

# تحلیل متغیرهای هیدرولوژیکی و خصوصیات خاک با روش‌های رگرسیون ریبج و بردار پشتیبان فضایی

مهدی روزبه<sup>۱\*</sup> و آرش عامری<sup>۲</sup>

<sup>۱</sup> گروه آمار، دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه سمنان، سمنان، ایران

<sup>۲</sup> گروه آمار، دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه سمنان، سمنان، ایران

تاریخ دریافت: ۱۴۰۴/۰۸/۲۱

تاریخ پذیرش: ۱۴۰۴/۱۰/۱۳

## چکیده

روش‌های سنتی آماری با گسترش داده‌های فضایی با ساختار مکانی-زمانی پیچیده با چالش‌های جدی مواجه شده‌اند. این داده‌ها به دلیل خودهمبستگی مکانی، ناهمسانی واریانس و وابستگی‌های جغرافیایی پیچیده، نیازمند روش‌های تخصصی هستند. در این پژوهش، رگرسیون بردار پشتیبان به‌عنوان یک روش نوین برای تحلیل و مدل‌سازی ساختار پیچیدی فضایی داده‌های زمین‌آماری مربوط به مقادیر کلسیم و منیزیم موجود در خاک معرفی می‌شود. این تحلیل بر اساس مختصات جغرافیایی مختلف (شرقی-غربی و شمالی-جنوبی)، در دو عمق ۰ تا ۲۰ و ۲۰ تا ۴۰ سانتی‌متر و در سه منطقه‌ی جغرافیایی متنوع انجام گرفته است. روش رگرسیون بردار پشتیبان با توانایی مدل‌سازی روابط غیرخطی پیچیده و در عین حال حفظ ساختار فضایی داده‌ها، امکان پیش‌بینی دقیق‌تر و واقع‌بینانه‌تر توزیع عناصر تغذیه‌ای خاک را فراهم می‌سازد. این رویکرد، با بهره‌گیری از توابع هسته، قابلیت تحلیل فضا‌های ویژگی با ابعاد بالا و پیچیدگی‌های ساختاری را فراهم می‌آورد و اثربخشی خود را در برابر جملات اختلال و داده‌های پرت اثبات می‌کند. برای سنجش دقیق کارایی رگرسیون بردار پشتیبان، عملکرد آن در برابر رگرسیون ریبج مورد ارزیابی مقایسه‌ای قرار می‌گیرد.

**واژه‌های کلیدی:** اعتبارسنجی متقابل تعمیم‌یافته، رگرسیون ریبج، رگرسیون بردار پشتیبان، ماشین‌های بردار پشتیبان  
رده‌بندی ریاضی: ۶۲H۲۵; ۶۲J۰۵.

## ۱ مقدمه

آن‌ها با چالش‌هایی مانند خودهمبستگی مکانی، ناهمسانی واریانس و ساختارهای وابستگی پیچیده همراه است. از آنجا که روش‌های آماری سنتی توانایی پاسخگویی کامل به این چالش‌ها را ندارند، توسعه و به‌کارگیری روش‌های نوین برای مدل‌سازی و تحلیل داده‌های فضایی ضرورت دارد. در این راستا، رگرسیون بردار پشتیبان<sup>۱</sup> (*SVR*) به‌عنوان یک روش قدرتمند در پردازش داده‌های مکانی مطرح است که توانایی

در عصر حاضر، پیشرفت‌های علمی و فناوری منجر به تولید حجم گسترده‌ای از داده‌های فضایی شده است که دارای ساختارهای پیچیده مکانی-زمانی هستند. این داده‌ها شامل اطلاعات جغرافیایی، موقعیت‌های مکانی و روابط میان پدیده‌ها در فضا می‌باشند و تحلیل

\*نویسنده مسئول مقاله، mahdi.roozbeh@semnan.ac.ir

<sup>۱</sup>Support Vectors Regression

<sup>۲</sup>Ridge Regression

جداگانه برای هر متغیر مستقل را فراهم می‌سازد که موجب بهبود بیشتر در دقت برآوردها می‌گردد. رگرسیون ریج به‌ویژه در تحلیل داده‌های فضایی که اغلب با چالش همخطی بین متغیرهای کمکی مواجه هستند، می‌تواند به‌عنوان یک روش کارآمد برای مدل‌سازی روابط پیچیده مورد استفاده قرار گیرد [۵].

## ۱.۲ جریمه در رگرسیون ریج

رگرسیون ریج با اعمال جریمه بر ضرایب رگرسیونی، آن‌ها را کاهش و مجموع توان‌های دوم باقیمانده‌ها را به حداقل می‌رساند. این روش، که به‌عنوان رگرسیون ریج شناخته می‌شود، پارامترها را به سمت صفر هدایت کرده و واریانس برآوردها را کاهش می‌دهد. بنابراین، می‌توان تعریف این رویکرد را مطابق رابطه (۱) به‌صورت زیر ارائه کرد:

$$\hat{\beta}^{Ridge} = \arg \min_{\beta} \left[ \sum_{i=1}^N \left( Y_i - \beta_0 - \sum_{j=1}^P X_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^P \beta_j^2 \right] \quad (1)$$

در رابطه (۱)  $\lambda \geq 0$  به‌عنوان پارامتر جریمه تعریف می‌شود که باید به‌طور جداگانه تعیین شود و بخش دوم معادله  $\lambda \sum_{j=1}^P \beta_j^2$  که به‌عنوان جریمه انقباض شناخته می‌شود، زمانی که ضرایب  $\beta_1, \beta_2, \dots, \beta_p$  نزدیک به صفر هستند، کوچک می‌شود و باعث کاهش برآوردهای  $\beta_j$  به سمت صفر می‌گردد.

همچنین پارامتر تنظیم‌کننده  $\lambda$  به کنترل تأثیر نسبی این دو بخش بر برآوردهای ضرایب رگرسیون کمک می‌کند. زمانی که  $\lambda = 0$  باشد، جریمه اثری ندارد و رگرسیون ریج مشابه برآوردهای کمترین توان‌های دوم معمولی خواهد بود. اما با افزایش  $\lambda$  به سمت بی‌نهایت، تأثیر جریمه انقباض بیشتر می‌شود و برآوردهای ضرایب رگرسیون ریج به سمت صفر سوق پیدا می‌کنند. همچنین می‌توان برآورد رگرسیون ریج را به شکل ماتریسی طبق رابطه (۲) نوشت:

$$\hat{\beta}^{Ridge} = \hat{\beta}(k) = (X^T X + \lambda I_p)^{-1} X^T Y, \quad \lambda \geq 0 \quad (2)$$

که  $\hat{\beta}^{Ridge}$  به‌عنوان برآوردگر ریج شناخته می‌شود. برخلاف بعضی از روش‌های رگرسیونی، در روش رگرسیون ریج، ماتریس  $X^T X$  دارای وارون است [۶].

مدل‌سازی روابط غیرخطی پیچیده در کنار حفظ ساختار فضایی داده‌ها را داراست [۹]. علاوه بر این، رگرسیون ریج<sup>۲</sup> نیز به‌عنوان یک روش مقایسه‌ای مورد استفاده قرار می‌گیرد. مزیت اصلی الگوریتم ماشین بردار پشتیبان در تحلیل‌های فضایی، قابلیت ترکیب اطلاعات مکانی با متغیرهای کمکی، مقاومت در برابر داده‌های پرت و اختلال، و توانایی پردازش فضاهای ویژگی پیچیده با استفاده از توابع هسته است. این ویژگی‌ها موجب می‌شود که این الگوریتم بتواند وابستگی‌های فضایی را بدون نیاز به فرضیات سخت‌گیرانه مدل‌سازی کند و در مسائل کاربردی نتایج دقیق‌تری ارائه دهد.

آمار فضایی شاخه‌ای از علم آمار است که به تحلیل داده‌هایی می‌پردازد که وابستگی مکانی یا ساختار جغرافیایی دارند. در این نوع داده‌ها، مقدار مشاهده شده در یک موقعیت نه‌تنها به ویژگی‌های همان نقطه بلکه به مقادیر نقاط مجاور نیز وابسته است. هدف اصلی آمار فضایی، شناسایی و مدل‌سازی الگوهای مکانی، بررسی روابط بین پدیده‌ها در فضا و پیش‌بینی رفتار آن‌ها در موقعیت‌های ناشناخته است. این حوزه با استفاده از ابزارهایی مانند مدل‌های خودرگرسیون فضایی، آزمون‌های وابستگی مکانی و روش‌های خوشه‌بندی فضایی، امکان درک بهتر ساختارهای جغرافیایی و تصمیم‌گیری علمی در زمینه‌هایی همچون محیط زیست، شهرسازی، سلامت عمومی و علوم زمین را فراهم می‌آورد [۱].

## ۲ رگرسیون ریج

رگرسیون ریج یکی از روش‌های تحلیل رگرسیونی است که برای مقابله با مشکل همخطی چندگانه<sup>۳</sup> در داده‌ها معرفی شده است. این روش نخستین بار توسط هورل و کনার در سال ۱۹۷۰ ارائه شد و به‌عنوان جایگزینی برای روش کمترین توان‌های دوم معمولی (OLS)<sup>۴</sup> در شرایطی که متغیرهای مستقل همبستگی بالایی دارند یا زمانی که تعداد متغیرها ( $p$ ) نزدیک به تعداد مشاهدات ( $n$ ) باشد، به کار می‌رود. در رگرسیون ریج، با افزودن یک پارامتر تنظیم ( $\lambda$ ) به ماتریس واریانس-کوواریانس، واریانس برآوردها کاهش می‌یابد [۸]. این روش با قبول اندکی اریبی در برآوردها، واریانس را به میزان قابل توجهی کاهش داده و منجر به برآوردهای پایدارتری می‌شود. همچنین، توسعه این روش به‌صورت رگرسیون ریج تعمیم‌یافته<sup>۵</sup>، امکان در نظر گرفتن پارامتر اریب

<sup>3</sup>Multicollinearity

<sup>4</sup>Ordinary Least Squares

<sup>5</sup>Generalized Ridge Regression

پارامتر تنظیم ( $C$ ) در این میان نقش مهمی دارد، به طوری که مقادیر بزرگتر  $C$  خطاهای طبقه‌بندی را کاهش داده اما حاشیه کوچک‌تری ایجاد می‌کند، در حالی که مقادیر کوچک‌تر  $C$  به مدل اجازه خطای بیشتر داده و حاشیه بزرگ‌تری فراهم می‌آورد.

### ۳.۳ روش‌های هسته در ماشین‌های بردار پشتیبان

برای مسائل غیرخطی، از توابع هسته برای نگاشت داده‌ها به فضایی با ابعاد بالاتر استفاده می‌شود، جایی که داده‌ها به صورت خطی جداپذیر می‌گردند. برخی از پرکاربردترین توابع هسته عبارتند از [۲]:

$$1: \text{ هسته خطی: } K(\mathbf{x}, \mathbf{u}) = \mathbf{x}^T \cdot \mathbf{u}$$

$$2: \text{ هسته چندجمله‌ای: } K(\mathbf{x}, \mathbf{u}) = (a\mathbf{x}^T \mathbf{u} + c)^q$$

$$3: \text{ هسته سیگموئید: } K(\mathbf{x}, \mathbf{u}) = \tanh(\beta \mathbf{x}^T \mathbf{u} + \gamma)$$

$$4: \text{ هسته گاوسی (شعاعی): } K(\mathbf{x}, \mathbf{u}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{u}\|^2}{\sigma^2}\right)$$

### ۴.۳ رگرسیون بردار پشتیبان

رگرسیون بردار پشتیبان تعمیمی از ماشین‌های بردار پشتیبان به مسائل رگرسیون است که به جای دسته‌بندی، خروجی‌های پیوسته را برآورد می‌کند. در این روش، یک ناحیه حساسیت‌ناپذیر ( $\varepsilon$ ) حول تابع برآورد تعریف می‌شود که خطاهای کوچک‌تر از  $\varepsilon$  را نادیده می‌گیرد و تنها خطاهای خارج از این ناحیه جریمه می‌شوند. بدین ترتیب، بردارهای پشتیبان نقاطی هستند که خارج از این حاشیه قرار دارند و ساختار مدل را تعیین می‌کنند. فرمول‌بندی بهینه‌سازی  $SVR$  به شکل رابطه (۵) ارائه می‌شود:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \xi^*} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \\ \text{s.t.} \quad & y_i - \mathbf{w}^T \mathbf{x}_i - b \leq \varepsilon + \xi_i, \quad i = 1, \dots, n \\ & \mathbf{w}^T \mathbf{x}_i + b - y_i \leq \varepsilon + \xi_i^*, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad \xi_i^* \geq 0, \quad i = 1, \dots, n \end{aligned} \quad (5)$$

که در آن  $C$  پارامتر تنظیم تعادل میان پیچیدگی مدل و دقت پیش‌بینی و همچنین  $\xi_i$  و  $\xi_i^*$  متغیرهای خطا برای نقاط خارج از ناحیه  $\varepsilon$  هستند. این فرمول‌بندی با استفاده از توابع زیان مختلف (خطی، درجه دوم، هوبر) و همچنین رویکرد حاشیه نرم، امکان کنترل خطا و مقابله با نقاط پرت را فراهم می‌کند.

## ۳ ماشین‌های بردار پشتیبان و رگرسیون بردار پشتیبان

در ادامه به بررسی ماشین‌های بردار پشتیبان<sup>۶</sup> ( $SVM$ ) و رگرسیون بردار پشتیبان می‌پردازیم.

### ۱.۳ ماشین‌های بردار پشتیبان با حاشیه سخت

ماشین‌های بردار پشتیبان یکی از روش‌های قدرتمند یادگیری ماشین هستند که توسط ولادیمیر وپنیک و الکسی چروننکیس در دهه ۱۹۹۰ میلادی توسعه یافتند. ماشین‌های بردار پشتیبان براساس مفاهیم ابرصفحه (مرز تصمیم‌گیری)، بردارهای پشتیبان (نقاط کلیدی نزدیک به مرز تصمیم‌گیری) و حاشیه (فاصله بین ابرصفحه و نزدیک‌ترین نقاط داده) عمل می‌کنند. هدف اصلی  $SVM$ ، بیشینه‌سازی حاشیه برای دستیابی به مدلی با تعمیم‌پذیری بهتر است [۴، ۷].

هنگامی که داده‌ها به طور خطی جداپذیر باشند، از حاشیه سخت استفاده می‌شود. در این حالت،  $SVM$  به دنبال یافتن ابرصفحه‌ای است که دو کلاس  $y_i \in \{-1, +1\}$  را با حداکثر فاصله از یکدیگر جدا کند، به طوری که هیچ نقطه‌ای درون حاشیه قرار نگیرد. این مسئله به صورت یک مسئله بهینه‌سازی برای مینیمم‌سازی تابع هدف موجود در رابطه (۳) فرمول‌بندی می‌شود [۲]:

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, n \end{aligned} \quad (3)$$

### ۲.۳ ماشین‌های بردار پشتیبان با حاشیه نرم

برای داده‌هایی که به طور خطی جداپذیر نیستند، حاشیه نرم معرفی شده است. در این روش با معرفی متغیرهای آستانه خطا ( $\xi_i$ )، به مدل اجازه داده می‌شود تا در طبقه‌بندی برخی از نقاط نزدیک به مرز مقداری خطا وجود داشته باشد. مسئله بهینه‌سازی مرتبط با مینیمم‌سازی تابع هدف در این شرایط مطابق رابطه (۴) تعریف می‌شود و پارامترها با توجه به این محدودیت‌ها تنظیم می‌گردند [۲]:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & J(\mathbf{w}, b, \xi) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad i = 1, \dots, n \end{aligned} \quad (4)$$

<sup>6</sup>Support Vector Machines

## ۴ پیاده‌سازی و تحلیل نتایج

خط کاهش یافته و پس از نقطه‌ای خاص دوباره افزایش می‌یابد که بیانگر وجود مقدار بهینه‌ای از  $\lambda$  برای هر منطقه است. به‌طور کلی، نتایج جدول ۱ و شکل ۱ نشان می‌دهند که تنظیم مناسب  $\lambda$  نقش مهمی در بهبود عملکرد مدل ریح دارد و حساسیت مدل نسبت به مقدار این پارامتر در مناطق مختلف متفاوت است.

بر اساس نتایج حاصل از ارزیابی مدل‌ها در سه منطقه مختلف موجود در جداول ۲ تا ۴، می‌توان دریافت که عملکرد بهینه مدل‌های رگرسیون به شدت تحت تأثیر شرایط خاص هر منطقه قرار دارد. در منطقه ۱، مدل بردار پشتیبان با هسته شعاعی با کسب ضریب تعیین  $0.86$  برای متغیر  $ca_{0.20}$  به‌عنوان بهترین مدل شناسایی شد. در منطقه ۲، همین مدل با ضریب تعیین  $0.42$  برای متغیر  $ca_{0.20}$  و مدل رگرسیون ریح با کمترین مقادیر خطا برای متغیر  $ca_{0.40}$  مناسب‌ترین نتایج را ارائه دادند. در منطقه ۳، مدل‌های هسته شعاعی و چندجمله‌ای به‌ترتیب با ضرایب تعیین  $0.51$  و  $0.37$  برای متغیرهای  $ca_{0.20}$  و  $ca_{0.40}$  بهترین عملکرد را از خود نشان دادند. در مقابل، مدل مبتنی بر هسته سیگموئید در تمامی مناطق و برای هر دو متغیر ضعیف‌ترین نتایج را ثبت کرد. این یافته‌ها به وضوح نشان می‌دهد که انتخاب مدل بهینه برای پیش‌بینی محتوای کلسیم خاک نیازمند در نظر گرفتن همزمان موقعیت جغرافیایی و عمق نمونه‌برداری است. نتایج این مطالعه نشان می‌دهد که تغییرپذیری مکانی خصوصیات خاک تأثیر قابل توجهی بر دقت مدل‌های پیش‌بینی دارد. به‌طور مشخص، مدل بردار پشتیبان با هسته شعاعی به دلیل توانایی بالا در مدل‌سازی روابط غیرخطی پیچیده بین متغیرهای مؤثر بر محتوای کلسیم خاک، در اکثر مناطق عملکرد مطلوبی از خود نشان داده است. همچنین مشاهده شد که با افزایش عمق نمونه‌برداری، دقت پیش‌بینی کلیه مدل‌ها کاهش می‌یابد که می‌تواند ناشی از پیچیدگی بیشتر روابط بین متغیرها در لایه‌های زیرین خاک باشد. تفاوت در عملکرد مدل‌ها در مناطق مختلف را می‌توان به عواملی همچون تغییرات در بافت خاک، شرایط اقلیمی، کاربری اراضی و مدیریت خاک در هر منطقه نسبت داد. این تنوع در نتایج بر اهمیت توسعه مدل‌های منطقه‌ای و عدم تعمیم یک مدل واحد به تمامی مناطق تأکید دارد. علاوه بر این، مقایسه معیارهای ارزیابی مختلف شامل  $R^2$ ،  $MSE$  و  $MAPE$  نشان داد که هر یک از این معیارها می‌توانند جنبه‌های متفاوتی از عملکرد مدل را مورد سنجش قرار دهند.

در این پژوهش از مجموعه داده *camg* موجود در بسته *geoR* استفاده شده است. این داده شامل اندازه‌گیری مقادیر کلسیم و منیزیم خاک در دو عمق ۲۰ تا ۲۰ و ۴۰ سانتی‌متر در ۱۷۸ موقعیت مکانی مختلف است. عناصر کلسیم و منیزیم از عناصر مغذی ضروری برای رشد گیاهان هستند و بررسی توزیع فضایی آن‌ها به درک بهتر حاصل‌خیزی خاک کمک می‌کند. در هر موقعیت، مختصات جغرافیایی (شرقی-غربی و شمالی-جنوبی)، ارتفاع زمین و شماره منطقه ثبت شده است. منطقه مورد مطالعه به سه بخش تقسیم می‌شود: منطقه اول که معمولاً در فصل بارندگی دچار سیلاب شده و دارای خاک طبیعی است، منطقه دوم که در گذشته کوددهی شده و عمدتاً شالیزار برنج است و منطقه سوم که به‌تازگی کوددهی شده و بیشتر برای آزمایش‌های کشاورزی به‌کار می‌رود. ساختار داده شامل یک *Data frame* با ۱۰ متغیر شامل مختصات مکانی<sup>۷</sup>، ارتفاع<sup>۸</sup>، شماره منطقه<sup>۹</sup>، مقادیر کلسیم، منیزیم و ظرفیت تبادل کاتیونی در دو عمق خاک ذکر شده است [۳]. این مجموعه داده برای تحلیل‌های زمین‌آماري با استفاده از نیم‌تغییرنگار<sup>۱۰</sup> و کریکینگ معمولی<sup>۱۱</sup> در بسته *geoR* به‌کار می‌رود و امکان بررسی وابستگی فضایی و پیش‌بینی مقادیر کلسیم و منیزیم در نقاط نمونه‌برداری نشده را فراهم می‌کند. پس از تقسیم داده‌ها به نسبت ۷۰ به ۳۰ برای آموزش و آزمون، عملکرد مدل‌های رگرسیون بردار پشتیبان (با چهار هسته مختلف) و رگرسیون ریح با استفاده از معیارهای  $R^2$ ،  $MSE$ ،  $RMSE$  و  $MAPE$  ارزیابی شد. پارامتر بهینه مدل ریح برای مناطق مختلف و متغیرهای  $ca_{0.20}$  و  $ca_{0.40}$  با اعتبارسنجی متقابل تعیین و در جدول ۱ ارائه شده است.

جدول ۱: مقادیر بهینه پارامتر ریح ( $\lambda$ )

منطقه	کلسیم در لایه‌های مختلف	
	$ca_{0.20}$	$ca_{0.40}$
منطقه ۱	۱۰۱۱۶۹۲	۱۹۱۷۹۸۱
منطقه ۲	۱۸۷۵۲۳۸	۲۰۷۱۷۷۱
منطقه ۳	۵۲۶۰۲۶۱	۵۲۳۹۵۳۴

شکل ۱ نتایج اعتبارسنجی متقابل برای رگرسیون ریح را در سه منطقه و دو متغیر نشان می‌دهد. در هر نمودار، محور افقی لگاریتم پارامتر ریح  $\lambda$  و محور عمودی میانگین توان‌های دوم خطا ( $MSE$ ) است. همان‌طور که منحنی‌های قرمز نشان می‌دهند، با افزایش  $\lambda$  ابتدا مقدار

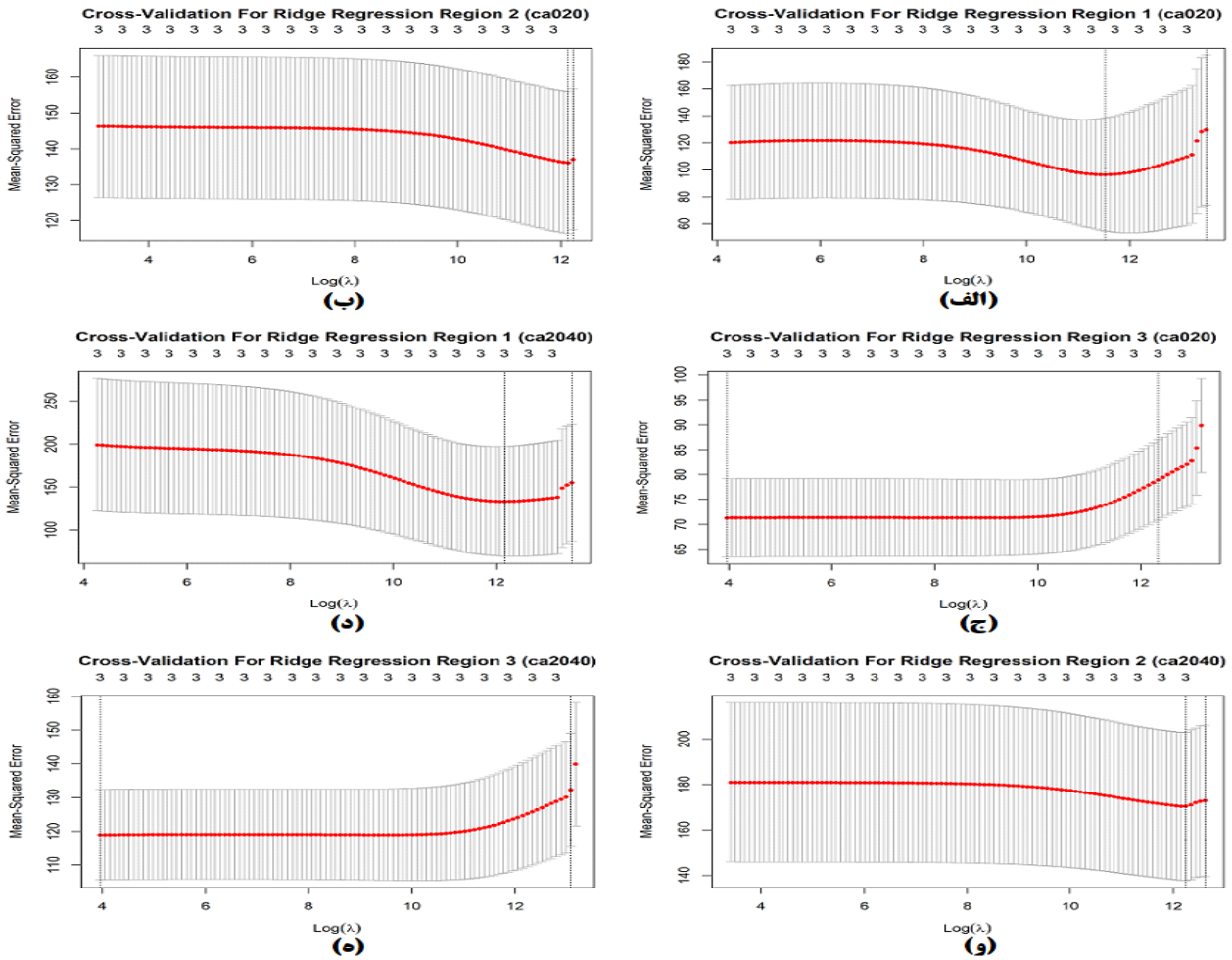
<sup>7</sup>East & North

<sup>8</sup>Elevation

<sup>9</sup>Region

<sup>10</sup>Semivariogram

<sup>11</sup>Ordinary Kriging



شکل ۱: نمودار اعتبارسنجی متقابل برای تعیین پارامتر بهینه  $\lambda$  در مناطق مختلف: الف) منطقه ۱ در  $ca_{20}$ ، ب) منطقه ۲ در  $ca_{20}$ ، ج) منطقه ۳ در  $ca_{20}$ ، د) منطقه ۱ در  $ca_{2040}$ ، و) منطقه ۲ در  $ca_{2040}$ ، ه) منطقه ۳ در  $ca_{2040}$ .

جدول ۲: مقایسه مدل‌های رگرسیون بر اساس معیارهای ارزیابی در منطقه ۱

معیارها								مدل رگرسیون
ca <sub>2040</sub> در منطقه ۱				ca <sub>20</sub> در منطقه ۱				
MAPE	RMSE	MSE	R <sup>2</sup>	MAPE	RMSE	MSE	R <sup>2</sup>	
۰/۴۸	۱۲/۷۶	۱۶۲/۸۴	۰/۰۰۲	۰/۲۳	۷/۳۷	۵۴/۲۸	۰/۸۲	بردار پشتیبان با هسته خطی
۰/۳۹	۱۰/۳۰	۱۰۶/۱۵	۰/۰۰۴	۰/۲۶	۸/۵۰	۷۲/۲۰	۰/۱۹	بردار پشتیبان با هسته چندجمله‌ای
۰/۳۷	۹/۴۶	۸۹/۴۲	۰/۰۲۲	۰/۱۶	۵/۵۴	۳۰/۶۹	۰/۸۶	بردار پشتیبان با هسته شعاعی
۰/۴۰	۱۰/۲۵	۱۰۵/۱۲	۰/۰۰۰۷	۰/۲۱	۶/۹۵	۴۸/۲۸	۰/۷۷	بردار پشتیبان با هسته سیگموئید
۰/۴۷	۱۲/۱۷	۱۴۸/۰۴	۰/۰۰۹	۰/۲۷	۹/۰۵	۸۱/۹۶	۰/۴۶	رگرسیون ریج

جدول ۳: مقایسه مدل‌های رگرسیون بر اساس معیارهای ارزیابی در منطقه ۲

معیارها								مدل رگرسیون
ca <sup>۲۰۴۰</sup> در منطقه ۲				ca <sup>۲۰۲۰</sup> در منطقه ۲				
MAPE	RMSE	MSE	R <sup>۲</sup>	MAPE	RMSE	MSE	R <sup>۲</sup>	
۰٫۳۱	۱۱٫۶۷	۱۳۶٫۲۵	۰٫۵۴	۰٫۱۳	۶٫۱۲	۳۷٫۴۳	۰٫۱۷	بردار پشتیبان با هسته خطی
۰٫۳۶	۱۳٫۳۹	۱۷۹٫۲۷	۰٫۵۷	۰٫۱۴	۷٫۰۰	۴۸٫۹۶	۰٫۰۰۰۲	بردار پشتیبان با هسته چندجمله‌ای
۰٫۳۳	۱۲٫۰۱	۱۴۴٫۲۸	۰٫۱۱	۰٫۱۳	۵٫۹۷	۳۵٫۶۶	۰٫۴۲	بردار پشتیبان با هسته شعاعی
۰٫۲۶	۱۰٫۶۷	۱۱۳٫۸۴	۰٫۵۰۴	۰٫۲۰	۱۰٫۳۴	۱۰۶٫۹۴	۰٫۵۳	بردار پشتیبان با هسته سیگموئید
۰٫۳۰	۱۰٫۱۹	۱۰۳٫۹۳	۰٫۵۲	۰٫۱۳	۶٫۶۸	۴۴٫۶۳	۰٫۵۲	رگرسیون ریح

جدول ۴: مقایسه مدل‌های رگرسیون بر اساس معیارهای ارزیابی در منطقه ۳

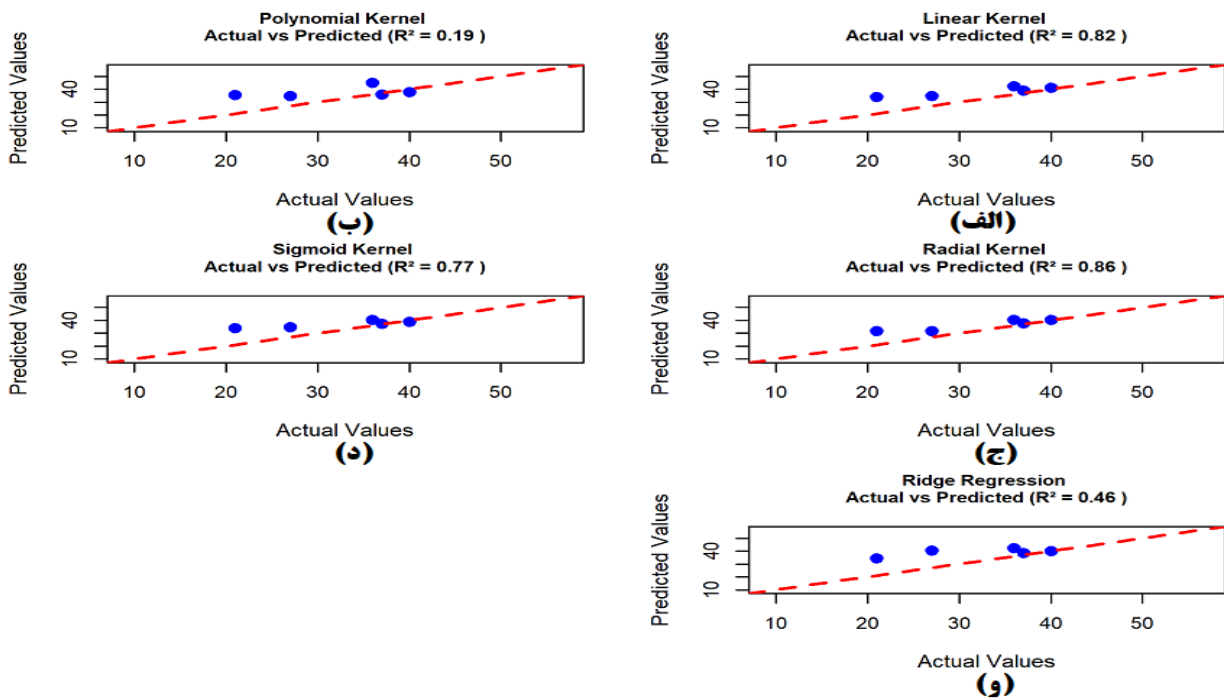
معیارها								مدل رگرسیون
ca <sup>۲۰۴۰</sup> در منطقه ۳				ca <sup>۲۰۲۰</sup> در منطقه ۳				
MAPE	RMSE	MSE	R <sup>۲</sup>	MAPE	RMSE	MSE	R <sup>۲</sup>	
۰٫۱۸	۱۰٫۷۹	۱۱۶٫۴۸	۰٫۲۳	۰٫۱۱	۷٫۷۱	۵۹٫۴۹	۰٫۲۸	بردار پشتیبان با هسته خطی
۰٫۱۷	۹٫۷۸	۹۵٫۵۸	۰٫۳۷	۰٫۱۰	۶٫۶۷	۴۴٫۴۷	۰٫۴۶	بردار پشتیبان با هسته چندجمله‌ای
۰٫۱۷۲	۹٫۹۹	۹۹٫۷۶	۰٫۳۶	۰٫۰۹	۶٫۶۶	۴۴٫۳۶	۰٫۵۱	بردار پشتیبان با هسته شعاعی
۰٫۴۳	۲۷٫۴۸	۷۵۵٫۴۵	۰٫۵۲	۰٫۳۲	۲۳٫۶۲	۵۵۷٫۷۳	۰٫۰۹	بردار پشتیبان با هسته سیگموئید
۰٫۱۹	۱۰٫۹۳	۱۱۹٫۴۹	۰٫۲۰	۰٫۱۲	۸٫۱۱	۶۵٫۸۴	۰٫۲۰	رگرسیون ریح

در منطقه ۲ و برای متغیر ca<sup>۲۰۲۰</sup> (شکل ۴)، مدل با هسته شعاعی با ضریب تعیین ۰٫۴۲ بهترین عملکرد را نشان داد و پس از آن مدل با هسته خطی با مقدار ۰٫۱۷ در رتبه دوم قرار گرفت. سایر مدل‌ها از جمله سیگموئید (۰٫۵۳)، ریح (۰٫۲۰) و چندجمله‌ای (۰) عملکرد قابل توجهی نداشتند. همچنین برای متغیر ca<sup>۲۰۴۰</sup> در همین منطقه (شکل ۵)، مدل‌های خطی (۰٫۵۴) و چندجمله‌ای (۰٫۵۷) نسبت به سایر مدل‌ها نتایج بهتری ارائه دادند، هرچند دقت کلی پایین ارزیابی شد.

در منطقه ۳ برای متغیر ca<sup>۲۰۲۰</sup> (شکل ۶)، مدل با هسته شعاعی با ضریب تعیین ۰٫۵۱ بهترین عملکرد را داشت و پس از آن مدل‌های چندجمله‌ای (۰٫۴۶) و خطی (۰٫۲۸) قرار گرفتند، مدل‌های ریح (۰٫۲) و سیگموئید (۰٫۰۹) عملکرد ضعیفی داشتند. برای متغیر ca<sup>۲۰۴۰</sup> (شکل ۷) نیز مدل چندجمله‌ای (۰٫۳۷) برتر بود و مدل‌های شعاعی (۰٫۳۶) و خطی (۰٫۲۳) در رتبه‌های بعدی قرار گرفتند، در حالی که ریح (۰٫۲) و سیگموئید (۰٫۵۲) نتایج مطلوبی ارائه نکردند.

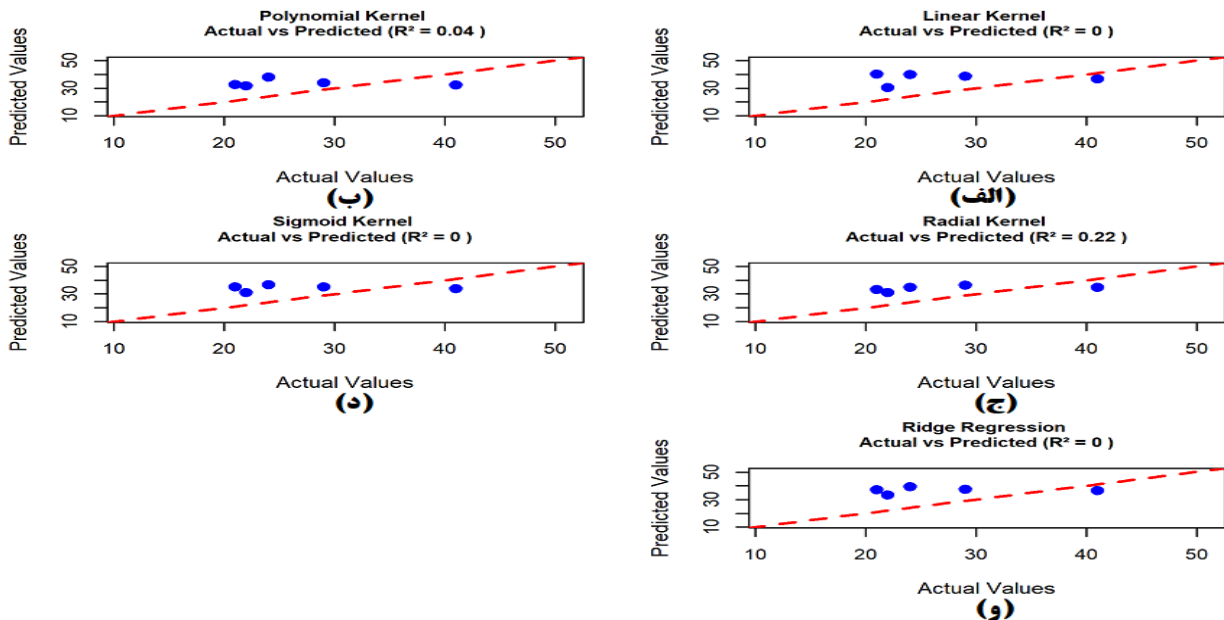
برای ارزیابی دقیق عملکرد مدل‌های رگرسیون در پیش‌بینی محتوای کلسیم خاک، نمودارهای مقادیر واقعی در مقابل مقادیر پیش‌بینی شده در سه منطقه و برای دو عمق نمونه‌برداری (ca<sup>۲۰۲۰</sup> و ca<sup>۲۰۴۰</sup>) در مجموعه نمودارهای موجود در شکل‌های ۲ تا ۷ مورد تحلیل قرار گرفت. در مجموعه نمودارهای موجود در شکل ۲ مربوط به منطقه ۱ و برای متغیر ca<sup>۲۰۲۰</sup>، مدل بردار پشتیبان با هسته شعاعی با ضریب تعیین ۰٫۸۶ بهترین عملکرد را نشان داد که بیانگر تطابق بسیار مطلوب بین مقادیر واقعی و پیش‌بینی شده و توزیع فشرده نقاط حول خط  $Y = X$  است. پس از آن، مدل‌های با هسته خطی (۰٫۸۲) و سیگموئید (۰٫۷۷) در رده‌های بعدی قرار گرفتند. در مقابل، مدل رگرسیون ریح (۰٫۴۶) و مدل با هسته چندجمله‌ای (۰٫۱۹) عملکرد ضعیف‌تری داشتند. برای متغیر ca<sup>۲۰۴۰</sup> در همین منطقه (شکل ۳)، اگرچه مقادیر ضریب تعیین به‌طور کلی کاهش یافت، مدل با هسته شعاعی (۰٫۲۲) و چندجمله‌ای (۰٫۴) بهتر از سایر مدل‌ها عمل کردند، در حالی که مدل‌های خطی، سیگموئید و ریح تقریباً فاقد توان پیش‌بینی معنادار بودند.

**Actual vs Predicted Values In Region 1 [ca020]**

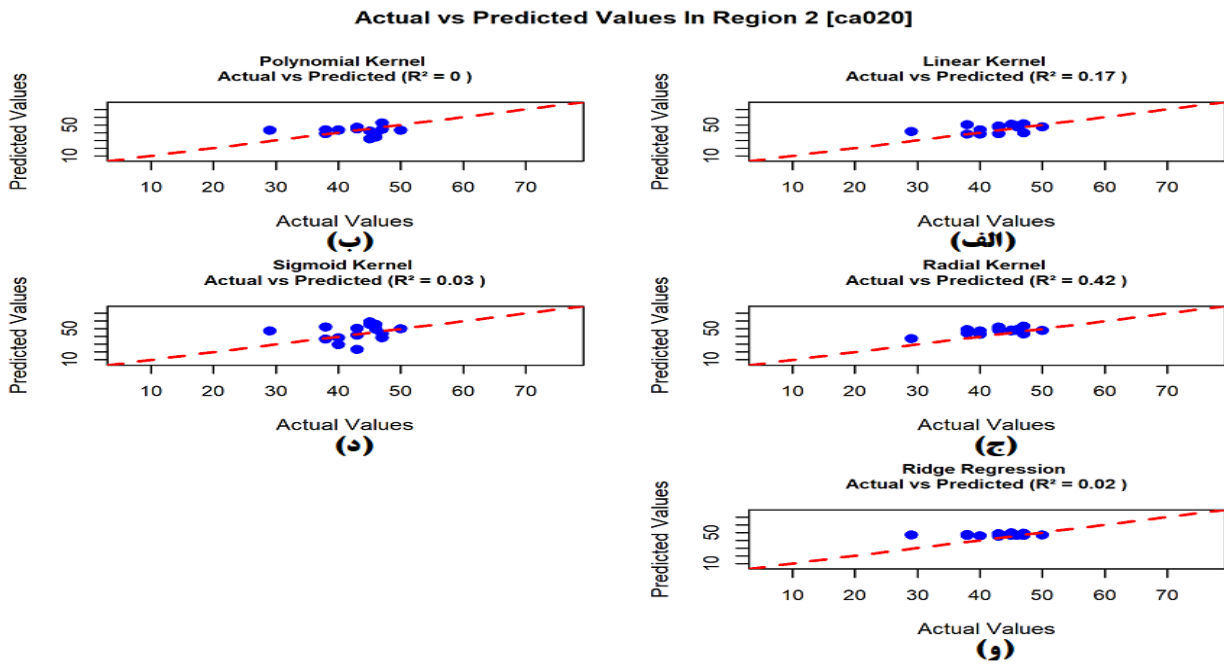


شکل ۲: نمودار مقادیر واقعی در برابر مقادیر پیش‌بینی شده برای چهار هسته بردار پشتیبان و رگرسیون ریج در منطقه ۱ در  $ca=020$ :  
 (الف) هسته خطی، (ب) هسته چندجمله‌ای، (ج) هسته شعاعی، (د) هسته سیگموئید، (و) رگرسیون ریج.

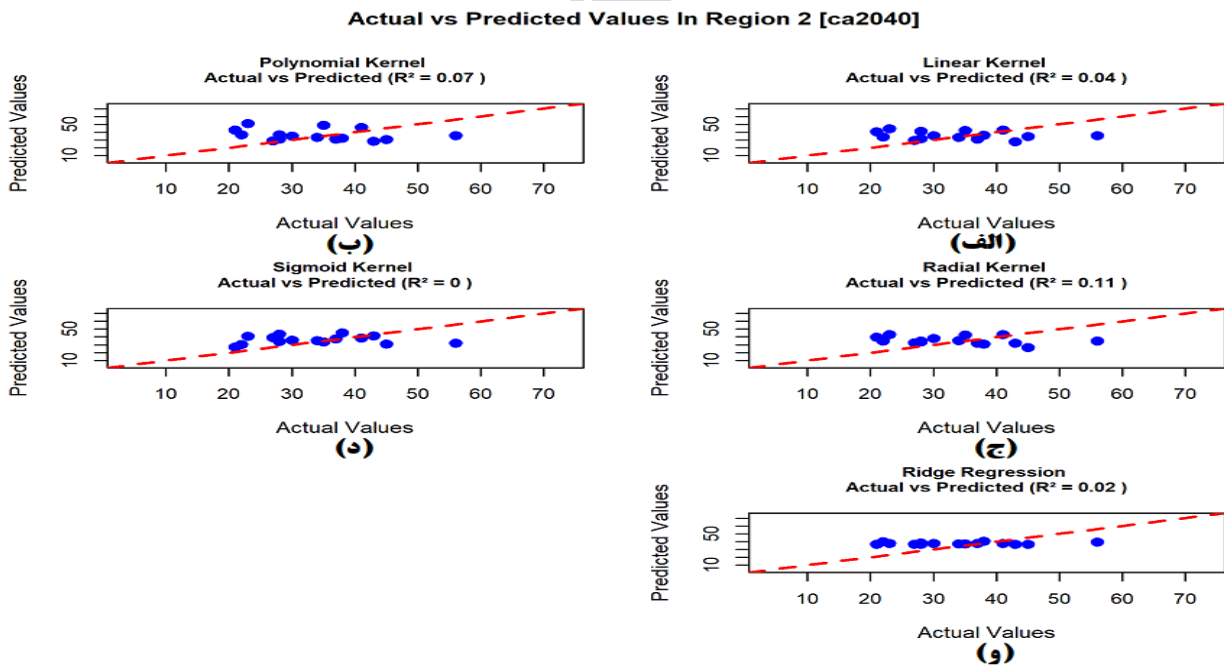
**Actual vs Predicted Values In Region 1 [ca2040]**



شکل ۳: نمودار مقادیر واقعی در برابر مقادیر پیش‌بینی شده برای چهار هسته بردار پشتیبان و رگرسیون ریج در منطقه ۱ در  $ca=2040$ :  
 (الف) هسته خطی، (ب) هسته چندجمله‌ای، (ج) هسته شعاعی، (د) هسته سیگموئید، (و) رگرسیون ریج.

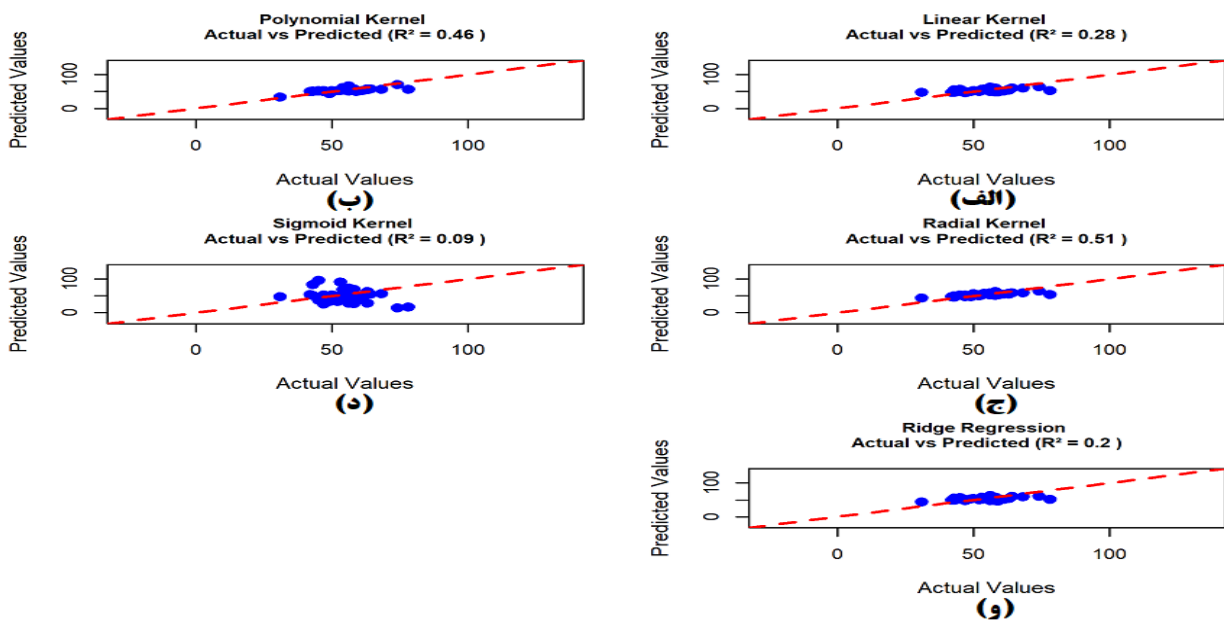


شکل ۴: نمودار مقادیر واقعی در برابر مقادیر پیش‌بینی شده برای چهار هسته بردار پشتیبان و رگرسیون ریج در منطقه ۲ در  $ca=0.20$ :  
 (الف) هسته خطی، (ب) هسته چندجمله‌ای، (ج) هسته شعاعی، (د) هسته سیگموئید، (و) رگرسیون ریج.



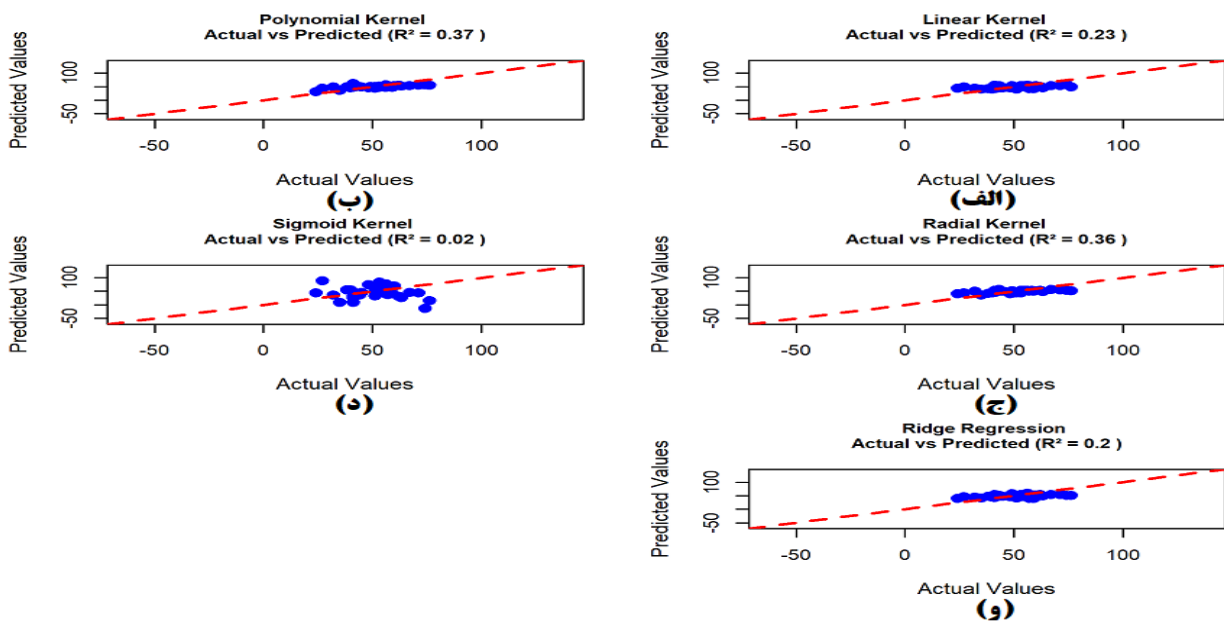
شکل ۵: نمودار مقادیر واقعی در برابر مقادیر پیش‌بینی شده برای چهار هسته بردار پشتیبان و رگرسیون ریج در منطقه ۲ در  $ca=0.40$ :  
 (الف) هسته خطی، (ب) هسته چندجمله‌ای، (ج) هسته شعاعی، (د) هسته سیگموئید، (و) رگرسیون ریج.

Actual vs Predicted Values In Region 3 [ca020]



شکل ۶: نمودار مقادیر واقعی در برابر مقادیر پیش‌بینی شده برای چهار هسته بردار پشتیبان و رگرسیون ریج در منطقه ۳ در  $ca=0.20$ :  
 (الف) هسته خطی، (ب) هسته چندجمله‌ای، (ج) هسته شعاعی، (د) هسته سیگموئید، (و) رگرسیون ریج.

Actual vs Predicted Values In Region 3 [ca2040]



شکل ۷: نمودار مقادیر واقعی در برابر مقادیر پیش‌بینی شده برای چهار هسته بردار پشتیبان و رگرسیون ریج در منطقه ۳ در  $ca=0.40$ :  
 (الف) هسته خطی، (ب) هسته چندجمله‌ای، (ج) هسته شعاعی، (د) هسته سیگموئید، (و) رگرسیون ریج.

می‌یابد. این امر حاکی از آن است که داده‌های مربوط به کلسیم در لایه زیرین خاک این منطقه نیاز به دقت پیش‌بینی بیشتری دارند. منطقه ۲ در هر دو لایه خاک دارای مقدار  $C$  بسیار پایین (۰/۱) است که نشان دهنده ساده‌سازی شدید مدل و اعمال جریمه سنگین برای خطاها می‌باشد. همچنین مقادیر بالای  $\varepsilon$  (۰/۸۸ و ۰/۸۲) بیانگر تحمل خطای گسترده است که احتمالاً به دلیل تغییرپذیری زیاد مقدار کلسیم یا وجود اختلال در داده‌های این منطقه می‌باشد. در منطقه ۳، مدل برای لایه سطحی خاک با مقادیر  $(C, \varepsilon) = (1, 0.02)$  تنظیم شده که حساسیت بسیار بالایی به خطاهای کوچک دارد، و در لایه زیرین این حساسیت تا حد  $\varepsilon$  صفر افزایش می‌یابد که نشان دهنده نیاز به دقت کامل در پیش‌بینی‌های مربوط به کلسیم است. این تنظیمات نشان می‌دهد که مدل  $SVR$  توانایی بالایی در تطبیق با شرایط مختلف لایه‌های خاک دارد و می‌تواند برای هر منطقه و عمق خاک به صورت بهینه تنظیم شود.

در منطقه ۳، حساسیت شدید به خطاهای کوچک در هر دو لایه مشاهده شد که بیانگر نیاز به دقت بالا در پیش‌بینی است. به‌طور کلی، نتایج بیانگر آن است که انتخاب مدل بهینه و تنظیم پارامترها باید متناسب با ویژگی‌های مکانی و عمقی داده‌ها صورت گیرد. استفاده از روش‌های مبتنی بر یادگیری ماشین مانند  $SVR$  این امکان را فراهم می‌کند که بدون نیاز به فرضیات سخت‌گیرانه، ساختارهای فضایی پیچیده مدل‌سازی شوند و پیش‌بینی‌های دقیق‌تری به دست آید. این یافته‌ها می‌تواند در مدیریت منابع خاک و بهینه‌سازی تصمیم‌گیری‌های کشاورزی و محیط زیستی مورد استفاده قرار گیرد.

## تقدیر و تشکر

نویسندگان مقاله کمال قدردانی و تشکر را از پیشنهادات ارزنده داوران، سردبیر و ویراستار محترم مجله که باعث ارائه بهتر و افزایش سطح کیفی مقاله شده است، دارند.

بر اساس جدول پارامترهای بهینه مدل  $SVR$  (جدول ۵)، می‌توان مشاهده کرد که هر منطقه دارای مشخصات منحصر به فردی است که منجر به تنظیمات پارامتری متفاوتی شده است. در منطقه ۱، مدل برای لایه سطحی خاک (۰-۲۰ سانتی‌متر) با مقادیر  $(C, \varepsilon) = (1, 0.066)$  تنظیم شده که نشان دهنده یک مدل نسبتاً ساده با تحمل خطای متوسط

جدول ۵: پارامترهای بهینه در مدل  $SVR$

منطقه	پارامترهای بهینه			
	ca2040		ca020	
	$\varepsilon$	$C$	$\varepsilon$	$C$
منطقه ۱	۰/۱۶	۱/۰	۰/۶۶	۱/۰
منطقه ۲	۰/۸۲	۰/۰۱	۰/۸۸	۰/۰۱
منطقه ۳	۰/۰	۱/۰	۰/۰۲	۱/۰

است، اما در لایه زیرین خاک (۲۰-۴۰ سانتی‌متر) با مقادیر  $(C, \varepsilon) = (10, 0.016)$  به مدلی با پیچیدگی بیشتر و حساسیت بالاتر به خطا تغییر

## ۵ نتیجه‌گیری

رشد حجم و پیچیدگی داده‌های فضایی، ضرورت بهره‌گیری از روش‌های نوین مدل‌سازی را آشکار ساخته است. در این مقاله، عملکرد رگرسیون بردار پشتیبان با هسته‌های مختلف و رگرسیون ریح برای پیش‌بینی محتوای کلسیم خاک در سه منطقه و دو لایه عمقی بررسی شد. نتایج نشان دادند که کارایی مدل‌ها به شدت به شرایط مکانی و عمق نمونه‌برداری وابسته است. مدل  $SVR$  با هسته شعاعی در اغلب موارد بهترین عملکرد را ارائه داد، در حالی که هسته سیگموئید ضعیف‌ترین نتایج را ثبت کرد. رگرسیون ریح نیز در برخی شرایط (مانند منطقه ۲ و متغیر  $ca2040$ ) توانست دقت بالاتری نسبت به مدل‌های غیرخطی داشته باشد. تحلیل پارامترهای بهینه نیز نشان داد که هر منطقه نیازمند تنظیمات متفاوتی است. در منطقه ۱ مدل برای لایه زیرین با  $C$  بالا و  $\varepsilon$  پایین سخت‌گیرانه‌تر تنظیم شد، در حالی که منطقه ۲ به دلیل اختلال و تغییرپذیری زیاد داده‌ها با  $C$  پایین و  $\varepsilon$  بالا به مدلی ساده‌تر گرایش یافت.

## مراجع

[۱] محمدزاده، م.، (۱۳۹۸)، آمار فضایی و کاربردهای آن، چاپ سوم، مرکز نشر آثار علمی دانشگاه تربیت مدرس، تهران،

[2] Awad, M., and Khanna, R. (2015), *Efficient Learning Machines: Theories, Concepts, and Applications for Engineers and System Designers, First Edition*, Apress Media, New York.

[3] Capeche, C.L., Macedo, J.R., Manzatto, H.R.H., and Silva, E.F. (1997), *Caracterização pedológica da fazenda Angra -*

*PESAGRO/RIO - Estação experimental de Campos (RJ)*, In: Congresso BRASILEIRO de Ciência do Solo, 26., Rio de Janeiro, Embrapa/SBCS.

- [4] Gutierrez, D. D. (2015), *Machine Learning and Data Science: An Introduction to Statistical Learning Methods with R, First Edition*, Technics Publications, NJ.
- [5] Izenman, A.J. (2008), *Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning*, Springer Science and Business Media, New York.
- [6] James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013), *An Introduction to Statistical Learning with Applications in R*, Springer Science and Business Media, New York.
- [7] Lauer, F., and Bloch, G. (2008), Incorporating Prior Knowledge in Support Vector Regression, *Machine Learning*, **70**, 89–118.
- [8] Roozbeh, M. (2025), Optimal ridge estimation in the restricted logistic semiparametric regression models using generalized cross-validation, *Journal of Applied Statistics*, 1–20.
- [9] Roozbeh, M., Rouhi, A., Mohamed, N.A., and Jahadi, F. (2023), Generalized support vector regression and symmetry functional regression approaches to model the high-dimensional data, *Symmetry*, **15**, 1262.

# Analysis of Hydrological Variables and Soil Properties Using Spatial Ridge and Support Vector Regression Methods

Mahdi Roozbeh<sup>1\*</sup>, Arash Ameri<sup>2</sup>

<sup>1</sup> *Department of Statistics, Faculty of Mathematics, Statistics and Computer Science, Semnan University, Semnan, Iran*

<sup>2</sup> *Department of Statistics, Faculty of Mathematics, Statistics and Computer Science, Semnan University, Semnan, Iran*

*Received: 2025/11/12*

*Accepted: 2026/01/03*

## Abstract

Traditional statistical methods have faced serious challenges due to the expansion of spatial data with complex spatiotemporal structure. These data require specialized methods due to spatial autocorrelation, variance heterogeneity, and complex geographical dependencies. In this study, support vector regression is introduced as a novel approach for analyzing and modeling the complex spatial structure of geostatistical data related to soil calcium and magnesium contents. This analysis is performed based on different geographical coordinates (east–west and north–south), at two depths of 0–20 cm and 20–40 cm, and across three distinct geographical regions. The support vector regression method, with its capability to model complex nonlinear relationships while preserving the spatial structure of the data, allows for more accurate and realistic prediction of the nutrient elements' distribution in soil. This approach, utilizing kernel functions, enables the analysis of high-dimensional feature spaces and structural complexities, proving its effectiveness against noise and outliers. To accurately measure the efficiency of support vector regression, its performance is compared against ridge regression.

**Keywords:** Generalized Cross-validation, Ridge Regression, Support Vector Machines, Support Vector Regression

**Mathematics Subject Classification:** 62J05; 62H25.