

## کاربرد عملگرهای وزنی در مدل رگرسیون قدرمطلق انحرافات مرتب شده

جلال چاچی<sup>۱</sup>، علیرضا چاچی<sup>۲</sup>

<sup>۱</sup> گروه آمار، دانشکده علوم ریاضی و کامپیوتر، دانشگاه شهید چمران اهواز

<sup>۲</sup> گروه برق، دانشکده فنی مهندسی، دانشگاه صنعتی شهدای هویزه

تاریخ دریافت: ۱۳۹۸/۱۱/۰۹ تاریخ پذیرش و انتشار: ۱۳۹۹/۱۰/۱۹

### چکیده:

در این مقاله رویکرد جدیدی در برآورد پارامترهای مدل رگرسیون خطی کمترین قدرمطلق انحرافات معرفی می‌شود که مبتنی بر مسائل بهینه‌سازی بر مبنای الحاق وزنی قدرمطلق انحرافات مرتب شده است. الحاق وزنی قدرمطلق انحرافات برازش مرتب شده در مسئله بهینه‌سازی در حالی که توابع نیکویی برازش مختلفی را بطور همزمان در مسئله مدل‌سازی در نظر می‌گیرد، توانایی تحلیل داده‌ها به منظور شناسایی نقاط دورافتاده را نیز فراهم می‌کند. بر این اساس این رویکرد تحت تاثیر مشاهدات دورافتاده قرار نمی‌گیرد و در هر مسئله متناسب با تعداد مشاهداتی که پتانسیل دورافتاده بودن را دارا هستند، به انتخاب بهترین برآوردر مدل با بهینه‌ترین مقدار نقطه شکست در بین مجموعه‌ای از برآوردرهای کاندید دیگر می‌پردازد. نیکویی برازش رویکرد پیشنهادی در مدل‌سازی داده‌های شبیه‌سازی شده و داده‌های واقعی در مهندسی آب با حضور مشاهدات دورافتاده تحلیل شده است. همچنین در انتها به تحلیل حساسیت برآوردرها شامل بررسی معیارهای ناریبی و کارایی برآوردرها پرداخته شده است.

واژه‌های کلیدی: رگرسیون وزنی، قدرمطلق انحرافات مرتب شده، نقطه شکست، توان‌های دوم خطا.

## ۱ مقدمه

تحلیل و برازش مدل‌های رگرسیونی (خطی و غیرخطی، استوار و غیر استوار) یکی از پرکاربردترین و متداول‌ترین موضوعات در مدل‌سازی داده‌ها است یکی از شیوه‌های عمده در تحلیل اینگونه مدل‌ها رسم نمودار پراکنش مقادیر متغیرهای تبیینی در مقابل مقادیر متغیر وابسته است که میزان اهمیت متغیرهای تبیینی در تبیین متغیر وابسته را آشکار می‌سازد (آکینسون و ریانی، ۲۰۰۰؛ هاوکینس و همکاران، ۱۹۸۴؛ ویس و لویس، ۱۹۹۴). بر این اساس پس از شناسایی متغیرها و میزان تاثیرگذاری آن‌ها بر یکدیگر، به مدل‌سازی رابطه بین متغیرها با بیشترین دقت و کارایی پرداخته می‌شود (سالیانی و همکاران، ۲۰۱۶). اغلب روش‌های برازش پرکاربرد متداول مانند روش مجموع کمترین توان‌های دوم خطا و روش مجموع کمترین قدرمطلق انحرافات به شدت تحت تاثیر حتی یک مشاهده دورافتاده (مخصوصاً نقاط اهرمی) هستند. اینگونه رویکردها قادر به شناسایی مشاهدات دورافتاده نیستند (بکر و قادر، ۱۹۹۹؛ هاوکینس، ۱۹۸۰؛ ویس و لویس، ۱۹۹۴) و در حضور مشاهدات دورافتاده نتایج غیر قابل اعتماد و دور از واقعیت دارند (درویلیس و همکاران، ۲۰۱۵). در این خصوص، روش‌های متنوعی در تحلیل مدل‌های رگرسیون استوار در حضور داده دورافتاده ارائه شده است (هوبرت و همکاران، ۲۰۰۸؛ گوین و ولش، ۲۰۱۰). این رویکردها عمدتاً بر مبنای استفاده از توابع هدفی هستند که تاثیرات نامناسب مشاهدات دورافتاده را در برآورد مدل کم اثر می‌سازند (بیلور و همکاران، ۲۰۰۶؛ روسو و لروی، ۱۹۸۷). به عنوان نمونه تابع هدف کمترین قدرمطلق انحرافات نسبت به مشاهدات دورافتاده عمودی استوار است ولی در مقابل نقاط دورافتاده اهرمی عملکرد ضعیفی دارد. در این خصوص رگرسیون‌های مبتنی بر  $M$ -برآوردها مشاهدات دورافتاده اهرمی را کم اثر می‌کنند (هوبر و رونچتی، ۲۰۰۹؛ ماروبینی و اورنتی، ۲۰۱۴؛ روسو و وانزومن، ۱۹۹۰). لذا بهتر است در مسئله بهینه‌سازی توابع هدف مناسبی در نظر گرفته شوند که متناسب با ماهیت و نوع مشاهدات دورافتاده عملکرد بهینه‌ای داشته باشند (ماروبینی و اورنتی، ۲۰۱۴). به این منظور رویکردهای وزنی در برآورد مدل‌های رگرسیون کمترین قدرمطلق انحرافات معرفی شد. محمودی (۱۳۸۴) به رگرسیون حداقل قدرمطلق انحرافات وزنی استوار و لینگ (۲۰۰۵) به بررسی اینگونه برآوردها در مدل‌های اتو-رگرسیو پرداخت. گائو و فنگ (۲۰۱۸) رگرسیون کمترین قدرمطلق انحرافات وزنی تاوانیده را با هدف شناسایی داده‌های دورافتاده و معرفی برآوردهای استوار معرفی کرد. لازم به ذکر است که در رویکردهای معرفی شده در برآورد مدل‌های رگرسیونی کمترین قدرمطلق انحرافات وزنی، پارامترهای مدل بر اساس بهینه‌سازی تابع هدف  $\sum_{i=1}^n w_i |e_i|$  برآورد می‌شوند. در این مقاله یک مدل رگرسیونی طبق شرایط زیر برآورد می‌شود:

۱- مسئله بهینه‌سازی با تابع هدفی بر مبنای الحاق وزنی قدرمطلق انحرافات مرتب شده  $\sum_{i=1}^n w_i |e_{(i)}|$

- در برازش مدل در نظر گرفته شود، که در آن  $e(i)$ ،  $i$ -امین بزرگترین مقدار خطا است.
- ۲- بهترین تابع هدف در برازش به داده‌ها در تناسب با تعداد و نوع مشاهدات دورافتاده به گونه‌ای تعیین شود که در برازش به مدل دارای بهینه‌ترین و بالاترین مقدار نقطه شکست<sup>۱</sup> باشد. نقطه شکست یک برآوردگر کمترین تعداد مشاهداتی است که اگر مقادیر نامتعارف یا دورافتاده اختیار کنند مقدار برآوردگر را به شدت تغییر دهند (روسو و هوبرت، ۲۰۱۳؛ هاوکنس و همکاران، ۱۹۸۴). به عنوان مثال نقطه شکست برآوردگرهای مجموع کمترین توان‌های دوم خطا و مجموع کمترین قدرمطلق انحرافات یک است، یعنی فقط یک مشاهده دورافتاده کافی است که مقدار تابع هدف مسئله بهینه‌سازی را به شدت تحت تاثیر قرار دهد (گوین و ولش، ۲۰۱۰). از طرفی نقطه شکست برآوردگر کمترین توان دوم میانه‌ی خطا  $\frac{n}{4}$  است یعنی  $\frac{n}{4}$  مشاهدات باید مقادیر دورافتاده اختیار کنند تا میانه تحت تاثیر قرار گیرد (روسو، ۱۹۸۴).
- ۳- به تحلیل داده‌ها و مقادیر خطاهای برازش به منظور شناسایی و تعیین نقاط دورافتاده پردازد.
- برای برقراری شرایط ۱ تا ۳، پس از مرتب کردن انحرافات، الحاق و جمع‌بندی آن‌ها با اعمال بردارهای وزنی صورت می‌گیرد. برای شناسایی مشاهدات دورافتاده وزن‌های بهینه در نظر گرفته می‌شوند تا مناسبترین تابع برازش مدل رگرسیونی با یک نقطه شکست بالا در مقابل دیگر برآوردگرهای کاندید حاصل شود (هوبرت و همکاران، ۲۰۰۸).
- در بخش ۲ رگرسیون کمترین قدرمطلق انحرافات مرتب شده وزنی معرفی می‌شود. در بخش ۳ الگوریتم برآوردیابی مدل بیان می‌شود. در بخش‌های ۴، ۵ و ۶ عملکرد مدل پیشنهادی در مثال‌های تحلیلی، شبیه‌سازی شده، کاربردی با داده‌های واقعی در مهندسی آب مورد بررسی و مطالعه قرار می‌گیرد. در بخش انتهایی به بحث و نتیجه‌گیری پرداخته می‌شود.

## ۲ رگرسیون کمترین قدرمطلق انحرافات مرتب شده وزنی

یک مدل آماری به صورت مجموع تابعی از متغیرهای تبیینی و مقادیر خطا به صورت

$$y = f_{\alpha}(x) + \varepsilon, \quad (1)$$

در نظر گرفته شود، که در آن  $y$  متغیر وابسته،  $x$  متغیر تبیینی،  $\varepsilon$  مقادیر خطا با شرایط زیربنایی هستند و  $f_{\alpha}(\cdot)$  با پارامترهای  $\alpha = (\alpha_0, \dots, \alpha_p)$  رابطه بین متغیرها را مشخص می‌کند. انتخاب مدل مناسب

<sup>1</sup>Break-down point

شامل تعیین نوع تابع  $f_{\alpha}(\cdot)$  و انتخاب روش برآوردیابی پارامترهای مدل است. پس از جمع‌آوری داده‌ها و انتخاب مدل رگرسیون خطی، به محاسبات آماری برای برآورد پارامترهای مدل و بررسی نکویی برازش مدل پرداخته می‌شود. فرض کنید برای مشاهدات  $(y_1, x_1), \dots, (y_n, x_n)$  رابطه خطی  $y = \alpha_0 + \alpha_1 x + \varepsilon$  برقرار باشد. برای برآورد پارامترهای  $\alpha_0$  و  $\alpha_1$  یک مسئله بهینه‌سازی با تابع هدف  $O = [O_1, \dots, O_n]$  در نظر گرفته می‌شود، که در آن  $i = 1, \dots, n, O_i = e_{(i)}$  به ترتیب هرکدام  $i$  امین قدرمطلق انحرافات مرتب شده را کمینه خواهند کرد، به طوری که

الف-  $e_i = |y_i - (\alpha + \beta \times x_i)|$  قدرمطلق مقادیر انحرافات برای  $i = 1, \dots, n$  است،

ب-  $e_{(1)} \leq \dots \leq e_{(n)}$  مقادیر مرتب شده قدرمطلق انحرافات هستند.

بر خلاف مسایل بهینه‌سازی معمولی که در آن‌ها یک تابع هدف منفرد به منظور یافتن جوابها بهینه می‌شود، بیش از یک تابع هدف در بهینه‌سازی وارد می‌شود (کوردوس و همکاران، ۲۰۱۹). در این مسئله بهینه‌سازی ابتدا کلیه توابع هدف توسط عملگرهایی (از جمله عملگرهای وزنی) به یک تابع تبدیل می‌شود (لیت و اسکریانس، ۲۰۱۹؛ یاری و چاجی، ۲۰۱۲b). در این مقاله با برقراری ارتباط بین وزن‌ها و نقاط دورافتاده به تعیین بهینه بردار وزن پرداخته می‌شود. در ادامه از بردارهای وزنی در الحاق و جمع‌بندی مقادیر خطا در یک تابع هدف به صورت

$$\min_W \min_{\alpha_0, \alpha_1} \sum_{i=1}^n \omega_i e_{(i)}, \quad \omega_i \in [0, 1], \quad \sum_{i=1}^n \omega_i = 1, \quad (2)$$

استفاده می‌شود، که در آن وزن‌ها با رویکردهای متنوعی (چاجی، ۲۰۱۷؛ چاجی و همکاران، ۲۰۱۸) قابل تعیین هستند. برای شناسایی مجموعه نقاط دورافتاده، مجموعه نقاط خوب و تحلیل حساسیت آستانه تمایز بین این دو مجموعه، بردارهای وزن به صورت زیر در نظر گرفته می‌شوند

$$W_k = (\underbrace{0, \dots, 0}_{\#k-1}, \underbrace{1}_{\#1}, \underbrace{0, \dots, 0}_{\#n-k}), \quad k = \lfloor \frac{n}{2} \rfloor + 1, \dots, n.$$

ملاحظه ۱. در رابطه (۲) خطاهای مرتب شده با اوزان مربوطه، تابع هدف مسئله بهینه‌سازی را تشکیل می‌دهند. منطقی است خطاهای بزرگ مربوط به مشاهداتی باشند که توسط مدل به خوبی برازش نشده‌اند. به این مشاهدات باید برچسب دورافتاده یا نامتعارف اختصاص داد. در نقطه مقابل مشاهداتی هستند که توسط مدل به خوبی برازش شده‌اند و دارای خطاهای برازشی کوچک هستند. به این مشاهدات برچسب

مشاهدات خوب اختصاص می‌دهیم. اکنون، باید مرز بین مشاهدات دورافتاده و خوب، تعیین شود. در اینجا مقدار  $k$  در بردار وزن  $\mathbf{W}_k$  نقش تعیین کننده مقدار آستانه‌ای بین مجموعه مشاهدات خوب و مجموعه مشاهدات دورافتاده را دارا است. مقدار بهینه  $k$  براساس معیاری که فاصله بزرگترین مقدار خطای برازش مدل بر مجموعه داده‌های خوب را با کوچکترین مقدار خطای برازش مدل بر مجموعه داده‌های دورافتاده بیشینه می‌کند، تعیین می‌شود. بدیهی است که با توجه به مرتب بودن خطاهای برازشی، هر چه این فاصله بزرگتر باشد، مدل بهینه تمایز آشکارتری بین دو مجموعه داده‌های خوب و دورافتاده ایجاد می‌کند.

**ملاحظه ۲.** در یک مجموعه مشاهدات حداکثر نیمی از آن‌ها را می‌توان به عنوان داده دورافتاده در نظر گرفت. لذا با تحلیل حساسیت مقادیر فاصله خطاهای برازش بین مجموعه داده‌های خوب و دورافتاده بر اساس کلیه مقادیر  $n, \dots, 1, \left\lceil \frac{n}{2} \right\rceil + 1$ ، بهترین مقدار  $k$  انتخاب می‌شود. جزییات نحوه محاسبه این فاصله و تعیین بهینه مقدار  $k$  در الگوریتمی که در ادامه بیان می‌شود، توضیح داده شده است.

با جایگذاری بردارهای  $\mathbf{W}_k$  در (۲)، مسئله بهینه‌سازی به صورت

$$\min_{\left\lceil \frac{n}{2} \right\rceil + 1 \leq k \leq n} \min_{(\alpha_0, \alpha_1)} e(k), \quad (3)$$

بازنویسی می‌شود و برای حل آن ابتدا مسئله‌های بهینه‌سازی

$$(\hat{\alpha}_{0k}, \hat{\alpha}_{1k}) = \operatorname{argmin} e(k), \quad k = \left\lceil \frac{n}{2} \right\rceil + 1, \dots, n. \quad (4)$$

در نظر گرفته می‌شود. نقطه شکست در هر یک از مسئله‌های بهینه‌سازی (۴) برابر  $n - k$  است، زیرا اگر  $n - k$  مشاهده، دورافتاده باشند، مقادیر خطاهای مرتب شده متناظر با آن‌ها به ترتیب در مکان‌های  $n, \dots, k + 1$  به صورت  $e_{(k+1)}, \dots, e_{(n)}$  قرار می‌گیرند. در نتیجه این مقادیر در مسئله بهینه‌سازی وارد نمی‌شوند و مقادیر پارامترهای برآورد شده را تحت تاثیر قرار نمی‌دهند. در این خصوص نیاز است با جستجو در  $n - \left\lceil \frac{n}{2} \right\rceil$  مسئله بهینه‌سازی (۴) مقدار بهینه  $k$ ، یعنی  $k^*$ ، به درستی تعیین شود. بر این اساس مناسبترین تابع نیکویی برازش برای برآوردیابی مدل به صورت  $\min_{(\alpha_0, \alpha_1)} e(k^*)$  نتیجه می‌شود. در نهایت با توجه به محاسبات مسئله‌های بهینه‌سازی (۴)، جواب مسئله بهینه‌سازی (۳)، به صورت  $(\hat{\alpha}_0, \hat{\alpha}_1) = (\hat{\alpha}_{0k^*}, \hat{\alpha}_{1k^*})$  حاصل می‌شود. توجه کنید که در هر مسئله مقدار بهینه  $k^*$  متناسب با تعداد مشاهدات دورافتاده در مجموعه داده‌ها مشخص می‌شود. این مقدار مناسبترین تابع نیکویی برازش

۴۴ ..... کاربرد عملگرهای وزنی در مدل رگرسیون

برای برآوردیابی مدل را در هر مسئله نتیجه می‌دهد و از طرفی تعیین کننده بهترین نقطه شکست در مسئله بهینه‌سازی است که برابر  $n - k^*$  است (هوبرت و همکاران، ۲۰۰۸؛ روسو و هوبرت، ۲۰۱۳).

ملاحظه ۳. بر اساس روش معرفی شده، هر انتخاب دلخواه و مناسب دیگری برای  $W_k$  ها بطور مشابه قابل بررسی است و مدل جدیدی با خصوصیات متفاوتی را معرفی خواهد نمود.

### ۳ برآورد پارامترهای مدل

در یک مدل رگرسیون مبتنی بر کمترین قدرمطلق انحرافات، پارامترهای بهینه مدل بر اساس یک مسئله بهینه‌سازی خطی محاسبه می‌شوند. در این گونه مدل‌ها اثبات می‌شود که مدل برازش شده حداقل از  $p + 1$  نقطه از مشاهدات عبور می‌کند، که  $p + 1$  تعداد پارامترهای برآورد شده در مدل است. لذا در یک مدل رگرسیون مبتنی بر مسئله بهینه‌سازی (۳) ابتدا کلیه ترکیب‌های مختلف  $p + 1$  تایی از  $n$  جفت مشاهدات  $(y, x)$  انتخاب می‌شود و معادله صفحه گذرا از این نقاط به همراه میزان خطاهایی که دیگر نقاط با صفحه گذرا بوجود می‌آورند، مشخص می‌شود (اوگوندل و همکاران، ۲۰۱۶). در انتها مدل بهینه در بین  $\binom{n}{p+1}$  مدل کاندید و بر اساس مجموعه‌ای از معیارهای برازش موجود، انتخاب می‌شود.

الگوریتم ۱. الگوریتم برآوردیابی ضرایب مدل:

گام ۱- قرار دهید  $k = \lfloor \frac{n}{p} \rfloor + 1$  و  $\mathcal{S} = \{[(y_m, x_m), (y_\ell, x_\ell)] | m, \ell = 1, \dots, n, m < \ell\}$   
 گام ۲- مدل‌های برازش کاندید را با انجام موارد زیر تعیین کنید:  
 ۱- معادله خط گذرا از دو جفت نقطه دلخواه  $[(y_m, x_m), (y_\ell, x_\ell)]$  در گردایه  $\mathcal{S}$  با شیب  $a$  و عرض از مبدا  $b$  را بدست آورید.

۲- قدرمطلق انحرافات کل مشاهدات را به صورت  $e_i = |y_i - a - bx_i|$  بدست آورید.

۳- مقادیر  $e_{(1)} \leq \dots \leq e_{(n)}$  را محاسبه کنید و مقدار  $e_{(k)}$  را تعیین نمایید.

۴- موارد ۱ تا ۳ را برای تمام جفت نقاط موجود در  $\mathcal{S}$  انجام دهید.

۵- معادله خط گذرا از جفت نقطه‌ای که کمترین مقدار  $e_{(k)}$  را در بین کلیه اعضای  $\mathcal{S}$  دارد را به عنوان بهترین خط برازش مرحله  $k$ ام انتخاب نموده، قرار دهید  $(\hat{a}_k, \hat{b}_k) = \operatorname{argmin}_{a,b} e_{(k)}$ .  
 توجه شود که در این مرحله  $(\hat{a}_k, \hat{b}_k)$  می‌تواند به عنوان یک برآورد بهینه برای مدل‌سازی مشاهدات  $(y_1, x_1), \dots, (y_n, x_n)$  استفاده شود.

۶- قدرمطلق انحرافات نسبت به مدل برآورد شده مرحله  $k$ ام و خطای نسبی هر مشاهده را به صورت

$$E_i = |y_i - (\hat{a}_k + \hat{b}_k \times x_i)|, \quad \delta_i = \frac{E(i)}{\sum_{i=1}^n e_i}, \quad i = 1, \dots, n,$$

محاسبه کنید، که در آن  $E_{(1)} \leq \dots \leq E_{(n)}$ . توجه کنید که مقادیر بزرگ  $\delta_i$  نشانگر برازش نامناسب مدل به مشاهده  $i$ ام است. مقدار  $\Delta_k = \delta_{k+1} - \delta_k$  آستانه بین نیکویی برازش مدل به داده‌های خوب و بد را اندازه‌گیری می‌کند. هر چه این مقدار برای یک مدل بزرگتر باشد، نشان می‌دهد که آن مدل توانایی بهتری در شناسایی و جداسازی مشاهدات دورافتاده دارد.

گام ۳- قرار دهید  $k = k + 1$  اگر  $k \leq n$  آنگاه گام ۲ را تکرار کنید.

گام ۴- (انتخاب مدل بهینه): در بین کلیه مدل‌های انتخاب شده در مرحله ۶ از گام ۲ با مقادیر  $\Delta_k$ ، مقدار بهینه اندیس  $k$  را به صورت  $k^* = \operatorname{argmax}_{k=n-\lfloor \frac{n}{4} \rfloor, \dots, n} \Delta_k$  انتخاب کنید. عدد  $k^*$  مربوط به مدلی است که در بین مدل‌های کاندید در مرحله ۶ از گام ۲ بیشترین مقدار اختلاف در مقادیر نیکویی برازش به مشاهدات خوب و دورافتاده را اعمال می‌کند. برآوردهای مدل مربوط به این مدل، یعنی  $(\hat{\alpha}_0, \hat{\alpha}_1) = (\hat{a}_{k^*}, \hat{b}_{k^*})$ ، به همراه تابع نیکویی برازش  $e_{(k^*)}$  مربوطه را به عنوان مدل بهینه و بهترین تابع نیکویی برازش انتخاب کنید و قرار دهید  $\hat{\alpha}_0 + \hat{\alpha}_1 x_i$ . مقدار  $\hat{y}_i = \hat{\alpha}_0 + \hat{\alpha}_1 x_i$  مقدار  $k^* - n$  بهترین نقطه شکست برآوردهای مدل بهینه را مشخص می‌کند. هرچه مقدار  $k^*$  بزرگتر باشد، تعداد مشاهداتی که توسط مدل به خوبی برازش می‌شوند، بیشتر هستند و به تبع آن مشاهدات کمتری توسط مدل به خوبی برازش نمی‌شوند و به عنوان مشاهدات دورافتاده شناسایی می‌شوند.

گام ۵- (تعیین نقاط دورافتاده): اگر مشاهدات شامل نقطه یا نقاط دورافتاده باشند، مقدار  $\delta_{k^*}$  تفاوت و پرش قابل ملاحظه‌ای با مقدار بعد از خود یعنی  $\delta_{k^*+1}$  دارد.  $k^*$  را به عنوان تعداد مشاهدات خوب مشخص کنید. اکنون قدرمطلق خطاهای برازش مدل بهینه یعنی  $e_i = |y_i - (\hat{\alpha}_0 + \hat{\alpha}_1 \times x_i)|$  را برای هر  $i = 1, \dots, n$  محاسبه نمایید. با جستجو در مجموعه مشاهدات، آن دسته از مشاهداتی که مقادیر  $\delta_{k^*+1}, \dots, \delta_n$  را تولید می‌کنند به عنوان مشاهدات دورافتاده گزارش کنید.

## ۴ مطالعه تحلیلی

برای بررسی و مقایسه رویکرد معرفی شده در مقایسه با برآوردهای متداول کمترین قدرمطلق انحرافات  $(LA)$  و کمترین توان‌های دوم خطا  $(LS)$  مشاهدات جدول ۱ را به صورت

$$x_i \sim \text{Uniform}(0, 3), \quad y_i = 2 + 4x_i + N(0, 1), \quad i = 1, \dots, 8,$$

$$x_i \sim \text{Uniform}(3, 4), \quad y_i = x_i - \text{Uniform}(1, 2), \quad i = 9, 10.$$

شبیه‌سازی شده‌اند. در جدول‌های ۱ و ۲ بر مبنای گام ۳ الگوریتم ۱، نتایج مدل‌های بهینه برای توابع هدف

جدول ۱. مشاهدات مطالعه تحلیلی										
$i$	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
$x_i$	۲/۵	۲/۸	۲/۱	۰/۸	۰/۳	۰/۸	۱/۱	۰/۶	۳/۸	۳/۶
$y_i$	۱۱	۱۳	۹	۷	۴	۴	۵	۵	۲	۲

$k = 6, \dots, 10, e(k)$  و دو مدل کمترین قدرمطلق انحرافات  $LA$  و کمترین توان‌های دوم خطا  $LS$  آورده شده است. توجه کنید که در این مرحله هر یک از مدل‌های موجود در این جداول، خود به تنهایی می‌تواند به عنوان یک مدل مناسب برای مدل‌سازی این مشاهدات بکار برده شود. اکنون مطابق گام ۴ الگوریتم، به انتخاب بهترین مدل بهینه در بین مدل‌های کاندید جدول ۱ پرداخته می‌شود. به این منظور نتایج جدول ۲ در جهت انتخاب و معرفی مدل بهینه برای این مشاهدات به صورت زیر تفسیر می‌شوند.

۱- در مدل با  $k = 6$  چهار مشاهده پتانسیل دورافتاده بودن را دارا خواهند بود. زیرا در این مدل  $82.6\%$  درصد از خطای کل  $(\delta_1 + \delta_2 + \delta_3 + \delta_4 = 0.826)$  مربوط به چهار مشاهده است (مشاهدات دورافتاده) و  $17.4\%$  درصد باقیمانده مربوط به شش مشاهده دیگر است (مشاهدات خوب). در این مدل مشاهداتی که مقادیر  $\delta_1, \delta_2, \delta_3, \delta_4$  را تولید می‌کند، مشاهده دورافتاده هستند.

۲- در مدل با  $k = 7$  سه مشاهده پتانسیل دورافتاده بودن را دارند. در این مدل  $82.1\%$  درصد از خطای کل مربوط به سه مشاهده است (مشاهدات دورافتاده) و  $17.9\%$  درصد باقیمانده مربوط به هفت مشاهده دیگر است (مشاهدات خوب). بنابراین باید سه مشاهده مربوط به  $\delta_1, \delta_2, \delta_3$  را به عنوان مشاهدات دورافتاده تعیین نمود. پس مشاهده‌ای که مقدار خطای نسبی  $\delta_4 = 0.07$  را تولید می‌کند به عنوان مشاهده دورافتاده معرفی می‌شود که مقدار خطای نسبی آن در مقایسه با  $\delta_2 = 0.37$  و  $\delta_1 = 0.38$  تفاوت زیادی



دارد. در صورت انتخاب این مدل یکی از مشاهداتی که برازش خوبی از مدل دریافت می‌کند باید به عنوان یک مشاهده دورافتاده معرفی شود. البته مقدار  $\Delta_7 = 0.02$  تاکید می‌کند که این مدل تمایز معنی‌داری در مقادیر نیکویی برازش بین مشاهدات دورافتاده و خوب ایجاد نمی‌کند.

۳- مدل  $k = 8$  که از کمینه کردن تابع برازش  $\min e_{(8)}$  به دست می‌آید، در مقایسه با مدل‌های دیگر مدل برتر است. در این مدل ۸۰ درصد از خطای کل ( $\delta_9 + \delta_{10} = 0.80$ ) مربوط به دو مشاهده است و ۲۰ درصد دیگر مربوط به هشت مشاهده باقیمانده است. از طرفی تفاوت معنی‌داری بین مقدار  $\delta_8 = 0.05$  و  $\delta_9 = 0.39$  است. لذا مشاهده مربوط به  $\delta_8$  و همچنین مشاهدات مربوط به  $\delta_1, \dots, \delta_7$  از برازش خوبی برخوردار هستند و نمی‌تواند به عنوان مشاهدات دورافتاده معرفی شوند. به عبارتی تفاوت معنی‌داری بین خطاهای برازش مدل به مشاهدات خوب و دورافتاده وجود دارد ( $\Delta_8 = 0.34$ ).

۴- مدل با  $k = 9$  نمی‌تواند به عنوان مدل بهینه انتخاب شود، زیرا در این مدل می‌توان یک مشاهده را قطعاً به عنوان مشاهده دورافتاده معرفی نمود که حداکثر خطای نسبی تولیدی آن ۲۹ درصد است. بدیهی است که ۷۱ درصد از خطای نسبی تولید شده توسط این مدل در ۹ مشاهده پخش شده است.

۵- در مدل‌های دیگر جدول ۱ تفاوت معنی‌داری بین بزرگترین مقادیر خطاهای نسبی وجود ندارد. به عنوان نمونه مدل با  $k = 10$  در مقایسه با دیگر مدل‌ها، اختلاف‌های معنی‌داری بین بزرگترین مقادیر خطاهای نسبی بوجود نمی‌آورد. همچنین با توجه به ماهیت بهینه‌سازی مدل‌های  $LS$  و  $LA$  مقدار کل خطا بطور متوسط در بین کلیه مشاهدات انتشار می‌یابد. ماهیت انتشار خطا بطور متوسط باعث می‌شود که بخشی از مشاهدات به خوبی برازش نشوند و دارای خطای برازش نسبتاً بزرگی باشند. این نتایج در شکل ۶ که نمودارهای میله‌ای مقادیر خطاهای نسبی مرتب شده مدل‌های جدول ۲ را نشان می‌دهد، قابل استنتاج است. لذا در این دو مدل نمی‌توان به تعیین مشاهدات دورافتاده پرداخت.

اکنون برای مدل بهینه  $\hat{y}_i = 3.05 + 3.18x_i$ ، مطابق گام ۵ الگوریتم ۱، مشاهدات دورافتاده تعیین می‌شوند. بردار خطاهای برازش این مدل در جدول ۴ بیانگر آن است که مشاهدات  $(y_9, x_9) = (2, 3.8)$  و  $(y_{10}, x_{10}) = (2, 3.6)$  به ترتیب بزرگترین مقادیر خطا ۱۳/۱ و ۱۲/۵ را تولید می‌کنند. این دو مقدار تفاوت زیادی با دیگر مقادیر خطا دارند و به عنوان مشاهدات دورافتاده شناسایی می‌شوند.

## ۵ مطالعه شبیه‌سازی و تحلیل حساسیت نتایج

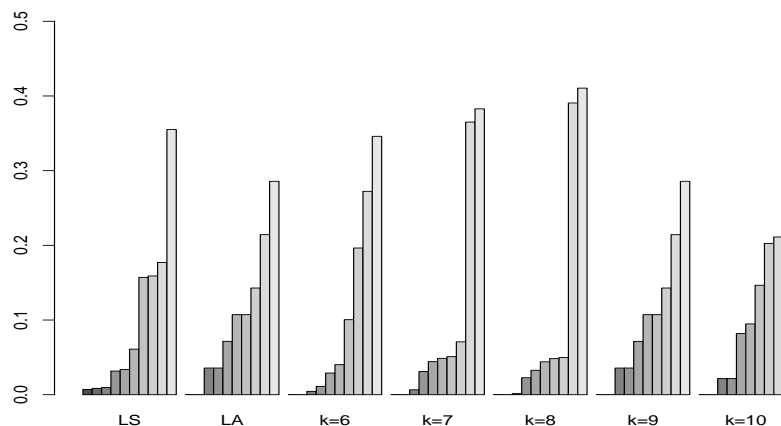
در این بخش عملکرد برآوردگرهای استوار پیشنهاد شده ( $\hat{\beta}_{rob}$ ) در مقابل برآوردگرهای کمترین توان دوم خطا ( $\hat{\beta}_{LS}$ ) و برآوردگرهای کمترین قدرمطلق انحرافات ( $\hat{\beta}_{LA}$ ) مقایسه می‌شوند. برای این منظور، از

جدول ۲. مدل‌های کاننید برآورد شده در گام ۳ الگوریتم به همراه مقادیر خطاهای مرتب شده

$k$	$\min e_{(k)}$	$(a_k, b_k)$	$\sum_{i=1}^n e_i$	$e_{(1)}$	$e_{(2)}$	$e_{(3)}$	$e_{(4)}$	$e_{(5)}$	$e_{(6)}$	$e_{(7)}$	$e_{(8)}$	$e_{(9)}$	$e_{(10)}$
۶	$\min e_{(6)}$	$(4,53, -0,67)$	۳۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
۷	$\min e_{(7)}$	$(3,40, 7,67)$	۳۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
۸	$\min e_{(8)}$	$(3,05, 3,18)$	۳۲	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
۹	$\min e_{(9)}$	$(5,00, 0,00)$	۲۸	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
۱۰	$\min e_{(10)}$	$(3,62, 1,25)$	۲۹	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
$L_A$	$\min \sum  e_i $	$(5,00, 0,00)$	۲۸	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
$L_S$	$\min \sum e_i^+$	$(5,59, 0,13)$	۱۲۵۳۳	۰,۸۸	۱,۰۷	۱,۲۱۰	۲,۹۶	۴,۲۴	۷,۶	۱۹,۶	۱۹,۹	۲۲,۲	۴۴,۵

جدول ۳. خطاهای نسبی مرتب شده مدل‌های کاننید بر اساس گام ۴ الگوریتم برای انتخاب مدل نهاییه

$k$	$\delta_1$	$\delta_2$	$\delta_3$	$\delta_4$	$\delta_5$	$\delta_6$	$\delta_7$	$\delta_8$	$\delta_9$	$\delta_{10}$	$\sum_{i=1}^k \delta_i$	$\sum_{i=k+1}^{10} \delta_i$	$\Delta_k$
۶	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰,۱۷۴	۰,۱۷۴	۰,۸۲۶
۷	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰,۱۷۹	۰,۱۷۹	۰,۸۲۱
۸	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰,۲۰	۰,۲۰	۰,۸۰
۹	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰,۲۱	۰,۲۱	۰,۷۹
۱۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۱,۰۰	۱,۰۰	۰,۰۰
$L_A$	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	—	—	—
$L_S$	۰,۰۰۷	۰,۰۰۸	۰,۰۰۹	۰,۰۱۰	۰,۰۱۱	۰,۰۱۲	۰,۰۱۳	۰,۰۱۴	۰,۰۱۵	۰,۰۱۸	—	—	—



شکل ۱. نمودارهای میله‌ای مقادیر خطاهای نسبی مدل‌های جدول ۲

جدول ۴. خطاهای برازش مدل بهینه

$i$	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
$e_i$	۰	۱/۰۵	۰/۷۳	۱/۴	۰	۱/۶	۱/۵	۰/۰۴۵	۱۳/۱	۱۲/۵

مقیاس‌های وسیع‌تری از مدل‌های شبیه‌سازی شده استفاده می‌شود، که هر یک از نقطه نظر حجم نمونه، انواع و تعداد مشاهدات دورافتاده و تعداد پارامترها یا تعداد متغیرهای ورودی متفاوت هستند. حال به تحلیل حساسیت نتایج در ارتباط با واریانس خطا  $\sigma_e^2$ ، حجم نمونه  $n$  (شامل  $n_1$  مشاهده خوب و  $n_2 = n - n_1$  مشاهده دورافتاده)، متوسط مقادیر برآورد شده برای پارامترها، انحراف معیار و کارایی برآوردها پرداخته می‌شود. برای ارایه نتایج کاملی پیرامون مقایسه بین رویکرد پیشنهادی با رویکردهای کمترین توان دوم خطا و کمترین قدرمطلق انحرافات، هر یک از این شبیه‌سازی‌ها به تعداد  $M = 10000$  بار تکرار شده‌اند. نتایج مربوط به مقادیر و شاخص‌های زیر برای هر مدل برازش شده محاسبه می‌شود:

$$\hat{\sigma}_e^2(m) = \frac{1}{n-1} \sum_{i=1}^n (e_i(m) - \bar{e}(m))^2, \quad m = 1, \dots, M,$$

$$\overline{\hat{\sigma}_e^2} = \frac{1}{M} \sum_{m=1}^M \hat{\sigma}_e^2(m), \quad s_{\hat{\sigma}_e^2}^2 = \frac{1}{M} \sum_{m=1}^M (\hat{\sigma}_e^2(m) - \overline{\hat{\sigma}_e^2})^2,$$

۵۰ ..... کاربرد عملگرهای وزنی در مدل رگرسیون

$$\begin{aligned}\overline{\hat{\beta}_j} &= \frac{1}{M} \sum_{m=1}^M \hat{\beta}_j(m), \quad s_{\hat{\beta}_j}^2 = \frac{1}{M} \sum_{m=1}^M (\hat{\beta}_j(m) - \overline{\hat{\beta}_j})^2, \\ RMSE(\hat{\beta}_j) &= \sqrt{\frac{1}{M} \sum_{m=1}^M (\hat{\beta}_j(m) - \beta_j)^2}, \quad j = 0, 1, \dots, p, \\ \text{eff}(\hat{\beta}_{\text{rob}}, \hat{\beta}_{\text{LA}}) &= \frac{\sum_{m=1}^M \|\hat{\beta}_{\text{LA}}(m) - \beta\|_2^2}{\sum_{m=1}^M \|\hat{\beta}_{\text{rob}}(m) - \beta\|_2^2}, \\ \text{eff}(\hat{\beta}_{\text{rob}}, \hat{\beta}_{\text{LS}}) &= \frac{\sum_{m=1}^M \|\hat{\beta}_{\text{LS}}(m) - \beta\|_2^2}{\sum_{m=1}^M \|\hat{\beta}_{\text{rob}}(m) - \beta\|_2^2}, \\ \text{eff}(\hat{\beta}_{\text{LS}}, \hat{\beta}_{\text{LA}}) &= \frac{\sum_{m=1}^M \|\hat{\beta}_{\text{LA}}(m) - \beta\|_2^2}{\sum_{m=1}^M \|\hat{\beta}_{\text{LS}}(m) - \beta\|_2^2},\end{aligned}$$

که در آن  $\beta = (\beta_0, \dots, \beta_p)$ ، نرم اقلیدسی،  $\|\cdot\|_2^2$ ،  $e(m) = (e_1(m), \dots, e_n(m))$  و  $\text{eff}(\cdot, \cdot)$  کارایی متقابل دو بردار از برآوردها است. مدل شبیه‌سازی شده و

## ۵.۱ شبیه‌سازی ۱

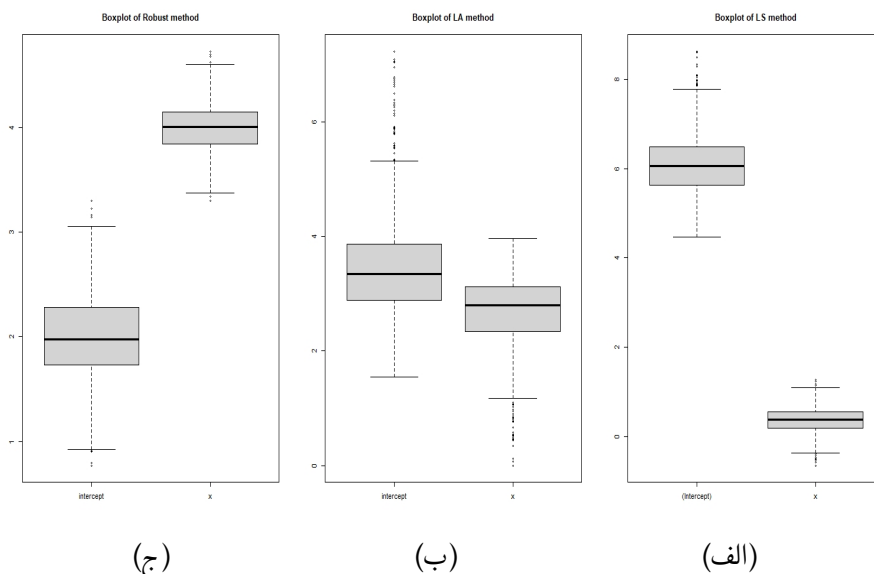
نتایج مقایسه رویکردهای موجود برای مدل

$$\begin{aligned}x_i &\sim \text{Uniform}(0, 3), \quad y_i = 2 + 4x_i + N(0, 1), \quad i = 1, \dots, n_1, \\ x_i &\sim \text{Uniform}(3, 4), \quad y_i = x_i - \text{Uniform}(1, 2), \quad i = n_1 + 1, \dots, n.\end{aligned}$$

در جدول ۵ ارائه شده است. همچنین نمودارهای جعبه‌ای مقادیر پارامترهای برآورد شده این مدل‌ها در شکل ۲ برای حالتی که  $n = 50$  و تعداد مشاهدات خوب  $n_1 = 40$  و تعداد مشاهدات دورافتاده  $n_2 = 10$  است، نشان داده شده است.

۱- به نظر می‌رسد برآوردهایی که روش پیشنهادی در این مقاله معرفی می‌کند، ویژگی نارایی را داشته باشند. اگرچه که این ویژگی باید در یک مطالعه نظری مورد بررسی قرار گیرد، اما بر اساس نتایج ارائه شده در سطر دوم جدول ۵ و جدول ۶ این موضوع دور از انتظار نیست. در این سطر میانگین مقادیر پارامترهای برآورد شده در ۱۰۰۰۰ تکرار برابر مقادیر اصلی پارامترها است.

۲- مقادیر  $RMSE$  برای برآوردهای استوار در مقایسه با برآوردهای کمترین توان دوم خطا و کمترین



شکل ۲. نمودار برآورد پارامترهای مدل الف- کمترین توان دوم خطا، ب- کمترین قدرمطلق انحرافات و ج- پیشنهاد شده

جدول ۵. نتایج مدل‌های برازش شده در ۱۰۰۰۰ تکرار شبیه‌سازی ( $n_1 = 40$  و  $n = 50$ )

مدل‌ها	پارامترها و شاخص‌ها		
$LS$	$LA$	$rob.$	
$(160, 19)$	$(229, 43)$	$(329, 26)$	$(\widehat{\sigma_e^2}, s_{\widehat{\sigma_e^2}})$
$(611, 0.35)$	$(348, 2.59)$	$(199, 3.99)$	$(\widehat{\beta_0}, \widehat{\beta_1})$
$(0.71, 0.28)$	$(0.88, 0.77)$	$(0.40, 0.23)$	$(s_{\widehat{\beta_0}}, s_{\widehat{\beta_1}})$
$(42, 36)$	$(17, 16)$	$(0.40, 0.23)$	$RMSE(\widehat{\beta})$
$eff(\widehat{\beta}_{LA}, \widehat{\beta}_{LS}) = 54$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LA}) = 242$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LS}) = 1393$	کارایی متقابل

قدرمطلق انحرافات کمتر است.

۳- مقادیر انحراف معیار برآوردگرهای استوار کمتر از مقادیر انحراف معیار برآوردگرهای حاصل شده از دو رویکرد دیگر است. این مقادیر نشان می‌دهند که برآوردگرهای استوار در مقایسه با دیگر برآوردگرها، کارا و از دقت بهتری در برآورد مقادیر پارامترها برخوردار است.

۴- بردار برآوردگرهای استوار در مقایسه با بردار برآوردگرهای کمترین توان دوم خطا و بردار کمترین قدرمطلق انحرافات کارا است.

جدول ۶. نتایج مقایسه مدل‌های برازش شده در ۱۰۰۰۰ تکرار شبیه‌سازی ( $n = ۲۰۰$  و  $n_1 = ۱۵۰$ )

پارامترها و شاخص‌ها	rob.	مدل‌ها LA	LS
$(\widehat{\sigma_e^2}, \widehat{s_{\sigma_e^2}})$	(۳۷۸, ۱,۴۴)	(۱۷,۴, ۲,۴۸)	(۱۶۵, ۱,۰۳)
$(\widehat{\beta_0}, \widehat{\beta_1})$	(۱,۹۹, ۳,۹۹)	(۵,۵۴, ۰,۴۳)	(۶,۵۱, -۰,۰۰۰۷)
$(s_{\widehat{\beta_0}}, s_{\widehat{\beta_1}})$	(۰,۲۰, ۰,۱۱)	(۰,۷۶, ۰,۷۴)	(۰,۳۷, ۰,۱۳)
$RMSE(\widehat{\beta})$	(۰,۲۰, ۰,۱۱)	(۳,۶۲, ۳,۶۴)	(۴,۵۲, ۴,۰۰)
کارایی متقابل	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LS}) = ۶۶۳,۸$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LA}) = ۴۸۰,۳$	$eff(\widehat{\beta}_{LA}, \widehat{\beta}_{LS}) = ۱,۳۸$

جدول ۷. نتایج مدل‌های ۱ در مطالعه شبیه‌سازی ۲ در ۱۰۰۰۰ بار تکرار ( $n = ۲۰۰$ )

$n_1$	پارامترها و شاخص‌ها	rob.	مدل‌ها LA	LS
۱۹۰	$(\widehat{\sigma_e^2}, \widehat{s_{\sigma_e^2}})$	(۵۱۶۸۶, ۴۱۸۳)	(۶۲۹۰, ۲۲۰,۱)	(۴۳۰۲, ۸۰,۱)
	$\widehat{\beta}$	(۲,۰۰, ۴,۰۲, ۶,۰۰, ۰,۹۹)	(۹۷,۰۳, ۴,۳۱, ۵,۴۷, ۰,۲۳)	(۱۱۲,۴۰, ۴,۰۳, ۵,۱۲, ۰,۱۰)
	$s_{\widehat{\beta}}$	(۰,۸۱, ۰,۲۲, ۰,۰۰۹, ۰,۰۰۶)	(۹۵,۵۱, ۳,۰۵, ۰,۵۴۵, ۰,۷۷۰)	(۸۰,۲, ۳,۵۲, ۰,۱۸۳, ۰,۰۶۲)
	$RMSE(\widehat{\beta})$	(۰,۸۱, ۰,۲۲, ۰,۰۰۹, ۰,۰۰۶)	(۱۰۳,۸, ۳,۸۳, ۱,۱۰, ۰,۸۴)	(۱۱۰,۶۹, ۳,۵۲, ۰,۸۹۵, ۰,۸۹۴)
	$Effi.$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LS}) = ۱۷۰,۲۱$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LA}) = ۱۲۶,۷۳$	$eff(\widehat{\beta}_{LA}, \widehat{\beta}_{LS}) = ۱,۳۴$
۱۶۰	$(\widehat{\sigma_e^2}, \widehat{s_{\sigma_e^2}})$	(۱۷۴۶۲, ۷۴۴)	(۹۹۹, ۱,۴۳)	(۸۴۶, ۹۱)
	$\widehat{\beta}$	(۲,۰۱, ۴,۰۰, ۵,۹۹, ۰,۹۹)	(۱۰۵,۶۴, ۴,۳۶, ۴,۹۰, ۰,۱۶)	(۸۱,۸۸, ۳,۴۵, ۳,۹۵, ۰,۳۴)
	$s_{\widehat{\beta}}$	(۰,۸۶, ۰,۲۴, ۰,۰۱, ۰,۰۰۷)	(۵۹, ۳,۸, ۰,۱۶, ۰,۴۴)	(۷۲, ۵,۱, ۰,۲۲, ۰,۰۵)
	$RMSE(\widehat{\beta})$	(۰,۸۶, ۰,۲۴, ۰,۰۱, ۰,۰۰۷)	(۱۰۳,۸, ۳,۸۳, ۱,۱۰, ۰,۸۴)	(۸۰,۲۱, ۵,۱۸, ۲,۰۵, ۰,۶۶)
	$Effi.$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LS}) = ۷۹,۶۳$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LA}) = ۱۳۲,۹۵$	$eff(\widehat{\beta}_{LA}, \widehat{\beta}_{LS}) = ۰,۵۹$
۱۳۰	$(\widehat{\sigma_e^2}, \widehat{s_{\sigma_e^2}})$	(۲۴۹۴۳, ۸۲۱,۵)	(۱۰۷۷, ۱۳۶,۷)	(۹۷۴, ۸۲,۴)
	$\widehat{\beta}$	(۲,۰۲, ۴,۰۰, ۵,۹۹, ۰,۹۹)	(۶۳,۴۲, ۴,۴۲, ۳,۳۶, ۰,۴۸)	(۵۶,۶۱, ۳,۰۰, ۳,۰۹, ۰,۵۳)
	$s_{\widehat{\beta}}$	(۰,۹۵, ۰,۲۷, ۰,۰۱, ۰,۰۰۷)	(۲۵,۵۸, ۵,۰۴, ۰,۸۷, ۰,۱۹۵)	(۶,۷۳, ۵,۶۹, ۰,۱۹, ۰,۰۴۵)
	$RMSE(\widehat{\beta})$	(۰,۹۵, ۰,۲۷, ۰,۰۱, ۰,۰۰۷)	(۶۶,۵۳, ۵,۰۶, ۲,۷۷, ۰,۵۵۴)	(۵۵,۰۳, ۵,۷۷, ۲,۹۱, ۰,۴۶۵)
	$Effi.$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LS}) = ۳۰۹,۶۰۸$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LA}) = ۴۴۹,۸۰۶$	$eff(\widehat{\beta}_{LA}, \widehat{\beta}_{LS}) = ۰,۶۸$
۱۱۰	$(\widehat{\sigma_e^2}, \widehat{s_{\sigma_e^2}})$	(۲۷۲۸۹, ۸۷۱,۷)	(۱۳۳۶, ۱۴۵,۸)	(۹۵۳, ۸۰,۱)
	$\widehat{\beta}$	(۱,۹۵, ۳,۹۹, ۶,۰۰, ۱,۰۰)	(۲,۲۷, ۱,۴۳, ۱,۱۹, ۰,۹۶)	(۴۳,۴۴, ۲,۶۵, ۲,۶۵, ۰,۶۳)
	$s_{\widehat{\beta}}$	(۱,۰۴, ۰,۲۸, ۰,۰۱, ۰,۰۰۸)	(۵,۵۲, ۱,۴۷, ۰,۲۰, ۰,۰۴۰)	(۶,۷۵, ۵,۵۲, ۰,۱۸, ۰,۰۴۳)
	$RMSE(\widehat{\beta})$	(۱,۰۴, ۰,۲۸, ۰,۰۱, ۰,۰۰۸)	(۵,۵۳, ۲,۹۶, ۴,۸۱, ۰,۰۵۲)	(۴۱,۹۸, ۵,۶۷, ۳,۳۵, ۰,۳۶۵)
	$Effi.$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LS}) = ۱۵۴,۱۲$	$eff(\widehat{\beta}_{rob.}, \widehat{\beta}_{LA}) = ۵۳,۳$	$eff(\widehat{\beta}_{LA}, \widehat{\beta}_{LS}) = ۲,۸۹$

۵- واریانس خطاهای به دست آمده در مدل پیشنهادی، بیشتر از واریانس خطای دو مدل دیگر است. این

موضوع در ابتدا به ماهیت برازش مدل‌های کمترین توان دوم خطا و کمترین قدرمطلق انحرافات که سعی بر کمینه‌سازی متوسط خطای هر مشاهده را دارند، مربوط است. توجه شود که در رویکرد پیشنهادی، هدف این است که تا حد امکان به مشاهدات خوب مدلی برازش شود که خطاهای اندکی داشته باشد. این موضوع باعث می‌شود که در مجموعه مقادیر خطاهای برازش، برخی از خطاها کوچک و با فاصله از خطاهایی باشند که از مشاهدات دورافتاده حاصل می‌شوند. هرچه این فاصله بیشتر باشد، تمایز آشکارتری در برازش مدل به مجموعه مشاهدات خوب و مجموعه مشاهدات دورافتاده حاصل می‌شود (به عنوان نمونه به نمودارهای میله‌ای در شکل ۶ برای  $k = 8$  و  $LA$  توجه کنید). لذا بدیهی است که پراکندگی چنین مقادیری نسبت به حالتی که خطاها بطور متوسط و میانگین توزیع شده‌اند، بیشتر باشند. از جهتی دیگر، باید در یک مطالعه نظری بررسی کرد که واریانس خطای مدل چه مقداری دارد تا بتوان از آن به عنوان یک نقطه مرجع استفاده نمود و مقدار واریانس خطاهای برازش شده را با آن مقایسه نمود.

## ۵.۲ شبیه‌سازی ۲

برای مدل

$$x_{i1} \sim \text{Uniform}(0, 1) - \text{Uniform}(0, 1), \quad i = 1, \dots, n_1,$$

$$x_{i2} \sim \text{Uniform}(5, 30) - \text{Uniform}(5, 30),$$

$$x_{i3} \sim \text{Uniform}(100, 150),$$

$$y_i = 2 + 4x_{i1} + 6x_{i2} + x_{i3} + N(0, 1),$$

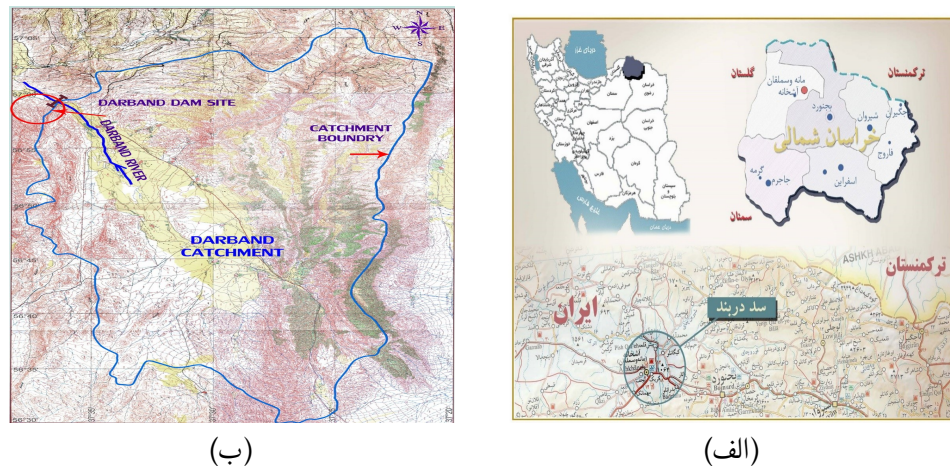
$$x_{i1} \sim \text{Uniform}(10, 11) - \text{Uniform}(10, 11), \quad i = n_1 + 1, \dots, n,$$

$$x_{i2} \sim \text{Uniform}(50, 80),$$

$$x_{i3} \sim \text{Uniform}(300, 500),$$

$$y_i = 1 + x_{i1} + x_{i2} + x_{i3} + N(0, 1),$$

داده‌ها ۱۰۰۰۰ بار شبیه‌سازی شده‌اند. در هر تکرار مدل‌های کمترین توان‌های دوم خطا، کمترین قدرمطلق انحرافات و مدل پیشنهاد شده در این مقاله برازش شده است. در جهت تحلیل و بررسی عملکرد مدل‌های مورد بحث، خلاصه‌ای از نتایج در جدول ۷ آورده شده است که مانند بخش ۵.۱ تحلیل می‌شوند.



شکل ۳. الف: سد دربند ب: موقعیت جغرافیایی محل گردآوری مشاهدات ناحیه آبخیزداری دربند

## ۶ تحلیل داده‌های بار معلق آب رودخانه

در مهندسی آب اندازه‌گیری و برآورد صحیح و دقیق بار معلق<sup>۱</sup> که توسط رودخانه‌ها حمل می‌شود در بسیاری از پروژه‌های منابع آب بسیار با اهمیت است (چاچی، ۲۰۱۹). بدین منظور فراهم نمودن نتایج و برآوردهای دقیقی از بار معلق بر مبنای متغیرهای مناسب موجود در ناحیه‌های آبخیزداری<sup>۲</sup> یکی از مسایل مهم در مهندسی آب است. در اینباره، جریان سالانه و داده‌های سری زمانی مربوط به روان آب<sup>۳</sup> و بار معلق ۱۲ ناحیه آبخیزداری مختلف، شامل سد دربند در رودخانه دربند در خراسان شمالی، ثبت شده است (شکل ۷-الف). سد دربند در رودخانه دربند، در خراسان شمالی و در طول جغرافیایی  $36^{\circ} 37'$  و عرض جغرافیایی  $48^{\circ} 58' 56''$  قرار دارد (شکل ۷-ب). مساحت این ناحیه آبخیزداری  $1117 \text{ km}^2$  است که بالاترین ارتفاع از سطح دریا آن  $2725 \text{ m}$  و پایین‌ترین ارتفاع از سطح دریا آن  $680 \text{ m}$  است. لذا بر اساس مطالعات یکسانی در ۱۲ ناحیه آبخیزداری مختلف، با استفاده از ابزارهای استاندارد برخی از خصوصیات آب‌شناسی این نواحی ثبت گردید. در این خصوص از ایستگاه‌های مختلف اندازه‌گیری مشاهدات جدول ۸ مربوط به متغیرهای زیر جمع‌آوری شدند (چاچی، ۲۰۱۹):

$y$ : میزان بار معلق که توسط رودخانه در یک سال حمل می‌شود ( $\frac{\text{ton}}{\text{year}}$ ).

<sup>1</sup>Suspended load

<sup>2</sup>Watershed areas

<sup>3</sup>Discharge



- $x_1$ : جریان سالیانه آب در رودخانه/ناحیه آبخیزداری ( $10^6 m^3$ ).
- $x_2$ : مساحت ناحیه آبخیزداری ( $km^2$ ).
- $x_3$ : تجمع سالیانه بار معلق در هر کیلومتر مربع از دلتای رودخانه ( $\frac{ton}{year} km^2$ ).
- $x_4$ : جریان آب مربوط به مقدار تجمع سالیانه بار معلق در هر کیلومتر مربع از دلتای رودخانه ( $\frac{1000 m^3}{km^2}$ ).
- در سال‌های پرباران که اغلب با سیل همراه است، مشاهداتی به عنوان دورافتاده یا دورافتاده ثبت

جدول ۸. داده‌های مهندسی آب

$x_4$	$x_3$	$x_2$	$x_1$	$y$	$i$
۲۲۸۵	۱۵۰	۱۱۱۷	۲۵۵	۱۶۳۶۵۶	۱
۹۰۷۶	۵۰۴	۴۸۷	۴۴۲	۲۴۵۳۲۷	۲
۱۱۱۲	۱۷۹	۵۹۲۵	۶۵۹	۱۰۶۱۳۵۳	۳
۳۸۴۶	۴۰۳	۴۰۳	۱۵۵	۱۶۲۳۶۶	۴
۷۶۳۳	۲۰۸	۲۰۷	۱۵۸	۴۳۱۳۳	۵
۳۸۶۹	۳۸۲	۱۲۵۶	۴۸۶	۴۸۰۲۷۰	۶
۲۴۸۷	۳۰۳	۱۳۵۵	۳۳۷	۴۱۰۵۲۴	۷
۴۸۳۶	۴۵۲	۵۴۸	۲۶۵	۲۴۸۰۰۰	۸
۲۰۴۴	۳۲۲	۱۲۰۳۷	۲۴۶۰	۳۸۷۷۱۱۲	۹
۱۷۳۷	۲۹۱	۱۷۲۶۶	۲۹۹۹	۵۰۳۳۴۲۵	۱۰
۱۶۳۴	۳۳۲	۲۷۲۴۱	۴۴۵۰	۹۰۴۱۴۸۰	۱۱
۶۷۸	۶۵۶	۸۸۴	۲۲۰	۵۷۹۶۷۰	۱۲

جدول ۹. مدل‌های کاندید برآورد شده بر اساس گام ۳ الگوریتم به همراه برآوردهای  $LS$  و  $LA$

$\sum_{i=1}^{n=12} e_i$	$(a_{0k}, a_{1k}, a_{2k}, a_{3k}, a_{4k})$	$\min e_{(k)}$	$k$
۲۳۲۴۴۰۵	$(-30542370, 128055, 674, 8566, -26804)$	$\min e_{(6)}$	۶
۲۰۱۷۳۰۵	$(-1958890, 46495, 2120, 7701, -28496)$	$\min e_{(7)}$	۷
۱۸۵۳۱۴۱	$(-2993938, 36604, 2737, 8682, -19020)$	$\min e_{(8)}$	۸
۱۹۰۵۹۱۱	$(-24249887, 104614, 1152, 7964, -47605)$	$\min e_{(9)}$	۹
۱۷۸۹۲۰۸	$(-2412375, -23368, 3699, 7094, 5747)$	$\min e_{(10)}$	۱۰
۱۹۶۱۵۱۲	$(-2898537, 101686, 1199, 8486, -24886)$	$\min e_{(11)}$	۱۱
۱۸۲۳۶۴۰	$(-3703418, 33092, 2620, 7417, 20003)$	$\min e_{(12)}$	۱۲
۱۶۷۳۷۹۷	$(-3278594, -2862, 2307, 9536, -5951)$	$\min \sum  e_i $	$LA$
۵۶۶۲۵۶۵۱۸۸۴۲	$(-6503148, 42265, 2573, 12674, 17306)$	$\min \sum e_i^2$	$LS$

می‌شوند. لذا در مدل‌سازی این مجموعه مشاهدات بهتر است که از روش‌های استوار در برآورد مدل رگرسیون  $\epsilon + \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \alpha_4 x_4$  استفاده شود. با استفاده از رویکرد پیشنهادی، نتایج جداول ۹ و ۱۰ برای مدل‌های مختلف کاندید حاصل شده است. شکل ۲ نمودار میله‌ای مقادیر خطاهای نسبی جدول ۱۰ را نشان می‌دهد. بر طبق این نتایج، در مدل  $k = 11$ ،  $\delta_i = 34$  درصد از خطای نسبی برآزش از یازده مشاهده به دست می‌آید و ۶۶ درصد دیگر از یک مشاهده به دست می‌آید. از اینجا می‌توان نتیجه گرفت که در مجموعه مشاهدات حداقل یک داده دورافتاده وجود دارد.

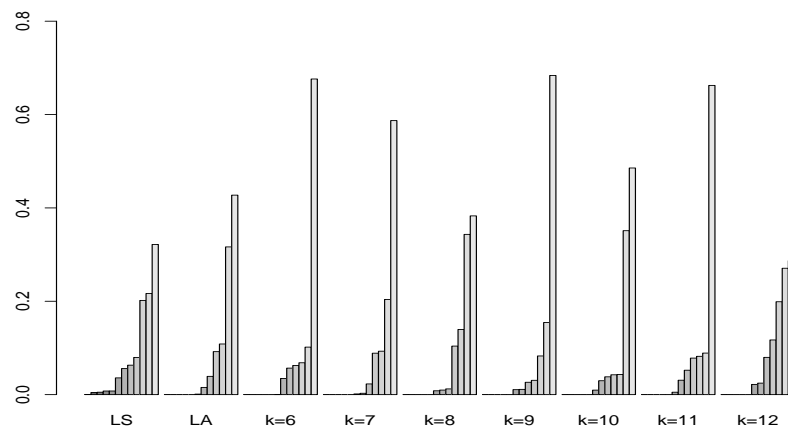
جدول ۱۰. مقادیر خطاهای نسبی مدل‌های کاندید بر اساس گام ۴ الگوریتم در تعیین بهترین مدل بهینه

$\Delta_k$	$\delta_{12}$	$\delta_{11}$	$\delta_{10}$	$\delta_9$	$\delta_8$	$\delta_7$	$\delta_6$	$\delta_5$	$\delta_4$	$\delta_3$	$\delta_2$	$\delta_1$	$k$
۰/۰۳۴	۰/۶۸	۰/۱۰	۰/۰۶۸	۰/۰۶۲	۰/۰۵۷	۰/۰۳۴	۰/۰۰۰	۰	۰	۰	۰	۰	۶
۰/۰۲۱	۰/۵۹	۰/۲۰	۰/۰۹۳	۰/۰۸۹	۰/۰۲۳	۰/۰۰۳	۰/۰۰۱	۰	۰	۰	۰	۰	۷
۰/۰۹۲	۰/۳۸	۰/۳۴	۰/۱۴۰	۰/۱۰۴	۰/۰۱۲	۰/۰۰۹	۰/۰۰۸	۰	۰	۰	۰	۰	۸
۰/۰۵۲	۰/۶۸	۰/۱۵	۰/۰۸۳	۰/۰۳۱	۰/۰۲۶	۰/۰۱۱	۰/۰۱۰	۰	۰	۰	۰	۰	۹
۰/۳۰۷	۰/۴۹	۰/۳۵	۰/۰۴۳	۰/۰۴۲	۰/۰۳۸	۰/۰۲۹	۰/۰۰۹	۰	۰	۰	۰	۰	۱۰
۰/۵۸۰	۰/۶۶	۰/۰۸	۰/۰۸۲	۰/۰۷۸	۰/۰۵۲	۰/۰۳۰	۰/۰۰۵	۰	۰	۰	۰	۰	۱۱
—	۰/۲۹	۰/۲۷	۰/۱۹۹	۰/۱۱۷	۰/۰۸۰	۰/۰۲۴	۰/۰۲۱	۰	۰	۰	۰	۰	۱۲
—	۰/۴۳	۰/۳۱	۰/۱۰۸	۰/۰۹۲	۰/۰۳۹	۰/۰۱۵	۰/۰۰۱	۰	۰	۰	۰	۰	LA
—	۰/۳۲	۰/۲۱	۰/۲۰۱	۰/۰۷۹	۰/۰۶۳	۰/۰۵۵	۰/۰۳۵	۰/۰۰۷	۰/۰۰۷	۰/۰۰۵	۰/۰۰۴	۰/۰۰۰	LS

همچنین تفاوت معنی‌داری بین مقدار  $\Delta_{11} = ۰/۵۸$  برای این مدل و دیگر مدل‌ها وجود دارد. در نتیجه، مدل برازش شده با  $k = ۱۱$  به صورت

$$\hat{y} = -۲۸۹۸۵۳/۷ + ۱۰۱۶۸/۶x_1 + ۱۱۹/۹x_2 + ۸۴۸/۶x_3 - ۲۴۸۸/۶x_4,$$

بهترین برازش به این مشاهدات است. بنابراین مشاهده مربوط به  $(y_3, x_3)$  یک داده دورافتاده است.



شکل ۴. نمودارهای میله‌ای مقادیر خطاهای نسبی مدل‌های جدول ۱۰

## بحث و نتیجه‌گیری

در این مقاله به معرفی و بررسی یک مدل رگرسیون وزنی-مبنا با انحرافات مرتب شده و مقایسه آن با حالات رگرسیون کمترین توان‌های دوم خطا و کمترین قدرمطلق انحرافات پرداخته شد. نتایج مثال‌های عددی به همراه نتایج تحلیل حساسیت برآوردهای ارایه شده در مطالعات شبیه‌سازی، نشانگر موثر بودن رویکرد پیشنهادی در تشخیص و بی‌اثرسازی مشاهدات دورافتاده دارد. رویکردی که در این مقاله مورد بحث و تحلیل قرار گرفت، جایگزین مناسبی برای مدل‌های رگرسیون مبتنی بر رویکردهای رگرسیون کمترین توان‌های دوم خطا و کمترین قدرمطلق انحرافات و حتی  $M$ -برآوردها است. بر طبق بردار وزن اختیار شده، حالات متنوعی از مدل‌های رگرسیون متداول از قبیل رگرسیون کمترین توان‌های دوم خطا و کمترین قدرمطلق انحرافات و حتی مدل‌های رگرسیون پایدار از قبیل  $LTS$ ،  $LMS$  و رگرسیون چندکی همراه با نقاط شکست بالا حاصل می‌شود. لذا کلاسی از فرمول‌ها و مدل‌های و مسایل بهینه‌سازی به وجود می‌آورد که روش‌های رگرسیون پایدار متداول را پوشش می‌دهد. در این کلاس همچنین می‌توان بر مبنای بردارهای وزن اختیار شده، الگوریتم‌ها و روش‌های حل عددی مختلفی معرفی نمود. در این مقاله روش بهینه‌سازی خاصی در کمینه‌سازی تابع برازش مورد استفاده قرار گرفت. این روش به درستی و با دقت قادر به شناسایی و کم اثر ساختن نقاط دورافتاده است. همچنین، بررسی نظری خواص برآوردهای ارایه شده به همراه بررسی خواص حدی چنین برآوردهایی می‌تواند در کارهای آتی مورد بررسی و تحقیق قرار گیرد.

## تقدیر و تشکر

نویسندگان مقاله از داوران گرامی و ویراستار محترم مجله که نظرات ارزشمند ایشان باعث بهبود مطالب ارایه شده در این مقاله گردید، کمال تشکر و قدردانی را دارند. مقاله حاضر برگرفته از طرح پژوهشی به شماره  $IR/21.MA/102$  در دانشگاه صنعتی شهدای هویزه است. نویسنده اول از اعتبارات پژوهشی به شماره  $SCU.MS98/38837$  در دانشگاه شهید چمران اهواز و نویسنده دوم از اعتبارات پژوهشی به شماره  $97/20/1073$  در دانشگاه صنعتی شهدای هویزه استفاده نموده است.

## مراجع

محمودی، م. (۱۳۸۴). رگرسیون حداقل قدرمطلق انحرافات وزنی استوار، پایان‌نامه کارشناسی ارشد، دانشکده ریاضی و کامپیوتر، دانشگاه صنعتی امیرکبیر.

- Atkinson, A. and Riani. M. (2000), *Robust Diagnostic Regression Analysis*, Springer, New York.
- Becker, C. and Gather, U. (1999), The Masking Breakdown Point of Multivariate Outlier Identification Rules, *Journal of the American Statistical Association*, **94**, 947-55.
- Billor, N., Chatterjee, S. and Hadi, A. S. (2006), A Re-weighted Least Squares Method for Robust Regression Estimation, *American Journal of Mathematical and Management Sciences*, **26**, 229-52.
- Chachi, J. (2019), A Weighted Least Squares Fuzzy Regression for Crisp Input Fuzzy Output Data, *IEEE Transactions on Fuzzy Systems*, **27**, 739-748.
- Chaji, A. R. (2017), Analytic Approach on Maximum Bayesian Entropy Ordered Weighted Averaging Operators, *Computers and Industrial Engineering*, **105**, 260-264.
- Chaji, A. R., Fukuyama, H. and Shiraz, R. K. (2018), Selecting a Model for Generating OWA Operator Weights in MGADM Problems by Maximum Entropy Membership Function, *Computers and Industrial Engineering*, **124**, 370-378.
- Dervilis, N., Worden, K. and Cross, E. J. (2015), On Robust Regression Analysis as a Mean of Exploring Environmental and Operational Conditions for SHM Data. *Journal of Sound and Vibration*, **347**, 279-296.
- Gao, X. and Feng, Y. (2018), Penalized Weighted Least Absolute Deviation Regression, *Statistical Interface*, **11**, 79-89.

Hawkins, D. M. (1980), *Identification of Outliers*, Chapman and Hall, London.

Hawkins, D. M., Bradu, D. and Kass, G. V. (1984), Location of Several Outliers in Multiple Regression Using Elemental Sets, *Technometrics*, **26**, 197-208.

Huber, P. and Ronchetti, E. M. (2009), *Robust Statistics*, 2ed., Wiley, Hoboken, NJ.

Hubert, M., Rousseeuw, P. J. and VanAelst, S. (2008), High-breakdown Robust Multivariate Methods, *Statistical Science*, **23**, 92-119.

Kordos, M., Arnaiz-González, Á. and García-Osorio, C. (2019), Evolutionary Prototype Selection for Multi-Output Regression, *Neurocomputing* **358**, 309-320.

Leite, D. and Škrjanc, I. (2019), Ensemble of Evolving Optimal Granular Experts, OWA Aggregation, and Time Series Prediction, *Information Sciences*, **504**, 95-112.

Ling, S. (2005), Self-Weighted Least Absolute Deviation Estimation for Infinite Variance Autoregressive Models, *Journal of the Royal Statistical Society, Series B*, **67**, 381-393.

Marubini, E. and A. Orenti. (2014), Detecting Outliers and/or Leverage Points: A Robust Two-Stage Procedure with Bootstrap Cut-Off Points, *Epidemiology, Biostatistics and Public Health*, **11**, 90-94.

Nguyen, T. D. and Welsch, R. (2010), Outlier Detection and Least Trimmed Squares Approximation Using Semi-Definite Programming, *Computational Statistics and Data Analysis*, **54**, 3212-3226.

- Ogundele, S. O., Mbegbu, J. I. and Nwosu, C. R. (2016), An Alternative Algorithm and R Programming Implementation for Least Absolute Deviation Estimator of the Linear Regression Models, *Journal of Modern Applied Statistical Methods*, **15**, 755-767.
- Rousseeuw, P. J. (1984), Least Median of Squares Regression, *Journal of the American Statistical Association*, **79**, 871-80.
- Rousseeuw, P. J. and Leroy. A. M. (1987), *Robust Regression and Outlier Detection*, John Wiley and Sons, New York.
- Rousseeuw, P. J. and Hubert, M. (2013), High-Breakdown Estimators of Multivariat Elocation and Scatter, In *Robustness and Complex Data Structures*, (Eds C. Becker, R. Fried and S. Kuhnt), Berlin: Springer, 49-66.
- Rousseeuw, P. J. and VanZomeren, B. C. (1990), Unmasking Multivariate Outliers and Leverage Points, *Journal of the American Statistical Association*, **85**, 633-639.
- Salini, S., Cerioli, A., Laurini, F. and Riani, M. (2016), Reliable Robust Regression Diagnostics, *International Statistical Review*, **84**, 99-127.
- Vic, B. and Lewis, T. (1994), *Outliers in Statistical Data*, Wiley, NewYork.
- Yari, G. and Chaji, A. R. (2012a), Maximum Bayesian Entropy Method for Determining Ordered Wighted Averaging Operator Weights, *Computers and Industrial Engineering*, **63**, 338-342.
- Yari, G. and Chaji, A. R. (2012b), Determination of Ordered Weighted Averaging Operator Weights Based on the M-Entropy Measures, *International Journal of Intelligent Systems*, **27**, 1020-1033.

## **Employing Weighted Operators in Ordered Least Deviations Regression Model**

Chachi, J.<sup>1</sup>, Chaji, A.<sup>2</sup>

<sup>1</sup>Department of Statistics, Shahid Chamran University of Ahvaz, Ahvaz, Iran.

<sup>2</sup>Department of Electrical Engineering, Shohadaye Hoveizeh University of Technology, Dashte-Azadegan, Iran.

**Abstract:** This article introduces a new method to estimate the least absolute linear regression model's parameters, which considers optimization problems based on the weighted aggregation operators of ordered least absolute deviations. In the optimization problem, weighted aggregation of ordered fitted least absolute deviations provides data analysis to identify the outliers while considering different fitting functions simultaneously in the modeling problem. Accordingly, this approach is not affected by outlier observations and in any problem proportional to the number of potential outliers selects the best model estimator with the optimal break-down point among a set of other candidate estimators. The performance and the goodness-of-fit of the proposed approach are investigated, analyzed and compared in modeling analytical dataset and a real value dataset in hydrology engineering at the presence of outliers. Based on the results of the sensitivity analysis, the properties of unbiasedness and efficiency of the estimators are obtained.

**Keywords:** Weighted regression, Ordered absolute-deviations, Breakdown point, Squared errors.

**Mathematics Subject Classification (2010):** 62J05, 62J02, 62G08.