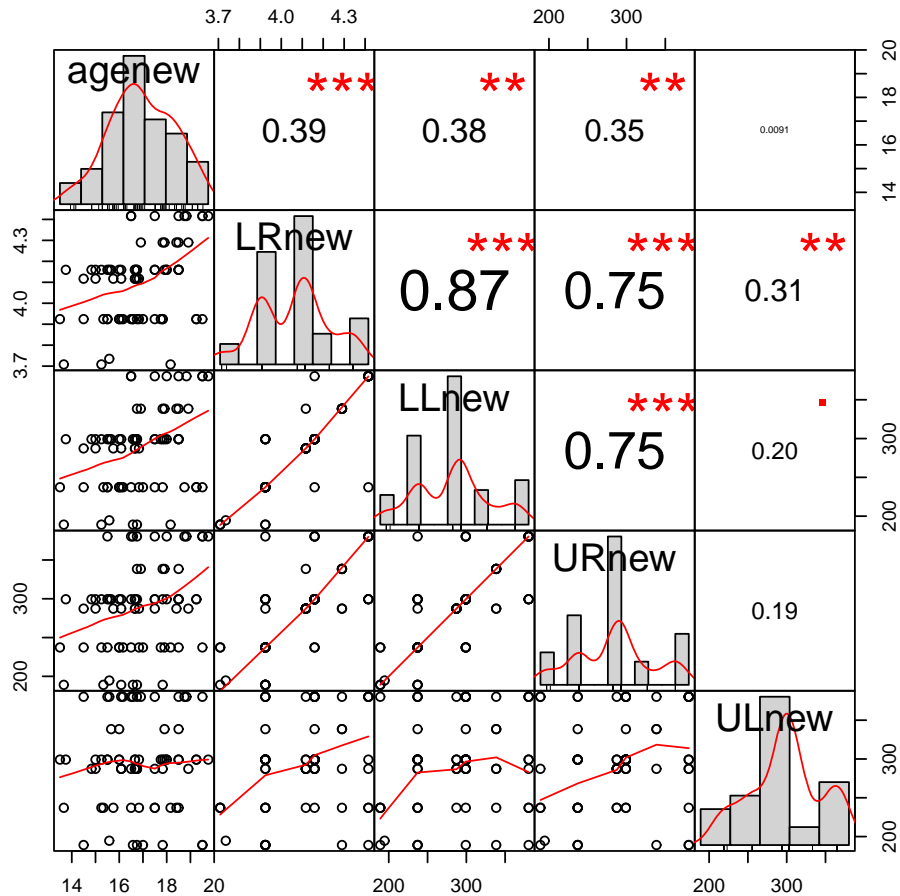
$$\text{kappa} = \sqrt{\frac{\max \lambda_j}{\min \lambda_j}} = ۱۴۶/۸۲۱۵, \quad j = ۱, \dots, p,$$

استفاده می‌شود، که در آن  $\lambda_i$  ها مقادیر ویژه ماتریس  $\mathbf{X}^T \mathbf{X}$  در داده‌های تبدیل یافته است. چون مقدار آماره کاپا بیشتر از ۱۰۰ است، شکل همخطی جدی وجود دارد.



شکل ۲. تبدیلات باکس-کاکس بر روی متغیرهای توضیحی.

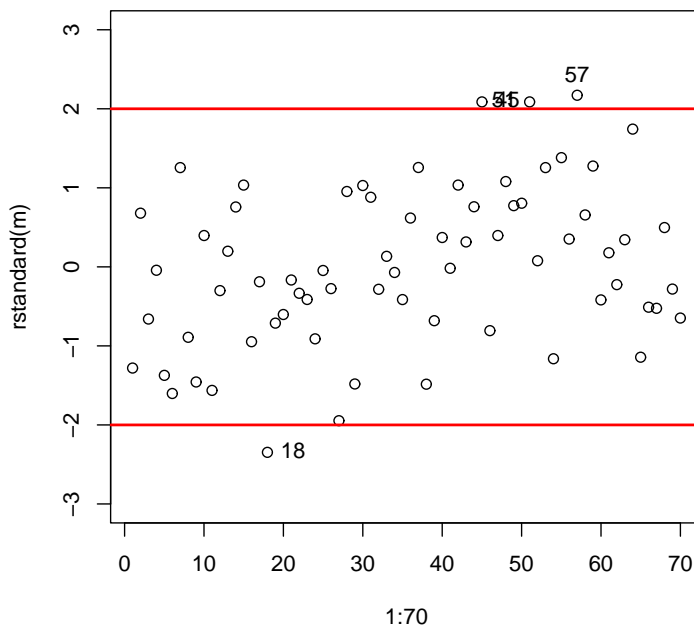
در ادامه، وجود نقاط دورافتاده در مدل مورد بررسی قرار می‌گیرد. اگر  $r_i$  مانده استاندارد شده مشاهده  $i$ ام را نشان دهد، به‌ازای  $|r_i| > 2$  باشد، مشاهده دورافتاده است. (شیدر، ۲۰۰۹). برای تشخیص



شکل ۳. نمودار همبستگی متغیرهای تبدیل‌یافته.

نقاط دورافتاده در مدل از شکل ۴ استفاده می‌شود. همان‌طور که لحظه می‌شود چهار مشاهده ۱۸ام، ۴۵ام، ۵۱ام و ۵۷ام دورافتاده‌اند. بنابراین به‌طور همزمان مشکل همخطی و نقاط دورافتاده در مدل رگرسیونی (۲) مشاهده می‌شود. در ادامه برای مقابله با این مشکلات از رگرسیون ستیغی استوار استفاده می‌کنیم. در این‌جا پارامترها به چهار روش کمترین توان‌های دوم، روش ستیغی، روش کمترین توان‌های دوم پیراسته، روش ستیغی کمترین توان‌های دوم پیراسته برآورد شده است و در نهایت این چهار روش را برای مدل نهایی

$$(agenew)_i = \beta_1(LRnew)_i + \beta_2(LLnew)_i + \beta_3(URnew)_i + \beta_4(ULnew)_i + \varepsilon_i$$



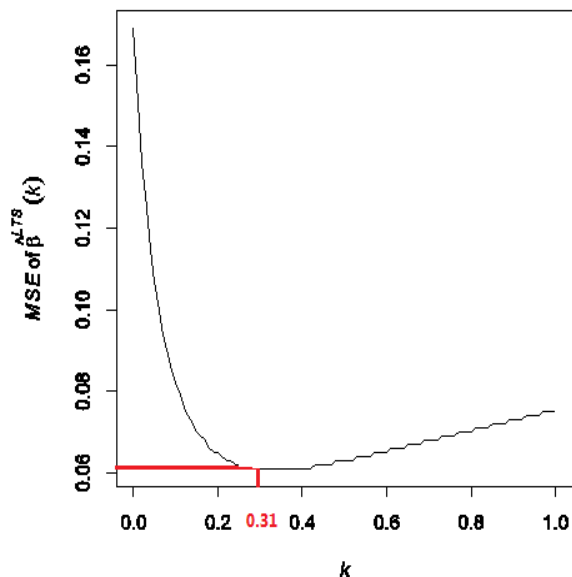
شکل ۴. نمودار مانده‌های استاندارد شده.

$$i = 1, \dots, 70, \quad (5)$$

با هم مقایسه می‌شوند، که در آن منظور از  $LLnew$  و  $URnew$ ،  $ULnew$ ،  $LRnew$ ،  $agenew$  متغیرهای تبدیل یافته بر پایه تبدیل باکس-کاکس شکل ۲ است. برای تعیین بهترین  $k$  از معیار اعتبارسنجی متقابل تعمیم یافته استفاده شده است.

جدول ۸. برآوردهای پیشنهادی به همراه انحراف استاندارد آن‌ها.

ستیفی استوار		ستیفی		کمترین توان دوم پیراسته		کمترین توان دوم		پارامتر
SD	برآورد	SD	برآورد	SD	برآورد	SD	برآورد	
۰/۰۵۹	-۰/۲۵۳۹	۰/۰۴۶۴	-۰/۱۹۱۵	۰/۰۹۷۴	۰/۳۸۴۵	۰/۰۹۳	۰/۲۸۴۸	$\hat{\beta}_1$
۰/۰۰۴	۰/۰۹۶۵	۰/۰۰۳	۰/۰۹۳۹	۰/۰۰۷	۰/۱۷۴۰	۰/۰۰۶	۰/۰۷۲۰	$\hat{\beta}_2$
۰/۰۰۱۹	-۰/۱۲۰۳	۰/۰۰۲	-۰/۰۲۰۰	۰/۰۰۳	-۰/۰۲۹۸	۰/۰۰۳	۰/۱۰۴۶	$\hat{\beta}_3$
۰/۰۰۰۲	-۰/۰۴۷۱	۰/۰۰۰۲	-۰/۰۲۸۸	۰/۰۰۰۴	-۰/۰۸۶۱	۰/۰۰۰۴	۰/۱۱۳۴	$\hat{\beta}_4$



شکل ۵. نمودار مقادیر  $k$  در مقابل MSE برای مدل رگرسیونی ستیغی استوار.

جدول ۰۹. MSE و سرعت محاسبه چهار برآوردگر داده‌های سن تقویمی.

روش برآورد	MSE	سرعت محاسبات
کمترین توان‌های دوم	۱/۱۷۰۹۱	۰۱ : ۰۰ : ۰۰
کمترین توان‌های دوم پبراسته	۰/۱۶۷۳۱	۰۶ : ۰۰ : ۰۰
ستیغی	۰/۰۸۱۰۲	۰۶ : ۰۰ : ۰۰
ستیغی استوار	۰/۰۶۳۰۰	۰۹ : ۰۰ : ۰۰

مقدار بهینه برای مدل رگرسیونی ستیغی استوار برابر ۰/۳۱ است که این موضوع در شکل ۵ به خوبی نشان داده شده است. همان‌طور که در جدول ۹ دیده می‌شود، کمترین میزان MSE مربوط به روش ستیغی استوار است، که نشان‌دهنده عملکرد بهتر این روش برای این مجموعه داده‌ها است.

## بحث و نتیجه‌گیری

با توجه به کاربرد گسترده مدل‌های رگرسیونی در علوم پزشکی، بررسی دقیق و عمیق داده‌ها و متغیرها از اهمیت بسیار برخوردار است. زیرا عدم بررسی دقیق داده‌ها منجر به تحلیل و مدل‌بندی نادقیق رگرسیونی



می‌شود. وجود داده‌های دور افتاده، وجود همبستگی بین متغیرهای توضیحی و همچنین فرضیات زیربنایی مدل رگرسیونی از مواردی است که همواره باید قبل از انجام رگرسیون باید به آن توجه داشت. در این مقاله وجود همزمان داده دورافتاده و همخطی بین متغیرهای توضیحی مورد بررسی قرار گرفت و روش رگرسیون ستیگی کمترین توان‌های دوم پیراسته برای مواجهه با آن پیشنهاد شد. به عنوان کاربردی از این روش، مطالعه‌ای با هدف برآورد سن براساس دندان‌های مولر مورد بررسی قرار گرفت. نتایج به‌دست آمده نشان می‌دهد که روش پیشنهادی در مقایسه با سایر روش‌ها MSE کمتری دارد. همچنین روش‌های ذکر شده در مقاله در داده‌های شبیه‌سازی شده نیز به کار برده شد. نتایج حاصل نشان داد که روش ستیگی استوار در شرایط وجود همخطی و نقاط پرت به طور همزمان بهتر از سایر برآوردها عمل می‌کند.

## تقدیر و تشکر

نویسندگان مقاله ضمن تشکر از اعضای مجترم 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

امامی، ه. (۱۳۹۷). روش‌های تشخیصی در مدل‌های خطی ریبج نیمه‌پارامتری با خطا در اندازه‌گیری، مجله مدل‌سازی پیشرفته ریاضی، ۸، ۱۹-۴۸.

راسخ، ع.، منصوری، ب. و نرگس هدایت‌پور، ن. (۱۳۹۸)، شناسایی مشاهدات پرت در مدل رگرسیونی ریبج تحت محدودیت‌های خطی تصادفی، مجله علوم آماری ایران، ۱۳، ۱۱۷-۱۳۷.

Amini, M. and Roozbeh, M. (2015), Optimal Partial Ridge Estimation in Restricted Semiparametric Regression Models, *Journal of Multivariate Analysis*, **136**, 26–40.

Amini, M. and Roozbeh, M. (2016), Least Trimmed Squares Ridge Estimation in Partially Linear Regression Model, *Journal of Statistical Computation and Simulation*, **86**, 2780–276.

Andrews, D. F. (1974), A Robust Method for Multiple Linear Regression, *Technometrics*, **16**, 523-531.

Bartlett, M. S. (1949), Fitting a Straight Line When Both Variables are Subject to Error, *Biometrics*, **5**, 207-212.

Bickel, P. J. (1973), On Some Analogues to Linear Combination of Order Statistics in Linear Model, *The Annals of Statistics*, **1**, 957-616.

Box, G.E., and Cox, D.R. (1964), An Analysis of Transformed Data, *Journal of Royal Statistical Society, Ser. B*, **39**, 211-252.

Brown, G. W. and Mood, A. M. (1951), *On Median Tests for Linear Hypotheses*, Proceedings of the 2nd Berkeley Symposium on Mathematical Statistics and Probability, edited by J. Neyman, University of California Press, Berkeley and Los Angeles, 159-166.

Demirjian, A. and Goldstein, H. (1976), New Systems for Dental Maturity Based on Seven and Four Teeth. *Annals Of Human Biology*, **3(5)**, 411-21.

Donho, D. L. and Huber, P. J. (1983), *The Notion of Breakdown Point*, in *A Festschrift for Erich Lehmann*, edited by P. Bickel, K. Doksum, and J. L. Hodges, Jr., Wadsworth, Belmont, CA.

Edgeworth, F. Y. (1887), On Observation Relating to Several Quantities, *Hermathena*, **6**, 279-28.

Ezoddini Ardakani, F., Navab, A., A., Bashardoost, N., Mansoorian, H., Ahmadih, M. H. and Sadat Correlation between Chronological, Skeletal, and Dental Age on Panoramic Radiography in Patients Referred to Yazd Dental Clinics on 2004-05. *Journal Dental School Shahid Beheshti University of Medical Sciences.*, **24(4)**, 474-484.

- Hampel, F. R. (1971), A General Qualitative Definition of Robustness, *Annals of Mathematical Statistics*, **42**, 1887-1896.
- Hampel, F. R. (1975), Beyond Location Parameters: Robust Concepts and Methods, *Bulletin International Statistical institute*, **46**, 375-382.
- Hodges, J. L. (1974), *Efficiency in normal Samples and Tolerance of Extreme Values for Some Estimates of Location*, Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, 163-168.
- Hoerl A. E. and Kennard R. W. (1970), Ridge Regression: Biased Estimation for Non-orthogonal Problems, *Thechnometrics*, **12**, 55-67.
- Jaeckel, L.A. (1972), Estimating Regression Coefficients by Minimizing the Dispersion of Residuals, *Annals of Mathematical Statistics*, **5**, 1449-1458.
- Jureckova, J. (1971), Nonparametric Estimate of Regression Coefficients, *Annals of Mathematical Statistics*, **42**, 1328-1338.
- Koenker, R. and Bassett, G.J. (1978), Regression Quantiles, *Econornetrica*, **46**, 33-50.
- Maronna R. A., Martin R. A. and Yohai V. J. (2006), *Robust Statistics: Theory and Methods*, John Wiley and Sons, New York.
- Nair, K. R. and Shrivastava, M. P. (1942), On a Simple Method of Curve Fitting, *Sankhya*, **6**, 121-132.
- Rousseeuw, P.J (1984), Least Median of Squares Regression, *Journal of the American Statistical Association*, **79**, 871-880.
- Rousseeuw, P.J. and Leroy, A.M. (1987), *Robust Regression and Outlier Detection*, John Wiley, New York.

Roozbeh M. (2016), Robust Ridge Estimator in Restricted Semiparametric Regression Models, *Journal of Multivariate Analysis* **147**, 127-144.

Roozbeh M. and Arashi, M. (2017), Least-Trimmed Squares: Asymptotic Normality of Robust Estimator in Semiparametric Regression Models, *Journal of Statistical Computation and Simulation* **147**, 1147-1130.

Siegel, A.F. (1982), Robust Regression Using Repeated Medians, *Biometrika*, **69**, 242-244.

Sheather, S. J. (2009), *A Modern Approach to Regression with R*, Springer, College Station, Texas.

Tukey, J. W. (1970), *Exploratory Data Analysis (Limited Preliminary Edition)*, Addison-Wesley, Reading, MA.

Velleman, P. F. and Hoaglin, D. C. (1981), *Applications, Basics, and Computing of Exploratory Data Analysis*, Duxbury-Press, Boston.

Wald, A. (1940), The Fitting of Straight Lines if Both Variables are Subject to Error, *Annals of Mathematical Statistics*, **11**, 284-300.

Wasserman, L. (2006), *All of Nonparametric Statistics*, Springer-Verlag, New York.

Journal of Statistical Sciences, Autumn and Winter, 2020  
Vol. 14, No. 2, pp 409-428  
DOI: 10.29252/jss.13.2.409

## **Modeling of Chronological Age Using Least Trimmed Squares Ridge Regression**

Roozbeh, M., Manavi, M.  
Department of statistics, Semnan University, Semnan, Iran.

**Abstract:** The popular method to estimation the parameters of a linear regression model is the ordinary least square method which, despite the simplicity of calculating and providing the BLUE estimator of parameters, in some situations leads to misleading solutions. For example, we can mention the problems of multi-collinearity and outliers in the data set. The least trimmed squares method which is one of the most popular of robust regression methods decreases the influence of outliers as much as possible. The main goal of this paper is to provide a robust ridge estimation in order to model dental age data. Among the methods used to determine age, the most popular method throughout the world is the modern modified Demirjian method that is based on the calcification of the permanent tooth in panoramic radiography. It has been shown that using the robust ridge estimator is leading to reduce the mean squared error in comparison with the OLS method. Also, the proposed estimators were evaluated in simulated data sets.

**Keywords:** Demirjian dental age, Least Trimmed Squares, Multicollinearity, Outliers, Ridge estimator.

**Mathematics Subject Classification (2010):** 62J07.