

رگرسیون خطی در حضور داده‌های بد تراز

عادله عصاره^۱ فیروزه ریواز^۲

چکیده:

در این مقاله چهار رویکرد به مسئله برازش یک مدل رگرسیون خطی در حضور داده‌های بد تراز فضایی ارائه می‌شود. این رویکردها عبارتند از روش باجایگذاری، شبیه‌سازی، رگرسیون کالبیدنی و ماکسیمم درستنمایی. در دو رویکرد اول، با مدل‌بندی همبستگی موجود در متغیر توضیحی، پیشگویی آن در موقعیت‌های متناظر با متغیر پاسخ تعیین می‌شود. سپس باجایگذاری پیشگویی‌های به دست آمده به جای مقادیر واقعی در مدل رگرسیونی، برازش مدل انجام می‌شود. نشان داده می‌شود این کار باعث ایجاد خطای برکسن شده و این خطای نیز منجر به ایجاد اریبی در برآورد شبیه مدل رگرسیونی می‌شود. برای تعديل این اریبی، رویکرد رگرسیون کالبیدنی ارائه می‌شود. در رویکرد ماکسیمم درستنمایی مستقیماً از داده‌های بد تراز استفاده شده و پارامترهای مدل رگرسیونی برآورد می‌شوند. در واقع، دیگر نیازی به پیشگویی متغیر توضیحی در مکان‌های متناظر با متغیر پاسخ نیست. اما متأسفانه بررسی دقیق خواص برآورده‌گر ماکسیمم درستنمایی به دلیل نداشتن فرم تحلیلی، امکان پذیر نیست. در یک مطالعه شبیه‌سازی، عملکرد کلیه رویکردها تحت چندین مدل فضایی برای متغیر توضیحی مورد بررسی قرار می‌گیرد. مشاهده می‌شود رگرسیون کالبیدنی می‌تواند به میزان قابل توجهی اریبی برآورده‌گر شبیه خط رگرسیونی را نسبت به روش‌های دیگر کاهش دهد. به علاوه، میزان پوشش اسمی بازه اطمینان شبیه خط رگرسیونی توسط این روش قابل توجه است.

واژه‌های کلیدی: داده‌های بد تراز فضایی، رویکرد باجایگذاری، رگرسیون کالبیدنی، خطای برکسن.

۱ مقدمه

t_m مشاهده شده‌اند. این نوع داده‌ها، به داده‌های بد تراز فضایی معروف هستند.

برای مثال فرض کنید ارتباط میان میزان نوعی آلاینده‌ها در شهر تهران و تعداد بیماران قلبی در بیمارستان‌های این شهر در مقطعی از زمان مورد توجه است. اما از آنجا که مکان ایستگاه‌های هواشناسی متفاوت با مکان بیمارستان‌ها است، مشاهدات مربوط به این دو متغیر از نوع داده‌های بد تراز فضایی هستند. برای درک بهتر موضوع، اگر در یک فضای مورد مطالعه مفروض، موقعیت‌های مربوط به مشاهدات متغیر توضیحی را با t و موقعیت‌های مربوط به مشاهدات متغیر پاسخ را با s نشان دهیم، تحقیق

پیشرفت‌های متعدد در تکنیک‌های جمع‌آوری داده، موجب وفور متغیرهای پیشگوی بالقوه برای توضیح یک متغیر پاسخ فضایی شده است. با این وجود هنگامی که داده‌ها از منابع مختلف می‌آیند، مکان‌ها و مقیاس‌های فضایی به ندرت باهم متناظر هستند. به عبارت دیگر، گاهی مجموعه مکان‌های مشاهده شده برای متغیرهای توضیحی با مجموعه مکان‌های مشاهده شده برای متغیر پاسخ متفاوت است. به بیان دیگر متغیر پاسخ در موقعیت‌های مانند s_1, s_2, \dots, s_n مشاهده شده است، در حالی که متغیرهای توضیحی در موقعیت‌های دیگری مانند t_1, t_2, \dots, t_m

^۱ گروه آمار، دانشگاه شهید بهشتی

^۲ گروه آمار، دانشگاه شهید بهشتی

ماکسیمم درستنمایی برآورد کرد. بنابراین در بخش ۲ ابتدا مدل‌بندی و نمادگذاری‌ها معرفی می‌شوند. سپس رویکردهای مختلف به مسئله برآش مدل رگرسیون خطی در حالت بد ترازی فضایی، برای متغیرهای کهی شامل رویکردهای باجایگذاری، شبیه‌سازی، رگرسیون کالبینی و ماکسیمم درستنمایی به ترتیب در بخش‌های ۱.۲، ۲.۲، ۳.۲ و ۴.۲ و ماکسیمم درستنمایی به ترتیب در بخش‌های ۱، ۲، ۳.۲ و ۴.۲ ارائه می‌شوند. در بخش ۳ نیز، براساس یک مطالعه شبیه‌سازی، عملکرد روش‌های مختلف مورد بررسی قرار می‌گیرد.

۲ مدل آماری

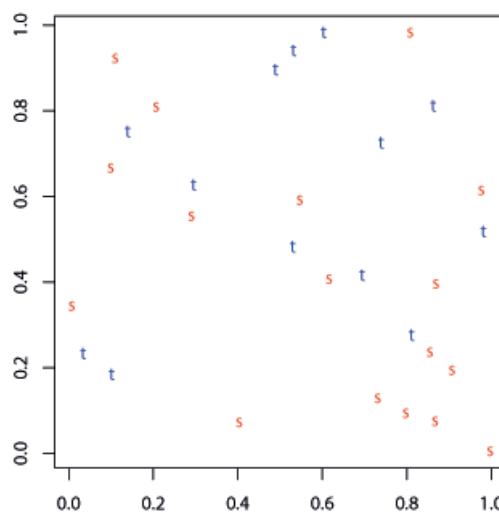
فرض کنید \mathbf{Y} و $\mathbf{X} = (Y(s_1), Y(s_2), \dots, Y(s_{n_y}))^T$ به ترتیب بردار مشاهدات متغیر پاسخ و بردار مشاهدات متغیر توضیحی در موقعیت‌های s_1, s_2, \dots, s_{n_y} باشند. همچنین Z ماتریسی از مقادیر دیگر پیشگویان اندازه‌گیری شده بدون خطا در نظر گرفته می‌شود. مدل رگرسیون خطی چندگانه

$$Y(s_i) = \beta_0 + \beta_1 X(s_i) + Z^T(s_i) \beta_z + \epsilon(s_i), \quad i = 1, 2, \dots, n_y \quad (1)$$

را در نظر بگیرید، که در آن $\sigma_\epsilon^2 \sim N(0)$. همچنین β_0 و β_z پارامترهای نامعلوم، $X(s_i)$ متغیر توضیحی در مکان s_i و $Z(s_i)$ بردار پیشگویان اندازه‌گیری شده بدون خطا در مکان s_i است. هنگامی که هر دو بردار \mathbf{Y} و \mathbf{X} مشاهده شده‌اند، می‌توان از روش‌های معمول برای برآورد پارامترهای مدل استفاده کرد. اما در مسائل بد ترازی، متغیر توضیحی X در مکان‌هایی متفاوت با مشاهدات متغیر پاسخ، مثلاً در مکان‌های t_1, t_2, \dots, t_{n_w} مشاهده شده است و با $\mathbf{X} = (X(t_1), X(t_2), \dots, X(t_{n_w}))^T$ نشان داده می‌شود. از این رو می‌توان $\mathbf{W} = (X(s_1), X(s_2), \dots, X(s_{n_y}))^T$ را با استفاده از بردار \mathbf{W} پیشگویی کرده و آن را با $\hat{\mathbf{X}} = (\hat{X}(s_1), \hat{X}(s_2), \dots, \hat{X}(s_{n_y}))^T$ نشان داد.

اگر متغیر پاسخ Y در n_y مکان و متغیر توضیحی X در n_w مکان مشاهده شده باشد، این اطلاعات در جدول ۱ خلاصه شده است. در ادامه چهار رویکرد برای برآش مدل رگرسیون خطی در مواجهه با داده‌های بد تراز فضایی ارائه می‌شود.

از داده‌های بد تراز فضایی را می‌توان در شکل ۱ مشاهده کرد.



شکل ۱: داده‌های بد تراز فضایی

نخستین بار ژو و همکاران (۲۰۰۳) مسئله تحلیل رگرسیونی را برای این نوع داده‌ها مطرح نمودند. پس از آن مادسن و همکاران (۲۰۰۸) با در نظر گرفتن یک مدل رگرسیون خطی ساده، با استفاده از روش کربیگین، متغیرهای توضیحی متناظر با پاسخ‌های مشاهده شده را پیشگویی و با دو روش کمترین توان‌های دوم و ماکسیمم درستنمایی برآش مدل را انجام دادند. همچنین گری‌پاریس و همکاران (۲۰۰۹) رویکردهای مختلفی برای تحلیل رگرسیونی داده‌های بد تراز فضایی با تکیه‌گاه نقطه‌ای ارائه دادند و به تعديل خطای اندازه‌گیری ناشی از پیشگویی در موقعیت‌های بد تراز فضایی پدست زدند.

سپس یانگ و همکاران (۲۰۰۹) برای تصحیح این نوع خطا در مواجهه با داده‌های بد تراز فضایی پرداختند و در یک مطالعه شبیه‌سازی این رویکردها را باهم مقایسه کردند. اشپیرو و همکاران (۲۰۱۱) به منظور تصحیح خطای اندازه‌گیری ناشی از پیشگویی در موقعیت‌های بد تراز فضایی، از روش‌های بوت استریوی استفاده کردند.

در این مقاله به مسئله برآش یک مدل خطی به داده‌های بد تراز فضایی می‌پردازیم. برای حل مسئله تحلیل رگرسیونی داده‌های بد تراز فضایی، در مکان‌هایی که متغیر پاسخ مشاهده شده است، متغیر توضیحی پیشگویی شده و از مقادیر حاصل به عنوان متغیر توضیحی در مدل رگرسیونی استفاده می‌شود. همچنین می‌توان بدون پیشگویی مقادیر متغیر توضیحی در مکان‌های متناظر با پاسخ، ضرایب رگرسیون را به روش

در مدل (۱) منجر به خطای برکسن^۴ می‌شود (گریپاریس و همکاران، ۲۰۰۹). با جایگذاری (۴) در مدل رگرسیون داریم:

$$\begin{aligned} Y &= \beta_0 + \beta_1 X + \beta_Z Z + \epsilon \\ &= \beta_0 + \beta_1 (\hat{X} + V) + \beta_Z Z + \epsilon \\ &= \beta_0 + \beta_1 \hat{X} + \beta_Z Z + \end{aligned}$$

که در آن $\epsilon = \beta_1 V + \eta$. ماتریس کوواریانس خطای جدید، η ، دیگر قطری نیست. زیرا واریانس ماندها برای مدل رگرسیون عبارت از:

$$var(\eta) = \beta_1^2 \Sigma_V + \sigma_\epsilon^2 I_{n_y}.$$

است. بنابراین اگرچه برآوردهای OLS برای β_1 ناریب است (با این فرض که \hat{X} برای X ناریب است)، اما واریانس آن دیگر واریانس واقعی $\hat{\beta}_1$ نخواهد بود. در این گونه موقع از روش کمترین توانهای دوم وزنی (WLS) استفاده می‌شود. برای این منظور عناصر قطری ماتریس کوواریانس خطای را به عنوان وزن‌هایی برای مدل رگرسیون در نظر می‌گیرند. اما چون در $var(\eta) = \beta_1^2 \Sigma_V + \sigma_\epsilon^2 I_{n_y}$ پارامترهای β_1 و σ_ϵ^2 نامعلوم هستند، لذا این ماتریس و در نتیجه عناصر قطری آن نامعلوم هستند و بنابراین نمی‌توان از روش WLS برای برآورد ضرایب رگرسیونی استفاده کرد.

لازم به ذکر است که حالتی خاص از روش باجایگذاری به روش کریگ و رگرس^۵ معروف است که در آن \hat{X} پیشگویی کریگی در نظر گرفته می‌شود و به عنوان متغیر توضیحی در مدل رگرسیون به کار می‌رود (مادسن و همکاران، ۲۰۰۸).

۲۰۲ رویکرد شبیه‌سازی

روش‌های پیشگویی فضایی معمولاً رویه‌ای هموار از فرایند مورد مطالعه ارائه می‌دهند. لذا تغییرپذیری مقادیر پیشگویی بسیار کمتر از تغییرپذیری واقعی متغیر موردنظر است. بنابراین در زمین آمار اغلب به جای روش‌های هموارسازی فضایی از شبیه‌سازی زمین آماری استفاده می‌شود. گریپاریس و همکاران (۲۰۰۹) رویکرد شبیه‌سازی را به این صورت ارائه دادند که از برآورد توزیع متغیر توضیحی به شرط داده‌ها، M نمونه s_1, s_2, \dots, s_M در موقعیت‌های $X_{(k)}$ را به این صورت

جدول ۱: جدول داده‌های بد تراز فضایی			
متغیرهای توضیحی	متغیر پاسخ	موقعیت	
s_i or t_j	Y	X	Z
s_1	$Y(s_1)$	-	$Z(s_1)$
s_2	$Y(s_2)$	-	$Z(s_2)$
:	:	:	:
s_{n_y}	$Y(s_{n_y})$	-	$Z(s_{n_y})$
t_1		$X(t_1)$	$Z(t_1)$
t_2		$X(t_2)$	$Z(t_2)$
:	:	:	:
t_{n_w}		$X(t_{n_w})$	$Z(t_{n_w})$

۱۰۲ رویکرد باجایگذاری

فرض کنید (s_i) در مدل (۱) نرمال باشد. برای برآذش مدل رگرسیون به روش باجایگذاری^۳ ابتدا X را با استفاده از یکی از روش‌های پیشگویی فضایی و مبتنی بر بردار W پیشگویی کرده و آن را با \hat{X} نشان می‌دهیم. سپس \hat{X} به عنوان متغیر توضیحی در مدل رگرسیون به کار برده می‌شود. در نهایت برآوردهای کمترین توانهای دوم معمولی برای β و واریانس آن به شکل زیر محاسبه می‌شوند:

$$\hat{\beta}_{\text{plug-in}} = (\hat{X}^T \hat{X})^{-1} \hat{X}^T \mathbf{Y} \quad (۲)$$

$$Var(\hat{\beta}_{\text{plug-in}}) = \sigma_\epsilon^2 (\hat{X}^T \hat{X})^{-1} \quad (۳)$$

که در آن $\hat{X} = [1_{n_y \times 1} \hat{X}]$. این روش برآوردهای ناریبی را برای پارامترهای خط رگرسیون نتیجه می‌دهد، اما استفاده از \hat{X} به جای X در مدل رگرسیونی، عدم حتمیتی را وارد مدل می‌کند و موجب ایجاد همبستگی در ساختار خطای می‌شود. برای توضیح این مطلب، توجه کنید که مقادیر حاصل از روش‌های پیشگویی فضایی تغییرپذیری کمتری نسبت به مقادیر واقعی دارند، لذا می‌توان نوشت:

$$X = \hat{X} + V \quad (۴)$$

که در آن $\hat{X} - X = V$ خطای مربوط به پیشگویی X و مستقل از ϵ است. بعلاوه $V \sim N(0, \Sigma_V)$. استفاده از \hat{X} به جای

³Plug-in Approach

⁴Berkson error

⁵Krige-and-Regress approach

دارند. لذا روش ارائه شده در (۵) نسبت به روش گری پاریس و همکاران (۲۰۰۹)، با آنچه در عمل اتفاق می‌افتد، سازگارتر است.

اکنون با به کار بردن هریک از $\mathbf{X}_{(k)}$ های تولیدشده به عنوان متغیر توضیحی در (۱)، مدل را برازش می‌دهیم و در نتیجه M برآورد برای β_1 و واریانس درون شبیه‌سازی^۶ حاصل می‌شود که آن‌ها را به ترتیب با $\hat{\beta}_{1(k)}$ و $k = 1, \dots, M$, $W_{(k)}$ نشان می‌دهیم. سپس برای به دست آوردن یک برآورد کلی از $\hat{\beta}_1$ ها میانگین گرفته می‌شود:

$$\hat{\beta}_{1M} = \frac{1}{M} \sum_{k=1}^M \hat{\beta}_{1(k)}.$$

تغییرپذیری این برآوردگر دو جزء دارد: میانگین واریانس درون شبیه‌سازی (\bar{W}_M) و واریانس بین شبیه‌سازی^۷ (B_M)، که به ترتیب عبارت از

$$\bar{W}_M = \frac{1}{M} \sum_{k=1}^M W_{(k)},$$

و

$$B_M = \frac{1}{M-1} \sum_{k=1}^M (\hat{\beta}_{1(k)} - \hat{\beta}_{1M})^2,$$

همستند. بنابراین تغییرپذیری کل به صورت

$$\widehat{Var}(\hat{\beta}_{1M}) = \bar{W}_M + \frac{M+1}{M} B_M,$$

خواهد بود که $(1 + \frac{1}{M})$ تعديلی برای M تعداد متناهی شبیه‌سازی است (یانگ و همکاران، ۲۰۰۹).

مقادیر $\mathbf{X}_{(k)}$ تولیدشده از شبیه‌سازی، تغییرپذیری یکسانی با مقادیر واقعی \mathbf{X} دارند. اختلاف این دو را باید به عنوان خطای اندازه‌گیری کلاسیک در نظر گرفت (گری پاریس و همکاران، ۲۰۰۹). بنابراین استفاده از $\mathbf{X}_{(k)}$ به جای \mathbf{X} در مدل (۱) را می‌توان در قالب یک مدل خطای اندازه‌گیری کلاسیک بررسی کرد و در نتیجه، برآوردگرهای حاصل به دلیل وجود این نوع خطای اریب خواهند بود.

۳۰۲ رویکرد رگرسیون کالبیدنی

رگرسیون کالبیدنی^۸ (RC) روشی آماری برای تعدیل آثار خطای اندازه‌گیری در برآورد پارامترهای یک مدل است. لذا برای تعدیل آثار این نوع خطای در برآوردهای حاصل از رویکردهای باجایگذاری و

⁶within-simulation variance

⁷between-simulation variance

⁸Regression Calibration Approach

s_{n_y} تولید می‌شود و هریک از این M نمونه به عنوان متغیر توضیحی در مدل رگرسیون (۱) مورد استفاده قرار می‌گیرد. در نتیجه M برآوردهای $\hat{\beta}_{1(k)}$ ، $k = 1, 2, \dots, M$ ، حاصل می‌شود و سپس برای به دست آوردن یک برآورد کلی از آن‌ها میانگین گرفته می‌شود. در این صورت واریانس $\hat{\beta}_1$ برابر

$$Var(\hat{\beta}_1) = Var(E(\hat{\beta}_{1(k)} | \mathbf{X}_{(k)})) + E(Var(\hat{\beta}_{1(k)} | \mathbf{X}_{(k)})).$$

است (گری پاریس و همکاران، ۲۰۰۹). یانگ و همکاران (۲۰۰۹) از روش تجزیه چولسکی برای تولید تحقیقاتی از متغیر توضیحی که خواص فضایی یکسانی با داده‌های اصلی دارند، استفاده کردند. برای توضیح این روش فرض کنید $\begin{pmatrix} \mathbf{W} \\ \mathbf{X}_{(k)} \end{pmatrix}$ دارای توزیع نرمال با ماتریس کوواریانس

$$\Sigma = \begin{pmatrix} \Sigma_W & \Sigma_{XW} \\ \Sigma_{XW}^T & \Sigma_X \end{pmatrix}$$

باشد. Σ را به صورت

$$\Sigma = LL^T$$

$$L = \begin{pmatrix} L_{11} & \circ \\ L_{21} & L_{22} \end{pmatrix} \quad \text{فرض کنید}$$

تجزیه می‌کنیم که در آن

$$N(0, \Sigma) \quad \mathbf{V} = \begin{pmatrix} \mathbf{V}_1 \\ \mathbf{V}_{2(k)} \end{pmatrix} \quad \text{برداری از تحقیقاتی مستقل و همتوزیع (۱)}$$

باشد. آنگاه با استفاده از

$$\begin{pmatrix} \mathbf{W} \\ \mathbf{X}_{(k)} \end{pmatrix} = \begin{pmatrix} L_{11} & \circ \\ L_{21} & L_{22} \end{pmatrix} \begin{pmatrix} L_{11}^{-1} \mathbf{W} \\ \mathbf{V}_{2(k)} \end{pmatrix}$$

$$= \begin{pmatrix} \mathbf{W} \\ L_{21} L_{11}^{-1} \mathbf{W} + L_{22} \mathbf{V}_{2(k)} \end{pmatrix},$$

$$k = 1, \dots, M, \quad (5)$$

می‌توان M تحقیق از فرایند \mathbf{X} در موقعیت‌های s_1, s_2, \dots, s_{n_y} تولید کرد. شایان ذکر است که مقادیر شبیه‌سازی شده بر اساس این روش در موقعیت‌هایی که مشاهده موجود است، مقادیر یکسانی با مشاهدات

در مدل رگرسیون (۱) خواهیم داشت:

$$\begin{aligned} \mathbf{Y} &= \beta_0 \mathbf{1}_{n_y \times 1} + \beta_1 \hat{\mathbf{X}}_{cal} + \beta_Z Z + \epsilon \\ &= \beta_0 \mathbf{1}_{n_y \times 1} + \beta_1 \hat{\mathbf{X}} \hat{\gamma} + \beta_Z Z + \epsilon \end{aligned} \quad (7)$$

ولذا برآوردهای کالبینی β به صورت

$$\begin{aligned} \hat{\beta}_{cal} &= \hat{\Gamma}^{-1} \hat{\beta}_{plug-in} \\ &= \hat{\Gamma}^{-1} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \mathbf{Y} \end{aligned} \quad (8)$$

به دست می‌آید که $\hat{\beta}_{plug-in}$ در رابطه (۲) ارائه شده است. همچنین واریانس $\hat{\beta}_{cal}$ با استفاده از بسط تیلور به شکل

$$Var(\hat{\beta}_{cal}) = \hat{\Gamma}^{-1} [\beta_1^2 \sigma_\delta^2 (D^T D)^{-1} + \sigma_\epsilon^2 (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1}] (\hat{\Gamma}^{-1})^T$$

به دست می‌آید که در آن $D = [1 \quad \hat{\mathbf{W}}_1 \quad Z_1]$

۴.۲ رویکرد ماکسیمم درستنمایی

تا اینجا روش‌هایی که ارائه شد همگی مبتنی بر پیشگویی مقادیر مشاهده نشده متغیر توضیحی در موقعیت‌های متناظر با متغیر پاسخ و استفاده از مقادیر حاصل در مدل رگرسیونی بودند. اکنون روش ماکسیمم درستنمایی ارائه می‌شود که در آن برای برآورد ضرایب رگرسیونی، مستقیماً از بردار W استفاده می‌شود.

مدل (۱) با توجه میانگین و کوواریانس

$$E \begin{bmatrix} \mathbf{Y} \\ \mathbf{X} \\ \mathbf{W} \end{bmatrix} = \begin{bmatrix} [\beta_0 + \beta_1 \mu_X + Z^T(s_i) \beta_Z]_{i=1}^{n_y} \\ \mu_X \mathbf{1}_{(n_y+n_w) \times 1}, \end{bmatrix} \quad (9)$$

و

$$cov \begin{bmatrix} \mathbf{Y} \\ \mathbf{X} \\ \mathbf{W} \end{bmatrix} = \begin{bmatrix} \beta_1^2 \Sigma_X + \Sigma_\epsilon & \beta_1 \Sigma_X & \beta_1 \Sigma_{XW} \\ \beta_1 \Sigma_X & \Sigma_X & \Sigma_{XW} \\ \beta_1 \Sigma_{XW}^T & \Sigma_{XW}^T & \Sigma_W, \end{bmatrix} \quad (10)$$

را در نظر بگیرید. با توجه به مفروضات مدل داریم:

$$\begin{bmatrix} \mathbf{Y} \\ \mathbf{W} \end{bmatrix} \sim N \left(\begin{bmatrix} [\beta_0 + \beta_1 \mu_X + Z^T(s_i) \beta_Z]_{i=1}^{n_y} \\ \mu_X \mathbf{1}_{n_w \times 1} \end{bmatrix}, \Sigma \right),$$

که در آن

$$\Sigma = \begin{bmatrix} \beta_1^2 \Sigma_X + \Sigma_\epsilon & \beta_1 \Sigma_{XW} \\ \beta_1 \Sigma_{XW}^T & \Sigma_W \end{bmatrix}.$$

شبیه‌سازی که در بخش‌های قبل ارائه شدند، می‌توان از رویکرد رگرسیون کالبینی که توسط گری‌پاریس و همکاران (۲۰۰۹) ارائه شده است، استفاده نمود.

در این روش، ابتدا بردار مشاهدات متغیر توضیحی یعنی $\mathbf{W} = (X(t_1), X(t_2), \dots, X(t_{n_w}))^T$ با مشاهدات متغیر پاسخ مشاهده شده‌اند، به دو بردار \mathbf{W}_1 با طول k و \mathbf{W}_2 با طول $n_w - k$ افزایش می‌شود. سپس \mathbf{W}_1 با استفاده از بردار Z پیشگویی شده و با $\hat{\mathbf{W}}$ نشان داده می‌شود. همچنین ماتریس \mathbf{W}_2 را به طور متناظر با \mathbf{W} ، به دو ماتریس Z_1 و Z_2 افزایش می‌کنیم. اکنون مبتنی بر \mathbf{W}_1 ، رگرسیون $\hat{\mathbf{W}}$ روی \mathbf{W}_1 برآورده می‌شود. برای این منظور، مدل

$$\begin{aligned} W_1(t_i) &= \gamma_0 + \gamma_1 \hat{W}_1(t_i) + \gamma_z^T Z_1(t_i) \\ &\quad + \delta(t_i), \quad i = 1, \dots, k, \end{aligned} \quad (9)$$

را با \circ $E(\delta(t_i)) = \sigma_\delta^2$ در نظر بگیرید. با برآوردهای مدل (۶) با روش کمترین توان‌های دوم، برآوردهای $\hat{\gamma} = (\hat{\gamma}_0, \hat{\gamma}_1, \hat{\gamma}_z^T)^T$ ، به دست می‌آید که از آن‌ها برای کالبینه کردن مقادیر پیشگویی $\hat{\mathbf{X}}$ استفاده می‌شود. با توجه به رابطه (۶) داریم:

$$E(X(s_i) | \hat{\mathbf{X}}(s_i)) = \gamma_0 + \gamma_1 \hat{\mathbf{X}}(s_i) + \gamma_z^T Z(s_i), \quad i = 1, \dots, n_y,$$

با جایگذاری $\hat{\gamma} = (\hat{\gamma}_0, \hat{\gamma}_1, \hat{\gamma}_z^T)$ به جای γ در عبارت فوق، مقادیر کالبینه شده $\hat{\mathbf{X}}$ به صورت

$$\hat{X}_{cal}(s_i) = \hat{\gamma}_0 + \hat{\gamma}_1 \hat{X}(s_i) + \hat{\gamma}_z^T Z(s_i), \quad i = 1, \dots, n_y,$$

به دست می‌آید و لذا می‌توان نوشت

$$\hat{X}_{cal} = \hat{\mathbf{X}} \hat{\gamma},$$

که در آن $\hat{\mathbf{X}} = [\mathbf{1}_{n_y \times 1} \quad \hat{\mathbf{X}} \quad Z]$ و برداری از یک‌هاست. ماتریس

تبديل Γ را به صورت زیر تعریف می‌کنیم:

$$\Gamma = \begin{bmatrix} 1 & \gamma_0 & \mathbf{0}_{1 \times q} \\ \mathbf{0} & \gamma_1 & \mathbf{0}_{1 \times q} \\ \mathbf{0}_{q \times 1} & \gamma_z & I_{q \times q} \end{bmatrix}.$$

لازم به ذکر است که برآورده ماتریس Γ ، ماتریس $[1_{n_y \times 1} | \hat{\mathbf{X}} | Z]$ را به ماتریس $[\hat{\mathbf{X}}_{cal} | Z]$ تبدیل می‌کند. حال با جایگذاری

برآورد کرد. لازم به ذکر است واریانس مجانی رابطه (۱۳)، به تحقق‌های مستقل و هم‌توزیع $\mathbf{W}^T \mathbf{Y}$ وابسته است. اما در اغلب کاربردها، فقط یک بردار \mathbf{W} مشاهده شده است. بنابراین برآورده‌گر واریانس (۱۳) ممکن است خیلی خوب عمل نکند.

۳ مطالعه شبیه‌سازی

به منظور مقایسه رویکردهای ارائه شده در این مقاله برای رگرسیون داده‌های بد تراز فضایی، از یک مطالعه شبیه‌سازی استفاده می‌کنیم. در این مطالعه چندین سناریوی مختلف با $N = 500$ مجموعه داده شبیه سازی شده برای هر سناریو استفاده می‌شود. مقادیر متغیر توضیحی بنابر رابطه

$$\mathbf{X} = \mathbf{g} + \boldsymbol{\delta}$$

تولید می‌شود که در آن

$$\mathbf{g} \sim N(\mu_1, R(\phi, \nu)),$$

و برای R تابع همبستگی مترن به صورت

$$\frac{1}{\Gamma(\nu)^{2\nu-1}} \left(\frac{2\sqrt{\nu}t}{\phi\pi} \right)^\nu K_\nu \left(\frac{2\sqrt{\nu}t}{\phi\pi} \right),$$

در نظر گرفته می‌شود که در آن t فاصله، ϕ دامنه فضایی، $\nu < 0$ پارامتر همواری و K_ν تابع بسل اصلاح شده نوع دوم از مرتبه ν است. همچنین δ تغییرات مقیاس کوچک است و فرض می‌شود

$$\delta \sim N(o, \sigma_\delta^2 I_{n_w}),$$

برای بردار متغیر پاسخ فرض می‌کنیم

$$\mathbf{Y} \sim N(\beta_0 + \beta_1 \mathbf{X}, \sigma^2 I),$$

و در کلیه سناریوها $\beta_0 = 1$ و $\beta_1 = 0$ قرار می‌دهیم. فرض استقلال خطاهای متغیر پاسخ بیان می‌کند که تنها عنصر مولد خودهمبستگی فضایی در متغیر پاسخ، متغیر توضیحی است.

معمولًا تعداد موقعیت‌هایی که متغیر پاسخ مشاهده می‌شود بیشتر از تعداد موقعیت‌هایی است که متغیر توضیحی مشاهده می‌شود. بنابراین برای تمام مجموعه داده‌ها فرض می‌کنیم $n_w = 80$ و $n_y = 200$ است. لازم به ذکر است که بردارهای \mathbf{X} و \mathbf{Y} را ابتدا در 280 در $n_w + n_y$ موقعیت به طور کامل شبیه‌سازی و سپس آن‌ها را بد تراز می‌کنیم. به این صورت که در هر مجموعه داده 80 مقدار اول بردار \mathbf{Y} و 200 مقدار

بنابراین چگالی توازن \mathbf{Y} و \mathbf{W} عبارت از

$$f_{YW}(y, w) = \frac{1}{\sqrt{(2\pi|\Sigma|)}} \times \exp\left(-\frac{1}{2}\mathbf{V}^T \Sigma^{-1} \mathbf{V}\right),$$

است که در آن

$$\mathbf{V} = \begin{bmatrix} \mathbf{V}_Y \\ \mathbf{V}_X \end{bmatrix} = \begin{bmatrix} \mathbf{Y} - [\beta_0 + \beta_1 \mu_X + Z^T(s_i) \boldsymbol{\beta}_Z]_{i=1}^{n_y} \\ \mathbf{W} - \mu_X \mathbf{1}_{m \times 1} \end{bmatrix}.$$

فرض کنید θ_X و θ_ϵ به ترتیب بردارهای پارامترهای نیم‌تغییرنگار فرایندهای پیشگو و خطأ باشند. برآورده‌گر ماکسیمم درستنمایی β با مینیمم کردن مقدار

$$l(\phi) = -\log(f_{YW}) = \frac{1}{2} \log|\Sigma| + \frac{1}{2} \mathbf{V}^T \Sigma^{-1} \mathbf{V}, \quad (11)$$

نسبت به پارامترهای $\theta_X^T \mu_X \theta_\epsilon^T$ به دست می‌آید. معمولاً θ_X و θ_ϵ بردارهای به طول ۳ هستند (زیرا شامل پارامترهای دامنه، اثر قطعه‌ای و ازاره هستند) و $\beta = [\beta_0, \beta_1, \beta_Z]^T$ نیز برداری به طول $q + 2$ است. بنابراین مینیمم کردن (۱۱) یک مسئله در فضای $(q + 9)$ بعدی است که به صورت مستقیم تحت قید مثبت بودن پارامترهای نیم‌تغییرنگار تعیین می‌شوند.

نتیجه بهینه سازی عددی، برداری از برآوردهای ML برای تمامی این پارامترهاست:

$$\hat{\phi}_{ML} = [\hat{\beta}_{ML}^T \ \hat{\mu}_{X,ML} \ \hat{\theta}_{X,ML}^T \ \hat{\theta}_{\epsilon,ML}^T]^T,$$

و برآورده ML پارامتر شیب β_1 ، عنصر دوم ϕ خواهد بود که با $\hat{\beta}_{1,ML}$ نشان داده می‌شود.

برآورده واریانس $\hat{\beta}_{1,ML}$

اگر توزیع $\mathbf{Y} \mathbf{W}^T$ شرایط نظم را داشته باشد، که در حالت گاوی سی بستگی به شکل توابع نیم‌تغییرنگار دارد، آنگاه برآورده‌گر ماکسیمم درستنمایی، نالاریب و بهطور مجانی نرم‌ال است (مادسون و همکاران، ۲۰۰۸). در این حالت ماتریس مجانی کوواریانس ϕ_{ML} معکوس ماتریس اطلاع فیشر $I(\phi)^{-1}$ است که $- (j, i)$ امین عنصر آن عبارت از

$$I(\phi)_{ij} = -E\left(\frac{\partial l^*}{\partial \phi_i \phi_j} l(\phi; \mathbf{Y}, \mathbf{W})\right) \quad (12)$$

است. واریانس $\hat{\beta}_{1,ML}$ را می‌توان به شکل زیر

$$\widehat{var}(\hat{\beta}_{1,ML}) = I^{-1}(\phi)_{22} \Big|_{\phi=\hat{\phi}_{ML}} \quad (13)$$

- سناریوی D : یک فرایند فضایی با همواری کم و دامنه ۱ و اثر قطعه‌ای بالا

$$g \sim N(0, R(1, 0/5)), \quad \delta \sim N(0, \sigma_\delta^2 I_{80}),$$

$$\sigma_\delta^2 = 0/42, \quad \sigma_\epsilon^2 = 0/82$$

- سناریوی E : یک فرایند فضایی با همواری کم و دامنه ۰/۳

$$g \sim N(0, R(0/3, 0/5)), \quad \delta \sim N(0, \sigma_\delta^2 I_{80}),$$

$$\sigma_\delta^2 = 0/22, \quad \sigma_\epsilon^2 = 0/82$$

بدیهی است هرچه فرایند ناهموارتر باشد، برآورد پارامترها چالش برانگیزتر خواهد بود. شکل ۳ تحقیقاتی از فرایند پیشگو را در این پنج سناریو نشان می‌دهد.

ابتدا در هر سناریو با استفاده از مجموعه داده‌های کامل (تراز)، مدل رگرسیونی خطی ساده برآش و پارامترهای مدل برآورده شوند. سپس داده‌های بد تراز در نظر گرفته می‌شوند و مدل رگرسیونی، با روش‌هایی که ذکر شد، برآش داده می‌شود.

در روش باجایگذاری، برای هر مجموعه داده ابتدا بهترین پیشگوی خطی نااریب (کریگی) متغیر توضیحی در مکان‌های متناظر با متغیر پاسخ را با استفاده از تابع `krige.conv` در بسته `geoR` در R محاسبه و سپس پارامترهای مدل رگرسیون خطی Y روی مقادیر کریگی حاصل برآورده شوند.

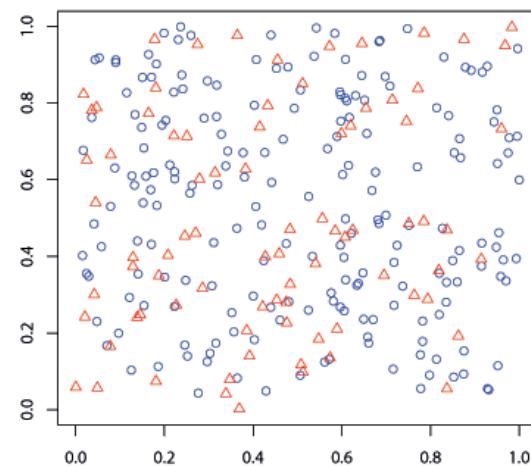
در رویکرد شبیه‌سازی، ابتدا ماتریس هم‌تغییرنگار مانا را با استفاده از تابع `cov.spatial` در بسته `geoR` و تجزیه چولسکی این ماتریس را با استفاده از تابع `chol` محاسبه می‌کنیم. سپس برای هر مجموعه داده به تعداد $M = 100$ نکرار عمل شبیه‌سازی متغیر توضیحی در موقعیت‌های متناظر با متغیر پاسخ را، همانگونه که در بخش ۲.۲ شرح داده شد، انجام می‌دهیم و در نتیجه ۱۰۰ مقدار برای پارامتر شیب خط رگرسیون حاصل می‌شود. میانگین این ۱۰۰ مقدار، برآورده β_1 برای آن مجموعه داده است.

به علاوه، داده‌های شبیه‌سازی شده را با عمل کالبیدن، کالیبره کرده و سپس با استفاده از مقادیر کالیبره شده، رگرسیون را انجام می‌دهیم.

در رویکرد رگرسیون کالبیدنی برای هر مجموعه داده تمامی عناصر بردار W را به نوبت با استفاده از $1 - n_w$ عنصر دیگر این بردار پیشگویی و سپس بردار W را روی بردار مقادیر کریگی آن مدل و بردار پارامترهای γ برآورده می‌کنیم. سپس با استفاده از بردار $\hat{\gamma}$ ، مقادیر کریگیدن در رویکرد

آخر بردار X را حذف می‌کنیم و به این ترتیب داده‌های فضایی بد تراز حاصل می‌شود.

فضای مطالعه را مریع واحد در نظر می‌گیریم و به تصادف ۲۸۰ موقعیت نقطه‌ای در آن انتخاب و فرض می‌کنیم در ۸۰ نقطه از این موقعیت‌ها، متغیر توضیحی و در ۲۰۰ نقطه دیگر، متغیر پاسخ مشاهده شده است. شکل ۲ فضای مطالعه و موقعیت‌های داده‌های بد تراز فضایی با Δ و می‌دهد که در آن ۸۰ موقعیت مربوط به مشاهدات متغیر توضیحی با \triangle و ۲۰۰ موقعیت مربوط به مشاهدات متغیر پاسخ با \circ نشان داده شده است.



شکل ۲: فضای مطالعه و موقعیت‌های داده‌های بد تراز فضایی در مطالعه شبیه‌سازی

اکنون پنج سناریو با پارامترهای همواری و دامنه‌های فضایی مختلف به صورت زیر در نظر می‌گیریم:

- سناریوی A : یک فرایند فضایی با همواری بالا و دامنه ۱/۶

$$g \sim N(0, R(1/6, 1)), \quad \delta \sim N(0, \sigma_\delta^2 I_{80}),$$

$$\sigma_\delta^2 = 0/12, \quad \sigma_\epsilon^2 = 0/82$$

- سناریوی B : یک فرایند فضایی با همواری بالا و دامنه ۱

$$g \sim N(0, R(1, 1)), \quad \delta \sim N(0, \sigma_\delta^2 I_{80}),$$

$$\sigma_\delta^2 = 0/22, \quad \sigma_\epsilon^2 = 0/82$$

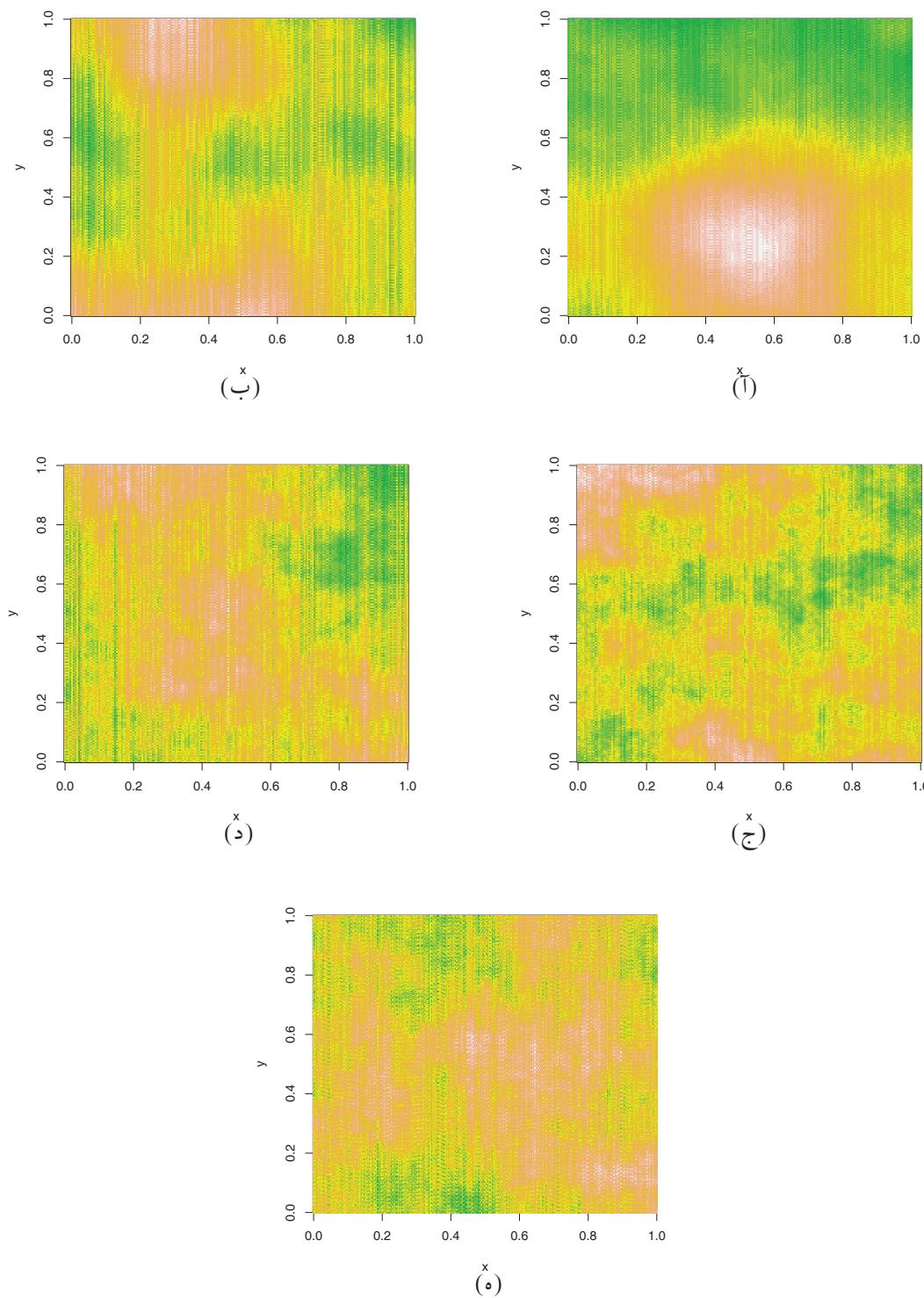
- سناریوی C : یک فرایند فضایی با همواری کم و دامنه ۱

$$g \sim N(0, R(1, 0/5)), \quad \delta \sim N(0, \sigma_\delta^2 I_{80}),$$

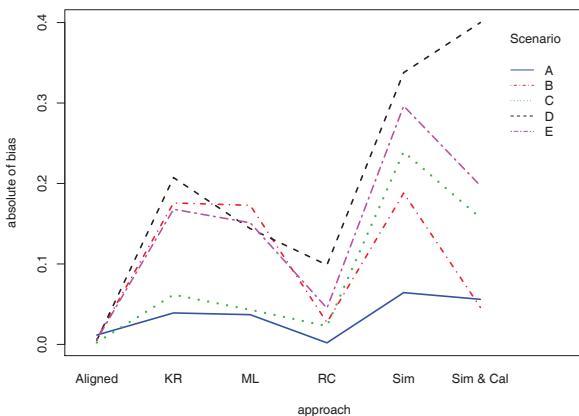
$$\sigma_\delta^2 = 0/22, \quad \sigma_\epsilon^2 = 0/82$$

سپس، برآوردهای ماکسیم درستمایی پارامترهای مدل رگرسیونی محاسبه می‌شوند.

با جایگذاری را، همانگونه که در بخش ۳.۲ شرح داده شد، کالیبره و در نهایت رگرسیون خطی \hat{Y} روی مقادیر کالیبره شده برآورد داده می‌شوند.

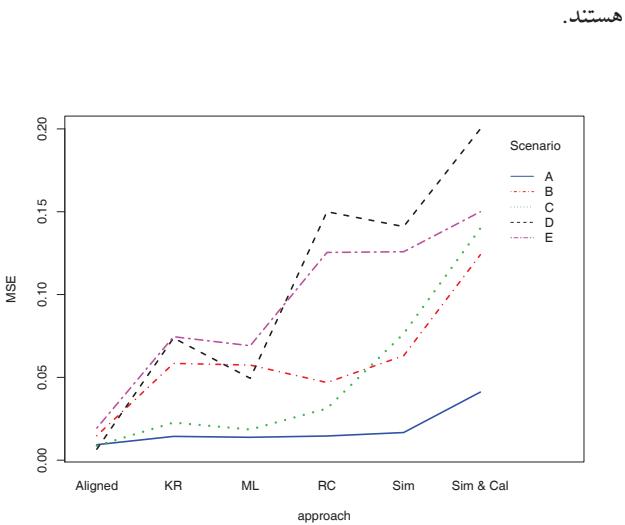


شکل ۳. تحقیقاتی از فرایند پیشگو (آ) سناریوی A ، (ب) سناریوی B ، (ج) سناریوی C ، (د) سناریوی D ، (ه) سناریوی E .



شکل ۵: قدر مطلق اریبی موجود در برآورد پارامتر شبیه‌سازی خط رگرسیون در رویکردهای مختلف

شکل ۶ مقدار MSE شبیه‌سازی خط رگرسیون را برای رویکردهای مختلف نشان می‌دهد. ملاحظه می‌شود که در سناریوی *A*، MSE تمامی رویکردها به جز رویکرد *Sim&Cal* تقریباً یکسان هستند. اما در فرایندهای ناهموارتر، MSE روش‌های مختلف با یکدیگر متفاوت هستند.



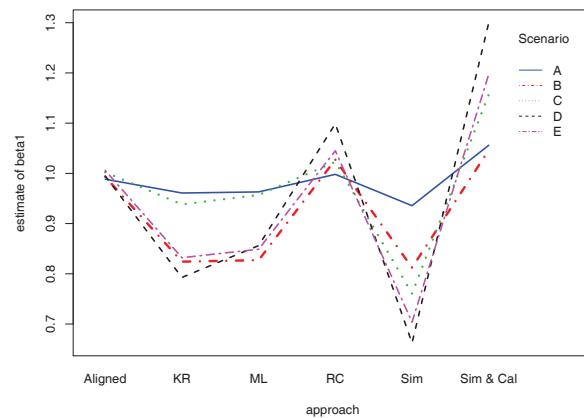
شکل ۶: $MSE(\hat{\beta}_1)$ در رویکردهای مختلف

بهطور کلی، روش *Sim&Cal* بیشترین MSE و روش *ML* کمترین MSE را دارد. بعلاوه مشاهده می‌شود که روش *RC* بیشتر از روش *KR* است. اما همانطور که در بخش‌های قبل توضیح داده شد و در این مطالعه شبیه‌سازی نیز مشاهده شد، روش *KR* برآوردهای اریب ایجاد می‌کند و از این روش *RC* نسبت به این روش برتری دارد.

نتایج حاصل از ۵۰۰ شبیه‌سازی در جدول ۶ خلاصه شده است. در این جدول، برآورد شبیه‌سازی خط رگرسیون، برآورد اریبی، متوسط انحراف معیار، میانگین توانهای دوم خط و میزان پوشش بازه‌های اطمینان ۹۵٪ گزارش شده است.

اکنون برای بررسی این نتایج از نمودارهای زیر استفاده می‌کنیم. توجه داریم که در این نمودارها مقادیر مربوط به رگرسیون خطی مجموعه داده‌های کامل (تراز) را با *Aligned*، رویکرد کریگ و رگرس را با *KR*، رویکرد رگرسیون کالبیدنی را با *RC*، رویکرد ماکسیمم درستنمایی را با *ML*، رویکرد شبیه‌سازی متغیر توضیحی را با *Sim* و رویکرد شبیه‌سازی و کالبیدن را با *Sim&Cal* نشان داده‌ایم.

شکل ۴ برآوردهای β_1 را تحت رویکردهای مختلف نشان می‌دهد. ملاحظه می‌شود که روش‌های رگرسیون کالبیدنی و شبیه‌سازی و کالبیدن، برای β_1 بیش برآورده و رویکردهای کریگ و رگرس، ماکسیمم درستنمایی و شبیه‌سازی کم برآورده دارند. از طرفی هرچه همواری فرایند کمتر می‌شود، انحراف برآوردها از مقدار واقعی پارامتر یعنی ۱ بیشتر می‌شود.

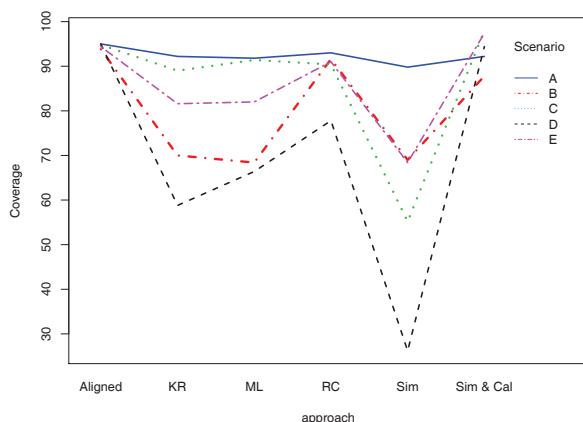


شکل ۴: برآورد شبیه‌سازی خط رگرسیون تحت رویکردهای مختلف

برای بررسی بهتر این موضوع، شکل ۵ را ببینید. در این شکل، قدر مطلق اریبی موجود در برآورد β_1 تحت رویکردهای مختلف ارائه شده است. مشاهده می‌شود به ویژه هنگامی که فرایند ناهموار می‌شود، اریبی روش‌های *KR*، *ML* و *Sim* به شدت افزایش می‌یابد. در مقایسه، روش *RC* در تمامی سناریوها اریبی بسیار کمی ایجاد کرده است و به عبارت دیگر اریبی موجود در روش *KR* را به خوبی تصحیح کرده است. همچنین مشاهده می‌شود که روش *Sim&Cal* اریبی روش *Sim* را به مقدار قابل توجهی تعديل کرده است.

جدول ۲: نتایج مطالعه شبیه‌سازی برای $\hat{\beta}_1$

پوشش (%)	MSE	انحراف معیار	اریبی	$\hat{\beta}_1$	رویکرد	سناریو
۹۵/۰	۰/۰۰۹	۰/۰۹۵	-۰/۰۱۲	۰/۹۸۸	داده‌های تراز	A
۹۲/۲	۰/۰۱۴	۰/۱۰۰	-۰/۰۳۹	۰/۹۶۱	کریگ و رگرس	
۹۳/۰	۰/۰۱۵	۰/۰۹۷	-۰/۰۰۲	۰/۹۹۸	رگرسیون کالبیدنی	
۹۱/۸	۰/۰۱۴	۰/۱۰۰	-۰/۰۳۷	۰/۹۶۳	ماکسیمم درستنایی	
۸۹/۸	۰/۰۱۷	۰/۱۰۳	-۰/۰۶۴	۰/۹۳۶	شبیه‌سازی	
۹۲/۲	۰/۰۴۱	۰/۱۵۲	۰/۰۵۶	۱/۰۵۶	شبیه‌سازی و کالبیدن	
۹۳/۸	۰/۰۱۵	۰/۱۱۵	-۰/۰۰۷	۰/۹۹۳	داده‌های تراز	B
۷۰/۰	۰/۰۵۸	۰/۱۲۶	-۰/۱۷۶	۰/۸۲۴	کریگ و رگرس	
۹۱/۶	۰/۰۴۶	۰/۱۴۵	۰/۰۲۷	۱/۰۲۷	رگرسیون کالبیدنی	
۶۸/۴	۰/۰۵۷	۰/۱۲۰	-۰/۱۷۳	۰/۸۲۷	ماکسیمم درستنایی	
۶۹/۰	۰/۰۶۳	۰/۱۲۷	-۰/۱۸۸	۰/۸۱۲	شبیه‌سازی	
۸۷/۸	۰/۱۲۴	۰/۴۱۸	۰/۰۴۶	۱/۰۴۶	شبیه‌سازی و کالبیدن	
۹۴/۸	۰/۰۰۹	۰/۰۹۲	۰/۰۰۲	۱/۰۰۲	داده‌های تراز	C
۸۹/۰	۰/۰۲۲۸	۰/۱۱۸	-۰/۰۶۱۸	۰/۹۳۸	کریگ و رگرس	
۹۰/۴	۰/۰۳۱	۰/۱۱۷	۰/۰۲۲	۱/۰۲۲	رگرسیون کالبیدنی	
۹۱/۴	۰/۰۱۹	۰/۱۱۷	-۰/۰۴۳	۰/۹۵۶	ماکسیمم درستنایی	
۵۵/۲	۰/۰۷۶	۰/۱۲۶	-۰/۲۳۹	۰/۷۶۱	شبیه‌سازی	
۹۸/۰	۱/۱۷۹	۱/۷۵۱	۰/۱۵۷	۱/۱۵۷	شبیه‌سازی و کالبیدن	
۹۵/۰	۰/۰۰۷	۰/۰۸۰	-۰/۰۰۵	۰/۹۹۵	داده‌های تراز	D
۵۸/۸	۰/۰۷۴	۰/۱۱۹	-۰/۲۰۷	۰/۷۹۳	کریگ و رگرس	
۷۷/۸	۰/۴۰۹	۰/۱۵۰	۰/۰۹۹	۱/۰۹۹	رگرسیون کالبیدنی	
۶۶/۴	۰/۰۴۹	۰/۱۱۱	-۰/۱۴۴	۰/۸۵۶	ماکسیمم درستنایی	
۲۶/۲	۰/۱۴۱	۰/۱۲۵	-۰/۳۳۸	۰/۶۶۲	شبیه‌سازی	
۹۴/۴	۴۱/۰	۵/۰۰۹	۰/۴۶۹	۱/۴۶۹	شبیه‌سازی و کالبیدن	
۹۴/۴	۰/۰۱۹	۰/۱۴۳	۰/۰۰۶	۱/۰۰۶	داده‌های تراز	E
۸۵/۸	۰/۰۷۰	۰/۱۸۲	-۰/۱۵۳	۰/۸۴۷	کریگ و رگرس	
۹۱/۸	۰/۱۰۹	۰/۲۲۱	۰/۰۶۸	۱/۰۶۸	رگرسیون کالبیدنی	
۴۶/۶	۰/۱۱۸	۰/۱۴۳	-۰/۳۰۴	۰/۸۹۶	ماکسیمم درستنایی	
۶۹/۶	۰/۱۱۹	۰/۱۸۵	-۰/۲۸۲	۰/۷۱۸	شبیه‌سازی	
۹۶/۸	۴/۲۰۳	۲/۲۲۱	۰/۱۷۱	۱/۱۷۱	شبیه‌سازی و کالبیدن	



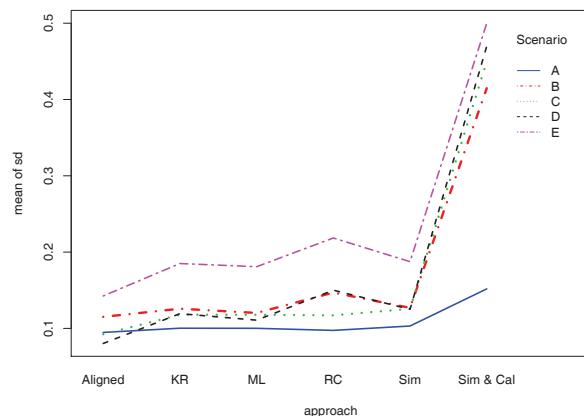
شکل ۸: درصد پوشش بازه‌های اطمینان ۹۵٪ در رویکردهای مختلف

بحث و نتیجه‌گیری

مسئله رگرسیون داده‌های بد تراز فضایی یکی از موضوعات مورد توجه محققان و تحلیل‌گران در بسیاری از زمینه‌های علوم است. در این مقاله رویکردهای مختلفی برای رگرسیون داده‌های بد تراز فضایی شامل رویکرد باجایگذاری، شبیه‌سازی، رگرسیون کالبیدنی و ماکسیمم درست‌نمایی ارائه شدند. در دو رویکرد باجایگذاری و شبیه‌سازی وجود خطای اندازه‌گیری منجر به ایجاد اریبی در برآوردهای پارامترهای مدل رگرسیونی شدند. لذا رویکرد رگرسیون کالبیدنی برای تعديل این اریبی پیشنهاد شد. نتایج مطالعه شبیه‌سازی نشان داد این رویکرد اگرچه MSE را مقداری افزایش می‌دهد اما اریبی موجود در روش‌های دیگر را تا حد زیادی تعديل می‌کند. همچنین در این روش، میزان پوشش اسمی بازه‌های اطمینان پارامتر شیب رگرسیون قابل توجه است.

لازم به ذکر است که این مقاله در راستای معرفی روش‌هایی برای کاهش اریبی موجود در برآوردهای ضرایب مدل رگرسیونی نگارش یافته است. ارزیابی پیشگویی توسط مدل‌های حاصل از این رویکردها از نکات شایان توجه است که توسط نگارندگان در حال بررسی است.

شکل ۷ میانگین انحراف معیار $\hat{\beta}_1$ را در هر سناریو برای هر رویکرد نشان می‌دهد. ملاحظه می‌شود که در تمامی سناریوها رویکرد *Sim&Cal* و پس از آن رویکرد *RC* بیشترین انحراف معیار را دارند. توجه داریم که میانگین انحراف معیار در روش *Sim&Cal* بسیار زیاد است اما میانگین انحراف معیار در روش *RC* تفاوت چندانی با سایر روش‌ها ندارد.



شکل ۷: میانگین انحراف معیار در رویکردهای مختلف

شکل ۸ درصد پوشش بازه‌های اطمینان ۹۵٪ را نشان می‌دهد. این درصد برای هر سناریو به این صورت محاسبه شده است که برای هر مجموعه داده، بازه اطمینانی به صورت $(\hat{\beta}_1 \pm 1/96sd)$ محاسبه شد و در نهایت ۵۰۰ بازه اطمینان بدست آمد. سپس بررسی شد که چند درصد از این بازه‌ها مقدار واقعی β_1 یعنی ۱ را در بر می‌گیرند. شکل ۸ نشان می‌دهد هنگامی که فرایند ناهموار می‌شود، درصد پوشش بازه‌های اطمینان ۹۵٪ در رویکردهای *KR*, *ML* و *Sim* به شدت کاهش می‌یابد. این درصد در رویکرد *Sim&Cal* بالا است و بدیهی است که دلیل آن بزرگی انحراف معیار این روش است. بازه‌های اطمینان ۹۵٪ در روش *RC* در تمامی سناریوها پوشش بالایی برای مقدار واقعی β_1 فراهم کرده‌اند و از این نظر نسبت به سایر روش‌های برآورد برتری دارد.

مراجع

- [1] Gryparis, A., C. J. Paciorek, A. Zeka, J. Schwartz, and B. Coull (2009). Measurement error caused by spatial misalignment in environmental epidemiology. *Biostatistics* **10**, 258–274.
- [2] Madsen, L., D. Ruppert, and N. S. Altman (2008). Regression with spatially misaligned data. *Environmetrics* **19**, 453–467.
- [3] Szpiro, A. A., L. Sheppard, and T. Lumley (2011). Efficient measurement error correction with spatially misaligned data. *Biostatistics* **0**, 1–14.
- [4] Young, L. J., C. A. Gotway, G. Kearney, and C. Duclos (2009). Assessing uncertainty in support-adjusted spatial misalignment problems. *Communications in Statistics— Theory and Methods* **38**, 3249–3264.
- [5] Zhu, L., B. Carlin, and A. Gelfand (2003). Hierarchical regression with misaligned spatial data: relating ambient ozone and pediatric asthma er visits in atlanta. *Environmetrics* **14**, 537–557.