

Modelling and Digitally Mapping of Surface Soil Organic Carbon in Semirrom County Employing Several Machine Learning Algorithms

F. Rahmati¹, S. Hojati², K. Rngzan³ and A. Landi²

1. Former PhD Student, Department of Soil Science, Faculty of Agriculture, Shahid Chamran University of Ahvaz, Ahvaz, Iran
2. Professor, Department of Soil Science, Faculty of Agriculture, Shahid Chamran University of Ahvaz, Ahvaz, Iran
3. Professor, Department of Remote Sensing and GIS, Faculty of Earth Sciences, Shahid Chamran University of Ahvaz, Ahvaz, Iran

Received: 22 August 2024 Accepted: 16 December 2024 *Corresponding Author: s.hojati@scu.ac.ir

Abstract

Introduction: Knowledge of the spatial distribution of soil organic carbon (SOC) is a practical tool for determining sustainable land management strategies. Estimating carbon content and stocks is essential for carbon sequestration, greenhouse gas emissions and national carbon balance inventories. Accurate mapping of SOC's spatial distribution is a key assumption for soil resource management and land use planning. Over the last two decades, the utilization of data mining approaches in the spatial modeling of SOC using machine learning algorithms has gained considerable attention. The digital environment requires continuous soil maps at local and regional scales; however, such information is not always available at the necessary scale. Therefore, the digital soil mapping (DSM) approach is a crucial solution for quantifying and assessing variations in soil properties such as SOC, using remotely sensed indices and digital elevation model (DEM) as the most commonly used ancillary data for SOC prediction. In this context, the data mining techniques serves as the pathway to create digital soil maps. This study aims to compare two common machine learning algorithms, random forest and multiple linear regression, in the digital mapping of surface SOC in Semirrom County, Isfahan Province. Digital maps of SOC using these two algorithms were created, and the most important variables affecting the spatial distribution of SOC in the study area were reported.

Materials and Methods: A total of 200 surface soil samples (0-10 cm) were collected from the Semirrom area (51° 17' - 52° 3' E; 30° 42' - 31° 51' N) in Isfahan, Iran. According to the reports from synoptic meteorological station, the annual average temperature ranged from 7.5 to 12.5 °C, and the annual precipitation varied between 350 and 450 mm. The soil moisture and temperature regimes are classified as Xeric and Mesic, respectively. Soil sampling was conducted using the Global Positioning System (GPS) at the surface layer (0-10 cm). The preparation of the soil samples involved air drying, pounding and softening, followed by sieving through a 2 mm mesh. The amount of organic carbon in the samples was determined using the Walkley-Black method. Additionally, in order to evaluate the effect of other soil properties on the organic content of the soils, laboratory analyses were conducted, including measurements of saturated soil moisture content, soil texture, soil pH in saturated pastes, electrical conductivity of the soil saturation extracts and the calcium carbonate equivalent of the soils, following standard laboratory protocols.

In this research, auxiliary variables including terrain parameters and vegetation indices were derived from DEM and the Landsat 8 OLI satellite images using ArcMap version 10.4.10 and SAGAGIS version 6.0.4. All auxiliary layers were then converted to raster format using the "raster" package and

merged using the “Covstack” function. The values of all environmental covariates at each sampling point were extracted into a single file using the “extract” function from the “sp” package in the RStudio environment. Subsequently, using SPSS software version 19 and the principal component analysis (PCA) method, the most important auxiliary covariates among the 29 variables used in this research were identified for the modeling process. The dataset was then split into two groups; calibration (80%) and validation (20%) subsets. Finally, SOC contents were predicted and mapped using multiple linear regression (MLR) and random forest (RF) algorithms in the RStudio environment. The MLR and RF algorithms were executed using the “lm” and “randomForest” packages, respectively. Five different statistics were utilized to evaluate the performance of each model, including the coefficient of determination (R^2), bias, root mean square error (RMSE), normalized RMSE (nRMSE), and mean bias error (MBE).

Results and Discussion: Based on the descriptive analysis of the soil samples, the soils in the study area were characterized as non-saline, alkaline, and calcareous. The SOC contents of the soils ranged from 0.3 % to 2.2% with a mean value of 0.89%. The coefficient of variation for the SOC contents was 21.7%, which classifies the soils of the study area as having moderate variability, according to the values proposed by Wilding (1985). The results of the PCA indicated that the most important auxiliary variables for the modeling process included slope aspect, channel network base level, catchment slope, total curvature, height, longitudinal curvature, mass balance index, modified catchment area, slope degree, slope length, topographic position index, vertical distance to the channel network, soil adjusted vegetation index, transformed vegetation index, difference vegetation index, ratio vegetation index, and general curvature. These variables explained 80% of the total variance in the study area. A comparison of two different SOC prediction models, demonstrated that the RF model ($n_{tree} = 1000$ and $m_{try} = 10$) outperformed the MLR model, with R^2 , RMSE, nRMSE, and bias values of 0.79, 0.12, 0.13, and 0.002 respectively. The five most important variables identified by the RF algorithm for predicting SOC contents in the study area were the transformed vegetation index, ratio vegetation index, soil adjusted vegetation index, and slope degree. The final map of surface SOC distribution reveals that, although the RF algorithm provided better predictions than the MLR model, it also resulted in overestimation and/or underestimation of the minimum and maximum values of surface SOC contents, respectively.

Conclusion: The results of this study demonstrated that the RF regression algorithms outperformed the MLR method, thanks to its ability to account for the nonlinear and complex relationships between SOC content and environmental covariates.

Key Words: *Environmental covariates, machine learning, performance, spatial distribution*

مدل سازی و نقشه برداری رقومی کربن آلی در خاک سطحی اراضی شهرستان سمیرم با استفاده از برخی روش های یادگیری ماشین

فاطمه رحمتی^۱، سعید حاجتی^{۲*}، کاظم رنگزن^۳ و احمد لندی^۴

۱- دانش آموخته دکتری گروه مهندسی علوم خاک، دانشکده کشاورزی، دانشگاه شهید چمران اهواز، اهواز، ایران

۲- استاد گروه مهندسی علوم خاک، دانشکده کشاورزی، دانشگاه شهید چمران اهواز، اهواز، ایران

۳- استاد گروه سنجش از دور و GIS، دانشکده علوم زمین، دانشگاه شهید چمران اهواز، اهواز، ایران

چکیده

آگاهی از توزیع مکانی کربن آلی خاک گامی موثر در دستیابی به استفاده پایدار از اراضی و تعیین استراژی های مدیریتی مربوط به آن است. از این رو، این مطالعه با هدف مدل سازی و نقشه برداری رقومی کربن آلی خاک سطحی (۰-۱۰ سانتی متری) شهرستان سمیرم با استفاده از روش های رگرسیون جنگل تصادفی و رگرسیون خطی چندمتغیره انجام شد. به این منظور ۲۰۰ نمونه خاک سطحی به صورت منظم و با فواصل نمونه برداری ۵ کیلومتر \times ۵ کیلومتر از سطح منطقه برداشت گردید و سپس کربن آلی نمونه ها با استفاده از روش واکلی- بلک اندازه گیری شد. در پایان، نقشه رقومی کربن آلی در خاک سطحی منطقه با روش های مزبور و به کمک متغیرهای کمکی استخراج شده از مدل رقومی ارتفاع و تصاویر ماهواره لندست ۸ در محیط نرم افزار RStudio تهیه شد. یافته های این مطالعه حاکی از آن است که الگوریتم جنگل تصادفی برای برآورد میزان کربن آلی خاک به ترتیب با مقادیر RMSE و R^2 معادل ۰/۱۲ و ۰/۷۹ نسبت به روش رگرسیون خطی چندمتغیره با RMSE و R^2 معادل ۰/۱۹۲ و ۰/۰۵۷ پیش بینی های بهتری ارائه داده است. نتایج نشان داد که مهم ترین متغیرهای محیطی مؤثر بر توزیع کربن آلی خاک در منطقه مطالعاتی در مدل های مورد استفاده یکسان نیستند. به گونه ای که در مدل جنگل تصادفی شاخص های استخراج از پوش گیاهی و در رگرسیون خطی چندمتغیره شاخص های توپوگرافی نقش بیشتری در توزیع کربن آلی داشته است. بررسی نقشه نهایی پراکنش کربن آلی خاک در منطقه مطالعاتی نشان داد که تخمین های انجام شده با روش جنگل تصادفی اگرچه در مقایسه با روش رگرسیون خطی چندمتغیره تخمین های بهتری را ارائه داده اما در تخمین مقادیر کمینه و بیشینه مقادیر کربن آلی سطحی خاکها موفق نبوده است.

تاریخچه مقاله

دریافت: ۱۴۰۳/۰۶/۰۱

پذیرش نهایی: ۱۴۰۳/۰۹/۲۶

کلمات کلیدی:

یادگیری ماشین،

متغیرهای کمکی،

عملکرد،

پراکنش مکانی

* عهده دار مکاتبات

Email: s.hojati@scu.ac.ir

مقدمه

رطوبت بر اساس انحنای اراضی، پوشش گیاهی، توپوگرافی و عمق رسوبات زمین‌شناسی و انسان در محاسبات در نظر گرفته شود (۱ و ۲۱).

روش‌های مدل‌سازی متعددی برای نشان دادن ارتباط متغیرهای محیطی و ویژگی‌های خاک در نقشه‌برداری رقومی خاک وجود دارد. از بین مدل‌ها می‌توان به رگرسیون خطی چند متغیره و جنگل تصادفی اشاره کرد. مدل‌های رگرسیون خطی یک مدل آماری است که برای تخمین ویژگی‌های خاک از همبستگی یا ارتباط خطی بین ویژگی خاک و متغیرهای کمکی استفاده می‌کند. در روش رگرسیون خطی چند متغیره، می‌توان رابطه بین متغیر هدف و بیش از یک متغیر کمکی را به صورت همزمان ارزیابی نمود. این روش در مقابل اطلاعات نادرست، حساسیت بالایی دارد و ورود چنین داده‌هایی ممکن است منجر به خطاهای بزرگی در نتایج به دست آمده شود. علاوه بر این، برای استفاده از این روش متغیرها باید توزیع نرمال داشته باشند و تغییر آنها از یک رابطه خطی پیروی کند (۳۶).

الگوریتم جنگل تصادفی یک روش قدرتمند در فرآیند مدل‌سازی است که پیشرفت قابل توجهی در روش‌های داده‌کاوی ارائه داده است (۲۱). الگوریتم جنگل تصادفی مدل توسعه یافته‌ای از مدل طبقه‌بندی و رگرسیون درختی است که داده‌های مشاهداتی و متغیرهای کمکی را به دفعات برای به دست آوردن ارتباط بهینه بین متغیر پاسخ و متغیرهای مستقل و انجام تخمین‌های بعدی تفکیک و دسته‌بندی می‌کند. در این روش، مجموعه‌ای از شرط‌ها به صورت یک الگوریتم درختی با ساختار منطقی (اگر-آنگاه^۲) برای طبقه‌بندی یا پیش‌بینی کمکی یک متغیر به کار می‌رود (۱).

روش جنگل تصادفی به واسطه مزایای فراوانی مانند کارایی بالا در پیش‌بینی، قابل اجرا بودن برای حجم فراوانی از داده‌ها، برآورد داده‌های گم شده، مقاوم بودن در برابر نویز و عدم تأثیرپذیری از واریانس متغیرهای محیطی و خاک در سال‌های اخیر مورد توجه بسیاری از محققان قرار گرفته است (۱ و ۲۱) و نتایج بسیاری از تحقیقات نشان داده است که مدل جنگل

مواد آلی بر بسیاری از ویژگی‌های خاک‌ها از جمله حاصلخیزی، ساختمان، نفوذ آب در خاک، ظرفیت نگهداشت آب، تراکم و فعالیت میکروبی خاک تأثیر دارد (۱۵ و ۲۷). به علاوه، کربن آلی نقشی حیاتی در پایداری محیط زیست دارد و از شاخص‌های مهم در ارزیابی کیفیت و سلامت خاک است؛ از این رو، مدل‌سازی و نمایش تغییرات توزیع مکانی کربن آلی به‌ویژه در خاک‌های مناطق خشک و نیمه خشک می‌تواند در مدیریت پایدار اراضی بسیار حائز اهمیت تلقی گردد. نقشه‌برداری ویژگی‌های خاک به روش سنتی زمان‌بر و پرهزینه است. این در حالی است که در سال‌های اخیر استفاده از تکنیک‌های نقشه‌برداری رقومی خاک با صرف حداقل هزینه و وقت برای پیش‌بینی مکانی خصوصیات خاک مورد توجه بسیاری از محققان قرار گرفته است. معادله ینی (۱۱) پایه و اساس نقشه‌برداری رقومی خاک است که در آن تشکیل خاک وابسته به پنج فاکتور اقلیم، موجودات زنده، پستی و بلندی، مواد مادری و زمان می‌باشد. مک براتنی و همکاران^۱ (۱۵) معادله ینی را بازبینی کردند و مدل اسکورپن ($S_{c,a}=f(s,c,o,r,p,a,n)$) را ارائه دادند که برای توصیف کمی روابط بین خاک و عوامل مکانی دیگر به کار می‌رود و شامل هفت عامل سایر خصوصیات خاک در یک نقطه (s) خصوصیات اقلیمی (c)، پوشش گیاهی، جانوری یا فعالیت انسانی (o)، توپوگرافی (r)، مواد مادری (p)، سن (a)، و موقعیت جغرافیایی (n) است. در مدل اسکورپن $S_{c,a}$ کلاس‌های خاک یا خصوصیات خاک است.

سه جزء اصلی معادله اسکورپن ویژگی‌های خاک، متغیرهای محیطی و مدل‌های (توابع) ارتباط دهنده بین ویژگی‌های خاک و متغیرهای محیطی می‌باشند. متغیرهای کمکی محیطی در واقع نماینده عوامل خاک‌سازی هستند. اکثر متغیرهای محیطی با استفاده از مدل‌های رقومی ارتفاع، داده‌های طیفی سنجش از دور و نقشه‌های حاصل از مطالعات گذشته به دست می‌آید (۱). استفاده از داده‌های سنجش از دور و مدل‌های رقومی ارتفاع این امکان را فراهم می‌سازد تا تأثیر ژئومورفولوژی و موقعیت خاک در زمین نما، تکامل خاک، پراکنش مکانی

متر بالاتر از سطح دریا قرار دارد و رژیم رطوبتی خاک، زیریک^۴ و رژیم حرارتی آن، مزیک^۵ می‌باشد. در این منطقه تشکیلات دوره دوم زمین‌شناسی تا عهد حاضر مشاهده می‌شود. از نظر چینه‌شناسی، سنگ‌های مربوط به دوره‌های مختلف در کل منطقه قابل مشاهده است. نهشته‌های کواترنری شامل تراس‌های جدید و قدیم رودخانه‌های فصلی و دائمی است که عمدتاً اراضی زراعی را شامل می‌باشد. در مجموع ۳۷/۶ درصد کل اراضی منطقه را اراضی مواج یا همان فلات‌ها یا دشت‌های بریده می‌پوشانند. خاک‌های منطقه اغلب متعلق به رده‌های اتی‌سولز، آلفی‌سولز و اینسپتی‌سولز می‌باشند (۱۰).

شهرستان سمیرم از نظر موقعیت قرارگیری در حاشیه جنوب شرقی رشته کوه‌های زاگرس قرار گرفته به طوری که کوه‌های دنا به صورت حصاری در غرب و جنوب غربی آن کشیده شده است. بخش مرکزی قسمت‌های مرکزی و شمالی شهرستان را شامل می‌شود و بخش پادنا قسمت‌های جنوبی شهرستان را در برمی‌گیرد. شهر سمیرم واقع در مرکز شهرستان با ارتفاع ۲/۴۶۰ از سطح دریا به صورت دشت نیمه همواری با ارتفاع زیاد و شیب فراوان توسط کوه‌ها محصور شده است. با توجه به نقشه کاربری اراضی (شکل ۲) بیشترین مساحت منطقه مراتع فقیر با پوشش گیاهی کم می‌باشد که با رنگ بنفش نشان داده شده است. بر این اساس، ۷۲ درصد مساحت منطقه مراتع کوهستان‌های نسبتاً سرد است و دارای گیاهان یکساله می‌باشد. نقشه ژئومرفولوژی منطقه نشان می‌دهد در این منطقه نهشته‌های کواترنری وجود دارند که کاربری کشاورزی دارند (۱۰). در بخش مرکزی در مناطق با ارتفاع زیاد و در منطقه جنوبی در مناطق با ارتفاع کم سنگ‌های کربناته وجود دارد. این سنگ‌ها در منطقه مرکزی، با چشم‌انداز کوه و تپه توأم می‌باشند. در حالی که در منطقه جنوبی این سازند با اراضی مسطح و کم ارتفاع توأم شده است. آب و هوای این شهرستان بر اساس تقسیم‌بندی کوپن جزء مناطق معتدل سرد با تابستان‌های گرم و خشک و در ارتفاعات جنوبی جزء مناطق معتدل سرد با تابستان‌های خنک و خشک می‌باشد (۱۰).

تصادفی روش مناسبی برای برآورد ویژگی‌های خاک می‌باشد. هنگل و هم‌کاران^۱ (۶) نقشه ر قومی خصوصیات خاک را در آفریقا با دو مدل رگرسیون چندگانه خطی و جنگل تصادفی تهیه کردند و نتیجه گرفتند که مدل جنگل تصادفی در مقایسه با مدل رگرسیون خطی چندگانه دقت بیشتری داشته است. راد یانتو و هم‌کاران^۲ (۲۵) از سه مدل شبکه عصبی مصنوعی، جنگل تصادفی و رگرسیون درختی برای تهیه نقشه ر قومی کربن آلی خاک در اندونزی استفاده کردند و بهترین ضرایب تبیین را برای این مدل‌ها به ترتیب ۰/۸۶-۰/۵۹، ۰/۹۸-۰/۹ و ۰/۹۹-۰/۹۵ گزارش کردند. هیونگ و هم‌کاران^۳ (۷) نیز در پژوهشی به این نکته اشاره کردند که مدل‌های جنگل تصادفی و رگرسیون درختی به دلیل سرعت پردازش پارامترها و تفسیر آسان داده‌های خروجی، مدل‌های مناسبی به منظور استفاده در اهداف طبقه‌بندی هستند

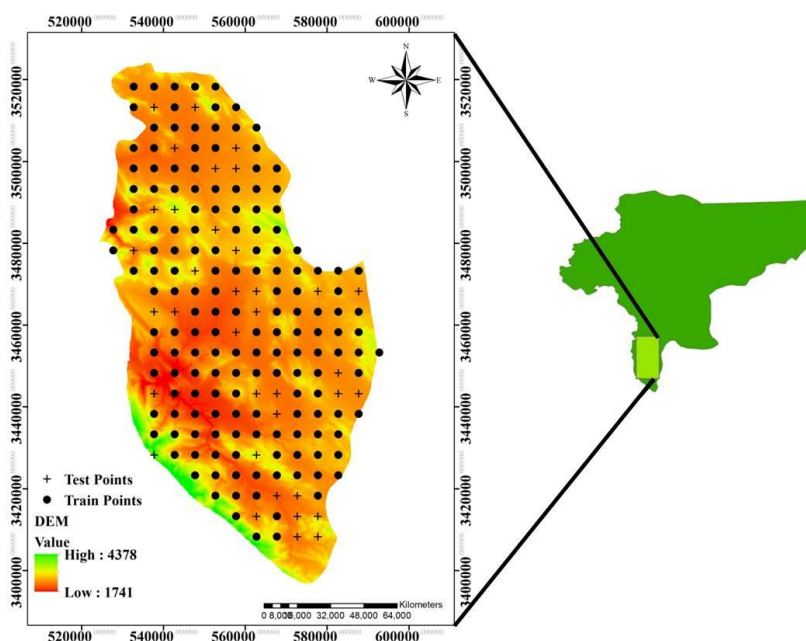
بررسی منابع انجام شده نشان می‌دهد تا کنون مطالعه‌ای در مورد مدل‌سازی تغییرات مکانی کربن آلی و شناسایی متغیرهای زودیافت‌تأثیرگذار بر پراکنش مکانی کربن آلی خاک در منطقه سمیرم استان اصفهان با استفاده از الگوریتم‌های مرسوم داده-کاوی انجام نشده است. از این رو، این پژوهش با اهداف مقایسه دو الگوریتم جنگل تصادفی و رگرسیون خطی چند متغیره، شناسایی بهترین عوامل محیطی کنترل‌کننده تغییرات کربن آلی و تهیه نقشه روند تغییرات کربن آلی در این منطقه انجام شد.

مواد و روش‌ها

موقعیت منطقه مورد مطالعه

شکل ۱ موقعیت منطقه مورد مطالعه در شهرستان سمیرم استان اصفهان را نشان می‌دهد. شهرستان سمیرم در جنوب غربی استان اصفهان با مساحت ۵۲۲۴ کیلومتر مربع بین طول‌های جغرافیایی ۱۷° ۵۱' تا ۳° ۵۲' شرقی و عرض ۳۰° ۴۲' تا ۳۱° ۵۱' شمالی واقع شده است.

میانگین دمای شهرستان ۷/۵ تا ۱۲/۵ درجه است و بارندگی سالانه ۳۵۰ تا ۴۵۰ میلی‌متر می‌باشد. ارتفاع منطقه ۲۰۰ تا ۳۰۰



شکل (۱) منطقه مورد مطالعه که نقاط آموزش (علامت دایره) و آزمون (به صورت +) را نشان می‌دهد

Figure (1) The study area showing training (Circle symbol) and testing (+) points

خاک با استفاده از دستگاه هدایت سنج الکتریکی (۲۷) و مقدار آهک به روش تیترا سیون برگشتی (۲۰) انجام گردید.

انتخاب متغیرهای محیطی و مدل‌سازی مکانی

در این پژوهش، تعداد ۲۸ متغیر کمی شامل شیب^۲، جهت شیب^۳، سطح مبنای شبکه زهکشی^۴، فاصله عمودی تا شبکه زهکشی^۵، شیب حوضه^۶، انحنا کلی^۷، انحنا طولی^۸، انحنا حداکثر^۹، ارتفاع از سطح دریا^{۱۰}، شاخص شاخص توازن جرم^{۱۱}، مساحت اصلاح شده حوضه^{۱۲}، طول شیب^{۱۳}، شاخص موقعیت توپوگرافی^{۱۴}، شاخص

نمونه‌برداری و تجزیه و تحلیل آزمایشگاهی

برای نمونه‌برداری از منطقه مورد مطالعه بعد از جمع‌آوری و بررسی اطلاعات اولیه و نقشه‌ها با استفاده از نرم افزار گوگل ارث موقعیت جغرافیایی تعداد ۲۰۰ نقطه با فواصل نمونه‌برداری ۵ کیلومتر در ۵ کیلومتر به صورت یک شبکه منظم تعیین گردید. سپس با استفاده از سیستم موقعیت یاب جهانی^۱ (GPS) از لایه سطحی خاک (۱۰ - ۰ سانتی متری) نمونه‌برداری انجام شد. آماده‌سازی نمونه‌های خاک شامل هوا خشک کردن، کوبیدن و نرم کردن نمونه‌های برداشت شده انجام و سپس نمونه‌ها از المک دو میلی‌متری عبور داده شدند. آنگاه میزان کربن آلی نمونه‌ها به روش واکلمی - بلک (۳۵) تعیین گردید. همچنین به منظور ارزیابی تأثیر ویژگی‌های خاک بر میزان کربن آلی خاک، تجزیه‌های آزمایشگاهی شامل تعیین درصد رطوبت اشباع خاک به روش وزنی (۸) بافت خاک به روش هیدرومتری (۲)، واکنش خاک (pH) در گل اشباع با استفاده از دستگاه pH متر (۱۶)، قابلیت هدایت الکتریکی عصاره اشباع

- 2- Slope
- 3- Aspect
- 4- Channel Network Base Level
- 5- Vertical distance to channel network
- 6- Catchment Slope
- 7- Total Curvature
- 8- Longitudinal Curvature
- 9 - Maximum Curvature
- 10- Height
- 11- Mass Balance Index
- 12- Modified catchment area
- 13- Slope length
- 14- Topographic position Index

- 1- Global Positioning System

ارائه شده توسط هنگل و روزیتر (۴ و ۵) تشریح گردیده است. لازم به ذکر است که پیش از استفاده از تابع extract برای استخراج مقادیر عددی متغیرهای محیطی در نقاط نمونه- برداری شده، ابتدا تمام لایه‌های اطلاعاتی کمکی با استفاده از بسته raster به فرمت رستری تبدیل شده و سپس با استفاده از تابع Covstack با یکدیگر ادغام شدند.

روش رگرسیون خطی چند متغیره

در این روش، بر اساس رگرسیون، رابطه بین متغیر هدف (کربن آلی خاک) و بیش از یک متغیر کمکی ارزیابی می‌شود (معادله ۱).

$$y_i = \beta_0 + \sum_{k=1}^P \beta_k X_{ik} \quad (1)$$

در معادله ۱، y_i مقدار متغیر هدف در یک مشاهده معین β_0 مقدار عرض از مبدا، β_k ضریب رگرسیون، X_{ik} داده‌های هر متغیر کمکی مورد نظر مربوط به نقطه مشاهداتی و P تعداد متغیرهای کمکی است. لازم به توضیح است که در روش رگرسیون خطی چند متغیره فرض بر خطی بودن روابط بین متغیرهای کمکی و متغیر هدف است. در این پژوهش برای اجرای الگوریتم رگرسیون خطی چند متغیره از بسته "lm" در محیط RStudio استفاده شد (۳). اما از آنجایی که ممکن است همیشه روابط بین متغیرهای کمکی و متغیر هدف خطی نباشد. لذا در این پژوهش از رگرسیون جنگل تصادفی نیز استفاده شد.

رگرسیون جنگل تصادفی

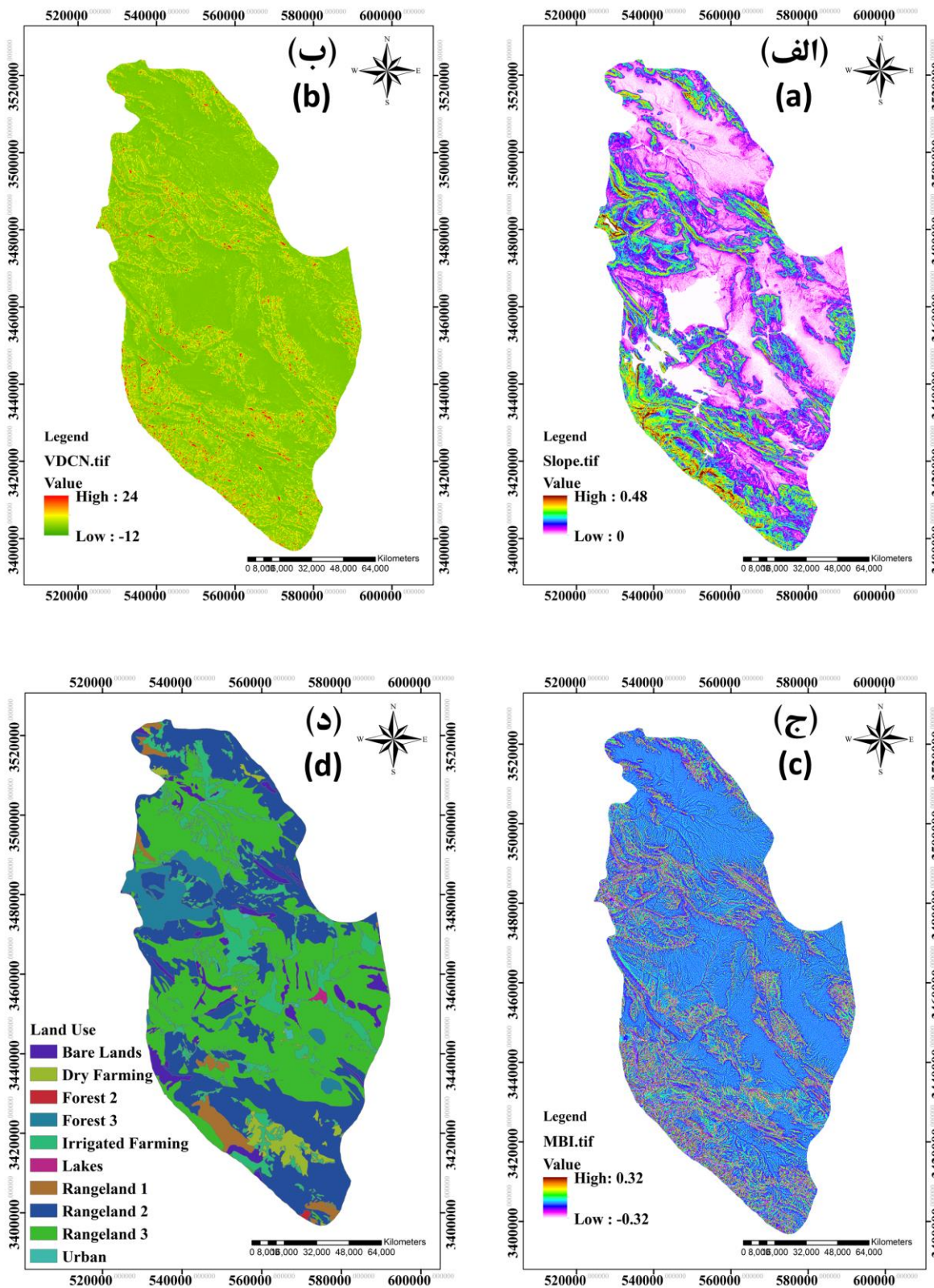
مدل جنگل تصادفی یکی از الگوریتم‌های پرکاربرد یادگیری ماشین محسوب می‌شود. رویکرد جنگل تصادفی مبتنی بر روش‌های جدید ترکیب اطلاعات است. این مدل دقیق و سریع است و به جای رشد دادن یک درخت، در این الگوریتم تعداد زیادی درخت مستقل تصمیم‌گیری (صدها یا هزارها) برای برقراری ارتباط بین متغیرهای کمکی و خصوصیات خاک رشد داده می‌شود (۱).

پوشش گیاهی نرمال شده^۱، شاخص پوشش گیاهی^۲، شاخص تعدیل شده برای خاک^۳، شاخص پوشش گیاهی تغییر یافته^۴، شاخص پوشش گیاهی تغییر یافته تصحیح شده^۵، شاخص تفاضل پوشش گیاهی^۶، شاخص نسبی^۷، انحنای عمومی^۸، مساحت ویژه^۹، شاخص خرسی توپوگرافی^{۱۰}، عمق دره^{۱۱}، شاخص ناهمواری توپوگرافی^{۱۲}، بافت سطح زمین^{۱۳}، شاخص همواری دره با درجه تفکیک بالا^{۱۴} و شاخص انباشت جریان^{۱۵} با استفاده از مدل رقومی ارتفاع با قدرت تفکیک مکانی ۹۰ متر و باندهای چندگانه تصاویر ماهواره لندست ۸ با استفاده از نرم‌افزارهای ArcMap نسخه ۱۰/۴/۱ و SAGAGIS نسخه ۶/۴/۰ (۱۳) تهیه شدند (۵). شکل ۲ به ضعی از متغیرهای محیطی تهیه شده در این پژوهش را نمایش می‌دهد.

سپس با استفاده از نرم‌افزار SPSS نسخه شماره ۱۹ و روش تجزیه به مؤلفه‌های اصلی (PCA)، از بین ۲۸ متغیر کمکی مورد استفاده در این پژوهش، مهم‌ترین متغیرهای کمکی برای فرآیند مدل سازی انتخاب گردیدند (۹). برای پیش‌بینی پراکنش ویژگی‌های خاک‌ها در منطقه مورد مطالعه، مقادیر متغیرهای محیطی و داده‌های مربوط به ویژگی‌های خاک در هر نقطه نمونه برداری شده به عنوان ورودی وارد محیط نرم‌افزار RStudio شدند و مقادیر متناظر آن‌ها در نقاط نمونه- برداری با استفاده از تابع extract در بسته نرم‌افزاری sp استخراج گردید. چگونگی استخراج این پارامترها در روش

- 1- Normalized Difference Vegetation Index
- 2- Soil Adjusted Vegetation Index (SAVI)
- 3- Transformed Vegetation Index (TVI)
- 4- Corrected Transformed Vegetation Index (CTVI)
- 5- Difference Vegetation Index (DVI)
- 6- Ratio Vegetation Index (RVI)
- 7- General Curvature
- 8- Surface Area
- 9- Topographic Wetness Index
- 10- Valley Depth
- 11- Topographic Ruggedness Index
- 12- Terrain Surface Texture
- 13- Multi-resolution Valley Bottom Flatness
- 14- Flow Accumulation

رحمتی و همکاران: مدل‌سازی و نقشه‌برداری رقمی کربن آلی در...



شکل (۲) نقشه‌های (الف) شیب بر حسب درجه (ب) نقشه فاصله عمودی تا شبکه زهکشی، (ج) نقشه شاخص توازن جرم و (د) کاربری اراضی در منطقه مطالعاتی

Figure (4) (a) The slope (%), (b) Vertical distance to channel network, (c) Mass Balance Index, and (d) Land Use maps of the study area

افزایش میانگین مربعات خطا برای تعیین اهمیت متغیرهای کمکی استفاده گردید. در این روش ابتدا برای هر درخت، خطای پیش بینی برای نمونه‌های خارج از کیسه (OOB) محاسبه می‌گردد. سپس این کار با جابجایی تمامی متغیرهای پیش‌بینی کننده تکرار می‌گردد و در پایان، میانگین اختلاف بین این دو مرحله در رابطه با تمامی درختان ر شد یافته در جنگل محاسبه می‌شود و با انحراف معیار تفاوت‌ها نرمال می‌شود. (۳۰). این کار برای تمامی متغیرهای مورد استفاده در مدل سازی متغیر هدف انجام و در نهایت نمودار اهمیت متغیرهای به کار برده شده تهیه شد.

در این پژوهش برای اجرای الگوریتم جنگل تصادفی از بسته نرم‌افزاری randomForest در محیط RStudio استفاده شد. برخی دیگر از بسته‌های استفاده شده در این مطالعه شامل readxl به منظور فراخوانی فایل‌های اکسل، GGally به منظور بررسی همزمان همبستگی بین متغیرها و توزیع داده‌ها و MASS به منظور شناسایی داده‌های گم شده است. همچنین از تابع importance در بسته randomForest برای تعیین اهمیت متغیرهای کمکی به کار رفته در فرآیند مدل‌سازی استفاده شد. از روش نمونه‌گیری بوت استرپ (نمونه‌گیری با جایگزینی به گونه‌ای که یک نمونه می‌تواند بیش از یک بار در فرآیند نمونه‌گیری انتخاب شود) و با استفاده از بسته caret به منظور تفکیک داده‌های آموزشی (حدوداً ۸۰ درصد) و آزمون (حدوداً ۲۰ درصد) و اعتبارسنجی نتایج استفاده گردید. (۱).

ارزیابی نتایج مدل‌ها

در این پژوهش برای ارزیابی نتایج مدل‌سازی کربن آلی خاک و بررسی دقت مدل از ضریب تبیین (معادله ۳)، شاخص میانگین انحراف خطا (معادله ۴)، میانگین ریشه دوم مربعات خطا (معادله ۵)، میانگین ریشه مربعات خطای نرمال شده (معادله ۶) و بایاس (معادله ۷) استفاده شد. در آمار تابع آریبی یا بایاس یک برآوردگر، اختلاف بین امید ریاضی مقادیر واقعی و تخمین زده شده را نشان می‌دهد.

ایجاد هر درخت تصمیم شامل دو مرحله است. مرحله اول، شامل ایجاد گره‌ها و انشعابات می‌باشد؛ گره در حقیقت اختصاص دادن متغیرهای مستقل به گروه‌های تعریف شده، براساس متغیر وابسته می‌باشد. انشعاب نیز عبارت است از انتخاب نقطه شکست از طریق انتخاب متغیری که بتواند بهترین شاخه‌های جدید را ایجاد کند. مرحله دوم، مرحله توقف و هرس است. در واقع مرحله توقف، شامل قوای است که میزان توسعه انشعاب در گره‌ها را محاسبه می‌کند. هرس نیز شامل حذف شاخه‌هایی است که اثر ناچیزی در مقادیر تخمینی حاصل از مدل نهایی دارند.

هدف از این مرحله به حداقل رساندن خطای پیش‌بینی است (۱). در این روش، یک نمونه به طور تصادفی از کل داده‌ها انتخاب شده و یک درخت روی این نمونه ساخته می‌شود. سپس در هر گره درخت، گروهی تخمین‌گر از کل تخمین‌گرها انتخاب می‌شوند و بهترین انشعاب با استفاده از این تخمین‌گرها تعیین می‌گردد و به این شکل خطای کلی مدل کاهش می‌یابد. سپس مدل سازی با روش جنگل تصادفی، ابتدا روی داده‌های آموزشی سپس روی داده‌های اعتبارسنجی اجرا می‌شود و در نهایت مدلی انتخاب می‌شود که دارای کمترین خطا است. در این مدل، تقریباً یک سوم داده‌های اصلی در ایجاد هر درخت استفاده نمی‌شوند. در این مدل، تقریباً یک سوم داده‌های اصلی در ایجاد هر درخت استفاده نمی‌شوند. چون این نمونه‌ها در آموزش استفاده نشده‌اند، و در مرحله آزمون مدل مورد استفاده قرار می‌گیرند (۱). برای اینکه مدل جنگل تصادفی یک مدل پیش‌بینی تولید کند، باید دو پارامتر مشخص شود. تعداد درختان رگرسیون در جنگل (ntree) و تعداد ویژگی‌های تصادفی انتخاب شده در هر گره (mtry) از جمله این عوامل هستند.

الگوریتم جنگل تصادفی برای تعیین اهمیت متغیرهای کمکی به طور پیش فرض از دو روش کاهش ضریب جینی^۱ یا درصد افزایش در مقدار میانگین مربعات خطا (IncMSE%) استفاده می‌کند که در این مطالعه از روش دوم یعنی درصد

رحمتی و همکاران: مدل‌سازی و نقشه‌برداری رقمی کربن آلی در...

$$R^2 = \left[\frac{\sum_{i=1}^n (Z(x_i) - Z(ave))(Z^*(x_i) - Z^*(ave))}{\sqrt{\sum_{i=1}^n (Z(x_i) - Z(ave))^2 (Z^*(x_i) - Z^*(ave))^2}} \right]^2 \quad (۳)$$

$$MBE = \frac{1}{n} \sum_{i=1}^n Z(x_i) - Z^*(x_i) \quad (۴)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n [(Z^*(x_i) - Z(x_i))^2]} \quad (۵)$$

$$nRMSE = RMSE / Mean \quad (۶)$$

$$Bias = Z(ave) - Z^*(ave) \quad (۷)$$

منطقه می‌باشند. بر این اساس، تغییرات کربن آلی، رس، سیلت و EC در منطقه در کلاس متوسط قرار دارد. تغییرات مقادیر کربنات کلسیم معادل و pH در منطقه کم و تغییرپذیری مقدار شن در منطقه زیاد است.

شکل ۳ نتایج همبستگی بین ویژگی‌های خاک را نشان می‌دهد. بر این اساس، درصد کربن آلی با واکنش خاک (pH) همبستگی منفی معنی‌دار (-۰/۲۲۳) در سطح ۰/۰۰۱ دارد. در این رابطه به نظر می‌رسد که با افزایش pH شرایط برای تجزیه مواد آلی خاک فراهم‌تر می‌شود و لذا محتوای کربن آلی خاک را تحت تاثیر قرار می‌دهد. نتایج اسکالبرگ (۲۹) پیرامون تغییرات pH در لایه‌های مختلف خاک نیز حاکی از آن است که محتوای pH خاک‌ها می‌تواند ارتباط معنی‌داری با کربن آلی خاک داشته باشد. همچنین یافته‌های این مطالعه حاکی از آن است که بین محتوای سیلت با درصد کربنات کلسیم معادل (۰/۸۵۴) و هدایت الکتریکی (۰/۶۲۸) همبستگی مثبت معنی‌داری در سطح ۰/۰۰۱ و با شن همبستگی منفی (۰/۷۵۹) معنی‌دار در سطح ۰/۰۰۱ وجود دارد. بر این اساس، می‌توان چنین نتیجه‌گیری نمود که عمده کربنات کلسیم در خاک‌های مورد مطالعه عمدتاً در اندازه سیلت قرار دارند و با افزایش مقدار سیلت خاک‌ها بر فراوانی آهک در خاک‌های مورد مطالعه افزوده می‌شود.

در معادله‌های ۳ تا ۷، n تعداد نمونه‌ها، $Z^*(x_i)$ مقدار پیش‌بینی شده، $Z^*(ave)$ میانگین مقادیر پیش‌بینی شده، $Z(x_i)$ مقدار اندازه‌گیری شده و $Z(ave)$ میانگین مقادیر اندازه‌گیری شده است.

تهیه نقشه پراکنش مکانی کربن آلی

نقشه پراکنش مکانی کربن آلی در خاک‌های سطحی منطقه مطالعاتی با استفاده از مقادیر پیش‌بینی شده از الگوریتم‌های جنگل تصادفی و رگرسیون خطی چندمتغیره استفاده گردید. برای درون‌یابی در نقشه‌های تولیدی نیز از روش وزن‌دهی معکوس فاصله استفاده شد که تمامی این فرآیندها در محیط نرم افزار ArcMap انجام شد.

نتایج و بحث

در جدول (۱) خلاصه آماری داده‌های خاکی ویژگی‌های خاک مورد بررسی ارائه شده است. بر این اساس، میانگین کربن آلی خاک در منطقه مطالعاتی برابر با ۰/۸۹ درصد به دست آمده است که بر اساس طبقه‌بندی سایس (۳۱) برای خاک‌های مناطق نیمه خشک در کلاس با مقادیر زیاد کربن آلی خاک طبقه‌بندی می‌شود. بر این اساس، مقدار حداقل و حداکثر کربن آلی در منطقه مطالعاتی به ترتیب ۰/۳ و ۲/۱۷ درصد می‌باشد. بر اساس طبقه‌بندی وایلدینگ (۳۹) ضریب تغییرات کمتر از ۱۵ درصد نشان دهنده تغییرپذیری کم، ضریب تغییرات ۱۵ تا ۳۵ درصد نشان دهنده تغییرپذیری متوسط و ضریب تغییرات بالاتر از ۳۵ درصد نشان دهنده تغییرات زیاد ویژگی‌های خاک در

جدول (۱) توصیف آماری ویژگی‌های خاک در منطقه مطالعاتی
 Table (1) Descriptive statistics of soil properties in the study area

رسانایی الکتریکی خاک (dS/m)	محتوای رطوبت وزنی خاک (%)	واکنش خاک pH	کربنات کلسیم معادل Calcium (%) Carbonate Equivalent (%)	کربن آلی (%) Organic Carbon (%)	رس (%) Clay (%)	شن (%) Sand (%)	سیلت (%) Silt (%)	ویژگی آماری Statistic
0.6	36.0	7.6	31.6	0.9	28	31	40	میانگین (Mean)
1.5	46.0	8.2	70.0	2.2	15	70	56	حداکثر (Max)
0.3	26.0	7.0	5.0	0.3	40	5	10	حداقل (Min)
1.2	2.0	1.2	65.0	1.9	25	65	46	دامنه (Range)
0.2	3.1	0.1	3.9	0.2	5	16	10	انحراف معیار (Standard Deviation)
28.7	8.6	1.3	12.4	21.7	19.0	50	26	ضریب تغییرات (%) (Coefficient of Variation (%))

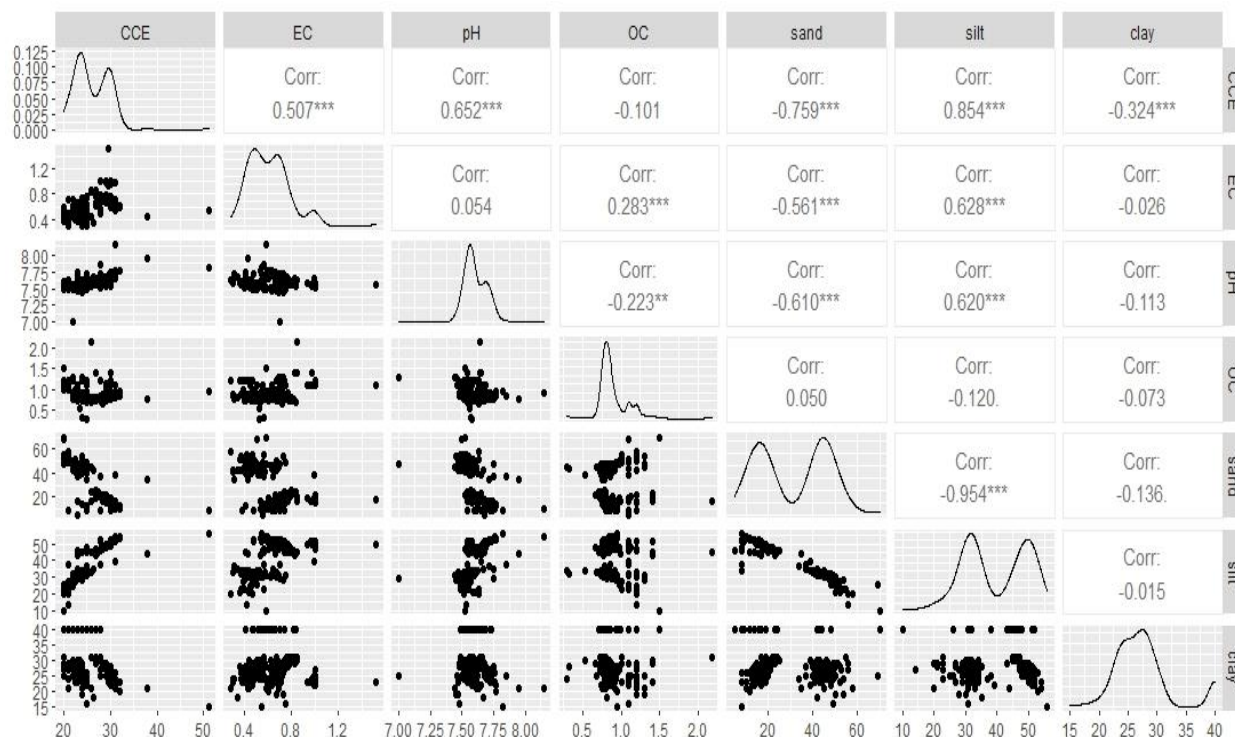
مطالعاتی توجیه می‌نمایند. درون هر مؤلفه اصلی متغیرهای کمکی که دارای بیشترین مقادیر همبستگی بودند، پس از انطباق با ضریب همبستگی بین آنها، در صورت وجود همبستگی بالا (بیشتر از ۰/۶)، متغیر با بالاترین مقدار همبستگی نگه داشته شد و بقیه حذف گردیدند (۴۰).
 بر این اساس، مهمترین متغیرهای کمکی مورد استفاده در فرآیند مدل‌سازی شامل جهت شیب، سطح مبنای شبکه زهکشی، شیب حوضه، انحنای کلی، ارتفاع، انحنای طولی، شاخص توازن جرم، مساحت اصلاح شده حوضه، شیب، طول شیب، شاخص موقعیت توپوگرافی، فاصله عمودی تا شبکه زهکشی، شاخص پوشش گیاهی تعدیل شده برای خاک، شاخص پوشش گیاهی تغییر یافته، شاخص تفاضل پوشش گیاهی، شاخص پوشش گیاهی نسبی، و انحنای عمومی بودند.

انتخاب متغیرهای کمکی برای فرآیند مدل‌سازی

جدول ۲ نتایج آزمون کرویت KMO و بارتلت را نشان می‌دهد. مقدار این شاخص همواره در محدوده صفر و ۱ قرار دارد. در صورتی که مقدار این شاخص بیشتر از ۰/۷ و نزدیک به ۱ باشد داده‌های مورد نظر برای تحلیل عاملی مناسب هستند. در مقابل، اگر مقدار این شاخص کمتر از ۰/۵ باشد برای تحلیل عاملی مناسب نخواهد بود. حال آن‌که، اگر مقدار این شاخص بین ۰/۵ تا ۰/۷ باشد، می‌توان با احتیاط به تحلیل عاملی پرداخت (۹). بر این اساس، در پژوهش حاضر داده‌های مورد نظر برای تحلیل عاملی مناسب هستند.

شکل ۴ نمودار اسکری (بازویی) و جدول ۳ نتایج تجزیه مؤلفه‌های اصلی را نشان می‌دهد. بر این اساس و با در نظر گرفتن مقادیر ویژه^۱ بزرگتر از ۱، تعداد ۷ مؤلفه اصلی برای ادامه مسیر و انتخاب متغیرهای موثر و کارآمد در فرآیند مدل‌سازی انتخاب شد که این مؤلفه‌ها در مجموع بیش از ۸۰ درصد از تغییرات کربن آلی را در منطقه

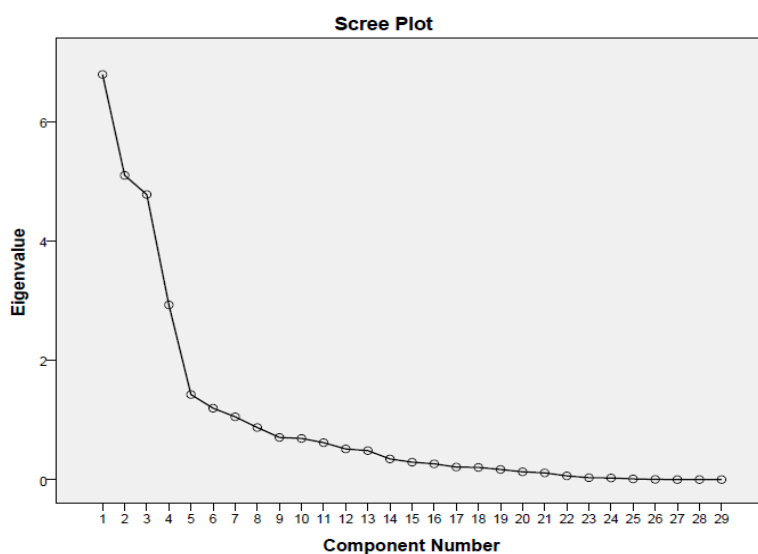
رحمتی و همکاران: مدل‌سازی و نقشه‌برداری رقمی کربن آلی در...



شکل (۳) ضریب همبستگی پیرسون بین ویژگی‌های خاک
Figure (3) Correlation coefficients between soil properties

جدول (۲) نتایج آزمون KMO و بارتلت
Table (2) Results of KMO and Bartlett's Test

پارامتر	مقدار عددی
Kaiser-Meyer-Olkin Measure of Sampling Adequacy	0.788
Bartlett's Test of Sphericity	Approx. Chi-Square 3182.879
	Df 120
	Sig 0.000



شکل (۴) نمودار اسکری (بازویی) برای تعیین تعداد مولفه‌های مناسب
Figure (4) The scree plot to determine suitable number of components

جدول (۳) واریانس تجزیه مؤلفه های اصلی متغیرهای محیطی
Table (3) The PCA variance of environmental variable

مؤلفه Component	مقادیر ویژه اولیه			مجموع مربعات مؤلفه های استخراج شده			مجموع مربعات مؤلفه ها پس از چرخش		
	کل Total	واریانس (%) variance (%)	درصد تجمعی Cumulative Percent	کل Total	واریانس (%) variance (%)	درصد تجمعی Cumulative Percent	کل Total	واریانس (%) Variance (%)	درصد تجمعی Cumulative Percent
PC1	6.7	23.4	23.4	6.7	23.4	23.4	6.5	22.5	22.5
PC2	5.09	17.5	40.9	5.09	17.6	40.9	3.7	13.1	35.6
PC3	4.7	16.4	57.4	4.7	16.4	57.4	3.7	12.9	48.6
PC4	2.9	10.1	67.5	2.9	10.1	67.5	3.2	11.3	59.9
PC5	1.4	4.9	72.4	1.4	4.9	72.4	3.1	10.7	70.6
PC6	1.19	4.1	76.5	1.19	4.1	76.5	1.6	5.5	76.2
PC7	1.05	3.6	80.2	1.05	3.6	80.2	1.1	3.9	80.2

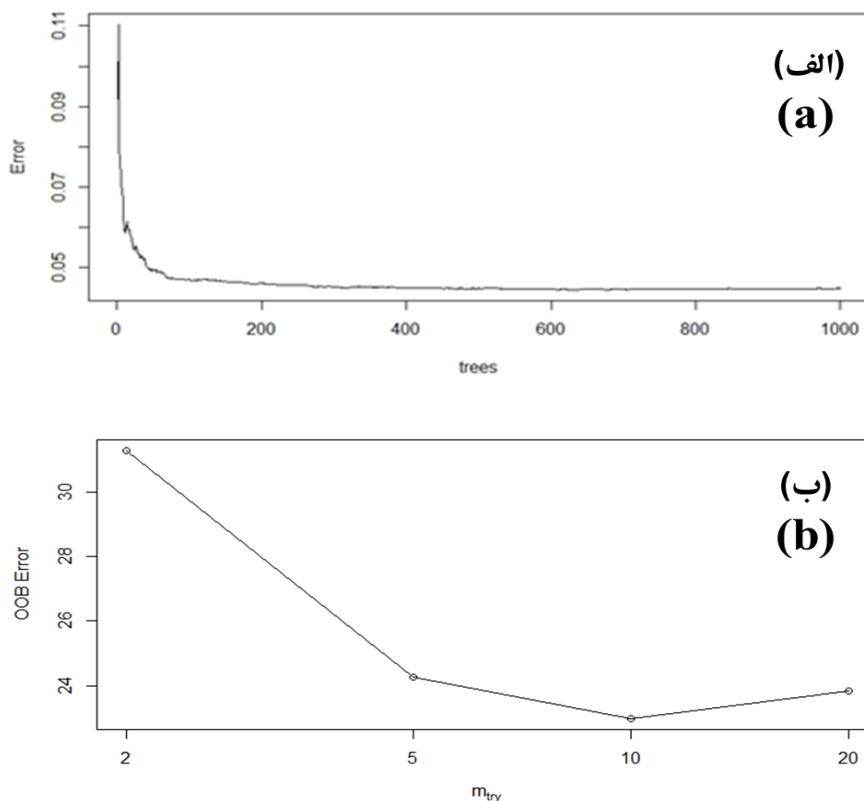
تخمین ویژگی های خاک استفاده می گردد. بنابراین مقدار mtry بهینه در این مطالعه ۱۰ انتخاب گردید.

تعیین اهمیت متغیرهای کمکی برای پیش بینی کربن آلی در مدل های مورد استفاده

شکل ۶ اهمیت متغیرهای کمکی مورد استفاده در مدل سازی کربن آلی خاک را با الگوریتم های جنگل تصادفی و رگرسیون خطی چندمتغیره نشان می دهد. همان گونه که مشاهده می گردد متغیرهای تاثیر گذار تشخیص داده شده در الگوریتم های مورد بررسی با یکدیگر متفاوت هستند. بر این اساس، در الگوریتم جنگل تصادفی شاخص های پوشش گیاهی تغییر یافته، پوشش گیاهی نسبی، شاخص پوشش گیاهی تعدیل کننده اثر خاک و شیب به عنوان پراهمیت ترین متغیرها در پیش بینی کربن آلی خاک شنا سایی شده اند. تقی زاده مهرجردی و همکاران^۱ (۳۲) برای تهیه نقشه رومی کلاس خاک و کربن آلی خاک در شهر بانه در استان کردستان، از متغیرهای کمکی داده های سنجش از دور و پارامترهای سرزمینی استفاده کردند و به این نتیجه رسیدند که مناسبترین متغیرهای کمکی جهت پهنه بندی کربن آلی خاک، شاخص پوشش گیاهی، شاخص رس، شیب، جهت شیب، انحنای سطح، شاخص همواری دره با درجه تفکیک بالا و شاخص خیزی بوده است.

تعیین پارامترهای الگوریتم جنگل تصادفی

در الگوریتم جنگل تصادفی تعیین دو پارامتر ntree و mtry اهمیت فراوانی دارد. بر این اساس، با افزایش تعداد مؤلفه های کمکی به عنوان ورودی مدل، مقدار خطای OOB کاهش می یابد. شکل ۵ الف نشان دهنده رابطه تعداد درختان و مقدار خطا در الگوریتم جنگل تصادفی است. بر این اساس، همانگونه که مشاهده می شود با افزایش تعداد درختان، مقدار خطا کاهش می یابد، به گونه ای که با افزایش تعداد درختان در جنگل تصادفی ساخته شده به ۱۰۰۰ درخت کمترین مقدار خطای تخمین به دست آمد. به همین دلیل برای افزایش دقت مدل، تعداد ntree (۱۰۰۰) به عنوان مقدار بهینه انتخاب گردید هر گره حاوی اطلاعاتی در مورد تعداد در آن گره و توزیع مقادیر متغیر وابسته است. نمونه ها موارد موجود در گره ریشه همه مشاهدات موجود در مجموعه آموزش است. در این پژوهش تعداد گره، ۲۴ در نظر گرفته شد. شکل ۵ ب ارتباط بین مقدار خطای تخمین خارج از کیسه (OOB) و تعداد متغیرهایی که به صورت تصادفی در هر بار نمونه گیری انتخاب می شوند (mtry) را نشان می دهد. با افزایش تعداد متغیرهای کمکی به ۱۰ متغیر به عنوان ورودی مدل، مقدار خطای OOB کاهش می یابد. بر اساس شکل ۵ب، همچنین با افزایش تعداد متغیرهای کمکی به بیش از ۱۰ متغیر مقدار خطای پیش بینی را به طور فزاینده ای بالا برده است. از این رو، در نقطه ای که مقدار خطای OOB کمینه مقدار خود را دارد، بهترین مقدار پارامتر mtry به دست می آید که از آن برای



شکل (۵) (الف) میزان خطا برای تعیین مقدار بهینه ntree (ب) میزان خطای OOB برای تعیین مقدار mtry
Figure (5) (a) The error to determining optimal ntree and (b) the OOB error to determining optimal mtry

تاثیرگذار در برآورد کربن آلی استفاده کرده است، حال آن‌که روش رگرسیون خطی چند متغیره تنها پارامترهایی را که با محتوای کربن آلی خاک رابطه مستقیم خطی دارند به عنوان متغیرهای تاثیرگذار معرفی نموده است. مصلح و همکاران^۱ (۱۸) اظهار کردند که پارامترهای مستخرج از مدل رقومی ارتفاع، حتی در مناطق دارای شدت پستی و بلندی کم، به عنوان متغیرهای محیطی مناسب در مدل‌سازی کلاس‌ها و خصوصیات خاک‌ها محسوب می‌شوند. وانگ و همکاران^۲ (۳۸) در بررسی تغییرات مکانی کربن آلی خاک دریافتند که شاخص پوشش گیاهی و جهت شیب از متغیرهای مهم در توزیع مکانی کربن آلی خاک هستند. ژو و همکاران^۳ (۴۳ و ۴۴) نیز در مطالعاتی گزارش کردند که متغیرهای توپوگرافی دارای

همانطور که در شکل ۶ مشاهده می‌شود، در الگوریتم جنگل تصادفی مهم‌ترین متغیرهای محیطی مؤثر بر توزیع کربن آلی خاک به ترتیب شامل شاخص‌های پوشش گیاهی (به دلیل برگشت بقایای گیاهی به خاک) و توپوگرافی می‌باشند؛ در صورتی که در روش رگرسیون خطی چندمتغیره، متغیرهای کمکی مرتبط با توپوگرافی اعم از درجه شیب، انحنای طولی، جهت شیب و انحنای کلی در مقایسه با شاخص‌های پوشش گیاهی به عنوان تاثیرگذارترین عوامل کنترل‌کننده پراکنش کربن آلی خاک در منطقه مطالعاتی تشخیص داده شده است. به اعتقاد نویسندگان تفاوت مشاهده شده بین دو روش در انتخاب متغیرهای تاثیرگذار در پیش‌بینی محتوای کربن آلی خاک تا حدودی به ماهیت عملکردی آنها برمی‌گردد؛ به گونه‌ای که در الگوریتم جنگل تصادفی علاوه بر روابط خطی از روابط غیرخطی نیز در شناسایی مولفه‌های

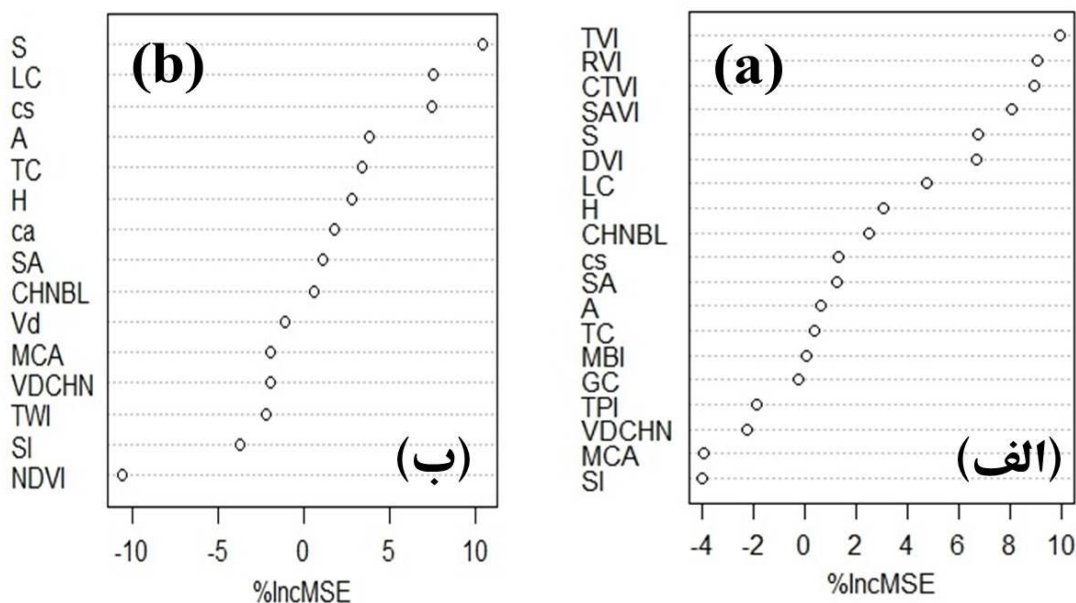
1- Mosleh *et al*
2- Wang *et al*
3- Zhou *et al*

۰/۰۳۷ و ۰/۰۰۵ برای مدل جنگل تصادفی با مقادیر RMSE، R²، nRMSE، MBE و بایاس ۰/۰۱۲، ۰/۰۷۹، ۰/۱۳، ۰/۱۳ و ۰/۰۰۲- نسبت به مدل رگرسیون خطی چند متغیره به ترتیب با مقادیر RMSE، R²، nRMSE، MBE و بایاس ۰/۱۹، ۰/۵۷، ۰/۲۲، ۰/۰۳۷ و ۰/۰۰۵ برای برآورد میزان کربن آلی خاک مناسب تر است. به طور کلی، هرچه میزان RMSE، MBE، RMSE و nRMSE مدل به صفر و هر چه میزان R² به یک نزدیک تر باشد، مدل دارای دقت بالاتری می باشد. با توجه به این که مقدار RMSE متناسب با واحد اندازه گیری متغیر است مقایسه مقدار آن بین مدل های ساخته شده با دو متغیر با واحدهای متفاوت درست نخواهد بود لذا مقدار RMSE را به دامنه داده های متغیر وابسته تقسیم کرده و آن را nRMSE یا RMSE نرمال شده می نامند. لازم به ذکر است nRMSE زیر ۱۰ درصد نشان دهنده دقیق بودن مدل، ۲۰-۱۰ درصد مناسب بودن مدل، ۳۰-۲۰ درصد دقت متوسط و بیش از ۳۰ درصد نشانه ضعیف بودن مدل است.

همبستگی بیشتری نسبت به سایر متغیرهای کمکی مورد استفاده با محتوای کربن آلی خاک و نیتروژن خاک بوده است. در واقع، متغیرهای توپوگرافی با کنترل جریان آب و رسوبگذاری بر توسعه خاک و توزیع مکانی ویژگی های خاک تأثیر می گذارند. میناسنی و مک برتنی (۱۹) در دو منطقه اندونزی نقشه رقوم ذخیره کربن آلی خاک را با استفاده از متغیرهای توپوگرافی شامل شیب، جهت شیب و شاخص رطوبت و به کارگیری مدل رگرسیون خطی، جنگل تصادفی و شبکه عصبی مصنوعی تهیه کردند و ضرایب تبیین را به ترتیب ۰/۹۹، ۰/۹۵ و ۰/۹۸ گزارش کردند.

ارزیابی نتایج مدل سازی پراکنش کربن آلی در منطقه

نتایج اعتبارسنجی و ارزیابی مدل های جنگل تصادفی و رگرسیون خطی چند متغیره در جدول ۴ نشان می دهد که مدل جنگل تصادفی با مقادیر RMSE، R²، nRMSE، MBE و بایاس ۰/۰۱۲، ۰/۰۷۹، ۰/۱۳، ۰/۱۳ و ۰/۰۰۲- نسبت به مدل رگرسیون خطی چند متغیره به ترتیب با مقادیر RMSE، R²، nRMSE، MBE و بایاس ۰/۱۹، ۰/۵۷، ۰/۲۲، ۰/۰۳۷ و ۰/۰۰۵، ۰/۱۳، ۰/۱۳ و ۰/۰۰۲- نسبت



شکل (۶) اهمیت متغیرهای محیطی مورد استفاده در مدل سازی با الف الگوریتم جنگل تصادفی و (ب) رگرسیون خطی چند متغیره
 Figure (6) The importance of environmental variable in (a) Random Forest and (b) Multiple Linear Regression modeling

محتوای کربن آلی خاک‌ها در سه کلاس کم (کمتر از ۰/۴ درصد)، متوسط (۰/۴ تا ۰/۸ درصد) و زیاد (بیش از ۰/۸ درصد) دسته‌بندی شده است (۳۰ و ۳۱). همان‌گونه که ملاحظه می‌گردد در نقشه به دست آمده به روش جنگل تصادفی، محتوای کربن آلی خاک سطحی منطقه از ۰/۵۹ تا ۱/۵۲ درصد تخمین زده شده است. در حالی که، محتوای کربن آلی تخمین زده شده با روش رگرسیون خطی چند متغیره در محدوده ۰/۸۳ تا ۰/۹۷ درصد می‌باشد که اختلاف نسبتاً زیادی را با مقادیر واقعی (اندازه‌گیری شده) کربن آلی در خاک سطحی منطقه نشان می‌دهد (جدول ۱). لذا به نظر می‌رسد روش جنگل تصادفی محتوای کربن آلی را در خاک‌های سطحی منطقه به نسبت بهتر پیش‌بینی کرده است، اگرچه به نوعی در پیش‌بینی مقادیر زیاد کربن آلی دچار کم‌تخمینی شده است. لازم به توضیح است که در روش جنگل تصادفی، ۹۹ درصد سطح منطقه از لحاظ محتوای کربن آلی در کلاس کربن آلی زیاد و تنها ۱ درصد سطح منطقه در کلاس کربن آلی متوسط قرار داشته است.

یافته‌های این مطالعه همچنین نشان داد که الگوریتم جنگل تصادفی به کار رفته در این پژوهش برای پیش‌بینی کربن آلی خاک منطقه مطالعاتی قادر به توجیه ۷۹ درصد از تغییرات مکانی کربن آلی است و متغیرهای به کار رفته در فرآیند مدل‌سازی قادر به پیش‌بینی حدود ۲۱ درصد از تغییرپذیری کربن آلی خاک در این منطقه نیست که به نظر می‌رسد تا اندازه زیادی به دلیل در نظر نگرفتن پارامترهای مدیریتی می‌باشد.

پایین بودن ضریب تبیین مدل رگرسیون چند متغیره احتمالاً به عدم استفاده از رابطه‌های غیرخطی و پیچیده بین مقدار کربن آلی خاک و پارامترهای کم‌کی مربوط می‌باشد. در الگوریتم جنگل تصادفی، مقادیر ویژگی‌های خاک با استفاده از یک فرمول ریاضی به دست نمی‌آید بلکه روابط بین متغیرهای محیطی و ویژگی خاک با استفاده از داده‌های آموزش کشف می‌شود (۲۶). ژانگ و همکاران^۱ (۴۱) نیز برتری روش جنگل تصادفی را نسبت به روش رگرسیون خطی در خصوص پیش‌بینی کربن آلی خاک نشان دادند. جیونگ و همکاران^۲ (۱۲) نیز برای پیش‌بینی توزیع کربن آلی خاک، مدل‌های جنگل تصادفی و رگرسیون خطی چندگانه را به کار برده و نشان دادند که روش جنگل تصادفی با مقادیر $R^2=0.97$ ، $RMSE=0.662$ ، $MSE=0.438$ و $ME=-0.248$ نتایج قابل قبولی در پیش‌بینی تغییرات کربن آلی خاک داشته است. سرینیواس و همکاران^۳ (۳۰) نیز نقشه رقمی کربن آلی و غیرآلی 1198 نمونه خاک را در هندوستان با استفاده از مدل جنگل تصادفی تهیه و گزارش کردند که مدل جنگل تصادفی مدل مناسبی بوده است.

نقشه نهایی پیش‌بینی کربن آلی خاک

شکل ۷ نقشه نهایی پیش‌بینی کربن آلی خاک را به روش‌های جنگل تصادفی و رگرسیون خطی چند متغیره نشان می‌دهد. لازم به توضیح است که در این مطالعه،

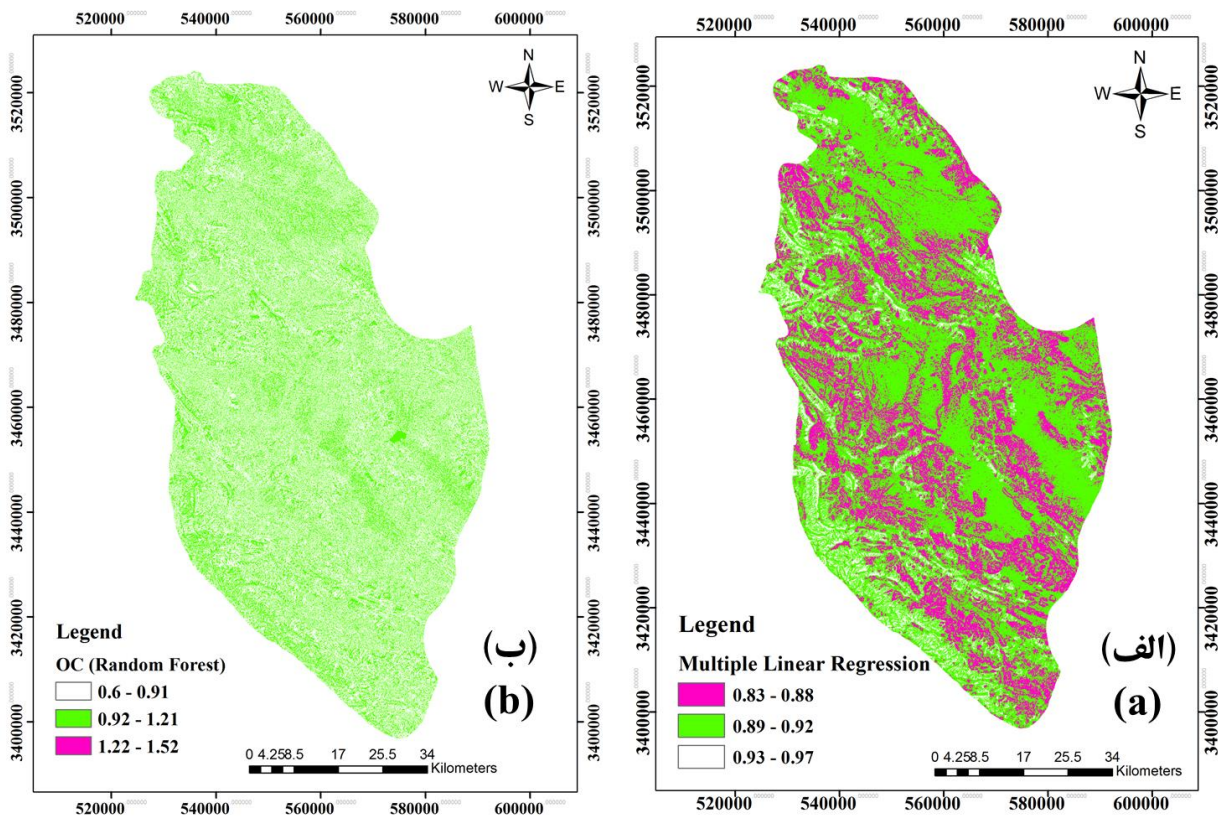
1- Zhou et al

2- Jeong et al

3- Sreenivas et al

جدول (۴) نتایج ارزیابی مدل‌های جنگل تصادفی و رگرسیون خطی چندمتغیره در پیش‌بینی کربن آلی خاک منطقه مطالعاتی
 Table (4) Evaluation of Random Forest and Multiple Linear Regression models in prediction of Soil Organic Carbon content of the study area

Bias	R ²	nRMSE	RMSE	MBE	مدل Model
-0.0023	0.79	0.13	0.12	0.013	رگرسیون جنگل تصادفی Random Forest Regression
0.0055	0.57	0.22	0.19	0.037	رگرسیون چند متغیره خطی Multiple Linear Regression



شکل (۷) نقشه توزیع کربن آلی در منطقه به روش (الف) رگرسیون خطی چند متغیره و (ب) جنگل تصادفی
 Figure (7) Map of the soil organic carbon produced by (a) Multiple Linear Regression, and (b) Random Forest algorithms

عملکرد بهتر روش رگرسیون جنگل تصادفی به دلیل توانایی آن در لحاظ نمودن روابط غیرخطی و پیچیده بین کربن آلی خاک با عوامل محیطی نسبت به روش رگرسیون خطی چندمتغیره در تخمین کربن آلی سطحی خاک بود. بر این اساس، در الگوریتم جنگل تصادفی

نتیجه‌گیری

این مطالعه با هدف مدل‌سازی و تهیه نقشه رقومی کربن آلی خاک سطحی منطقه سمیرم استان اصفهان با استفاده از مدل‌های جنگل تصادفی و رگرسیون خطی چندمتغیره انجام شد. نتایج حاصل از مدل‌سازی‌های انجام گرفته بیانگر

روش‌های مدل‌سازی مانند ماشین‌های بردار پشتیبان، جنگل تصادفی چندکی، الگوریتم ژنتیک، یادگیری عمیق نیز بررسی و نتایج این روش‌ها با نتایج مطالعه حاضر مقایسه و بهترین روش برای نقشه‌سازی ویژگی‌های خاک منطقه اعم از کربن آلی مورد استفاده قرار بگیرد.

سپاس‌گزاری

هزینه‌های اجرای این پژوهش توسط دانشگاه شهید چمران اهواز (پژوهانه شماره SCU.AS99.365) پرداخت شده است.

ساخته شده با ۱۰۰۰ درخت تصمیم (ntree) و mtry معادل ۱۰، مهمترین متغیرهای موثر بر تخمین کربن آلی در خاک سطحی منطقه مطالعاتی، به ترتیب شامل شاخص‌های مختلف پوشش گیاهی و درجه شیب می‌باشند. با این حال بررسی نقشه نهایی پراکنش کربن آلی در منطقه مطالعاتی نشان می‌دهد تخمین‌های انجام شده با روش جنگل تصادفی اگرچه در مقایسه با روش رگرسیون خطی چندمتغیره تخمین‌های بهتری را ارائه داده است اما در پیش‌بینی کمینه و بیشینه مقادیر کربن آلی سطحی خاک‌ها دچار کم‌تخمینی و/یا بیش‌تخمینی شده است. لذا پیشنهاد می‌شود در تهیه نقشه پراکنش کربن آلی خاک در منطقه مطالعاتی سایر

References

1. Breiman, L., 2001. Random forests. *Journal of Machine learning*, 45(1): 5-32. DOI: 10.1023/A:1010950718922.
2. Bouyoucos, G. J. 1951. A recalibration of hydrometer method for making mechanical analysis of soil. *Agronomy*, 43: 434-438. DOI: 10.2134/agronj1951.00021962004300090005x
3. Cutler, D.R., Edwards, J.T.C Beard, A. Cutler. K. H., and Hess, K.T. 2007. Random forests for classification in ecology. *Journal of Ecology*, 88(11): 2783-2792. DOI: 10.1890/07-0539.1
4. Hengl T., Rossiter D. G., and Stein A. 2003. Soil sampling strategies for spatial prediction by correlation with auxiliary maps. *Geoderma*, 120: 75–93. DOI: 10.1071/SR03005.
5. Hengl, T., Heuvelink, G.B.M., and Stein, A. 2004. A generic framework for spatial prediction of soil variables based on regression kriging. *Geoderma*, 120: 75–93. DOI: 10.1016 / j.geoderma. 2003.08.018.
6. Hengl, T., Heuvelink, B. M., Kempen, B., Leenaars, J.G. B., Walsh., M. G., Shepherd, K.D., Sila, A., MacMillan, R.A., Jesus, J. M., Tamene, L., and Tondoh, J.E. 2015. Mapping soil properties of Africa at 250 m resolution: random forests significantly improve current predictions, *PLOS ONE* 10 (6): e0125814. DOI: 10.1371/journal.pone.0125814.
7. Heung, B., H.C. Ho., J. Zhang., A. Knudby., C.E. Bulmer., and M.G. Schmidt. 2016. An overview and comparison of machine-learning techniques for classification purposes in digital soil mapping. *Geoderma*, 265: 62-77. DOI: 10.1016 / j.geoderma.2015.11.014.
8. Fatholouloumi, S., Vaezi, A. R., Alavipanah, S. K., Ghorbani, A., Saurette and D., Biswas, A. (2021). Effect of multi-temporal satellite images on soil moisture prediction using a digital soil mapping approach. *Geoderma*, 385, 114901. DOI: 10.1016/j.geoderma.2020.114901.
9. IBM Corp (2010). IBM SPSS Statistics for Windows, version 19, SPSS Inc., Chicago, Ill., USA.
10. Jalalian A. 1997. The studies of land resources and capability determination in Semirrom area. The Ministry of Jahad Sazandegi, Isfahan Province. (in Persian)
11. Jenny, H. 1941. *Factors of soil formation: A system of quantitative pedology*. McGraw-Hill, New York.
12. Jeong, G.H., Oeverdieck, S.J., Park, Huwe, B. and Ließ, M. 2017. Spatial soil nutrients prediction using three supervised learning methods for assessment of land potentials in complex terrain. *Catena*, 154: 73-84. DOI: 10.1016/j.catena.2017.02.006.
13. Lanyon, L. E. and Heald, W. R. 1982. Magnesium, calcium, strontium and barium. In: Page, A.L., et al. (Eds.), *Methods of Soil Analysis. Part II, Agronomy. Monograph*, American Society of Agronomy and Soil Science Society of America, Madison, Wisconsin, pp. 247-260.
14. Mahmoudzadeh, H., Matinfar, H. R., Taghizadeh-Mehrjardi, R., and Kerry, R. (2020). Spatial prediction of soil organic carbon using machine learning techniques in western Iran. *Geoderma Regional*, 21, e00260. DOI: 10.1016/j.geodrs.2020.e00260.
15. McBratney, A.B., Mendonça Santos, and M.L., Minasny, B. 2003. On digital soil mapping. *Geoderma*, 117: 3-52. DOI: 10.1016/S0016-7061(03)00223-4.
16. McBratney, A.B., Stockmann, U., Angers, D., Minasny, B., and Field, D. 2014. Challenges for soil organic carbon research. In: Hartemink, A.E., McSweeney, K., (Eds.), *Soil Carbon*, Cham: Springer, pp. 3-16. DOI: 10.1007/978-3-319-04084-4-1.

17. Mondal, A., Khare, D., Kundu, S., Mondal, S., Mukherjee, S. and A. Mukhopadhyaya. 2016. Spatial soil organic carbon (SOC) prediction by regression kriging using remote sensing data. *Egypt. Journal of Remote Sensing and Space Science*, 20 (1): 61-70. DOI: 10.1016/j.ejrs.2016.06.004.
18. Mosleh, Z., Salehi, M.H., Jafari, A., Borujeni, I.E., and Mehnatkesh, A. 2016. The effectiveness of digital soil mapping to predict soil properties over low-relief areas. *Environmental Monitoring and Assessment*, 188 (3): 195. DOI: 10.1007/s10661-016-5204-8.
19. Minasny, B., and McBratney, A.B. (2016). Digital soil mapping: A brief history and some lessons. *Geoderma*, 264: 301-311. DOI: 10.1016/j.geoderma.2015.07.017
20. Nelson, R. E. 1982. Carbonate and gypsum. In: Page, A.L., et al. (Eds.), *Methods of Soil Analysis. Part 2: Chemical Methods*, 2nd Ed., Agronomy Monograph, No. 9, American Society of Agronomy and Soil Science Society of America, Madison, WI. pp. 181-196.
21. Olaya, V. 2004. A gentle introduction to SAGA GIS. The SAGA User Group eV, Gottingen, Germany. 208 pp.
22. Osat, M., Heidari, A., Karimian Eghbal, M., and Mahmoodi, Sh. (2016). Spatial variability of soil development indices and their compatibility with soil taxonomic classes in a hilly landscape: a case study at Bandar village, Northern Iran. *Journal of Mountain Science*, 13(10): 1746-1759. DOI: 10.1007/s11629-016-3952-0.
23. R Development Core Team. 2015. R: a language and environment for statistical computing. R. Foundation for Statistical Computing, Vienna, Austria. <http://www.Rproject.org>.
24. Rossel, R.A.V., and McBratney, A.B. 2009. Diffuse reflectance spectroscopy as a tool for digital soil mapping. In: Hartemink, A.E., et al., (Eds.), *Digital Soil Mapping with Limited Data*. Springer, Dordrecht, pp. 165-172. DOI: 10.1007/978-1-4020-8592-5.
25. Rudiyanto, R., Minasny, B., Setiawan, B.I., Arif, C., Saptomo, S.K., and Chadirin, Y. (2016). Digital mapping for cost-effective and accurate prediction of the depth and carbon stocks in Indonesian peatlands. *Geoderma*, 272: 20–31. DOI: 10.1016/j.geoderma.2017.10.018.
26. Richards, L.A. 1954. *Diagnosis and Improvement of Saline-Alkali Soils*. USDA Hand book, No. 60. Washington, D.C., U.S.A.
27. Rhoades, J.D. 1996. Salinity: electrical conductivity and total dissolved solids. In: Sparks, D.L. (Ed.), *Methods of Soils Analysis, Part 3: Chemical Methods*. Soil Science Society of America Book series Number 5, Soil Science Society of America, Madison, Wisconsin, pp. 417-435.
28. Tarkalson, D.D., Brown, B., Kok, H., and Bjerneberg, D.L. 2009. Irrigated small-grain residue management effects on soil chemical and physical properties and nutrient cycling. *Soil Science*, 174:303-311. DOI: 10.1097/SS.0b013e3181a82a5f
29. Skullberg, U. 1991. Seasonal Variation of pH H₂O and pH CaCl₂ in centimeter- layers of mor Humus in a *Picea Abies* (L.) Karst stand. Sweden University of Agri Science, Department of Forest Site Research.
30. Sreenivas, K., Dadhwal, V.K., Kumar, S., Harsha, G.S., Mitran, T., Sujatha., G., Janaki Rama, S., Fyzee, M.A, and Ravisankar, T. 2016. Digital mapping of soil organic and inorganic carbon status in India. *Geoderma*, 269: 160-173. DOI: 10.1016/j.geoderma.2016.02.002
31. Sys, C., Van Ranst, E. and Debaveye, J. 1991. *Land Evaluation*. Agricultural Publication No. 7, General Administration for Development Cooperation, Brussels.
32. Taghizadeh-Mehrjardi, R., Nabiollahi, K. and Kerry, R. 2016. Digital mapping of soil organic carbon at multiple depths using different data mining techniques in Baneh region, Iran. *Geoderma*, 266: 98–110. DOI: 10.1016/j.geoderma.2015.12.003

33. Vaysse, K. and Lagacherie, K. 2015. Evaluating digital soil mapping approaches for mapping Global Soil Map soil properties from legacy data in Languedoc-Roussillon (France). *Geoderma Regional*, 4: 20-30. DOI: 10.1016/j.geodrs.2014.11.003
34. Venables, W.N., and B.D. Ripley. 2013. *Modern applied statistics with S-PLUS*. Springer, Dordrecht, 498 p.
35. Walkley A. and Black, I.A. 1934. An examination of the Degtjareff method for determining organic carbon in soils: effect of variations in digestion conditions and of inorganic soil constituents. *Soil Science*, 63: 251-263. DOI: 10.1097/00010694-194704000-00001
36. Wallach, D., Makowski, D., Jones, J.W., and Brun, F. 2006. *Working with dynamic crop models: Evaluation, analysis, parameterization, and applications*. Elsevier.
37. Wang Sh., Jin X., Adhikari, K., Li, W., Yu, M., Bian, Zh. and Wang, Q. 2018. Mapping total soil nitrogen from a site in northeastern China. *Catena*, 166: 134-146. DOI: 10.1016/j.catena.2018.03.023.
38. Wang, S., Wang, Q., Adhikari, K., Jia, S., Jin, X. and Liu, H. 2016. Spatial-temporal changes of soil organic carbon content in Wafangdian, China. *Sustainability*, 8: 1154. DOI: 10.3390/su8111154.
39. Wilding, L. 1985. Soil spatial variability: Its documentation, accommodation, and implication to soil surveys. In: Nielson, D.R., Bouma, J. (Eds.), *Wagenigen, Netherland*, pp. 166-194.
40. Yarali, J., Esmaili, A. and Esmaili, G.H. 2013. *Statistical Analyze with SPSS 20*. Kankash Publication, pp. 220-234.
41. Zhang, H., Wu, P., Yin, A., Yang, X., Zhang, M. and Gao, C. 2017. Prediction of soil organic carbon in an intensively managed reclamation zone of eastern China: A comparison of multiple linear regressions and the random forest model. *Science of the Total Environment*, 592: 704-713. DOI: 10.1016/j.scitotenv.2017.02.146
42. Zhao, Z., Yang, Q., Benoy, G., Chow, T.L., Xing, Z., Rees, H.W., and F.R. Meng. 2010. Using artificial neural network models to produce soil organic carbon content distribution maps across landscapes. *Canadian Journal of Soil Science*, 90 (1): 75–87. DOI: 10.4141/CJSS08057
43. Zhou, T., Geng, Y., Chen, J., Pan, J., Haase, D. and Lausch, A. 2020. High-resolution digital mapping of soil organic carbon and soil total nitrogen using DEM derivatives, Sentinel-1 and Sentinel-2 data based on machine learning algorithms. *Science of the Total Environment*, 729: 138244. DOI: 10.1016/j.scitotenv.2020.138244
44. Zhou, Y., Hartemink, A.E., Shi, Z., Liang, Z., and Lu, Y. 2019. Land use and climate change effects on soil organic carbon in North and Northeast China. *Science of the Total Environment*, 647: 1230-1238. DOI: 10.1016/j.scitotenv.2018.08.016.