

مجله علوم آماری، بهار و تابستان ۱۳۹۵

جلد ۱۰، شماره ۱، ص ۹۵-۱۱۲

DOI: 10.7508/jss.2016.01.006

برآورد تفاضلی استوار مدل‌های خطی جزئی

مهدی روزبه، جلال چاچی

گروه آمار، دانشگاه سمنان

تاریخ دریافت: ۱۳۹۳/۵/۴ تاریخ آخرین بازنگری: ۱۳۹۵/۳/۳۰

چکیده: رگرسیون خطی استوار یکی از متداول‌ترین رویکردها در روش‌های آماری استوار است. پارامترهای این روش اغلب از طریق کمترین توان‌های دوم پیراسته برآورد می‌شوند که در آن تابع هدف به گونه‌ای صورت‌بندی می‌شود که مجموع k تا از کوچکترین توان دوم باقیمانده‌ها (خطاها) کمینه شود. لذا این روش در مقایسه با روش متداول کمترین توان دوم خطا از محاسبات پیچیده‌تری برخوردار است. هدف اصلی این مقاله ارائه یک روش جدید برآورد مدل‌های خطی جزئی با رویکرد تشخیص داده‌های پرت و معرفی برآوردگرهای استوار بر مبنای کمترین توان‌های دوم پیراسته است. در این راستا ابتدا روش تفاضلی در برآورد پارامترهای مدل خطی جزئی بیان می‌شود. سپس روش به‌دست آوردن برآوردگرهای تفاضلی استوار در مدل‌های خطی جزئی بر اساس یک مسئله بهینه‌سازی مبتنی بر کمینه‌سازی مجموع k تا از کوچکترین توان دوم باقیمانده‌ها معرفی می‌شود. این رویکرد توانایی تشخیص داده‌های پرت را دارد. نتایج عددی مطالعه شبیه‌سازی و مطالعه کاربردی با داده‌های واقعی نشان‌دهنده دقت بسیار زیاد برآوردگرهای تفاضلی استوار معرفی شده در این مقاله در مقایسه با برآوردگرهای کلاسیک و متداول مدل‌های خطی

آدرس الکترونیک مسئول مقاله: جلال چاچی، jchachi@profs.semnan.ac.ir

کد موضوع بندی ریاضی (۲۰۱۰): ۶۲J۲۰، ۶۲G۳۵، ۶۲G۰۸

جزئی هستند.

واژه‌های کلیدی: برآوردگر تفاضلی استوار، کمترین توان‌های دوم پیراسته، مدل خطی جزئی استوار، داده‌های پرت.

۱ مقدمه

مدل‌های خطی جزئی انعطاف‌پذیرتر از مدل‌های خطی متداول هستند، که از دو قسمت خطی (بخش پارامتری مدل) و غیر خطی (بخش ناپارامتری مدل) تشکیل شده‌اند. در واقع استفاده از این مدل زمانی مفید است که متغیر پاسخ y به طور خطی با متغیر توضیحی x و به طور غیر خطی با متغیر توضیحی t (از طریق تابع $f(\cdot)$) رابطه داشته باشد. این‌گونه مدل‌ها نخستین بار توسط انگل و همکاران (۱۹۸۶) معرفی شدند و تا کنون و در طول دهه‌های اخیر بسیار مورد توجه محققان قرار گرفته‌اند و کاربردهای فراوانی در زمینه‌های علمی مختلف پیدا کرده‌اند، از قبیل اقتصاد (ویلیس، ۱۹۸۶؛ بلانچ‌فلاور و اسوالد، ۱۹۹۴)، مصرف سوخت خانوارها (اشمالنسی و استوکر، ۱۹۹۹). همچنین در مطالعات اخیر می‌توان به کاربردهای زیر اشاره نمود. انگل و همکاران (۱۹۸۶) برای بررسی رابطه بین مصرف برق ماهیانه خانوارها (متغیر وابسته y) با متغیرهای مستقل درآمد ماهیانه (X_1)، قیمت برق ماهیانه (X_2) و دمای هوا (t) از مدل خطی جزئی استفاده کردند. آنها حدس زدند که متغیر وابسته y رابطه خطی با متغیرهای X_1 و X_2 و رابطه غیر خطی با متغیر t دارد. یو و چن (۲۰۰۷) مدل‌های خطی جزئی را برای بررسی رابطه بین مصرف مشروبات الکلی (y) با متغیرهای درآمد (X_1)، قیمت مشروبات الکلی (X_2) و سال (t) به کار بردند. یو و همکاران (۲۰۰۷) از این مدل‌ها برای بررسی رابطه بین دستمزد ساعتی افراد شاغل (y) با متغیرهای مدت زمان اشتغال فرد در یک سال (X_1)، توانایی شغلی افراد (X_2)، آموزش‌های اضافی (X_3)، شهری و روستایی بودن افراد (X_4) و تعداد سال‌های تحصیل (t) استفاده نمودند. آکدیز و همکاران (۲۰۱۲) مدل‌های خطی جزئی را برای به دست آوردن رابطه بین مصرف ماهیانه برق (y) با متوسط درآمد افراد خانوار (X_1)، نسبت نرخ برق به نرخ گاز (X_2) و متوسط دمای هوا (t) به کار بردند.

کاربردهای گسترده اینگونه مدل‌ها در شرایط مختلف باعث ارائه روش‌ها و رویکردهای متنوع و متعددی در برآورد مدل‌های خطی جزئی شده است، از جمله: روش تفاضلی (یاتچو، ۱۹۹۷)، روش کمترین مربعات جریمه‌ای (انگل و همکاران، ۱۹۸۶)، روش مانده‌های جزئی (کوزیک، ۱۹۹۲)، روش درست‌نمایی نیم‌رخ (کارول و همکاران، ۱۹۹۷) و همچنین

روش‌هایی که در دو مرحله به برآورد قسمت خطی مدل (بخش پارامتری) و قسمت غیرخطی مدل (بخش ناپارامتری) می‌پردازند (فن و همکاران، ۱۹۹۸) (برای مطالعه بیشتر در این زمینه به روزبه (۱۳۹۰) و روزبه همکاران (۲۰۱۱) مراجعه شود).

اما باید در نظر داشت که اعتبار روش‌های برآوردیابی در آمار کلاسیک (به‌خصوص برآوردگرهای مدل‌های خطی از جمله رگرسیون کمترین توان دوم خطا) مبتنی بر فرضیه‌های زیربنایی هستند. یک برآوردگر یا یک روش برآوردیابی آماری استوار نامیده می‌شود هرگاه با به کار بردن آن به توان نتایج قابل اطمینانی حداقل در شرایط زیر گرفت:

(۱) برخی فرضیه‌های زیربنایی روش برقرار نباشد؛

(۲) در داده‌ها، مشاهدات پرت وجود داشته باشد.

تعیین داده‌های پرت و معرفی روش‌های آماری استوار در هر تجزیه و تحلیل آماری به خصوص برازش مدل‌های رگرسیون آماری بسیار مهم و قابل توجه است (حاجی باقری و همکاران، ۱۳۹۳). در رویکردهای مبتنی بر کمترین توان دوم خطا یک داده پرت تاثیر مخرب (بسیار) زیادی در شیوه برآورد پارامترها و در مدل برازش شده به داده‌ها دارد. داده‌های پرت نباید صرفاً فقط به خاطر اینکه خارج از محدوده دیگر داده‌ها قرار می‌گیرند از مطالعه حذف شوند. بلکه بسیار مهم است که تاثیر این داده‌ها در برازش مدل شناخته شود و رویکردهایی در این زمینه معرفی شوند که با وجود داده‌های پرت در مشاهدات نمونه به توانند استنباط قابل قبولی ارائه دهند. در این راستا می‌توان به برازش مدل‌های رگرسیون استوار مانند برآورد پارامترهای مدل بر اساس کمینه‌سازی مجموع توان دوم خطاهای پیراسته که اثر مخرب داده‌های پرت در برازش مدل را تا حدود (بسیار) زیادی کاهش می‌دهد، اشاره نمود. لذا در این راستا و در مقاله حاضر رویکردی پیشنهاد می‌شود که با تعیین و شناسایی مشاهدات پرت به برآورد استوار پارامترهای مدل رگرسیون خطی جزئی بر مبنای روش تفاضلی و بر اساس یک مسئله بهینه‌سازی مبتنی بر کمینه‌سازی مجموع k تا از کوچکترین توان دوم باقیمانده‌ها می‌پردازد. نتایج عددی مطالعه‌های شبیه‌سازی و کاربردی با داده‌های واقعی نشان‌دهنده دقت بسیار زیاد برآوردگرهای استوار معرفی شده در این مقاله در مقایسه با برآوردگرهای کلاسیک و متداول موجود هستند.

در بخش ۲ روش تفاضلی در برآورد مدل‌های خطی جزئی تشریح می‌شود و نحوه برآورد استوار پارامترهای مدل رگرسیون خطی جزئی بر مبنای روش تفاضلی در بخش ۳ بیان می‌شود. در بخش‌های ۴ و ۵ نتایج مثال‌های عددی (مطالعه شبیه‌سازی و تحلیل با داده‌های واقعی) آورده می‌شود و در انتها به بحث و نتیجه‌گیری پرداخته می‌شود.

۲ روش تفاضلی برآورد مدل‌های رگرسیون خطی جزئی

مدل رگرسیون خطی جزئی^۱ به صورت

$$y = X\beta + f(t) + \epsilon, \quad (1)$$

تعریف می‌شود، که در آن $y = (y_1, \dots, y_n)'$ مشاهدات متغیر پاسخ، $X = (x_1, \dots, x_n)'$ ماتریس $n \times p$ از مشاهدات متغیر توضیحی، $f(\cdot)$ یک تابع حقیقی-مقدار (نامعلوم) از مقادیر $t_1 \leq t_2 \leq \dots \leq t_n$ ، $\beta = (\beta_1, \dots, \beta_p)'$ بردار ضرایب مدل، و ϵ بردار n تایی خطاها از توزیعی با شرایط $E(\epsilon\epsilon') = \sigma^2 V$ و $E(\epsilon) = 0$ (ماتریس واریانس-کواریانس است) می‌باشند. در مدل فوق متغیر وابسته y به طور خطی با متغیر توضیحی X رابطه دارد و به طور غیرخطی و از طریق تابع $f(\cdot)$ با متغیر t در ارتباط است (هاردل و همکاران، ۲۰۰۰).

در روش تفاضلی مرتبه m برای برآورد پارامترهای مدل، ابتدا ماتریس تفاضلی

به صورت $D_{(n-m) \times n}$

$$D = \begin{bmatrix} d_0 & d_1 & \dots & d_m & 0 & 0 & \dots & 0 \\ 0 & d_0 & d_1 & \dots & d_m & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & d_0 & d_1 & \dots & d_m \end{bmatrix},$$

تشکیل می‌شود، که در آن d_j ها وزن‌های تفاضلی مدل با شرایط $\sum_{j=0}^m d_j = 0$ و

$$\sum_{j=0}^m d_j^2 = 1$$

به عنوان مثال وزن‌ها را می‌توان به صورت

$$d_j = \begin{cases} \sqrt{\frac{m}{m+1}} & j = 0, \\ -\sqrt{\frac{1}{m(m+1)}} & j = 1, \dots, m, \end{cases}$$

در نظر گرفت. هال و همکاران (۱۹۹۰) مقادیر بهینه وزن‌های تفاضلی مرتبه m را به صورت عددی تا مرتبه $m = 10$ مطابق جدول ۱ به دست آوردند.

با ضرب ماتریس تفاضلی D در طرفین معادله (۱) اثر تابع $f(\cdot)$ به صورت

$$\begin{aligned} Dy &= DX\beta + Df(t) + D\epsilon \\ &\cong DX\beta + D\epsilon, \quad (Df(t) \cong 0), \end{aligned}$$

حذف می‌شود. اکنون با تغییر متغیرهای $\tilde{y} = Dy$ ، $\tilde{X} = DX$ و $\tilde{\epsilon} = D\epsilon$ ، مدل (۱) به مدل خطی

$$\tilde{y}_{(n-m) \times 1} \cong \tilde{X}_{(n-m) \times p} \beta_{p \times 1} + \tilde{\epsilon}_{(n-m) \times 1},$$

^۱ Partial linear model

جدول ۱: ضرایب تفاضلی بهینه تا مرتبه $m = 10$

m	۱	۲	۳	۴	۵
d_0	۰/۷۰۷۱	۰/۸۰۹۰	۰/۸۵۸۲	۰/۸۸۷۳	۰/۹۰۶۴
d_1	-۰/۷۰۷۱	-۰/۵۰۰۰	-۰/۳۸۳۲	-۰/۳۰۹۹	-۰/۲۶۰۰
d_2		-۰/۳۰۹۰	-۰/۲۸۰۹	-۰/۲۴۶۴	-۰/۲۱۹۷
d_3			-۰/۱۹۴۲	-۰/۱۹۰۱	-۰/۱۷۷۴
d_4				-۰/۱۴۰۹	-۰/۱۴۲۰
d_5					-۰/۱۱۰۳
m	۶	۷	۸	۹	۱۰
d_6	۰/۹۲۰۰	۰/۹۳۰۲	۰/۹۳۸۰	۰/۹۴۴۳	۰/۹۴۹۴
d_7	-۰/۲۲۳۸	-۰/۱۹۶۵	-۰/۱۷۵۱	-۰/۱۵۷۸	-۰/۱۴۳۷
d_8	-۰/۱۹۲۵	-۰/۱۷۲۸	-۰/۱۵۶۵	-۰/۱۴۲۹	-۰/۱۳۱۴
d_9	-۰/۱۶۳۵	-۰/۱۵۰۶	-۰/۱۳۸۹	-۰/۱۲۸۷	-۰/۱۱۹۷
d_{10}	-۰/۱۳۶۹	-۰/۱۲۹۹	-۰/۱۲۲۴	-۰/۱۱۵۲	-۰/۱۰۸۵
d_{11}	-۰/۱۱۲۶	-۰/۱۱۰۷	-۰/۱۰۶۹	-۰/۱۰۲۵	-۰/۰۹۷۸
d_{12}	-۰/۰۹۰۶	-۰/۰۹۳۰	-۰/۰۹۲۵	-۰/۰۹۰۵	-۰/۰۸۷۷
d_{13}		-۰/۰۷۶۸	-۰/۰۷۹۱	-۰/۰۷۹۲	-۰/۰۷۸۲
d_{14}			-۰/۰۶۶۶	-۰/۰۶۸۷	-۰/۰۶۹۱
d_{15}				-۰/۰۵۸۸	-۰/۰۶۰۶
d_{16}					-۰/۰۵۲۷

تبدیل می‌شود، که در آن $\bar{\epsilon}_{(n-m) \times 1} \in \mathbb{R}^{n-m}$ بردار خطاها از توزیعی با شرایط $E(\bar{\epsilon}) = 0$ و ماتریس واریانس-کواریانس $E(\bar{\epsilon}\bar{\epsilon}') = \sigma^2 V_D$ است، به طوری که $V_D = DV'D' \neq I_{n-m}$ و I_{n-m} ماتریس همبندی از مرتبه $n-m$ است. مدل جدید یک مدل رگرسیون خطی متداول است، لذا پارامتر β را می‌توان با روش کمترین توان‌های دوم خطا برآورد کرد (بلسلی و همکاران، ۱۹۸۰). در روش کمترین توان‌های دوم تعمیم یافته، β با کمینه کردن مجموع توان‌های دوم وزنی خطاها به صورت

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta} (\tilde{y} - \tilde{X}\beta)' V_D^{-1} (\tilde{y} - \tilde{X}\beta) \\ &= (\tilde{X}' V_D^{-1} \tilde{X})^{-1} \tilde{X}' V_D^{-1} \tilde{y}, \end{aligned} \quad (2)$$

برآورد می‌شود. $\hat{\beta}$ برآوردگر کلاسیک مدل رگرسیون خطی جزئی است. با جایگذاری برآورد β در مدل اصلی، تابع f با روش هموارسازی رگرسیون خطی موضعی در مدل

$$f(t) = y - X\hat{\beta} - \epsilon,$$

تخمین زده می‌شود (تاکزاوا، ۲۰۰۶؛ واسرمن، ۲۰۰۵).

۳ برآورد استوار مدل‌های رگرسیون خطی جزئی

در این بخش رویکردی پیشنهاد می‌شود که پس از تعیین مشاهدات پرت به برآورد استوار پارامترهای مدل رگرسیون خطی جزئی می‌پردازد. پارامترهای این روش از طریق کمینه کردن مجموع k تا از کوچکترین توان‌های دوم باقیمانده‌ها (خطاها) برآورد می‌شوند. در این مقاله فرض می‌شود که نسبت مشاهداتی که پرت نیستند، $\frac{k}{n}$ ، معلوم است (روسو و لروی، ۱۹۸۷؛ گوین و والش، ۲۰۱۰).

اکنون به منظور تبدیل مسئله بهینه‌سازی کمترین توان‌های دوم خطاها (۲) به مسئله بهینه‌سازی کمترین توان‌های دوم خطاهای پیراسته، شاخص z_i به صورت زیر تعریف می‌شود. اگر برای مشاهده k ام $z_i = 0$ باشد، این مشاهده پرت است و اگر $z_i = 1$ باشد، این مشاهده پرت نیست. اکنون با اعمال این شاخص در مسئله بهینه‌سازی کمترین توان‌های دوم خطاها (۲)، مسئله بهینه‌سازی کمترین توان‌های دوم پیراسته به صورت

$$\min_{\beta, Z} \{(\tilde{y} - \tilde{X}\beta)' V_D^{-1} Z V_D^{-1} (\tilde{y} - \tilde{X}\beta)\}; \sum_{i=1}^n z_i = k, z_i \in \{0, 1\}$$

نوشته می‌شود، که در آن ماتریس قطری با مقادیر $z = [z_1, \dots, z_n]'$ در قطر اصلی است. مسئله بهینه‌سازی بالا با در نظر گرفتن $Z = I_n$ و $k = n$ ، به مسئله بهینه‌سازی کمترین توان‌های دوم خطای تعمیم یافته (۲) تبدیل می‌شود. در ادامه ساختار مسئله بهینه‌سازی فوق بیشتر مورد بررسی قرار می‌گیرد. برای این منظور فرض کنید

$$g(z, \beta) = (\tilde{y} - \tilde{X}\beta)' D(z) (\tilde{y} - \tilde{X}\beta),$$

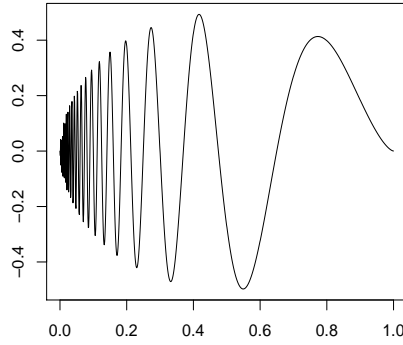
که در آن $D(z) = V_D^{-1} Z V_D^{-1}$. مسئله کمینه‌سازی $\min_{\beta, z} g(z, \beta)$ را می‌توان به صورت

$$\min_{\beta, z} g(z, \beta) = \min_z (\min_{\beta} g(z, \beta)),$$

در نظر گرفت. اکنون می‌توان مسئله بهینه‌سازی $\min_{\beta} g(z, \beta)$ را برای هر z ثابتی حل کرد و سپس جواب بهینه به دست آمده را روی کلیه مقادیر z کمینه کرد. با فرض ثابت بودن z بهترین جواب مسئله بهینه‌سازی $\min_{\beta} g(z, \beta)$ عبارت است از

$$\hat{\beta}(z) = (\tilde{X} D(z) \tilde{X})^{-1} \tilde{X}' D(z) \tilde{y}.$$

با در نظر گرفتن $\hat{\beta}(z)$ ، مجموع مانده‌های توان دوم پیراسته (یا به عبارتی $g(z, \hat{\beta}(z))$) فقط تابعی از z است که به صورت



شکل ۱: نمودار موجی شکل و هموار تابع داپلر

$$\begin{aligned} h(z) &= (\tilde{y} - \tilde{X}\hat{\beta}(z))' D(z) (\tilde{y} - \tilde{X}\hat{\beta}(z)) \\ &= \tilde{y}' D(z) \tilde{y} - \tilde{y}' D(z) \tilde{X} (\tilde{X}' D(z) \tilde{X})^{-1} \tilde{X}' D(z) \tilde{y} \end{aligned}$$

نوشته می‌شود. در نتیجه مسئله بهینه‌سازی در برآورد بردار z عبارت است از

$$\min_z h(z); \quad \sum_{i=1}^n z_i = k, \quad z_i \in \{0, 1\}.$$

اکنون با حل مسئله بهینه‌سازی بالا و پس از به‌دست آوردن مقدار \hat{z} برآوردگر استوار پارامتر β در مدل خطی جزئی (۱) به صورت

$$\begin{aligned} \hat{\beta}_R &= \hat{\beta}(\hat{z}) \\ &= (\tilde{X} D(\hat{z}) \tilde{X})^{-1} \tilde{X}' D(\hat{z}) \tilde{y} \end{aligned}$$

به‌دست می‌آید. با جایگذاری برآورد β در مدل اصلی، تابع f با روش هموارسازی رگرسیون خطی موضعی در مدل

$$f(t) = y - X\hat{\beta}_R - \epsilon,$$

تخمین زده می‌شود.

۴ مطالعه شبیه‌سازی

در این بخش، در یک مطالعه شبیه‌سازی به بررسی و مقایسه برآوردگر استوار مدل‌های خطی جزئی و برآوردگر کلاسیک متداول مدل‌های خطی جزئی پرداخته می‌شود. در این مطالعه متغیر پاسخ از مدل

$$y_i = \beta_0 + \beta_1 x_i + f(t_i) + \epsilon_i, \quad i = 1, \dots, 200,$$

به حجم $n = 200$ شبیه‌سازی می‌شود، که در آن

$$\beta = [\beta_0, \beta_1]' = [1/5, -3]',$$

$$X_i \sim N(1, 1),$$

$$f(t_i) = \sqrt{t_i(1-t_i)} \sin\left(\frac{2/1\pi}{t_i + 0/05}\right), \quad t_i = \frac{i - 0/5}{200},$$

$$\epsilon = [\epsilon_1, \epsilon_2, \dots, \epsilon_{150}, \epsilon_{151}, \epsilon_{152}, \dots, \epsilon_{200}],$$

$$\epsilon_i \sim \begin{cases} N(0, 0/01) & i = 1, \dots, 150, \\ t_{2,8} & i = 151, \dots, 200, \end{cases}$$

که در آن‌ها $f(\cdot)$ تابع معروف داپلر نام دارد و در شکل ۱ نمودار آن رسم شده است (واسرمن، ۲۰۰۵). باید توجه داشت که توابع سینوسی و موجی شکل هموار انتخاب‌های مناسبی برای تابع $f(\cdot)$ به منظور آزمون کارایی برآوردگرهای معرفی شده می‌باشند. اکنون برای اینکه داده‌ها شامل مشاهدات پرت نیز باشند، مقادیر بردار خطا از دو توزیع متفاوت شبیه‌سازی می‌شوند. در اینجا ۱۵۰ تا از مقادیر خطاها به طور مستقل از توزیع نرمال با میانگین صفر و واریانس $0/01$ تولید و شبیه‌سازی می‌شوند (مشاهدات خوب) و بقیه آنها، یعنی ۵۰ تا باقیمانده به طور مستقل از توزیع t -استیودنت غیرمرکزی با ۲ درجه آزادی و پارامتر غیرمرکزی ۸ تولید و شبیه‌سازی می‌شوند (این تعداد داده نقش مشاهدات پرت را دارند). در این شبیه‌سازی از آنجا که حجم نمونه بزرگ است ($n = 200$) و دیگر اینکه همگرایی برآوردگرها مورد توجه نمی‌باشد، تعداد تکرارها برابر یک در نظر گرفته شده است (واسرمن، ۲۰۰۵).

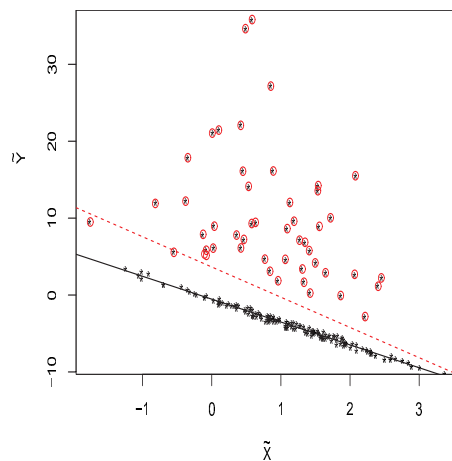
در ادامه در دو مرحله پارامتر β و تابع $f(\cdot)$ برآورد می‌شوند. ابتدا با روش تفاضلی مرتبه $m = 3$ پارامتر β برآورد می‌شود و سپس با جایگذاری $\hat{\beta}$ در مدل اصلی، تابع f با روش هموارسازی رگرسیون خطی موضعی تخمین زده می‌شود.

مرحله ۱. برآورد پارامتر β

با توجه به ضرایب تفاضلی بهینه مرتبه $m = 3$ (جدول ۱) ماتریس تفاضلی

$$D = \begin{pmatrix} d_0 & d_1 & d_2 & d_3 & 0 & 0 & \dots & 0 \\ 0 & d_0 & d_1 & d_2 & d_3 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & d_0 & d_1 & d_2 & d_3 \end{pmatrix}_{197 \times 200}$$

به دست می‌آید، که در آن $d_0 = 0/8582$ ، $d_1 = -0/2832$ ، $d_2 = -0/2809$ و $d_3 = -0/1942$. اکنون با توجه به روش معرفی شده در بخش‌های ۲ و ۳ پارامترهای مدل خطی جزئی به روش استوار (شناسایی داده‌های پرت و بی‌تاثیر ساختن آنها در برآورد پارامترهای مدل) و روش کلاسیک یا متداول برآورد می‌شوند. در شکل ۲ نمودار پراکنش داده‌ها به همراه نمودار برآوردهای تفاضلی تعمیم یافته استوار و کلاسیک داده‌ها رسم شده است. در این شکل نقاط پرت شناسایی شده توسط مدل استوار با دایره نشان داده شده‌اند.



شکل ۲: برآوردهای خطی استوار و غیراستوار

بدیهی و مشهود است که نقاط پرت شناسایی شده تاثیر زیادی در برآوردگر کلاسیک گذاشته و باعث انحراف نمودار به سمت این نقاط گردیده است. در جدول ۲ برآوردگر تفاضلی تعمیم یافته استوار ($\hat{\beta}_R$) و برآوردگر تعمیم یافته تفاضلی کلاسیک یا غیر استوار ($\hat{\beta}$) به همراه معیارهایی برای مقایسه آنها نشان داده شده است که در ادامه بیان می‌شوند:

(۱) برآورد استوار $\hat{\beta}_R = [1/28, -2/96]'$ در مقایسه با برآورد کلاسیک

۱۰۴ برآورد تفاضلی استوار مدل‌های خطی جزئی

$\hat{\beta} = [5/89, -3/80]'$ بسیار نزدیکتر به مقدار واقعی پارامتر $\beta = [1/5, -3]'$ است؛
 (۲) ضریب تعیین R^2 برای مدل برآورد شده به شیوه استوار برابر $0/90$ است که این مقدار بسیار بیشتر از $0/16$ ضریب تعیین مدلی است که پارامترهای آن به شیوه کلاسیک برآورد شده است. لازم به ذکر است که ضریب R^2 در اینجا برابر است با تغییرات کل بیان شده توسط مدل خطی جزئی به تغییرات کل، یعنی

$$R^2 = \frac{\sum_{i=1}^{200} (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^{200} (y_i - \bar{y})^2},$$

که در آن

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{f}(t_i), \quad i = 1, \dots, 200.$$

(۳) s^2 خطای کل در برآورد مدل خطی جزئی به شیوه استوار برابر $0/09$ است که این مقدار بسیار کمتر از $7/95$ خطای کل مدلی است که پارامترهای آن به شیوه کلاسیک برآورد شده است. s^2 همچنین برآوردگر مناسبی برای σ^2 است. لذا طبق نتایج فوق s^2 برآورد دقیق‌تر و نزدیکتری از مقدار $\sigma^2 = 0/01$ در مدل برآورد شده به شیوه استوار ارائه می‌دهد؛
 (۴) خطاهای استاندارد (SE) برآوردگر $\hat{\beta}_R$ کمتر از خطاهای استاندارد برآورد کلاسیک $\hat{\beta}$ است.

محاسبات این قسمت در نرم‌افزار R صورت گرفته است (فوکس و ویزبرگ، ۲۰۱۱).

جدول ۲: مقایسه برآوردگرهای استوار و غیر استوار در مطالعه شبیه‌سازی

استوار		غیر استوار		
SE_R	$\hat{\beta}_R$	SE	$\hat{\beta}$	
0/11	1/38	0/86	5/89	$\hat{\beta}_0$
0/08	-2/96	0/62	-3/80	$\hat{\beta}_1$
0/90		0/16		R^2
0/09		7/95		s^2

مرحله ۲. برآورد تابع $f(\cdot)$

در این مرحله پس از جایگذاری برآورد پارامتر β (به عنوان مثال برآورد استوار یعنی $\hat{\beta}_R = [1/38, -2/96]'$)

$$f(t_i) = y_i - (1/38 - 2/96 x_i) - \epsilon_i, \quad i = 1, \dots, 200,$$

جدول ۳: مقایسه کارایی نسبی برآوردهای استوار و غیر استوار در مطالعه شبیه سازی با حجم های نمونه ای مختلف و درصدهای متغیری از داده های پرت

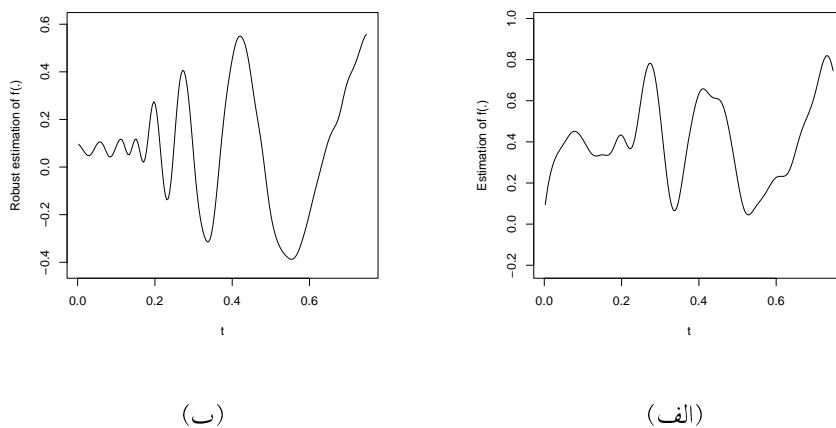
کارایی نسبی $\hat{\beta}$ و $\hat{\beta}_R$	حجم نمونه	در صد داده های پرت
۳/۹۲	۵۰	۲۵
۲/۲۰	۱۰۰	۲۵
۱/۹۸	۱۵۰	۲۵
۱/۶۹	۲۰۰	۲۵
۷/۴۰	۵۰	۳۳
۵/۰۲	۱۰۰	۳۳
۴/۹۵	۱۵۰	۳۳
۲/۲۵	۲۰۰	۳۳
۱۲/۶۰	۵۰	۵۰
۸/۴۶	۱۰۰	۵۰
۵/۵۳	۱۵۰	۵۰
۵/۶۴	۲۰۰	۵۰

به برآورد تابع $f(\cdot)$ با روش هموارسازی رگرسیون خطی موضعی پرداخته می شود. به طور مشابه با جایگذاری برآورد غیر استوار کلاسیک پارامتر β (یعنی $\hat{\beta} = [5/89, -3/80]$) در مدل اصلی تابع $f(\cdot)$ برآورد می شود. در شکل های ۱ و ۳ تابع f و برآوردهای آن با استفاده از روش های استوار $\hat{\beta}_R$ و غیر استوار کلاسیک $\hat{\beta}$ رسم شده اند. با توجه به مقایسه تابعی این برآوردها با نمودار اصلی تابع $f(\cdot)$ بدیهی است که برآورد تابع $f(\cdot)$ با استفاده از برآورد استوار قسمت خطی مدل رگرسیون خطی جزئی (یعنی $\hat{\beta}_R$) نزدیکتر به تابع اصلی است.

در انتها با محاسبه خطای کل مدل در برآورد قسمت های خطی و غیر خطی مدل، یعنی $SSE = \sum_{i=1}^{20} (y_i - \hat{y}_i)^2$ ، برای دو مدل برآورد شده به شیوه استوار و غیر استوار به مقایسه آنها پرداخته می شود. این مقدار برای مدلی که پارامترهای آن به شیوه استوار برآورد شده است برابر $SSE_R = 13/55$ است، که کمتر از $SSE = 1589/43$ خطای کل مدلی است که پارامترهای آن به شیوه معمولی برآورد شده اند.

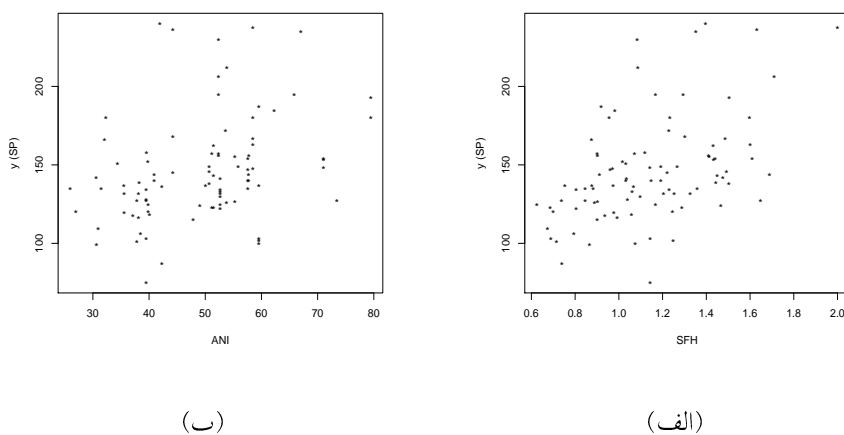
برای مقایسه کارایی نسبی برآوردهای استوار ($\hat{\beta}_R$) و برآوردهای کمترین توان های دوم خطا ($\hat{\beta}$) شبیه سازی به تعداد $M = 10^3$ بار تکرار می شود و مقادیر $\hat{\beta}(m)$ و $\hat{\beta}_R(m)$ برای $m = 1, \dots, M$ به دست آورده می شوند. کارایی نسبی این دو برآوردها به صورت

$$eff(\hat{\beta}_R, \hat{\beta}) = \frac{\frac{1}{M} \sum_{m=1}^M \|\hat{\beta}(m) - \beta\|_2^2}{\frac{1}{M} \sum_{m=1}^M \|\hat{\beta}_R(m) - \beta\|_2^2}$$



شکل ۳: برآورد تابع داپلر به روش‌های الف: غیر استوار و ب: استوار

محاسبه می‌شود، که در آن $\|x\|_2^2 = x'x$ نرم اقلیدسی بردار x است. در جدول ۳ مقادیر کارایی نسبی این دو برآوردگر گزارش شده است. از طرفی برای این که نشان داده شود برآوردگرهای استوار معرفی شده در این مقاله برتری قابل ملاحظه‌ای بر برآوردگرهای کمترین توان‌های دوم خطا دارند، در این جدول کارایی نسبی برآوردگرها بر طبق حجم‌های نمونه‌ای مختلف و با درصدهای متغیری از داده‌های پرت محاسبه شده است. در کلیه حالت‌ها نتایج به دست آمده نشان‌دهنده این است که برآوردگرهای استوار کارایی بیشتری نسبت به برآوردگرهای کمترین توان‌های دوم خطا دارند. در انتها توجه کنید که روش کمترین توان‌های دوم پیراسته، روشی بسیار استوار در برآورد مدل‌های رگرسیون خطی است. نشان داده شده است که کارایی روش کمترین توان‌های دوم پیراسته وقتی تعداد داده‌های پرت از ۵ تا $\frac{2}{3}$ تغییر می‌کند در مقایسه با روش‌های کمترین توان‌های دوم خطا و کمترین میانگین خطا بسیار قابل توجه است (روسو، ۱۹۸۴). از طرفی توجه کنید که وقتی در داده‌ها هیچ مشاهده پرتی وجود نداشته باشد این روش معادل روش کمترین توان‌های دوم خطا است. لذا می‌توان گفت که کارایی برآوردگرهای کمترین توان‌های دوم پیراسته حداقل به اندازه کارایی برآوردگرهای کمترین توان‌های دوم خطا است.



شکل ۴: نمودار پراکنش متغیر وابسته قیمت مسکن در مقابل متغیرهای مستقل الف: مساحت مسکن و ب: متوسط درآمد خانوار

۵ تحلیل با داده‌های واقعی

در ادامه به بررسی و مقایسه برآوردگرهای استوار پیشنهاد شده در این مقاله با برآوردگرهای غیر استوار در مدل‌های خطی جزئی برای داده‌های قیمت مسکن در اوتاوا پرداخته می‌شود (هو، ۱۹۹۵). در این داده‌ها قیمت مسکن فروخته شده برحسب هزار دلار متغیر وابسته (SP) است و متغیرهای مستقل عبارتند از SFH مساحت خانه برحسب فوت مربع و ANI متوسط درآمد خانوار برحسب دلار. هو (۱۹۹۵) پس از بررسی متغیرهای موجود دریافت که قیمت مسکن رابطه‌ای خطی با مساحت آن و رابطه غیرخطی با متوسط درآمد خانوار دارد. در شکل ۴ متغیر قیمت مسکن در مقابل متغیرهای مستقل مساحت خانه برحسب فوت مربع (SFH) و متوسط درآمد خانوار برحسب دلار (ANI) رسم شده است. با توجه به رویکرد معرفی شده در بخش ۳، مدل خطی جزئی

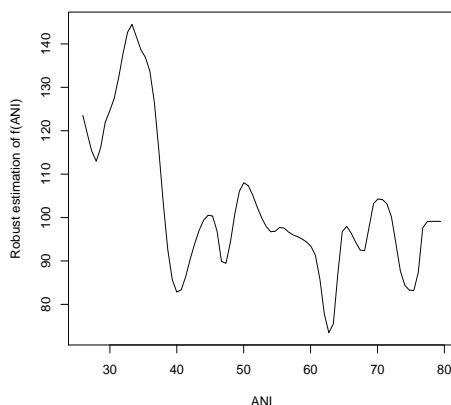
$$SP_i = \beta_0 + \beta_1 SFH_i + f(ANI_i) + \epsilon_i, \quad i = 1, \dots, 92,$$

بر اساس $n = 92$ مشاهده موجود در نظر گرفته می‌شود. برآوردگرهای استوار و غیر استوار قسمت خطی مدل در جدول ۴ به دست آمده‌اند و در شکل ۵ نشان داده شده‌اند. همچنین برآوردگرهای استوار و غیر استوار قسمت غیر خطی مدل (تابع f) در شکل ۶ نشان داده

شده‌اند.

جدول ۴: مقایسه برآوردهای استوار و غیر استوار برای داده‌های واقعی

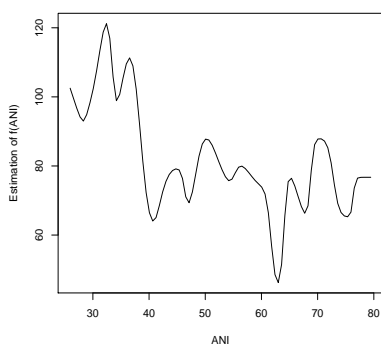
استوار		غیر استوار		
SE_R	β_R	SE	β	
۱/۹۲	-۱/۱۷	۲/۸۷	۰/۴۲	$\hat{\beta}_0$
۷/۲۷	۳۶/۹۷	۱۰/۳۷	۵۴/۹۵	$\hat{\beta}_1$
۰/۸۷۰۴		۰/۵۲۶۴		R^2
۵۸/۸		۵۲۳/۵۵		s^2



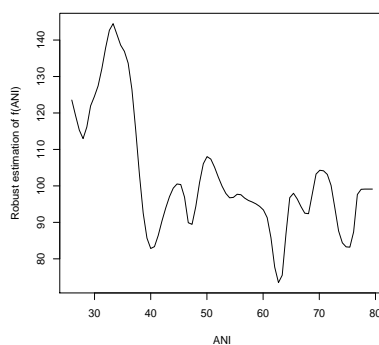
شکل ۵: برآوردهای استوار $\hat{y} = -1/17 + 35/97x$ و غیراستوار $\hat{y} = 0/42 + 54/95x$ برای قسمت خطی مدل

در جدول ۴ برآوردگر استوار ($\hat{\beta}_R$) و برآوردگر غیر استوار ($\hat{\beta}$) به همراه معیارهایی برای مقایسه آنها نشان داده شده‌اند. بر اساس معیارهای به‌دست آمده در این جدول برای مدل‌های خطی جزئی استوار و غیر استوار، به‌وضوح استدلال می‌شود (مشابه بحث و استدلال ارائه شده در مقایسه دو مدل در مطالعه شبیه‌سازی شده) که مدل استوار در برازش به داده‌ها بهتر عمل می‌کند. به عبارتی این معیارها را می‌توان به صورت زیر تفسیر نمود:

(۱) ضریب تعیین R^2 برای مدل برآورد شده به شیوه استوار برابر ۰/۸۷۰۴ است که این مقدار بسیار بیشتر از ۰/۵۲۶۴ ضریب تعیین مدلی است که پارامترهای آن به شیوه کلاسیک



(ب)



(الف)

شکل ۶: نمودار برآوردگر الف: استوار برای قسمت غیرخطی مدل \hat{f}_R و ب: غیراستوار برای قسمت غیرخطی مدل \hat{f}

برآورد شده است.

(۲) s^2 خطای کل در برآورد مدل خطی جزئی به شیوه استوار برابر $58/81$ است که این مقدار بسیار کمتر از $523/55$ خطای کل مدلی است که پارامترهای آن به شیوه کلاسیک برآورد شده است.

(۳) خطاهای استاندارد (SE) برآوردگر $\hat{\beta}_R$ برابر $2/87$ است که کمتر از خطاهای استاندارد برآورد کلاسیک $\hat{\beta}$ ، یعنی $10/37$ ، است.

با توجه به استدلال مطرح شده در بالا بر طبق نتایج عددی به دست آمده در جدول ۴، داده‌های پرت شناسایی شده در شکل ۵ (نقاط دایره‌ای شکل) تاثیر مخربی در برآورد پارامترهای مدل خطی جزئی به روش استوار وارد نمی‌کنند.

بحث و نتیجه‌گیری

در این مقاله روش جدیدی در برآورد مدل‌های خطی جزئی با رویکرد تشخیص داده‌های پرت و معرفی برآوردگرهای استوار ارائه شد. در این روش داده‌های پرت صرفاً فقط به خاطر اینکه خارج از محدوده دیگر داده‌ها قرار می‌گیرند از مطالعه حذف نمی‌شوند. بلکه رویکردی معرفی می‌شود که تاثیر مخرب داده‌های پرت را در برازش مدل بی‌اثر می‌کند. رویکرد استوار

معرفی شده در این مقاله با رویکرد غیر استوار متداول کمترین توان‌های دوم خطا در مدل‌های خطی جزئی مقایسه شد. نتایج عددی مطالعه‌های شبیه‌سازی و کاربردی با داده‌های واقعی نشان می‌دهد که تشخیص داده‌های پرت همچنان که برآوردگرهای دقیق‌تری برای قسمت خطی مدل رگرسیون خطی جزئی $(X\beta)$ ارائه می‌دهد، تاثیر فراوانی در به دست آوردن برآوردگری بهتر و با دقت بیشتر برای قسمت غیرخطی (یا تابع f) نیز دارد. بررسی خصوصیات حدی، سازگاری و کارایی برآوردگرهای استوار مدل‌های خطی جزئی و یا همچنین کاربرد روش‌های بوت استرپ در تعیین خواص توزیعی این برآوردگرها (از قبیل میانگین، واریانس و ...) می‌توانند در تحقیقات آتی مورد توجه قرار گیرند (ویزک، ۲۰۰۶a؛ ۲۰۰۶b؛ ۲۰۰۶c).

تقدیر و تشکر

نویسندگان مقاله از داوران محترم که نظرات ارزشمند ایشان باعث بهبود مطالب ارائه شده در این مقاله گردید، کمال تشکر و قدردانی را دارند.

مراجع

حاجی باقری، ف.، راسخ، ع. و آخوند، م. (۱۳۹۳)، تشخیص نقاط پرت در مدل رگرسیونی لئو، مجله علوم آماری، ۸، ۳۶-۱۹.

روزبه، م. (۱۳۹۰)، برآورد در مدل‌های خطی جزئی، رساله دکترای آمار ریاضی، دانشگاه فردوسی مشهد.

Akdeniz Duran, E., Härdle, W. K. and Osipenko, M. (2012), Difference Based Ridge and Liu Type Estimators in Semiparametric Regression Models, *Journal of Multivariate Analysis*, **105**, 164-175.

Belsley, D. A., Kuh, E. and Welsch, R. E. (1980), *Regression Diagnostics*, John Wiley, New York.

Blanchfower, D. G. and Oswald, A. J. (1994), *The Wage Curve*, MIT Press Cambridge, MA.

- Carroll, R. J., Fan, J., Gijbels, I. and Wand, M. P. (1997), Generalized Partially Single-Index Models, *Journal of American Statistical Associations*, **92**, 477-489.
- Cuzick, J. (1992), Semiparametric Additive Regression, *Journal of Royal Statistical Society, Series B*, **54**, 831-843.
- Engle, R. F., Granger, C. W. J., Rice, J. and Weiss, A. (1986), Semiparametric Estimates of the Relation Between Weather and Electricity Sales, *Journal of American Statistical Associations*, **81**, 310-320.
- Fan, J., Hardle, W. and Mammen, E. (1998), Direct Estimation of Additive and Linear Components for High-Dimensional Data, *Annals of Statistics*, **26**, 943-971.
- Fox, J. and Weisberg, S. (2011), *An R Companion to Applied Regression*, 2nd Ed., Sage Publications, Thousand Oaks, CA.
- Hall, P., Kay, J. W. and Titterton, D. M. (1990), On Estimation of Noise Variance in Two-Dimensional Signal Processing, *Advanced in Applied Probability*, **23**, 476-495.
- Hardle, W., Liang, H. and Gao, J. (2000), *Partially Linear Models*, Physika Verlag, Heidelberg.
- Ho, M. (1995), *Essay on the Housing Market*, Unpublished Ph.D. Dissertation, University of Toronto.
- 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.
- Roosbeh, M., Arashi, M. and Niroumand, H. A. (2011), Ridge Regression Methodology in Partial Linear Models with Correlated Errors, *Journal of Statistical Computation and Simulation*, **81**, 517-528.
- Rousseeuw, P. J. (1984), Least Median of Squares Regression, *Journal of American Statistical Associations*, **79**, 871-880.

- Rousseeuw, P. J. and Leroy, A. M. (1987), *Robust Regression and Outlier Detection*, John Wiley, New York.
- Schmalensee, R. and Stoker, T. M. (1999), Household Gasoline Demand in the United States, *Econometrica*, **67**, 645-662.
- Takezawa, K. (2006), *Introduction to Nonparametric Regression*, John Wiley, New Jersey.
- Visek, J. A. (2006a), The Least Trimmed Squares, Part I: Consistency, *Kybernetika*, **42**, 1-36.
- Visek, J. A. (2006b), The Least Trimmed Squares, Part II: \sqrt{n} -Consistency, *Kybernetika*, **42**, 181-202.
- Visek, J. A. (2006c), The Least Trimmed Squares, Part III: Asymptotic Normality, *Kybernetika*, **42**, 203-224.
- Wasserman, L. (2005), *All of Nonparametric Statistics*, Springer, New York.
- Willis, R. J. (1986), Wage Determinants: A Survey and Reinterpretation of Human Capital Earnings Functions, *in: Ashenfelter, O. and Layard, R. The Handbook of Labor Economics*, North Holland-Elsevier Science Publishers, Amsterdam, **1**, 525-602.
- Yatchew, A. (1997), An Elementary Estimator of the Partial Linear Model, *Economic Letters*, **57**, 135-143.
- You, J. and Chen, G. (2007), Semiparametric Generalized Least Squares Estimation in Partially Linear Regression Models with Correlated Errors, *Journal of Statistical Planning and Inference*, **137**, 117-132.
- You, J., Chen, G. and Zhou Y. (2007), Statistical Inference of Partially Linear Regression Models with Heteroscedastic Errors, *Journal of Multivariate Analysis*, **98**, 1539-1557.