



## تحلیل اثر داده‌های نامتعادل بر دقت و خودهمبستگی مکانی مدل‌سازی کربن آلی خاک

اعظم جعفری<sup>۱</sup>، فریدون سرمیدیان<sup>۲\*</sup>، زهرا رسائی<sup>۳</sup>

۱- گروه علوم خاک، دانشکده کشاورزی، دانشگاه شهید باهنر کرمان، کرمان، ایران

۲- گروه علوم و مهندسی خاک، دانشکده کشاورزی، دانشکده‌گان کشاورزی و منابع طبیعی، دانشگاه تهران، کرج، ایران،

fsarmad@ut.ac.ir

۳- گروه علوم و مهندسی خاک، دانشکده کشاورزی، دانشکده‌گان کشاورزی و منابع طبیعی، دانشگاه تهران، کرج، ایران

### چکیده

توزیع نامتعادل داده‌ها یکی از چالش‌های مهم در مدل‌سازی مکانی متغیرهای خاکی، به‌ویژه کربن آلی خاک (SOC)، محسوب می‌شود. این پژوهش با هدف بررسی تأثیر متعادل‌سازی داده‌ها بر دقت مدل و ساختار مکانی باقیمانده‌ها، در بخشی از دشت قزوین با مساحت حدود ۶۰۰۰۰ هکتار انجام شد. در این راستا، ابتدا میزان SOC در ۲۸۰ نمونه سطحی (عمق ۰-۳۰ سانتی‌متر) اندازه‌گیری و مجموعه‌ای از متغیرهای محیطی مانند پارامترهای توپوگرافی، شاخص‌های پوشش گیاهی، رطوبت و اقلیم جمع‌آوری شد. پس از انتخاب متغیرهای مهم با استفاده از روش RFE، مدل جنگل تصادفی (RF) در دو سناریوی استفاده از داده‌های نامتعادل و متعادل‌شده با روش SMOTE اجرا گردید. عملکرد مدل‌ها با شاخص‌های Moran's I و R<sup>2</sup>، RMSE مورد ارزیابی قرار گرفت. نتایج نشان داد که مدل RF با داده‌های متعادل‌شده (RMSE=0.3)، در مقایسه با داده‌های نامتعادل (R<sup>2</sup>=0.60، RMSE=0.4) عملکرد بهتری از خود نشان داد. همچنین شاخص خودهمبستگی مکانی موران برای باقیمانده‌ها در حالت متعادل‌شده از ۰/۳۰ به ۰/۰۳ کاهش یافت. بررسی کورلوگرام نیز حاکی از کاهش ساختار مکانی در خطاهای مدل پس از متعادل‌سازی داده‌ها بود. این یافته‌ها نشان می‌دهد که متعادل‌سازی آماری داده‌ها، علاوه بر بهبود دقت عددی مدل، به کاهش خودهمبستگی مکانی در خطاها و افزایش کیفیت فضایی مدل منجر می‌شود.

واژگان کلیدی: شاخص موران، SMOTE، نقشه‌برداری رقمی خاک

### مقدمه

کربن آلی خاک (Soil Organic Carbon – SOC) یکی از مهم‌ترین شاخص‌های کیفیت خاک و عاملی کلیدی در فرآیندهای زیست‌محیطی از جمله چرخه جهانی کربن، تولید محصولات کشاورزی، پایداری اکوسیستم‌ها و مقابله با تغییر اقلیم به شمار می‌آید (Lal, 2004; Minasny et al., 2017). برآورد دقیق و مکانی SOC نه‌تنها در مدیریت پایدار منابع خاک و آب ضروری است، بلکه به عنوان ابزاری مؤثر در برنامه‌ریزی‌های اقلیمی، سیاست‌گذاری‌های کشاورزی و پیش سلامت خاک در مقیاس‌های محلی تا جهانی مطرح شده است (Adhikari & Hartemink, 2016; Hengl et al., 2017). با توجه به هزینه و زمان‌بر بودن نمونه‌برداری میدانی و آزمایشگاهی، بهره‌گیری از روش‌های مدل‌سازی رقمی خاک (Digital Soil Mapping – DSM) و الگوریتم‌های یادگیری ماشین به‌منظور پیش‌بینی مکانی SOC در سال‌های اخیر به‌طور گسترده مورد توجه قرار گرفته است (Heung et al., 2016; Wadoux et al., 2021). یکی از چالش‌های اساسی در فرآیند مدل‌سازی، توزیع نامتوازن داده‌ها (Data Imbalance) است که به‌ویژه در متغیرهایی مانند SOC که دارای توزیع چوله و غیریکنواخت در سطح منطقه هستند، به‌وضوح مشاهده می‌شود. در بسیاری از موارد،

درصد زیادی از نمونه‌ها مقادیر SOC پایینی دارند و تنها تعداد معدودی دارای مقادیر بالا هستند. این نابرابری در توزیع داده می‌تواند موجب تعصب (bias) در فرآیند یادگیری مدل شده، دقت پیش‌بینی در مقادیر اقلیت را کاهش داده و ساختارهای مکانی واقعی متغیر را مخدوش نماید (Zhao et al., 2019).

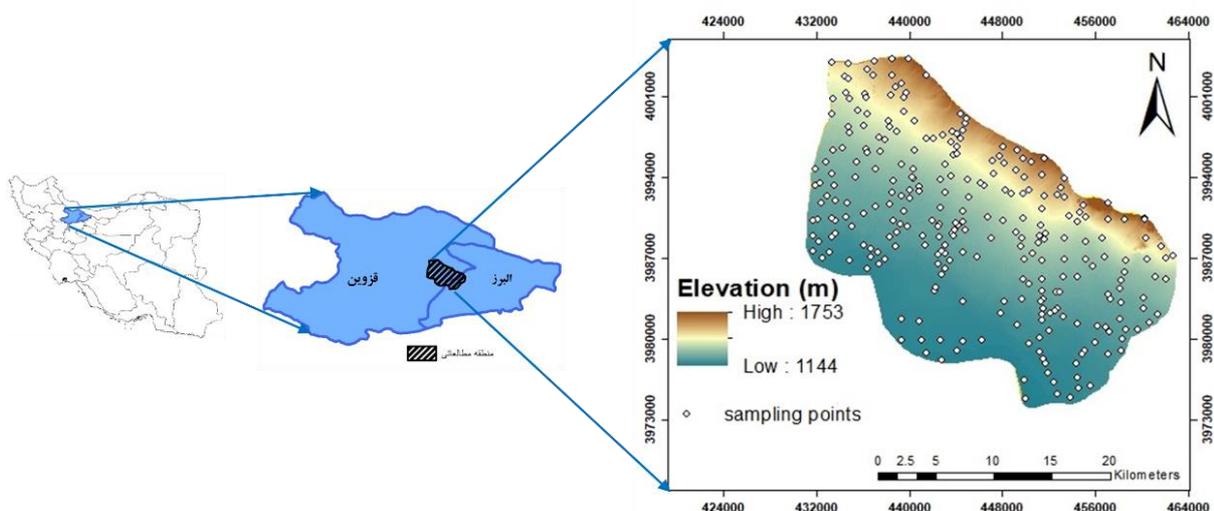
از سوی دیگر، یکی از مفاهیم کلیدی در تحلیل داده‌های مکانی، خودهمبستگی مکانی (Spatial Autocorrelation) است که نشان‌دهنده میزان شباهت بین مقادیر یک متغیر در موقعیت‌های مکانی مجاور می‌باشد. شاخص موران (Moran's I) به‌طور گسترده برای سنجش خودهمبستگی مکانی باقیمانده‌های مدل‌ها استفاده می‌شود. اگر باقیمانده‌های مدل دارای خودهمبستگی مکانی بالا باشند، نشان‌دهنده آن است که مدل قادر به بازنمایی صحیح روابط فضایی متغیر نبوده است. از این رو، تحلیل خودهمبستگی مکانی نه تنها ابزاری آماری، بلکه ابزاری تشخیصی برای ارزیابی و بهبود مدل‌سازی مکانی محسوب می‌شود (Hengl et al., 2007; Brus et al., 2011). با وجود پیشرفت‌های صورت گرفته، پژوهش‌های اندکی به بررسی مستقیم تأثیر عدم تعادل آماری داده‌ها بر ساختار مکانی باقیمانده‌ها در مدل‌سازی متغیرهای خاکی پرداخته‌اند. بنابراین، ضروری است رابطه میان عدم تعادل توزیع داده‌ها و خودهمبستگی مکانی در فرآیند مدل‌سازی، به‌ویژه برای متغیرهای محیطی مهمی مانند SOC، مورد بررسی دقیق قرار گیرد.

از این رو، پرسش اساسی این پژوهش آن است که چگونه توزیع نامتوازن داده‌ها می‌تواند باعث افزایش خودهمبستگی مکانی در باقیمانده‌های مدل گردد؟ همچنین آیا اصلاح عدم تعادل داده‌ها می‌تواند منجر به کاهش خودهمبستگی مکانی و افزایش دقت مکانی پیش‌بینی شود؟ بنابراین، هدف این مطالعه بررسی تأثیر توزیع داده‌ها (متعادل یا نامتعادل) بر میزان دقت و خودهمبستگی مکانی باقیمانده‌های مدل در مدل‌سازی کربن آلی خاک با استفاده از مدل جنگل تصادفی است. برای این منظور، ابتدا توزیع مقدار کربن آلی خاک در مجموعه‌ای از ۲۸۰ نمونه سطحی متعادل‌سازی و اصلاح شده و سپس عملکرد مکانی مدل در دو حالت با و بدون تعادل توزیع داده‌ها مقایسه شده است.

## مواد و روش‌ها

### تشریح منطقه مطالعاتی و داده‌های خاک

این پژوهش در بخشی از دشت قزوین با مساحت تقریبی ۶۰۰۰۰ هکتار انجام شد (شکل ۱). منطقه مورد مطالعه بین عرض‌های جغرافیایی ۳۵ درجه و ۵۵ دقیقه تا ۳۶ درجه و ۱۰ دقیقه شمالی و طول‌های جغرافیایی ۵۰ درجه و ۱۵ دقیقه تا ۵۰ درجه و ۳۵ دقیقه شرقی واقع شده و ارتفاع آن از ۱۱۴۴ متر در پایین‌ترین نقطه تا ۱۷۵۳ متر در بلندترین نقطه متغیر است. مواد مادری خاک در بخش‌های شمالی و مرکزی منطقه عمدتاً شامل رسوبات آبرفتی و کوهرفتی مربوط به دوره کواترنر و در نواحی جنوبی کفه‌های گلی و نمکی می‌باشند. از نظر توپوگرافی، منطقه دارای تنوع فیزیوگرافی قابل توجهی است که شامل تپه‌ها، واریزه‌های سنگریزه‌ای، دشت‌های دامنه‌ای، دشت‌های آبرفتی و اراضی پست می‌شود. کاربری اراضی عمدتاً شامل اراضی زراعی آبی و دیم، مراتع شور و غیرشور است.



شکل ۱- منطقه مورد مطالعه و پراکنش مکانی نقاط نمونه‌برداری شده

به منظور انجام این مطالعه، ۲۸۰ نمونه سطحی (عمق ۰-۳۰ سانتی متری) طبق الگوی نمونه برداری تصادفی طبقه بندی شده از قسمت های مختلف منطقه برداشت شد. بعد از کوبیدن و عبور نمونه های خاک از الک ۲ میلی متری، میزان درصد کربن آلی نمونه ها به روش سوزاندن تر (Walkely and Black, 1934) اندازه گیری شد.

### مدل سازی با الگوریتم جنگل تصادفی

در این پژوهش، از الگوریتم جنگل تصادفی (Random Forest) به منظور مدل سازی مکانی کربن آلی خاک (SOC) استفاده شد. این الگوریتم به عنوان یکی از روش های قدرتمند یادگیری ماشین، با بهره گیری از مجموعه ای از درخت های تصمیم گیری و فرآیند نمونه گیری بوت استرپ، قادر است روابط غیرخطی پیچیده بین متغیر هدف و مجموعه ای از متغیرهای مستقل را به خوبی شناسایی نماید. پیاده سازی این مدل با استفاده از بسته randomForest در محیط نرم افزار R انجام شد.

متغیرهای محیطی مورد استفاده در مدل سازی، دارای وضوح مکانی ۱۲/۵ متر بودند. داده های توپوگرافی از مدل ارتفاعی دیجیتال ALOS استخراج شده و برای محاسبه شاخص های مرتبط با پوشش گیاهی، شوری، کانی ها، رطوبت و نسبت های طیفی، تصاویر ماهواره ای Sentinel-2 به کار گرفته شد. علاوه بر این، میانگین دمای سالانه و میانگین بارندگی سالانه به عنوان شاخص های اقلیمی لحاظ شدند. به منظور انتخاب مؤثرترین متغیرها، از روش حذف بازگشتی ویژگی ها (Recursive Feature Elimination - RFE) استفاده شد (Chen & Jeong, 2007). بر این اساس، متغیرهای مهم شامل ارتفاع، عمق دره، شاخص سبزیگی، شاخص رطوبت ساگا و میانگین بارندگی سالانه برای مدل سازی نهایی برگزیده شدند.

### متعادل سازی داده ها

بررسی توزیع آماری کربن آلی خاک نشان داد که داده ها به شدت نامتوازن هستند، به طوری که بیش از ۷۰ درصد نمونه ها دارای مقادیر کمتر از ۱ درصد بودند. برای کاهش اثر این عدم تعادل در فرآیند یادگیری مدل، از روش SMOTE (Synthetic Minority Oversampling Technique) در محیط R استفاده شد. این روش با استفاده از نمونه گیری مصنوعی از مقادیر حداقل، موجب متعادل سازی توزیع SOC گردید.

### اعتبارسنجی مدل

به منظور ارزیابی دقت مدل ها، از روش اعتبارسنجی متقابل ۱۰-تایی (10-fold cross-validation) استفاده شد. برای اعتبارسنجی مدل ها، از شاخص های آماری ضریب تعیین ( $R^2$ ) و RMSE، استفاده شد.

$$R^2 = 1 - \frac{\sum_{i=1}^n (O_i - P_i)^2}{\sum_{i=1}^n (O_i - \bar{O})^2} \quad (1)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (P_i - O_i)^2} \quad (2)$$

در این معادلات،  $O_i$  و  $\bar{O}$  به ترتیب مقادیر مشاهده شده و میانگین آن ها و  $P_i$  مقادیر پیش بینی شده است. مقدار شاخص موران با استفاده از ماتریس وزنی همسایگی ایجاد شده بر مبنای فاصله جغرافیایی بین نمونه ها و تابع moran.test از بسته spdep در محیط R محاسبه گردید. کاهش مقدار Moran's I در باقیمانده ها نشانه ای از بهبود دقت مکانی مدل است.

### نتایج و بحث

بررسی مشخصات آماری کربن آلی خاک (SOC) در عمق ۰ تا ۳۰ سانتی متری نشان می دهد (جدول ۱) که این متغیر از تنوع قابل توجهی در منطقه برخوردار است. مقدار SOC بین ۰ تا ۲/۸ درصد متغیر بوده و میانگین آن برابر با ۰/۸۷ درصد محاسبه شد. انحراف معیار نسبتاً بالا (۰/۴۸) و ضریب تغییرات (CV) برابر با ۰/۵۶ نشان دهنده پراکندگی متوسط تا زیاد داده ها در سطح منطقه است. توزیع داده ها چولهی مثبت ( $Skewness = ۱/۱۴$ ) و کشیدگی بالاتر از توزیع نرمال ( $Kurtosis = ۱/۴۳$ ) دارد که بیانگر توزیع نامتقارن با غلبه مقادیر پایین SOC در نمونه ها می باشد. این ویژگی آماری، احتمال ایجاد عدم تعادل (Imbalance) در داده ها را افزایش می دهد و بر عملکرد مدل های پیش بینی اثرگذار است، لذا در ادامه تحلیل، تأثیر متعادل سازی داده ها بر دقت و ساختار مکانی مدل بررسی می گردد.

جدول ۱- شاخص های آماری کربن آلی خاک (SOC) در عمق ۰-۳۰ سانتی متر

حداقل	حداکثر	میانگین	انحراف معیار	واریانس	چولگی	کشیدگی	ضریب تغییرات
-------	--------	---------	--------------	---------	-------	--------	--------------

۰/۵۶	۱/۴۳	۱/۱۴	۰/۲۳	۰/۴۸	۰/۸۷	۲/۸۰	۰/۰۰
------	------	------	------	------	------	------	------

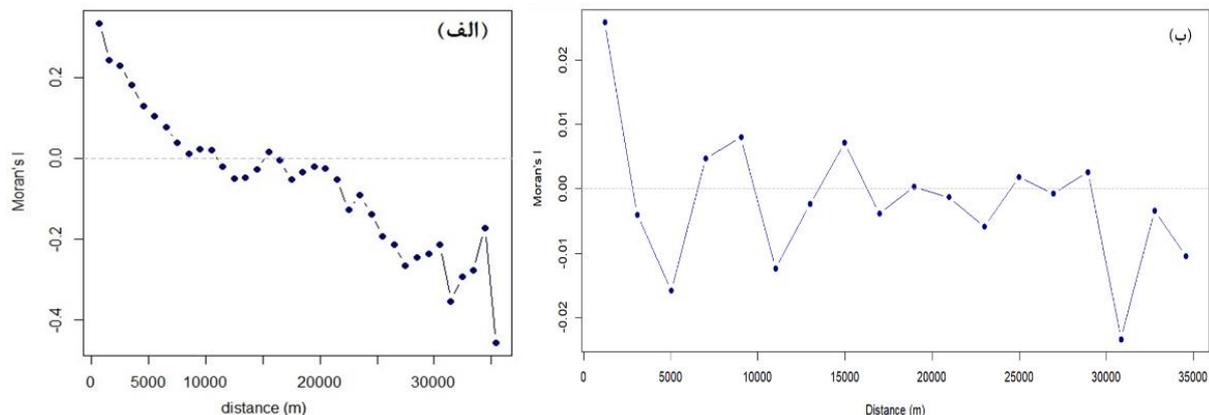
### متعادل سازی داده ها

مقایسه عملکرد مدل جنگل تصادفی در دو سناریوی استفاده از داده های نامتعادل و متعادل شده، نشان می دهد که متعادل سازی آماری داده ها نه تنها موجب بهبود دقت پیش بینی مدل می شود، بلکه می تواند به طور مؤثری ساختار مکانی خطاها را کاهش دهد. در مطالعات اخیر نیز تأکید شده است که عدم تعادل در توزیع داده ها، به ویژه در متغیرهای محیطی با توزیع چوله مانند کربن آلی خاک، می تواند منجر به سوگیری مدل نسبت به مقادیر غالب شده و موجب بروز الگوهای فضایی در باقیمانده ها گردد (Santos et al., 2020; Li et al., 2021; Ebrahimi et al., 2023). در مقابل، به کارگیری روش هایی مانند SMOTE، ضمن بهبود عملکرد عددی مدل، موجب کاهش خودهمبستگی مکانی در باقیمانده ها و افزایش قابلیت تعمیم پذیری مدل می شود (Chen et al., 2022; Wang et al., 2024).

در حالت استفاده از داده های نامتعادل، مقدار ریشه میانگین مربعات خطا (RMSE) برابر با ۰/۴ و ضریب تعیین ( $R^2$ ) معادل ۰/۳۹ به دست آمد (جدول ۲). همچنین مقدار شاخص خودهمبستگی مکانی Moran's I برای باقیمانده های مدل ۰/۳۰ محاسبه شد، که نشان دهنده وجود ساختار مکانی نسبتاً قوی در خطاهای مدل است. تحلیل کورلوگرام (correlogram) نیز بیانگر آن بود که مقدار شاخص موران با افزایش فاصله بین نقاط کاهش می یابد (شکل ۲)، که دلالت بر خودهمبستگی مکانی مثبت در فواصل کوتاه دارد. چنین رفتاری در باقیمانده ها نشان می دهد که مدل در بازنمایی صحیح الگوهای فضایی SOC ناکارآمد عمل کرده و اطلاعات مکانی به طور کامل از طریق مدل سازی دریافت نشده اند (Li et al., 2021; Nussbaum et al., 2018). در مقابل، پس از متعادل سازی داده ها با استفاده از روش SMOTE و اصلاح نسبت نمونه ها بین بیشترین و کمترین مقادیر، عملکرد مدل به طور محسوسی بهبود یافت. در این شرایط، مقدار RMSE به ۰/۳ کاهش یافته و  $R^2$  به ۰/۶ افزایش پیدا کرد. همچنین مقدار Moran's I به ۰/۰۳ رسید که بیانگر کاهش شدید خودهمبستگی مکانی در باقیمانده ها است. افزون بر این، روند تغییرات Moran's I در کورلوگرام برای داده های متعادل شده، نوساناتی اندک و حول صفر نشان داد که دلالت بر عدم وجود ساختار مکانی معنادار در خطاها دارد (Ebrahimi et al., 2023; Chen et al., 2022).

این نتایج با یافته های مطالعات مشابه هم راستا است که نشان داده اند عدم تعادل در داده ها منجر به تعصب مدل نسبت به مقادیر غالب شده و دقت پیش بینی در مقادیر با کمترین فراوانی را به شدت کاهش می دهد (Scholten, et al. 2023). علاوه بر این، این تعصب می تواند ساختارهای فضایی پنهان را در باقیمانده ها حفظ کند. در حالی که متعادل سازی داده ها به کمک روش هایی مانند SMOTE باعث توزیع متوازن تر خطاها در فضای جغرافیایی شده و کیفیت فضایی مدل را افزایش می دهد (Santos et al., 2020; Wang et al., 2024).

در مجموع، یافته های این پژوهش نشان می دهد که متعادل سازی آماری داده ها نه تنها باعث بهبود شاخص های عددی عملکرد مدل می شود، بلکه به طور معناداری خودهمبستگی مکانی باقیمانده ها را کاهش می دهد و مدل را از نظر مکانی پایاتر و قابل اطمینان تر می سازد.

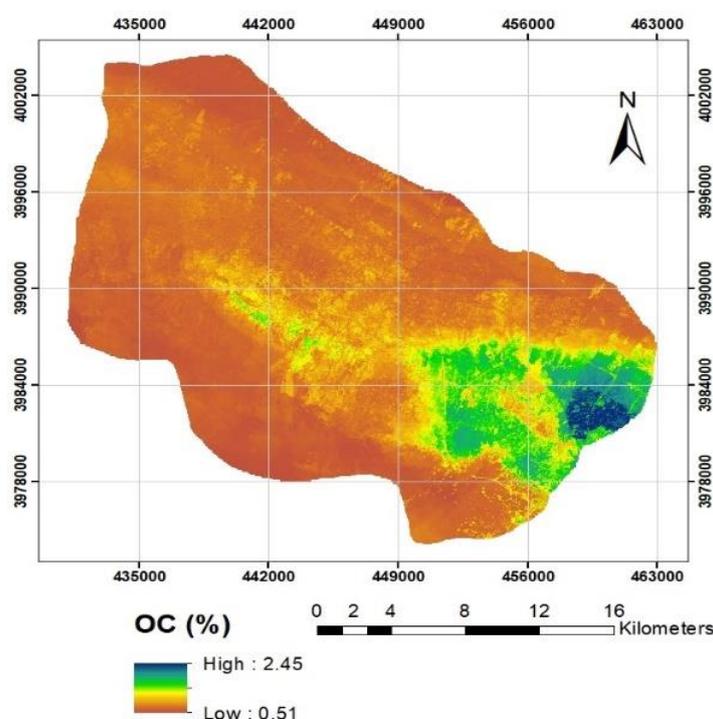


شکل ۲- روند تغییرات شاخص موران در داده های نامتعادل (الف) و متعادل شده (ب)

جدول ۲- مقایسه شاخص‌های ارزیابی مدل جنگل تصادفی در دو حالت داده‌های نامتعادل و متعادل شده

شاخص ارزیابی	داده‌های نامتعادل	داده‌های متعادل شده
ریشه میانگین مربعات خطا (RMSE)	۰/۴۰	۰/۳۰
ضریب تعیین ( $R^2$ )	۰/۳۹	۰/۶۰
شاخص خودهمبستگی مکانی (Moran's I)	۰/۳۰	۰/۰۳
روند شاخص موران با فاصله	کاهش تدریجی	ثابت و نزدیک به صفر
توزیع فضایی باقیمانده‌ها	خودهمبسته	تصادفی

در نهایت، نقشه توزیع مکانی کربن آلی خاک براساس داده‌های متعادل شده تهیه شد (شکل ۳). در قسمت شمال شرق، شرق و مرکز منطقه مورد مطالعه بیشترین مقدار کربن آلی مشاهده شد که به علت وجود کاربری کشت آبی، دیم و همچنین وجود مراتع غیر شور می‌باشد در حالیکه در قسمت جنوب و جنوب شرق، به علت وجود دشت‌های شور و عدم وجود پوشش گیاهی، مقدار کربن آلی کاهش یافته است.



شکل ۳- توزیع مکانی کربن آلی خاک براساس داده‌های متعادل شده

### نتیجه‌گیری

نتایج این مطالعه نشان داد که عدم تعادل در توزیع داده‌ها می‌تواند منجر به کاهش دقت پیش‌بینی و حفظ ساختارهای مکانی در خطاهای مدل گردد. استفاده از روش SMOTE برای متعادل‌سازی آماری داده‌ها، نه تنها عملکرد مدل جنگل تصادفی را از نظر شاخص‌های RMSE و  $R^2$  بهبود بخشید، بلکه به‌طور معناداری خودهمبستگی مکانی باقیمانده‌ها را کاهش داد. مقدار شاخص موران و پایداری کورلوگرام حول صفر در مدل متعادل‌شده، حاکی از کارایی بیشتر مدل در بازنمایی الگوهای فضایی SOC بود. بنابراین، در مطالعات آینده توصیه می‌شود پیش از مدل‌سازی، وضعیت تعادل داده‌ها بررسی و در صورت نیاز، اصلاح گردد تا نتایج دقیق‌تر و معتبرتری حاصل شود. ترکیب روش‌های یادگیری ماشین با تکنیک‌های آماری متعادل‌سازی، می‌تواند به عنوان رویکردی مؤثر در مدل‌سازی مکانی متغیرهای خاکی مورد استفاده قرار گیرد.

### فهرست منابع

- Adhikari, K., & Hartemink, A. E. (2016). Linking soils to ecosystem services — A global review. *Geoderma*, 262, 101–111.
- Brus, D. J., Kempen, B., & Heuvelink, G. B. M. (2011). Sampling for validation of digital soil maps. *European Journal of Soil Science*, 62(3), 394–407.
- Chen, Z., et al. (2022). Mitigating spatial bias in random forest models with balanced sampling. *Remote Sensing*, 14(12), 2809.
- Ebrahimi, A., et al. (2023). Improving prediction of soil properties by handling imbalanced data in machine learning models. *Environmental Modelling & Software*, 161, 105594.
- Hengl, T., et al. (2017). SoilGrids250m: Global gridded soil information based on machine learning. *PLOS ONE*, 12(2), e0169748.
- Hengl, T., Heuvelink, G. B. M., & Stein, A. (2007). Comparison of kriging with external drift and regression-kriging. *Geoderma*, 140(1–2), 480–496.
- Heung, B., Bulmer, C. E., & Schmidt, M. G. (2016). Predictive soil parent material mapping at a regional-scale: A Random Forest approach. *Geoderma*, 279, 68–83.
- Lal, R. (2004). Soil carbon sequestration impacts on global climate change and food security. *Science*, 304(5677), 1623–1627.
- Li, W., et al. (2021). Effect of sample imbalance on digital soil mapping accuracy. *Catena*, 203, 105309.
- Minasny, B., et al. (2017). Soil carbon 4 per mille. *Geoderma*, 292, 59–86.
- Nussbaum, M., Papritz, A., Baltensweiler, A., & Walthert, L. (2018). Evaluating digital soil mapping models for soil organic carbon at the national scale. *Geoderma*, 318, 45–58.
- Santos, H. G., et al. (2020). Addressing data imbalance in soil mapping using synthetic oversampling. *Geoderma*, 373, 114399.
- Scholten, T., et al. (2023). Coping with imbalanced data problem in digital mapping of soil functional classes using SMOTE. *European Journal of Soil Science*, 74(2), e13368.
- Wadoux, A. M., et al. (2021). Machine learning for digital soil mapping: Applications, challenges, and perspectives. *Earth-Science Reviews*, 223, 103824.
- Walkley, A., & Black, I. A. (1934). An examination of the Degtjareff method for determining soil organic matter and a proposed modification of the chromic acid titration method. *Soil Science*, 37(1), 29–38.
- Wang, Y., et al. (2024). Spatial implications of data distribution in machine learning-based soil mapping. *Soil Systems*, 8(1), 21.

## Evaluating the Influence of Data Imbalance on Accuracy and Spatial Autocorrelation in Soil Organic Carbon Modeling

Azam Jafari <sup>1</sup>, Fereydoon Sarmadian <sup>2\*</sup>, Zahra Rasaei <sup>2</sup>

<sup>1</sup>Soil Science Department, Faculty of Agriculture, Shahid Bahonar University of Kerman, Kerman, Iran

<sup>2</sup>Soil Science Department, Faculty of Agricultural, University College of Agriculture and Natural Resources, University of Tehran, Karaj, Iran

### Abstract

Data imbalance in environmental datasets presents a significant challenge in the spatial modeling of soil properties, particularly soil organic carbon (SOC). This study explores the effect of data balancing on model accuracy and the spatial structure of residuals across a 60,000-hectare area in the Qazvin Plain, Iran. A total of 280 surface samples (0–30 cm depth) were analyzed for SOC content, alongside a suite of environmental covariates including topographic metrics, vegetation indices, soil moisture parameters, and climatic variables. Key predictors were identified using recursive feature elimination (RFE), and random forest (RF) models were implemented under two scenarios: imbalanced data and balanced data via the synthetic minority oversampling technique (SMOTE). Model performance was evaluated using root mean square error (RMSE), coefficient of determination ( $R^2$ ), and Moran's I index of spatial autocorrelation. The SMOTE-adjusted RF model exhibited improved predictive accuracy (RMSE = 0.3,  $R^2$  = 0.60) relative to the imbalanced case (RMSE = 0.4,  $R^2$  = 0.40). Moreover, Moran's I for model residuals dropped substantially from 0.30 to 0.03 following data balancing. Correlogram analysis corroborated this reduction, indicating weakened spatial structure in residuals. These results highlight the value of statistical data balancing in enhancing both numerical precision and spatial consistency in SOC prediction models.

**Keywords:** Moran's I, SMOTE, Spatial modeling