



Machine Learning–Based Downscaling of GLDAS Surface Soil Moisture and Performance Evaluation Across Climatic Regions of Iran

Seyedeh Mahsa Mousavi Reineh,¹ Saman Javadi,^{2*} Hossein Yousefi,³ Aminreza Neshat,⁴ Alireza Massah Bavani,⁵

1. PhD candidate, Department of Water Engineering, Aburaihan Campus, University of Tehran, Iran. Email: mahsa.moosavi.rei@ut.ac.ir
2. Professor, Department of Water Engineering, Aburaihan Campus, University of Tehran, Iran. Email: javadis@ut.ac.ir
3. Professor, Faculty of Interdisciplinary Science and Technology, University of Tehran, Iran. Email: hosseinyousefi@ut.ac.ir
4. Assistant professor, Department of Civil Engineering, SR.C., Islamic Azad University, Tehran, Iran. Email: neshat.aminreza@srbiau.ac.ir
5. Associate Professor, Department of Water Engineering, Aburaihan Campus, University of Tehran. Email: armassah@ut.ac.ir

Article Info

Article type:
Research Article

Article history:

Received 14 October 2025
Revised 14 November 2025
Accepted 16 December 2025
Available online 22 December 2025

Keywords:

Downscaling,
Soil Moisture,
GLDAS Data,
Machine Learning,
CatBoost.

ABSTRACT

Research Topic: Due to scarce in-situ soil moisture observations and coarse reanalysis resolution, downscaling GLDAS surface soil moisture (SSM) using machine learning is essential.

Objective: This study evaluates the performance of machine learning algorithms for downscaling GLDAS SSM across Iran's diverse climatic zones.

Method: In this study, Random Forest (RF), XGBoost, CatBoost, and LightGBM algorithms were employed for downscaling process. Model inputs comprised station-based climatic variables (minimum and maximum temperature, precipitation, and evaporation) and spatial attributes on a monthly scale over a 31-year period. The models were trained using GLDAS surface soil moisture values, producing downscaled soil moisture at a higher spatial resolution as output. In-situ soil moisture observations were utilized exclusively for independent validation. The dataset was partitioned into training (80%) and testing (20%) sets based on a chronological order.

Results: CatBoost demonstrated a strong ability to capture nonlinear soil moisture patterns, achieving coefficients of determination exceeding 0.73 across all climatic zones considered in this study. Although model accuracy varied depending on climatic characteristics and the spatial distribution of reference data, CatBoost was identified as an efficient algorithm for soil moisture prediction over the study area due to its high generalization capability and satisfactory performance in independent national-scale validation ($R^2 = 0.607$, RMSE = 4.286, Bias = 2.131).

Conclusions: The findings indicate that machine learning frameworks, particularly CatBoost, offer a reliable approach for downscaling GLDAS SSM. Given its high generalization capability across varying hydro-climatic conditions, CatBoost is recommended for enhancing drought monitoring and water resource management in data-scarce regions.

Cite this article: Mousavi Reineh, S., Javadi, S., Yousefi, H., Neshat, A., Massah Bavani, A. (2026). Machine Learning–Based Downscaling of GLDAS Surface Soil Moisture and Performance Evaluation Across Climatic Regions of Iran. *Journal of ECOHYDROLOGY*, 12(4), 982-1004. <http://doi.org/10.22059/ije.2026.404312.1891>



© Seyyedeh Mahsa Mousavi Reineh, Saman Javadi, Hossein Yousefi, Aminreza Neshat, Alireza Massah Bavani.
Publisher: University of Tehran Press.

DOI: <http://doi.org/10.22059/ije.2026.404312.1891>

Introduction

Soil moisture is a pivotal variable in eco-hydrological processes. Despite its critical role in drought monitoring, flood prediction, and agricultural productivity, the scarcity of in situ measurements and the limited spatio-temporal coverage of monitoring stations have led researchers to rely on reanalysis products, such as GLDAS. However, the coarse spatial resolution of GLDAS data (ranging from 0.25 to 1) remains a significant bottleneck for local-scale applications and precision water resource management. To address this limitation, spatial downscaling has emerged as an essential tool. While machine learning (ML) algorithms offer robust capabilities for capturing complex non-linear relationships between climatic variables and soil moisture, their performance across Iran's diverse climatic zones requires rigorous evaluation. This study addresses the challenge of model-based downscaling of GLDAS surface soil moisture at a regional scale. The primary objective is to generate a high-resolution, refined version of GLDAS data. By leveraging innovative ML algorithms, this research provides a scalable framework to compensate for the lack of ground-based monitoring in many parts of Iran, offering a reliable substitute for enhanced hydrological and agricultural decision-making.

Method

This study employs a model-based downscaling framework across Iran's diverse climates to enhance the spatial resolution of GLDAS surface soil moisture (0–10 cm). The methodology integrates three datasets: station-based climatic variables (minimum and maximum air temperature, mean air temperature at 2 m, precipitation, evaporation, solar radiation, and soil surface temperature) and spatial attributes on a monthly scale over a 31-year period (1994–2024) from national stations, GLDAS soil moisture as the target variable, and in-situ measurements reserved strictly for independent validation. Four machine learning algorithms, namely Random Forest, XGBoost, CatBoost, and LightGBM, were trained for downscaling process. To ensure robust temporal generalization and prevent data leakage, a chronological 80/20 split was implemented for training and testing. Model interpretability was further analyzed using the SHAP approach to quantify variable contributions, while performance was rigorously evaluated through R², RMSE, and Bias metrics, ensuring the models applicability in regions with sparse observational data.

Results

Data preprocessing revealed that climatic extremes in Iran's arid regions are inherent physical features rather than outliers. While missing values for temperature and precipitation (<5%) were successfully reconstructed using regression, solar radiation was excluded due to a high missing rate (54%). To mitigate multi-collinearity ($r > 0.90$), T_{max} and T_{min} were selected as primary temperature predictors. Among the evaluated algorithms, CatBoost demonstrated superior national performance (R² = 0.854, RMSE = 2.345), outperforming Random Forest, LightGBM, and XGBoost—the latter showing significant overfitting. SHAP analysis identified T_{max} (35.07%) and precipitation (24.66%) as the dominant drivers, with station elevation (13.71%) playing a critical role in local-scale soil moisture distribution. Although internal validation indicated a systematic negative bias (underestimation), independent validation against ground-truth station data yielded an R² of 0.607 and a positive bias (2.131 kg/m²). This discrepancy is attributed to spatial scale mismatch between satellite-derived grid data and point-based measurements. Regional analysis confirmed CatBoost's stability across diverse climates, particularly in desert (BWk) and semi-arid (BSk) zones. Overall, the CatBoost model proved highly capable of capturing non-linear soil moisture dynamics and seasonal fluctuations across Iran's varied hydro-climatic regions.

Conclusions

This study demonstrates that the CatBoost algorithm is the most robust and stable tool for downscaling GLDAS soil moisture data across Iran's diverse hydro-climatic regions, followed by Random Forest and LightGBM. While XGBoost exhibited superior training performance, it failed to generalize due to significant overfitting. Independent validation against in-situ observations yielded a moderate national R^2 of 0.607, with regional accuracy peaking at 0.89 in areas with higher data density and climatic homogeneity. The observed systematic positive bias (2.131 kg/m²) and reduced sensitivity to extremes are attributed to the inherent tendency of machine learning models to converge toward mean values and the spatial scale mismatch between grid and point data. Although the current model requires local refinement for field-scale applications, it provides a powerful framework for regional drought monitoring and water resource management, effectively reproducing complex spatio-temporal patterns of soil moisture at a high spatial resolution.

Author Contributions

All authors participated in the design and development of the research project. The initial version of the article was written by Mahsa Mousavi, and all authors reviewed subsequent versions of the article and provided their views. Finally, all authors read and approved the final version of the article.

Data Availability Statement

It is a part of Ph.D thesis which is undergoing in the Aburaihan Campus, University of Tehran.

Acknowledgements

We gratefully acknowledge the Iran Meteorological Organization for making meteorological data available to the researchers of this study.

Ethical considerations

The authors avoided from data fabrication and falsification.

Funding

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

Conflict of interest

The authors declare no conflict of interest.

ریزمقیاس‌نمایی رطوبت سطحی خاک GLDAS با استفاده از الگوریتم‌های یادگیری ماشین و

ارزیابی عملکرد آن‌ها در اقلیم‌های مختلف ایران

سیده مهسا موسوی رینه^۱، سامان جوادی^{۲*}، حسین یوسفی^۳، امین رضا نشاط^۴، علیرضا مساح بوانی^۵

۱. دانشجوی دکتری مهندسی منابع آب، گروه مهندسی آب، پردیس ابوریحان، دانشگاه تهران، ایران. ایمیل: mahsa.moosavi.rei@ut.ac.ir

۲. استاد، گروه مهندسی آب، پردیس ابوریحان، دانشگاه تهران، تهران، ایران. ایمیل: javadis@ut.ac.ir

۳. استاد، گروه انرژی‌های نو و محیط زیست، دانشکده علوم و فناوری‌های میان‌رشته‌ای، دانشگاه تهران، تهران، ایران. ایمیل: hosseinyousefi@ut.ac.ir

۴. استادیار، گروه مهندسی عمران، واحد علوم و تحقیقات، دانشگاه آزاد اسلامی، تهران، ایران. ایمیل: neshat.aminreza@srbiau.ac.ir

۵. دانشیار، گروه مهندسی آب، پردیس ابوریحان، دانشگاه تهران، تهران، ایران. ایمیل: armassah@ut.ac.ir

چکیده

اطلاعات مقاله

موضوع: ریزمقیاس‌نمایی داده‌های رطوبت سطحی خاک پایگاه داده GLDAS به دلیل محدودیت داده‌های مشاهده‌ای رطوبت

خاک و دقت مکانی پایین داده‌های بازتحلیل، نیازمند بهره‌گیری از الگوریتم‌های یادگیری ماشین است.

هدف: آموزش و ارزیابی عملکرد الگوریتم‌های یادگیری ماشین در ریزمقیاس‌نمایی رطوبت سطحی خاک GLDAS در نواحی اقلیمی مختلف ایران.

روش تحقیق: الگوریتم‌های Random Forest, XGBoost, CatBoost و LightGBM برای ریزمقیاس‌نمایی به کار گرفته شدند. ورودی مدل‌ها شامل متغیرهای اقلیمی ایستگاهی (دما، حد اقل و حداکثر، بارش و تبخیر) و ویژگی‌های مکانی ایستگاه‌ها در مقیاس ماهانه و طی دوره‌ای ۳۱ ساله بود. مدل‌ها با استفاده از مقادیر رطوبت سطحی خاک GLDAS آموزش داده شدند و خروجی آن‌ها رطوبت خاک ریزمقیاس‌شده در تفکیک مکانی بالاتر است. داده‌های رطوبت خاک مشاهده‌ای فقط برای اعتبارسنجی مستقل استفاده شدند. تقسیم داده‌ها به مجموعه‌های آموزش (۸۰ درصد) و آزمون (۲۰ درصد) بر اساس ترتیب زمانی انجام شد.

یافته‌ها: الگوریتم CatBoost با ارائه ضریب تعیین فراتر از ۰/۷۳ در تمامی نواحی اقلیمی مورد مطالعه، توانمندی بالایی در استخراج الگوهای غیرخطی رطوبت خاک دارد. علی‌رغم وابستگی دقت به ویژگی‌های اقلیمی و پراکندگی مکانی داده‌ها، CatBoost به دلیل قابلیت تعمیم‌پذیری بالا و عملکرد مطلوب در اعتبارسنجی مستقل در سطح ملی ($R^2 = 0.607$, RMSE = 4.286, Bias = 2.131)، الگوریتمی کارآمد برای پیش‌بینی رطوبت خاک در تمامی پهنه‌های مورد مطالعه شناخته شد.

نتیجه‌گیری: نتایج نشان می‌دهد الگوریتم‌های یادگیری ماشین، به‌ویژه CatBoost، چارچوبی کارآمدی برای ریزمقیاس‌نمایی مدل‌محور رطوبت سطحی خاک GLDAS در مقیاس منطقه‌ای و ملی فراهم می‌کنند و می‌توانند در پایش خشکسالی و مدیریت منابع آب در مناطق فاقد داده‌های مشاهده‌ای کافی مورد استفاده قرار گیرند.

نوع مقاله:

مقاله پژوهشی

تاریخ دریافت: ۱۴۰۴/۰۷/۲۲

تاریخ بازنگری: ۱۴۰۴/۰۸/۱۳

تاریخ پذیرش: ۱۴۰۴/۰۹/۲۵

تاریخ انتشار: ۱۴۰۴/۱۰/۰۱

کلیدواژه‌ها:

ریزمقیاس‌نمایی،

رطوبت خاک،

GLDAS

یادگیری ماشین،

CatBoost

استناد: موسوی رینه، سیده مهسا؛ جوادی، سامان؛ یوسفی، حسین؛ نشاط، امین رضا؛ مساح بوانی، علیرضا. ریزمقیاس‌نمایی رطوبت سطحی خاک GLDAS با استفاده از

الگوریتم‌های یادگیری ماشین و ارزیابی عملکرد آن‌ها در اقلیم‌های مختلف ایران. مجله اکوهیدرولوژی، ۱۲(۴)، ۹۸۲-۱۰۰۴.

<http://doi.org/10.22059/ije.2026.404312.1891>

ناشر: انتشارات دانشگاه تهران. © سیده مهسا موسوی رینه، سامان جوادی، حسین یوسفی، امین رضا نشاط، علیرضا مساح بوانی.



مقدمه

رطوبت خاک یکی از عوامل کلیدی در فرایندهای بوم‌هیدرولوژیکی محسوب می‌شود که با تأثیرگذاری بر آلودگی سطح، رشد پوشش گیاهی و میزان تبخیر، نقش مهمی در تبادل آب و انرژی میان سطح زمین و جو ایفا می‌کند (Chahine, 1992; Shang et al., 2003; ZHANG et al., 2007). علی‌رغم اهمیت این پارامتر در پایش خشکسالی، سیلاب و بهره‌وری کشاورزی، محدودیت اطلاعات میدانی، هزینه‌بر بودن اندازه‌گیری و عدم پوشش مکانی - زمانی کافی ایستگاه‌ها، پژوهشگران را به استفاده از داده‌های بازتحلیل مانند سیستم جهانی یکپارچه‌سازی داده‌های زمین^۱ (GLDAS) سوق داده است. با این حال، چالش اساسی داده‌های GLDAS، تفکیک مکانی پایین آن‌ها (۰/۲۵ تا ۱ درجه) است که برای کاربردهای مقیاس محلی و مدیریت دقیق منابع آب کافی نیست (Batchu et al., 2023; Senanayake et al., 2024). در این راستا، ریزمقیاس‌نمایی به عنوان ابزاری ضروری برای بهبود دقت مکانی این داده‌ها مطرح می‌شود. بهره‌گیری از الگوریتم‌های یادگیری ماشین می‌تواند امکان استخراج الگوهای پیچیده بین متغیرهای اقلیمی و رطوبت خاک را فراهم کند و دقت برآوردها را افزایش دهد. با توجه به تنوع اقلیمی گسترده ایران، بررسی کارایی این روش‌ها در شرایط اقلیمی متفاوت اهمیت زیادی دارد (Park et al., 2017).

مسئله اصلی این پژوهش، ریزمقیاس‌نمایی مدل محور رطوبت سطحی خاک پایگاه داده GLDAS در مقیاس منطقه‌ای با استفاده از الگوریتم‌های یادگیری ماشین است. داده‌های GLDAS اگرچه پوشش زمانی بلندمدت و پیوسته‌ای دارند، اما به دلیل تفکیک مکانی درشت، برای بسیاری از کاربردهای منطقه‌ای و محلی در مطالعات هیدرولوژیکی و کشاورزی به صورت مستقیم قابل استفاده نیستند. در مقابل، داده‌های رطوبت خاک مشاهده‌ای دارای تفکیک مکانی بالا هستند، اما به دلیل محدودیت تعداد ایستگاه‌ها، ناپیوستگی زمانی و پوشش مکانی نامتوازن، امکان استفاده از آن‌ها به عنوان منبع اصلی مدل‌سازی در مقیاس ملی را ندارند. بر این اساس، در این پژوهش هدف، پیش‌بینی مستقیم رطوبت خاک اندازه‌گیری شده در ایستگاه‌ها نیست، بلکه یادگیری رابطه بین متغیرهای اقلیمی ایستگاهی و ویژگی‌های مکانی با مقادیر رطوبت سطحی خاک استخراج شده از داده‌های GLDAS و تولید نسخه‌ای ریزمقیاس‌شده و اصلاح‌شده از این داده‌ها در مقیاس منطقه‌ای است.

پیشینه پژوهش

GLDAS یکی از محصولات برجسته ناسا و سازمان ملی اقیانوسی و جوی آمریکا^۲ است که با هدف ادغام داده‌های رصدی زمینی و داده‌های سنجش از دور و ورودی‌های مختلف مدل‌های فیزیکی سطح زمین و همچنین، داده‌های مشاهداتی ایستگاه‌های زمینی (San Liang & Zhang, 2018) به منظور تخمین و ارائه وضعیت متغیرهای سطحی از جمله رطوبت خاک، دما، برف، رواناب، تبخیر و تعرق و غیره در مقیاس جهانی توسعه یافته است (Rodell et al., 2004; Xia et al., 2012). مأموریت اصلی GLDAS تولید داده‌های پیوسته و سازگار مکانی - زمانی با پوشش جهانی از متغیرهای سطحی خاک و اراضی، با هدف پشتیبانی از مدل‌سازی هیدرولوژیکی و اقلیمی، هشدار سیل، پایش خشکسالی، مدیریت منابع آب و تغذیه مدل‌های عددی هواشناسی است. از این رو، در مناطقی که داده‌های مشاهده‌ای رطوبت خاک در دسترس نیست، بهره‌گیری از داده‌های پایگاه GLDAS به عنوان منبعی معتبر به طور گسترده رایج است (Hasan et al., 2024; Koster & Suarez, 1992; Liang, 1994; Rodell et al., 2009). مجموعه داده‌های مرتبط با محتوای آب خاک GLDAS در مطالعات متعددی برای تخمین ذخیره آب ناحیه غیر اشباع استفاده شده است، در واقع جایی که در بسیاری مناطق مختلف جهان اندازه‌گیری میدانی آن انجام نمی‌شود (Rodell et al., 2007; Strassberg et al., 2007).

1. Global Land Data Assimilation System

2 NOAA

(2007; Tiwari et al., 2009). به منظور حل مشکل کمبود داده‌های مشاهداتی رطوبت خاک، GLDAS اطلاعاتی قابل اتکا برای مطالعه رطوبت خاک فراهم می‌آورد (Rodell et al., 2004). جدول ۱ اطلاعات مربوط به داده‌های مستخرج از پایگاه داده GLDAS را نشان می‌دهد.

جدول ۱. ویژگی‌های کلیدی داده‌های رطوبت خاک مستخرج از پایگاه داده (GIOVANNI, 2024) GLDAS

ویژگی	جزئیات
تفکیک‌پذیری مکانی	معمولاً از ۰/۲۵ در ۰/۲۵ درجه قوسی (۲۵×۲۵ کیلومتر) تا ۱ در ۱ درجه قوسی (۱۰۰×۱۰۰ کیلومتر)
نمایه عمقی	رطوبت در اعماق مختلف (۰-۱۰، ۱۰-۴۰، ۴۰-۱۰۰، ۱۰۰-۲۰۰ سانتی‌متر)
تفکیک زمانی	۳ ساعته، روزانه و ماهانه
نوع محصول	مدل‌سازی شبیه‌سازی شده بر اساس مدل‌های سطح زمین (مانند NOAH, VIC, MOSAIC, CLSM)
محدوده زمانی	از سال ۱۹۴۸ تا به امروز، بسته به مدل
واحد اندازه‌گیری	kg/m ²

مدل‌های سطح زمین در GLDAS از مدل‌های مختلفی (مانند NOAH, VIC, MOSAIC و CLSM) برای شبیه‌سازی دینامیک انتقال رطوبت در خاک استفاده می‌کنند. این مدل‌ها اگرچه در نحوه محاسبه پارامترهایی مثل معادل آبی برف یا آب ذخیره‌شده در پوشش گیاهی تفاوت‌هایی دارند، اما همگی بر پایه اصول فیزیکی انتقال جرم و انرژی استوار هستند (Rodell et al., 2004). در این میان، مطالعات متعددی نشان داده‌اند مدل NOAH (به‌ویژه نسخه ۲.۱) به دلیل دقت زیاد در شبیه‌سازی فرایندهای هیدرولوژیکی، کارایی بهتری نسبت به سایر مدل‌ها در این سامانه دارد (Ali et al., 2023; L. Chen et al., 2019). بر همین اساس، این پژوهش بر خروجی‌های این مدل متمرکز شده است (Y. Chen et al., 2013; Bi et al., 2016). از جمله این مطالعات می‌توان به مطالعات (Cai et al., 2017; Fatolazadeh et al., 2020; Jiménez et al., 2011) اشاره کرد.

چالش اصلی داده‌های GLDAS تفکیک مکانی درشت (۲۵ تا ۱۰۰ کیلومتر) است که استفاده از آن‌ها را در مقیاس‌های خرد مانند مدیریت حوضه‌های آبریز یا کشاورزی دقیق با محدودیت مواجه می‌کند (Hasan et al., 2024; Yin et al., 2018). برای رفع این چالش، تکنیک‌های ریزمقیاس‌نمایی آماری و یادگیری ماشین به عنوان روش‌هایی به‌صرفه مطرح شده‌اند. بررسی مرور منابع نشان داد روش یادگیری ماشین جنگل تصادفی^۱ با دقت بالایی توانایی انجام این کار را دارد (L. Chen et al., 2019; Park et al., 2017). با مقایسه مدل‌های مبتنی بر تقویت‌گرادیان مانند کت بوست^۲، ایکس جی بوست^۳ و لایت جی بی ام^۴ برای ریزمقیاس‌نمایی داده‌های رطوبت خاک در چین، دریافتند که ترکیب این مدل‌ها می‌تواند دقت برآورد را در مناطق با شرایط زمینی متنوع از وضوح مکانی اولیه ۹ کیلومتر به وضوح مکانی ۱ کیلومتر به طور محسوسی افزایش دهد (J. Xu et al., 2024). در میان مطالعات داخلی سلطانی و همکاران (۱۴۰۴) الگوریتم یادگیری ماشین جنگل تصادفی را برای ریزمقیاس‌نمایی و شناسایی مؤثرترین مدل‌های رطوبت خاک CMIP6 (با مجموع ۹۸ درصد دقت) استفاده کردند. نیک‌داد، محمدی قلعه‌نی و مقدسی (۱۴۰۲) به ارزیابی داده‌های رطوبت خاک از پایگاه‌های مختلف جهانی در اقلیم‌های مختلف ایران پرداختند. عبده کلاهچی و همکاران (۱۴۰۲) در استان لرستان، دقت داده‌های رطوبت خاک حاصل از پایگاه‌های GLDAS، ESA و سنجنده SMAP در مقایسه با داده‌های زمینی را ارزیابی کردند، نتایج نشان داد این محصولات رطوبت خاک از دقت قابل قبولی برخوردارند، این ارزیابی در صورتی انجام شد که روی این داده‌ها فرایند

1 Random Forest

2 CatBoost

3 XGBoost

4 LightGBM

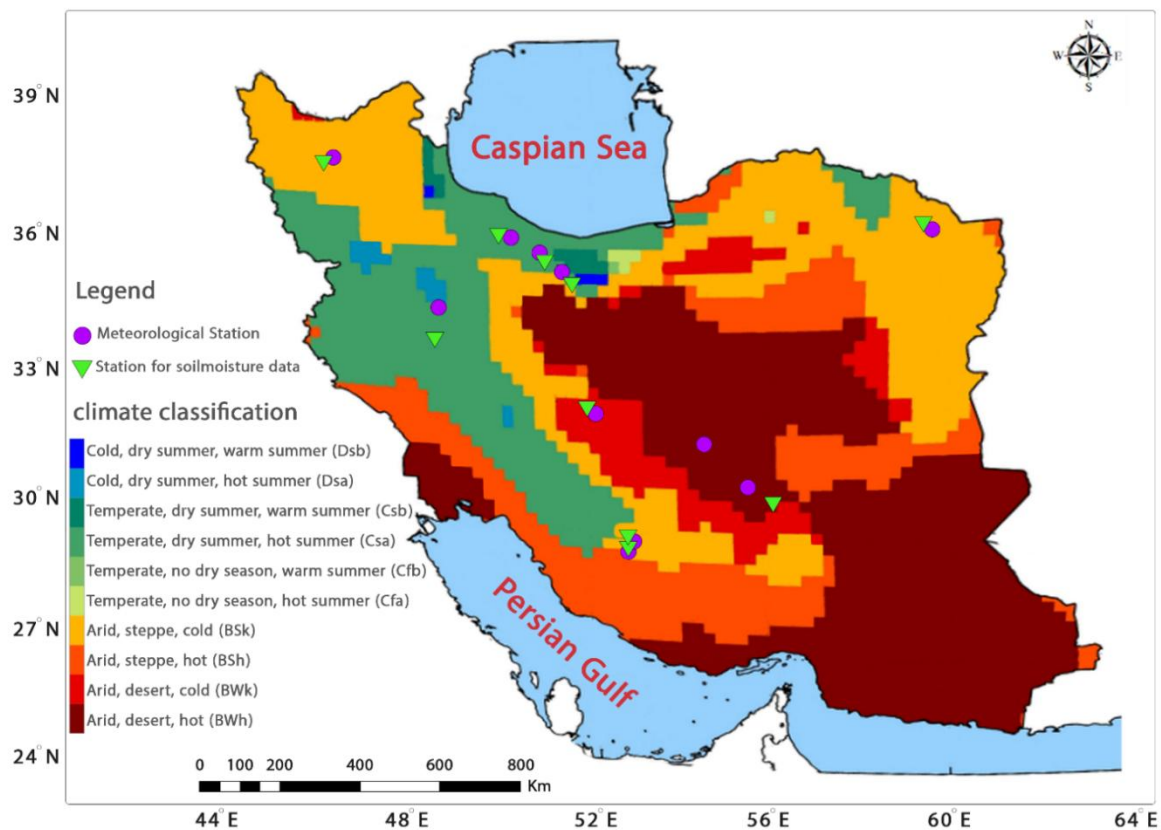
ریزمقیاس‌نمایی انجام نگرفت و یکی از پیشنهاد‌های نویسندگان برای تحقیقات آتی ریزمقیاس‌نمایی این داده‌ها در سطح کشور ایران بوده است. در مطالعات داخلی و خارجی متعددی از داده‌های رطوبت خاک حاصل از پایگاه داده GLDAS در برآورد داده‌های مفقود در چاه‌های مشاهداتی (Evans et al., 2020)، ارزیابی و تخمین تغییرات ذخیره آب زیرزمینی در گذشته (Wu et al., 2019)، ارزیابی تغییرات کمی رطوبت سطحی خاک (Zuo et al., 2019) و ارزیابی خشکسالی (Amini et al., 2023; Nouraki et al., 2024) ارتباط رطوبت خاک با میزان تولید کشاورزی و امنیت غذایی (Park et al., 2017; Senanayake et al., 2024) استفاده شده است، اما تحقیقی که به طور مشخص ریزمقیاس‌نمایی داده‌های رطوبت خاک پایگاه داده GLDAS را با استفاده از الگوریتم‌های یادگیری ماشین در ایران انجام دهد یافت نشد. از آنجا که در بیشتر مناطق ایران داده رطوبت خاک اندازه‌گیری نمی‌شود، یافتن جایگزینی مناسب برای رفع این نقیصه در مدیریت آب بسیار راه‌گشا و کمک‌کننده خواهد بود. این تحقیق با هدف ریزمقیاس‌نمایی داده‌های رطوبت خاک بزرگ‌مقیاس GLDAS با استفاده از الگوریتم‌های نوین یادگیری ماشین انجام شده است. استفاده از الگوریتم‌های مختلف یادگیری ماشین برای ریزمقیاس‌نمایی رطوبت خاک مبتنی بر پایگاه داده GLDAS و مقایسه این الگوریتم‌ها به لحاظ دقت و کارایی گامی نوآورانه در ریزمقیاس‌نمایی داده‌های رطوبت خاک پایگاه داده GLDAS در نواحی اقلیمی مختلف ایران به حساب می‌آید.

روش‌شناسی پژوهش

در چارچوب این مسئله، الگوریتم‌های یادگیری ماشین با استفاده از متغیرهای هواشناسی و ویژگی‌های مکانی ایستگاه‌ها آموزش داده شدند تا بتوانند الگوی مکانی - زمانی تغییرات رطوبت سطحی خاک GLDAS را در اقلیم‌های مختلف کشور بازتولید و در مقیاس‌های مکانی کوچک‌تر تعمیم دهند. به این ترتیب، خروجی مدل‌ها نمایانگر ریزمقیاس‌نمایی مدل‌محور داده‌های GLDAS است و نه بازسازی مقادیر نقطه‌ای واقعی رطوبت خاک در محل ایستگاه‌ها. این رویکرد امکان بهره‌گیری از پوشش زمانی بلندمدت داده‌های بازتحلیل را در کنار اطلاعات اقلیمی ایستگاهی فراهم کرده و بستری برای تولید محصولات رطوبت خاک با تفکیک مکانی بالاتر در مناطق فاقد داده‌های مشاهده‌ای کافی ایجاد می‌کند.

منطقه مورد مطالعه

بر اساس طبقه‌بندی کوپن - گایگر، ایران عمدتاً دارای اقلیم‌های خشک و نیمه‌خشک است. با توجه به این طبقه‌بندی کشور ایران ده تیپ اقلیمی دارد که اقلیم‌های BWh (بیابانی گرم)، BSk (نیمه‌خشک سرد)، Csa (مدیترانه‌ای با تابستان گرم)، BSh (نیمه‌خشک گرم) و BWk (بیابانی سرد) در مجموع بیش از ۹۸ درصد مساحت کشور را پوشش می‌دهند. این نشان می‌دهد الگوی اقلیمی ایران تحت سلطه شرایط خشک و مدیترانه‌ای قرار دارد (Raziei, 2022).



شکل ۱. موقعیت ایستگاه‌های هواشناسی و هواشناسی کشاورزی مورد مطالعه در پژوهش حاضر با توجه به طبقه‌بندی اقلیمی کوپن - گایگر (منبع طبقات اقلیمی Raziei, 2022)

اطلاعات کامل ایستگاه‌های مورد بررسی در این مطالعه در جدول ۲ آورده شده است.

جدول ۲. مشخصات جغرافیایی ایستگاه‌های هواشناسی مورد استفاده در تحقیق (منبع طبقات اقلیمی کوپن-گایگر (Raziei, 2022))

ایستگاه هواشناسی	استان	نوع داده دریافتی	طول جغرافیایی	عرض جغرافیایی	ارتفاع	طبقه اقلیمی
تهران فرودگاه مهرآباد	تهران	هواشناسی	۵۱/۳۰۹	۳۵/۶۹۳	۱۱۹۱	Bsk
ورامین	تهران	مرجع رطوبت خاک	۵۱/۶۵	۳۵/۳۱۷	۹۷۳	Bsk
اصفهان فرودگاه	اصفهان	هواشناسی	۵۱/۸۶۳	۳۲/۴۹۳	۱۵۵۱/۹	BWk
کبوترآباد	اصفهان	مرجع رطوبت خاک	۵۱/۸۳۴	۳۲/۵۱۵	۱۵۴۲/۵	BWk
هشتگرد	البرز	هواشناسی	۵۰/۷۴۷	۳۶/۰۰۶	۱۶۱۲/۹	Csa
کرج	البرز	مرجع رطوبت خاک	۵۰/۹۵۴	۳۵/۸۰۶	۱۲۹۲/۹	Csa
انار	کرمان	هواشناسی	۵۵/۲۵۰	۳۰/۸۸۰	۱۴۰۹	BWh
رفسنجان	کرمان	مرجع رطوبت خاک	۵۵/۹۳۰	۳۰/۳۸۰	۱۵۲۴	BWh
قزوین	قزوین	هواشناسی	۵۰/۰۴۵	۳۶/۲۴۵	۱۲۷۹/۱	Csa

ریزمقیاس‌نمایی رطوبت سطحی خاک GLDAS با استفاده از الگوریتم‌های یادگیری ماشین و... / موسوی رینه و همکاران ۹۹۰						
Csa	۱۲۴۹	۳۶/۲۵۰	۴۹/۸۹۹	مرجع رطوبت خاک	قزوین	اسماعیل‌آباد-اقبالیه
BSk	۱۵۹۶	۲۹/۷۷۸	۵۲/۷۰۴	هواشناسی و مرجع رطوبت خاک	فارس	زرقان
BSk	۱۴۸۸	۲۹/۵۴۵	۵۲/۶۰۴	هواشناسی و مرجع رطوبت خاک	فارس	شیراز
BSk	۱۳۶۱	۳۸/۱۲۲	۴۶/۲۳۴	هواشناسی	آذربایجان شرقی	تبریز
BSk	۱۳۳۸	۳۷/۹۷۶	۴۶/۰۴۷	مرجع رطوبت خاک	آذربایجان شرقی	خسروشاه
BSk	۱۱۷۶	۳۶/۴۸۳	۵۹/۲۸۳	هواشناسی و مرجع رطوبت خاک	خراسان	گلمکان
Dsa	۱۷۴۰/۸	۳۴/۸۵۰	۴۸/۵۳۰	هواشناسی	همدان	همدان فرودگاه
Dsa	۱۶۷۷/۸	۳۴/۱۵۰	۴۸/۴۰	مرجع رطوبت خاک	همدان	نپاوند
BWh	۱۲۳۰/۲	۳۱/۹۰	۵۴/۲۸۰	هواشناسی	یزد	یزد

داده‌های هواشناسی ایستگاه‌ها از سازمان هواشناسی کشور^۱ دریافت شدند. از آنجا که بسیاری از این ایستگاه‌ها مقادیر رطوبت خاک را ثبت نکرده بودند، برای هر ایستگاه هواشناسی نزدیک‌ترین ایستگاهی که داده‌های رطوبت خاک در آن موجود بود انتخاب و داده‌های مربوط به آن دریافت شد. تنها در ایستگاه‌های زرکان، شیراز و گلمکان، هر دو نوع داده (هواشناسی و رطوبت خاک) به صورت هم‌زمان در دسترس بود. درخور یادآوری است که برای ایستگاه یزد، هیچ داده‌ای از رطوبت خاک (در خود ایستگاه یا ایستگاه‌های اطراف) موجود نبود.

داده‌ها

در این پژوهش، سه دسته داده مستقل شامل داده‌های اقلیمی ایستگاهی، داده‌های رطوبت سطحی خاک استخراج‌شده از پایگاه داده GLDAS، و داده‌های رطوبت خاک مشاهده‌ای مورد استفاده قرار گرفتند که هر یک نقش مشخصی در فرایند آموزش مدل، ریزمقیاس‌نمایی و اعتبارسنجی مستقل^۲ ایفا می‌کنند. داده‌های اقلیمی ایستگاهی و داده‌های رطوبت سطحی خاک GLDAS دارای پوشش زمانی بلندمدت و پیوسته با تفکیک زمانی ماهانه برای دوره ژانویه ۱۹۹۴ تا دسامبر ۲۰۲۴ (۳۱ سال) هستند، در حالی که داده‌های رطوبت خاک مشاهده‌ای دارای پوشش مکانی و زمانی محدود، ناپیوسته و ناهمگن بوده و فقط برای اعتبارسنجی مستقل نتایج مورد استفاده قرار گرفتند.

دسته اول شامل سری‌های زمانی داده‌های اقلیمی ایستگاهی مورد استفاده به عنوان متغیرهای ورودی مدل‌های یادگیری ماشین است. این داده‌ها شامل ۴۰۹۲ سطر (۳۱ سال داده ماهانه برای ایستگاه‌ها) و ستون‌هایی با نام‌های تاریخ، نام منطقه، نام ایستگاه، اقلیم منطقه، طول جغرافیایی، عرض جغرافیایی، ارتفاع ایستگاه از سطح دریا، تابش خورشیدی^۳، دمای هوای میانگین^۴، دمای سطح

1 <https://data.irimo.ir>

2 independent validation

3 radglo24

4 tm_m

خاک^۱، دمای حداقل^۲ و حداکثر^۳، بارش^۴، تبخیر^۵ هستند که از سازمان هواشناسی کشور برای ایستگاه‌های منتخب در نواحی اقلیمی مختلف ایران گردآوری شدند. انتخاب تعداد ایستگاه‌های مورد استفاده در این تحقیق بر اساس دسترسی به سری‌های زمانی پیوسته، کامل و با کیفیت در دوره ۳۱ساله انجام شده است. با وجود ضرورت این معیار به منظور تضمین پایداری و روایی آماری نتایج، در برخی نواحی باعث محدود شدن تعداد ایستگاه‌های منتخب شده است.

دسته دوم شامل داده‌های رطوبت سطحی خاک^۶ (لایه ۰ تا ۱۰ سانتی‌متر) با تفکیک زمانی ماهانه استخراج شده از پایگاه داده GLDAS است که از طریق سامانه آنلاین ناسا^۷ دریافت شدند. شایان یادآوری است که واحد رطوبت سطحی در محصولات GLDAS کیلوگرم بر متر مربع است که نشان‌دهنده عمق معادل آب (میلی‌متر) در لایه مورد نظر است. این داده‌ها به عنوان متغیر وابسته در مدل‌های یادگیری ماشین مورد استفاده قرار گرفته و خروجی آن‌ها رطوبت خاک ریزمقیاس شده در تفکیک مکانی بالاتر است.

دسته سوم شامل داده‌های رطوبت خاک سطحی مشاهده‌ای در ایستگاه‌های هواشناسی منتخب است. این داده‌ها دارای پوشش مکانی محدود، دوره‌های زمانی کوتاه‌تر و ناپیوسته، و پراکنش نامتوازن در سطح کشور هستند و از این رو برای استفاده در فرایند آموزش یا آزمون داخلی مدل‌های یادگیری ماشین مناسب نبوده‌اند. در این مطالعه، داده‌های رطوبت خاک مشاهده‌ای فقط به عنوان داده‌های مرجع برای انجام اعتبارسنجی مستقل خروجی‌های ریزمقیاس‌نمایی شده حاصل از الگوریتم‌ها به کار گرفته شدند. این رویکرد امکان ارزیابی مستقل قابلیت کاربرد عملی داده‌های ریزمقیاس شده در مقیاس محلی را فراهم می‌کند، بدون آن‌که داده‌های رطوبت خاک مشاهده‌ای در فرایند آموزش مدل‌ها دخالتی داشته باشند؛ در نتیجه، مدل‌ها فقط بر داده‌های در دسترس گسترده (GLDAS) و متغیرهای هواشناسی) متکی بوده و قابلیت به‌کارگیری در مناطق فاقد داده‌های مشاهده‌ای متراکم را حفظ می‌کنند.

روش کار

مراحل انجام پژوهش حاضر در شکل ۲ نشان داده شده است. پس از دریافت داده‌ها از سازمان هواشناسی، پیش‌پردازش داده‌ها از جمله تحلیل اکتشافی، بررسی ساختار داده‌ها، آمار توصیفی و بررسی روندها، بررسی داده‌های پرت و گمشده، همبستگی متغیرها و در نهایت انتخاب متغیرهای ورودی الگوریتم‌ها انجام شد. سپس الگوریتم‌های مختلف یادگیری ماشین (Random Forest، XGBoost، CatBoost و LightGBM) برای ریزمقیاس‌نمایی داده‌های رطوبت خاک GLDAS استفاده شدند.

1 tsoilm_m

2 tmin_m

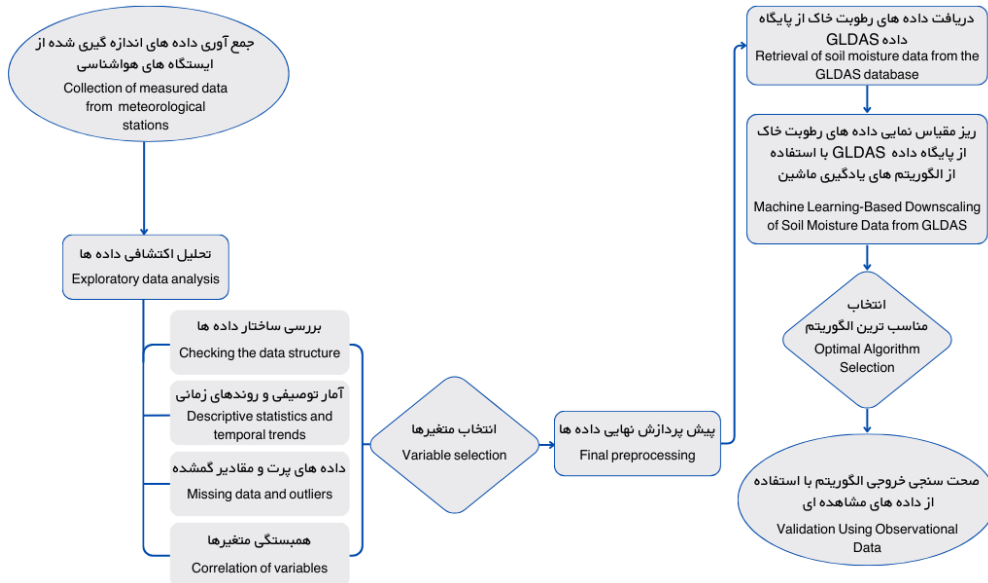
3 tmax_m

4 rrr24

5 evt_s

6 mean_GLDAS_SoilMoi

7 <https://giovanni.gsfc.nasa.gov/giovanni/>



شکل ۲. مراحل انجام پژوهش حاضر

برای هر ایستگاه، میانگین رطوبت خاک GLDAS در پیکسل مکانی متناظر با مختصات ایستگاه و نواحی اطراف آن تعیین شد. این هم‌مکانی‌سازی به مدل امکان یادگیری ارتباط بین داده‌های اقلیمی ایستگاهی و رطوبت خاک در سطح پیکسل GLDAS را می‌دهد. در این مطالعه، الگوریتم‌های یادگیری ماشین برای مدل‌سازی و ریزمقیاس‌نمایی مقادیر رطوبت سطحی خاک پایگاه GLDAS آموزش داده شدند.

به منظور آموزش و ارزیابی مدل‌های یادگیری ماشین، داده‌ها با در نظر گرفتن ماهیت زمانی سری‌های زمانی رطوبت خاک به دو دسته آموزش و آزمون^۱ تقسیم‌بندی شدند. از آنجا که اختلاط تصادفی داده‌های زمانی می‌تواند به نشت اطلاعات و برآورد خوش‌بینانه عملکرد مدل منجر شود، تقسیم داده‌ها بر اساس ترتیب زمانی انجام گرفت. به این منظور، برای هر طبقه اقلیمی، سری زمانی ماهانه داده‌ها به صورت مستقل مرتب‌سازی شده و ۸۰ درصد ابتدایی به عنوان مجموعه آموزش و ۲۰ درصد انتهایی به عنوان مجموعه آزمون در نظر گرفته شد. این رویکرد امکان ارزیابی توانایی مدل در تعمیم به دوره‌های زمانی آینده را فراهم کرده و برآورد واقع‌بینانه‌تری از عملکرد مدل ارائه می‌دهد.

Random Forest با ترکیب چندین درخت تصمیم‌گیری و توانایی ارزیابی اهمیت ویژگی‌ها، و الگوریتم‌های تقویت گرادینتی (XGBoost، CatBoost و LightGBM) با دقت و سرعت زیاد، گزینه‌های مناسبی برای مدل‌سازی داده‌های پیچیده هستند (S. Chen et al., 2019; T. Chen & Guestrin, 2016; Ke et al., 2017). مدل‌ها در محیط پایتون برنامه‌نویسی و آموزش داده شدند.

همچنین به منظور ارتقای تفسیرپذیری مدل‌های توسعه‌یافته و درک سهم هر یک از متغیرهای ورودی در بازتولید رطوبت خاک، از رویکرد SHAP استفاده شد. این روش مبتنی بر نظریه بازی‌های مشارکتی است و به خلاف روش‌های سنتی اهمیت ویژگی، می‌تواند نقش مثبت یا منفی هر متغیر را در خروجی نهایی مدل به صورت کمی تعیین کند (Lundberg & Lee, 2017).

ارزیابی الگوریتم‌ها

به منظور ارزیابی الگوریتم‌های ریزمقیاس نمایی رطوبت خاک، از شاخص‌های ضریب تعیین R^2 ، خطای میانگین، ریشه میانگین مربعات خطا RMSE و انحراف پیش بینی Bias استفاده شده است. در ادامه، فرمول‌های هر یک از این شاخص‌ها آورده شده است (Arah et al., 2008; Chicco et al., 2021).

$$R^2 = 1 - \frac{\sum_{i=1}^n (SM_p - SM_o)^2}{\sum_{i=1}^n (SM_o - \bar{SM}_o)^2} \quad (1)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^n |SM_p - SM_o|}{n}} \quad (2)$$

$$Bias = \frac{\sum_{i=1}^n (SM_o - SM_p)}{n} \quad (3)$$

در روابط ۱ تا ۳ SM_p مقادیر رطوبت خاک ریزمقیاس شده و SM_o مقادیر رطوبت خاک مرجع یا اندازه‌گیری شده و n تعداد هستند.

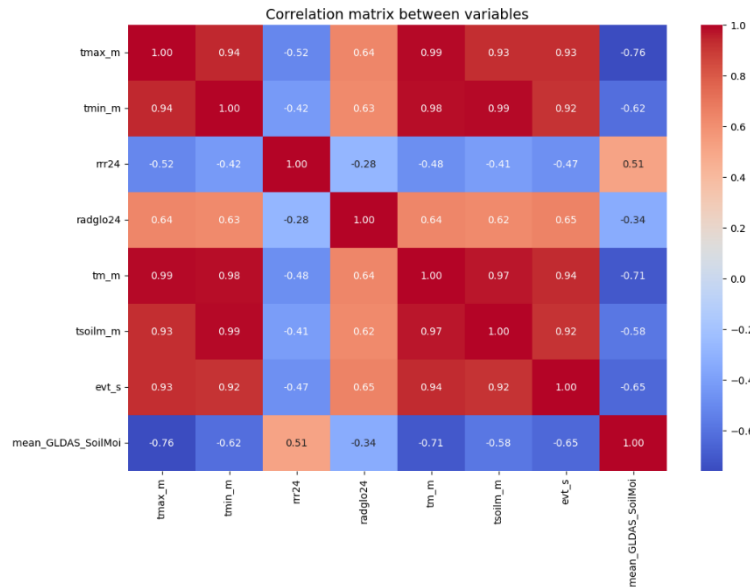
یافته‌های پژوهش و بحث

تحلیل اکتشافی و پیش‌پردازش داده‌ها

بررسی سری‌های زمانی نشان‌دهنده نوسانات فصلی منظم و انطباق متغیرها با چرخه‌های سالانه اقلیمی ایران است (Gedara et al., 2025). تحلیل توزیع داده‌های اقلیمی (هیستوگرام و نمودار جعبه‌ای) (Mehmood et al., 2025) تأیید کرد که مقادیر حدی در بارش و تبخیر، ناشی از طبیعت اقلیم‌های خشک و نیمه‌خشک ایران (مانند بارش‌های سیلابی ناگهانی یا تبخیر شدید در تابستان) بوده و داده‌پرت خطا محسوب نمی‌شوند (Z. Xu et al., 2024). در بخش داده‌های گمشده، متغیرهای دمایی، بارش و تبخیر با نقصانی کمتر از ۵ درصد با روش رگرسیونی بازسازی شدند (Little & Rubin, 2019)، اما متغیر تابش خورشیدی به دلیل ۵۴ درصد داده گمشده از روند مدل‌سازی حذف شد؛ وجود چنین مواردی ممکن است به بایاس شدید در مدل‌سازی منجر شود یا الگوهای نادرست را در داده‌ها تقویت کند (Azur et al., 2011).

تحلیل همبستگی و انتخاب ویژگی‌ها

در مدل‌سازی، هدف آن است که متغیرهای مستقل هر یک اطلاعات متمایز و غیرتکراری در ارتباط با متغیر وابسته فراهم کنند. مطابق با شکل ۳، تحلیل همبستگی پیرسون نشان‌دهنده هم‌خطی شدید ($r > 0.90$) میان متغیرهای دمایی t_{soil_m} ، t_{m_m} ، t_{min_m} و t_{max_m} بود. به منظور جلوگیری از بیش‌برازش و کاهش پایداری مدل (Dormann et al., 2013)، متغیرهای t_{min_m} و t_{max_m} به عنوان نمایندگان دما حفظ شدند، این دو متغیر نماینده حداقل و حداکثر دما هستند و بسیاری از مدل‌های معتبر مدل‌سازی رطوبت خاک به‌ویژه در مطالعات هیدرولوژیکی مانند (Gaona et al., 2023) از آن‌ها استفاده کرده‌اند؛ همچنین علاوه بر اثر متمایز این دو متغیر بر تبخیر، همبستگی کمتری با یکدیگر نسبت به سایر جفت‌ها نشان دادند. در نهایت در این مرحله از داده‌های هواشناسی t_{min_m} ، t_{max_m} ، $rrr24$ و evt_s به همراه اطلاعات مکانی ایستگاه‌ها و داده‌های رطوبت خاک پایگاه داده GLDAS در ۴ الگوریتم یاد شده استفاده شد.



شکل ۳. همبستگی بین متغیرها

ارزیابی عملکرد الگوریتم‌ها در سطح ملی

به منظور ارزیابی عملکرد الگوریتم‌های مختلف یادگیری ماشین در ریزمقیاس‌نمایی داده‌های رطوبت خاک GLDAS نتایج با استفاده از معیارهای R^2 ، RMSE و Bias در دو مقیاس کل کشور و منطقه‌ای با در نظر گرفتن اقلیم‌های مختلف مورد بررسی قرار گرفتند.

مطابق با جدول ۳ نتایج نشان داد CatBoost در مقیاس کل کشور بهترین عملکرد را داشته است.

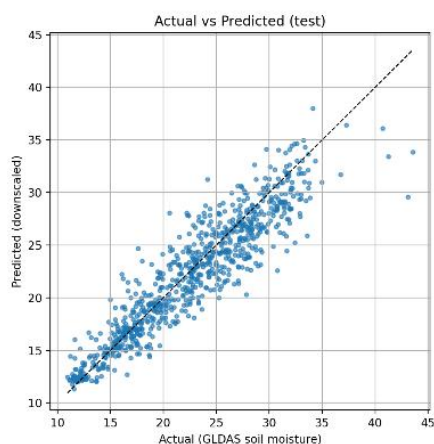
جدول ۳. نتایج ارزیابی الگوریتم‌های یادگیری ماشین در ریزمقیاس‌نمایی داده‌های رطوبت خاک GLDAS در سطح ملی

Model	R2		RMSE (kg/m ²)		Bias (kg/m ²)	
	Train	Test	Train	Test	Train	Test
RF	0.969	0.847	1.114	2.4	0.008	-0.363
XGBoost	0.975	0.823	1.006	2.579	0	-0.407
CatBoost	0.91	0.854	1.892	2.345	-0.025	-0.452
LightGBM	0.909	0.843	1.902	2.426	0	-0.442

این الگوریتم بالاترین مقدار ضریب تعیین ($R^2 = 0.854$) روی داده‌های آزمون و همچنین کمترین مقادیر خطا ($RMSE=2.345$) را ارائه داد. پس از آن، الگوریتم‌های Random Forest و LightGBM در رتبه‌های بعدی قرار گرفتند که نتایج آن‌ها بسیار نزدیک به هم بود ($R^2 = 0.847$) و ($RMSE = 2.4$)، در مقابل، الگوریتم XGBoost با وجود عملکرد بسیار خوب روی داده‌های آموزش ($R^2 = 0.975$)، در داده‌های آزمون دچار افت محسوسی شد ($R^2 = 0.823$) و ($RMSE=2.579$)، که نشان‌دهنده بروز بیش برآزش در این مدل است.

تحلیل بایاس در داده‌های آموزش و آزمون داخلی نشان داد تمامی مدل‌ها تمایل سیستماتیک به کم‌برآورد رطوبت خاک دارند، به طوری که میانگین بایاس در داده‌های آزمون برای همه الگوریتم‌ها منفی است. این در حالی است که مقادیر نسبتاً بالای R^2 و RMSE قابل قبول نشان می‌دهد مدل‌ها توانسته‌اند الگوی کلی تغییرات رطوبت خاک را به خوبی بازتولید کنند. وجود بایاس منفی

هم‌زمان با عملکرد مناسب R^2 و RMSE بیانگر آن است که خطای مدل‌ها بیشتر از جنس جهت‌گیری سیستماتیک^۱ است تا ناتوانی در مدل‌سازی تغییرپذیری داده‌ها. این رفتار می‌تواند ناشی از تمایل الگوریتم‌های یادگیری ماشین به همگرایی به مقادیر میانگین و کاهش حساسیت به مقادیر حدی رطوبت خاک باشد، که پدیده‌ای رایج در الگوریتم‌های رگرسیون (اعم از ساده و پیچیده) در ادبیات یادگیری ماشین محسوب می‌شود (Ting, 2024).



شکل ۴. رابطه بین مقادیر واقعی و پیش‌بینی‌شده رطوبت خاک GLDAS

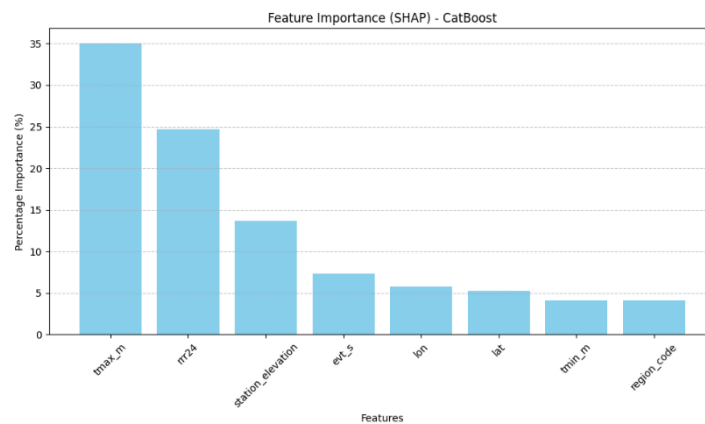
در شکل ۴، رابطه بین مقادیر واقعی رطوبت خاک GLDAS و مقادیر ریزمقیاس‌شده توسط مدل CatBoost در دوره آزمون نمایش داده شده است. تمرکز اغلب نقاط داده در اطراف خط ۱:۱ بیانگر توافق مناسب بین مقادیر مشاهده‌ای و پیش‌بینی‌شده است. همان‌طور که مشاهده می‌شود مدل در بازه مقادیر کمتر از حدود 35 kg/m^2 عملکرد مطلوبی از خود نشان داده و اختلاف معناداری به طور میانگین بین مقادیر واقعی و مقادیر بازتولیدشده مشاهده نمی‌شود. در مقابل، در مقادیر بالاتر (تقریباً بین ۳۵ تا ۴۵)، پراکندگی نقاط افزایش یافته و در برخی موارد انحراف از خط ۱:۱ قابل مشاهده است که نشان‌دهنده کاهش نسبی دقت مدل در مقادیر بالای رطوبت خاک است. با این حال، روند کلی داده‌ها حاکی از توانایی مناسب مدل در بازتولید الگوی تغییرات رطوبت خاک GLDAS است.

تحلیل اهمیت متغیرها

تحلیل تفسیرپذیری مدل با استفاده از رویکرد SHAP (شکل ۵) نشان داد متغیرهای اقلیمی و فیزیکی به طور ترکیبی فرایند ریزمقیاس‌نمایی را هدایت می‌کنند. دمای حداکثر با $35/07$ درصد بیشترین سهم را در پیش‌بینی رطوبت خاک داراست که بیانگر نقش کلیدی توان تبخیر اتمسفر در تخلیه رطوبت سطحی است. پس از آن، بارش با $24/66$ درصد به عنوان پیش‌ران اصلی ورود رطوبت قرار دارد. نکته قابل توجه، اهمیت بالای ارتفاع ایستگاه ($13/71$ درصد) است که فراتر از متغیرهایی نظیر تبخیر و دمای حداقل قرار گرفته است. این امر تأییدکننده تأثیر قابل ملاحظه توپوگرافی بر الگوهای بازتولید شده رطوبت خاک در مقیاس محلی است. همچنین، سهم در مجموع ۱۱ درصدی مختصات جغرافیایی (طول و عرض) در کنار متغیر کد منطقه (که به منظور درک تمایزهای اقلیمی - ساختاری به صورت یک ویژگی دسته‌ای^۲ به مدل معرفی شد)، نشان می‌دهد مدل علاوه بر روابط فیزیکی، تفاوت‌های ساختاری مکان‌مند را نیز در فرایند یادگیری لحاظ کرده است.

1 systematic regression bias or attenuation bias

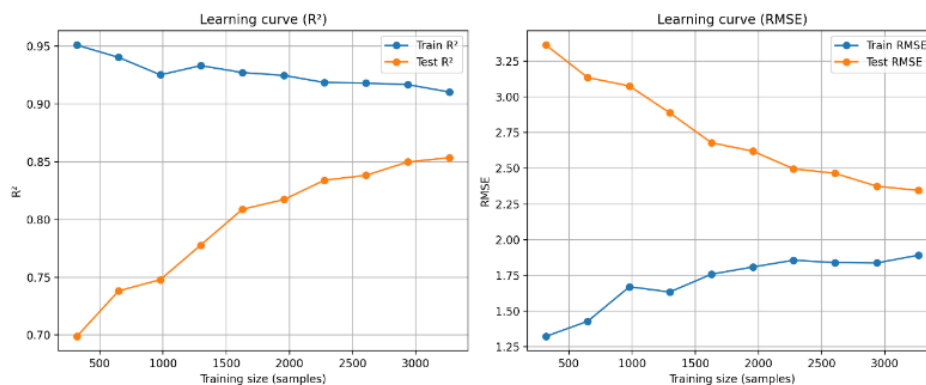
2 Categorical



شکل ۵. سهم نسبی متغیرهای ورودی در ریزمقیاس‌نمایی رطوبت خاک با مدل Catboost با استفاده از رویکرد تفسیرپذیر SHAP

تحلیل منحنی یادگیری

منحنی‌های یادگیری (شکل ۶) مدل CatBoost نشان دادند با افزایش حجم داده‌های آموزشی، عملکرد مدل در داده‌های آزمون بهبود یافته است، به طوری که در داده‌های آزمون $(R^2 = 0.70)$ به $(R^2 = 0.85)$ ، افزایش یافته و $(RMSE=3.3)$ به $(RMSE=2.2)$ ، کاهش یافته است. هم‌زمان، منحنی‌های آموزش و آزمون به تدریج به هم نزدیک شده و اختلاف اولیه میان آن‌ها کاهش یافته است، که نشان‌دهنده کاهش بیش‌برازش و همگرایی مناسب مدل است. به بیان دیگر با افزایش حجم داده، فاصله بین منحنی آموزش و آزمون به طور پیوسته کاهش و در نتیجه، دقت افزایش و خطای آزمون کم شده که نشان از تعمیم‌پذیری مناسب مدل دارد. این روند بیانگر همگرایی مناسب مدل و توانایی تعمیم قابل قبول آن در پیش‌بینی رطوبت خاک است.



شکل ۶. منحنی یادگیری (Learning curve) الگوریتم catboost

تحلیل منطقه‌ای و اقلیمی

تحلیل منطقه‌ای نشان داد (جدول ۴) که عملکرد در بین الگوریتم‌ها متفاوت است.

جدول ۴. نتایج ارزیابی الگوریتم‌های یادگیری ماشین در ریزمقیاس‌نمایی داده‌های رطوبت خاک GLDAS به تفکیک منطقه‌ای و طبقه اقلیمی

Region	climate	Random Forest			XGBoost		
		R ²	RMSE (kg/m ²)	Bias (kg/m ²)	R ²	RMSE (kg/m ²)	Bias (kg/m ²)
Esfahan	BWk	0.86	2.95	-1.40	0.83	3.22	-1.37
Hamedan	Dsa	0.81	2.84	-0.24	0.76	3.14	-0.35
Hashtgerd	Csa	0.79	2.75	-0.92	0.73	3.10	-0.67

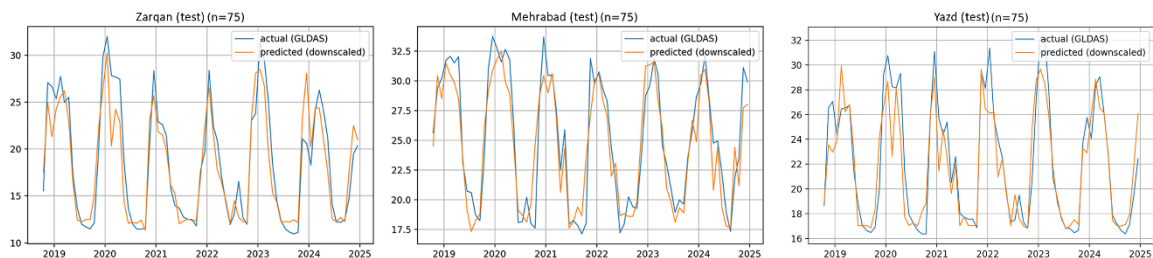
Zarqan	BSk	0.86	2.37	-0.13	0.85	2.44	-0.45
Golmakan	BSk	0.73	2.01	-0.39	0.67	2.23	-0.54
Qazvin	Csa	0.75	2.39	0.12	0.74	2.45	-0.13
Anar	BWh	0.69	2.23	0.12	0.68	2.28	0.48
Shiraz	BSk	0.83	2.09	-0.47	0.82	2.13	-0.75
Tabriz	BSk	0.74	2.37	-0.09	0.72	2.49	0.01
Mehrabad	BSk	0.85	2.07	-0.20	0.81	2.33	-0.48
Yazd	BWh	0.82	2.11	-0.40	0.80	2.23	-0.23

Region	climate	CatBoost			LGBM		
		R ²	RMSE (kg/m ²)	Bias (kg/m ²)	R ²	RMSE (kg/m ²)	Bias (kg/m ²)
Esfahan	BWk	0.87	2.84	-1.30	0.87	2.84	-1.26
Hamedan	Dsa	0.82	2.77	-0.50	0.79	2.96	-0.32
Hashtgerd	Csa	0.81	2.61	-0.94	0.80	2.67	-1.16
Zarqan	BSk	0.85	2.44	-0.50	0.84	2.47	-0.32
Golmakan	BSk	0.75	1.94	-0.42	0.72	2.03	-0.55
Qazvin	Csa	0.76	2.34	0.16	0.75	2.37	0.10
Anar	BWh	0.74	2.06	0.48	0.67	2.30	0.45
Shiraz	BSk	0.83	2.12	-0.85	0.81	2.22	-0.81
Tabriz	BSk	0.76	2.31	-0.16	0.74	2.40	-0.13
Mehrabad	BSk	0.85	2.09	-0.53	0.85	2.09	-0.35
Yazd	BWh	0.82	2.07	-0.41	0.81	2.14	-0.51

- اقلیم‌های بیابانی (BWh و BWk)، CatBoost بهترین عملکرد را از نظر ضریب تعیین (R^2) و معیارهای خطا ارائه داد، به‌ویژه در ایستگاه‌های انار و یزد. در اصفهان نیز CatBoost و LightGBM عملکرد مشابه و بهتری نسبت به RF و XGBoost داشتند. بنابراین، در اقلیم‌های بیابانی CatBoost پایدارترین و دقیق‌ترین الگوریتم شناخته شد.
- اقلیم‌های نیمه‌خشک سرد (BSk): ایستگاه‌های زرکان، گل‌مکان، تبریز و شیراز در این گروه قرار می‌گیرند. رقابت نزدیکی میان CatBoost و RF مشاهده شد، هرچند CatBoost در مناطق سردتر (تبریز و گل‌مکان) برتری نسبی داشت. بنابراین، CatBoost همچنان برتری نسبی دارد، اما فاصله آن با RF و LGBM نیز کم است. بررسی عملکرد الگوریتم‌ها در ایستگاه مهرآباد نشان می‌دهد الگوریتم Random Forest بهترین عملکرد کلی را با توجه هم‌زمان به مقدار بالای R^2 ، کمترین RMSE و بایاس محدود ارائه می‌دهد. با این حال، عملکرد مدل‌های CatBoost و LGBM نیز بسیار نزدیک بوده و می‌توان آن‌ها را به عنوان گزینه‌های قابل اعتماد جایگزین در نظر گرفت، در حالی که مدل XGBoost در این ایستگاه نسبت به سایر الگوریتم‌ها عملکرد ضعیف‌تری از خود نشان می‌دهد.
- اقلیم‌های مدیترانه‌ای و قاره‌ای (Csa, Dsa): ایستگاه‌های هشتگرد و قزوین در دسته (Csa)، قرار دارند. در هشتگرد و قزوین، CatBoost اندکی بالاتر از RF و LGBM عمل کرد. بنابراین در این اقلیم، هر دو الگوریتم RF و CatBoost با برتری جزئی قابل توصیه‌اند. در ایستگاه همدان (Dsa)، CatBoost بهترین نتایج را ارائه کرد ($R^2=0.82$) همچنین RF با ($R^2=0.81$)، با عملکرد ضعیف‌تر در رتبه بعدی قرار گرفت. بنابراین، در اقلیم Dsa نیز CatBoost پایدارترین و دقیق‌ترین

الگوریتم بود. در مقابل، XGBoost در اغلب مناطق عملکرد ضعیف‌تری نسبت به سایر الگوریتم‌ها نشان داد و تنها در برخی مناطق نتایج قابل قبول ارائه داد.

اگرچه عملکرد الگوریتم‌ها تحت تأثیر ویژگی‌های اقلیمی با نوسانات جزئی همراه بود، اما الگوریتم CatBoost به دلیل حفظ برتری نسبی در اکثر ایستگاه‌ها و ارائه ضریب تعیین فراتر از 0.73 در تمامی مناطق، به عنوان کارآمدترین و پایدارترین مدل منتخب برای پیش‌بینی رطوبت خاک در تمامی اقلیم‌های مورد مطالعه شناخته شد.



شکل ۷. مقایسه سری زمانی رطوبت خاک GLDAS و مقادیر ریزمقیاس شده توسط الگوریتم CatBoost در ایستگاه‌های زرقان، مهرآباد و یزد برای دوره آزمون

ارزیابی عملکرد الگوریتم CatBoost در مناطق مورد مطالعه نشان‌دهنده توانمندی بالای این مدل در استخراج الگوهای غیرخطی است، به طوری که در تمامی مناطق، $R^2 > 0.73$ گزارش شد. با این حال، تحلیل تطبیقی حاکی از آن است که عملکرد مدل در تمامی ایستگاه‌ها یکنواخت نبوده و ایستگاه‌های زرقان (اقلیم نیمه‌خشک)، مهرآباد (اقلیم نیمه‌خشک) و یزد (اقلیم بیابانی) به دلیل همگرایی بهینه میان بیشترین میزان دقت تبیین و کمترین مقدار خطا نسبتاً پایین و بایاس محدود، به عنوان مناطق برتر در پیش‌بینی شناخته شدند. تطابق بالای مقادیر پیش‌بینی شده با داده‌های واقعی در این ایستگاه‌ها که در جدول ۴ از منظر شاخص‌های آماری و در شکل ۷ به صورت بصری نمایش داده شده است، بیانگر پایداری الگوریتم در بازسازی نوسانات فصلی رطوبت خاک در این اقلیم‌ها است. به عبارتی مدل در ایستگاه زرقان با برخورداری از بهترین انطباق الگو^۱ در درک پیچیدگی‌ها و رفتارهای متغیر وابسته بسیار موفق‌تر عمل کرده است. در مقابل، علی‌رغم بالاتر بودن ضریب تعیین ($R^2=0.86$) در ایستگاه اصفهان، وجود سوگیری (Bias) منفی و خطا بالاتر در این منطقه نشان‌دهنده حساسیت مدل به داده‌های پرت یا نوسانات شدید متغیرهای ورودی است. در مجموع، می‌توان استنباط کرد که الگوریتم یادشده از قابلیت تعمیم‌پذیری مناسبی برخوردار است، اما دقت عملیاتی آن تحت تأثیر ویژگی‌های محلی هر منطقه دستخوش تغییرات جزئی می‌شود.

اعتبارسنجی با داده‌های رطوبت خاک مشاهده‌ای (مرجع)

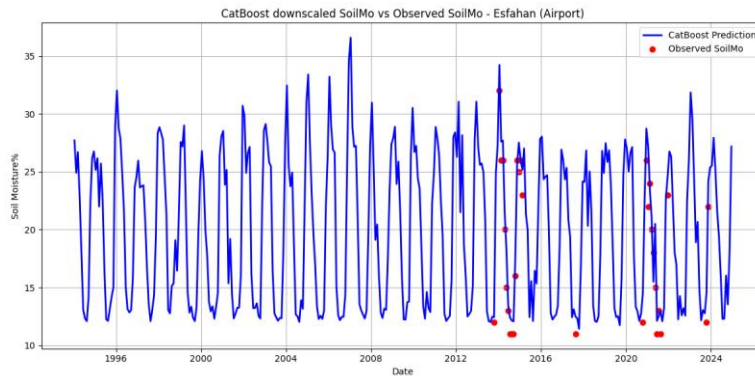
به منظور ارزیابی کارایی عملی مدل در مقیاس محلی، خروجی‌های ریزمقیاس شده با داده‌های مشاهده‌ای ایستگاهی که در فرایند آموزش مدل دخالتی نداشتند، مقایسه شدند. شکل ۵ مقایسه زمانی مقادیر پیش‌بینی شده و مشاهده‌ای را برای ایستگاه اصفهان به عنوان نمونه‌ای از عملکرد مطلوب مدل نشان می‌دهد.

نتایج اعتبارسنجی مستقل در سطح ملی، ضریب تعیین $R^2 = 0.607$ و $RMSE = 4.286 \text{ kg/m}^2$ را نشان داد. کاهش شاخص R^2 نسبت به ارزیابی داخلی^۱ ($R^2 = 0.854$)، پدیده‌ای متداول در مدل‌سازی‌های هیدرولوژیکی است که عمدتاً از عدم تطابق مقیاس^۲ ناشی می‌شود (Senanayake et al., 2024). با این حال، تبیین بیش از ۶۰ درصد از تغییرپذیری رطوبت خاک توسط مدل، بیانگر موفقیت فرایند ریزمقیاس‌نمایی در بازتولید الگوهای محلی است. بسیاری از پژوهش‌های پیشین نیز بر این موضوع تأکید داشته‌اند که عملکرد مدل‌ها در داده‌های آزمون داخلی معمولاً بهتر از عملکرد آن‌ها در اعتبارسنجی با داده‌های مشاهده‌ای است (Senanayake et al., 2024; Wang & Gao, 2025). این مرحله از تحلیل به منظور ارزیابی عملکرد کاربردی خروجی‌های ریزمقیاس‌نمایی شده انجام شده است، نه برای آموزش یا بهینه‌سازی مدل. به این ترتیب، این رویکرد امکان تفکیک روشن بین عملکرد داخلی مدل و دقت واقعی ریزمقیاس‌نمایی در مواجهه با داده‌های مستقل مشاهده‌ای را فراهم می‌کند.

در مورد بایاس، اگرچه در ارزیابی داخلی مدل‌ها (بر اساس داده‌های آموزشی و آزمون) میانگین بایاس منفی مشاهده شد که این منفی بودن نشان‌دهنده تمایل مدل‌ها به کاهش دامنه پیش‌بینی و همگرایی به مقادیر میانگین است، نتایج اعتبارسنجی مستقل با استفاده از داده‌های رطوبت خاک مشاهده‌ای رفتار متفاوتی را نشان داد. در این مرحله، بایاس مثبت به دست آمده ($Bias = 2.131 \text{ kg/m}^2$) بیانگر بیش‌برآورد نسبی رطوبت خاک توسط مدل در مقایسه با مشاهدات نقطه‌ای ایستگاهی است. این بیش‌برآورد می‌تواند به کم‌برآورد نیاز آبی واقعی گیاه، تشخیص دیرهنگام خشکسالی، برآورد نادرست رواناب و خطا در مدل‌سازی‌های هیدرولوژیکی و اکولوژیکی منجر شود. این تأثیرات کاربردی از اهمیت ویژه‌ای برخوردارند، زیرا رطوبت خاک به عنوان متغیری کلیدی در شاخص‌های خشکسالی، مدل‌های هیدرولوژیکی و فرایندهای تبخیر - تعرق شناخته می‌شود و خطای سیستماتیک در آن می‌تواند باعث عدم انطباق خروجی‌های مدل با شرایط واقعی محیطی شود. این تفاوت ناشی از ماهیت متفاوت داده‌های مرجع، ناهمگنی مقیاس مکانی، و عدم قطعیت ذاتی داده‌های مشاهده‌ای بوده و لزوماً به معنای تناقض در عملکرد مدل نیست. لازم به تأکید است که اعتبارسنجی کمی مدل با استفاده از شاخص‌های آماری فقط در ایستگاه‌هایی انجام شده است که داده‌های رطوبت خاک اندازه‌گیری شده و قابل اعتماد در دسترس بوده است. در سایر مناطق، به دلیل نبود داده‌های رطوبت خاک مشاهده‌ای، نتایج ارائه شده به عنوان تحلیل عملکرد منطقه‌ای و ارزیابی نسبی مدل بر اساس الگوهای مکانی و روندهای زمانی تفسیر می‌شوند و به معنای اعتبارسنجی مستقیم مبتنی بر داده‌های مشاهده‌ای نیستند.

تحلیل منطقه‌ای نشان داد که برخی مناطق مانند اصفهان ($R^2 = 0.89$)، هشتگرد ($R^2 = 0.805$)، قزوین ($R^2 = 0.636$) و شیراز ($R^2 = 0.747$) عملکرد خوبی دارند و ریزمقیاس‌نمایی‌ها معتبر هستند، اما در مناطقی مانند گل‌مکان و تبریز دقت مدل بسیار پایین است، چنین نتیجه‌ای معمولاً زمانی رخ می‌دهد که داده‌ها ناهمگن یا پراکنده باشند، تعداد داده‌های مرجع کم باشد، یا ویژگی‌های منطقه‌ای پیچیده بوده و مدل نتواند الگوهای آن را یاد بگیرد. در رفسنجان، همدان و مهرآباد و یزد داده‌های مرجع کافی برای ارزیابی وجود ندارد. نمودار مقایسه‌ای ایستگاه اصفهان به عنوان نماینده منطقه‌ای که الگوریتم CatBoost عملکرد خوبی داشته در شکل ۸ آورده شده است.

1 internal validation
2 scale mismatch



شکل ۸. مقایسه مقادیر رطوبت خاک واقعی و خروجی الگوریتم CatBoost در ایستگاه اصفهان (اعتبارسنجی مستقل)

این نتایج نشان می‌دهد مدل در سطح کل کشور عملکرد متوسط دارد، اما برای تحلیل دقیقه منطقه‌ای نیازمند بهبود مجموعه داده‌های مرجع است.

نتیجه‌گیری

در این تحقیق، ریزمقیاس‌نمایی رطوبت خاک با استفاده از الگوریتم‌های یادگیری ماشین شامل Random Forest، XGBoost، CatBoost و LightGBM انجام شد. بر اساس نتایج ملی و منطقه‌ای، الگوریتم CatBoost مناسب‌ترین روش برای مدل‌سازی رطوبت خاک در این مطالعه تشخیص داده شد، زیرا علاوه بر دقت بالا در مقیاس ملی، پایداری مناسبی در مناطق مختلف با شرایط اقلیمی متفاوت نشان داد. الگوریتم Random Forest نیز به عنوان گزینه دوم عملکرد قابل قبولی داشت، در حالی که LightGBM با دقتی نزدیک به این دو الگوریتم، می‌تواند در کاربردهایی که سرعت و کارایی محاسباتی اهمیت بیشتری دارد، مورد استفاده قرار گیرد. در مقابل، عملکرد XGBoost به دلیل بروز بیش‌برازش و افزایش خطای پیش‌بینی نسبت به سایر الگوریتم‌ها ضعیف‌تر بود و این الگوریتم گزینه مناسبی برای مسئله حاضر محسوب نشد؛ رفتاری که در مطالعات انجام‌شده در مناطق وسیع و اقلیمی ناهمگن، از جمله شمال غرب چین، نیز گزارش شده است (J. Xu et al., 2024). اعتبارسنجی مستقل با استفاده از داده‌های رطوبت خاک مشاهده‌ای نشان داد عملکرد CatBoost در مقیاس ملی در سطح متوسط قرار دارد ($R^2 = 0.607$)، با خطای پیش‌بینی متوسط ($RMSE = 4.286$) و تمایل به بیش‌برآورد مقادیر رطوبت خاک ($Bias = 2.131$). این نتایج بیانگر آن است که مدل قادر به بازتولید الگوهای کلی و روندهای مکانی - زمانی رطوبت خاک در سطح کشور است، اما دقت نقطه‌ای آن در همه مناطق یکسان نیست. با این حال، نتایج تحلیل منطقه‌ای اعتبارسنجی مستقل نشان داد عملکرد مدل در برخی مناطق با داده‌های مرجع بیشتر و شرایط اقلیمی همگن‌تر به مراتب بهتر بوده است، به طوری که مقادیر R^2 تا حدود ۰/۸۹ نیز حاصل شد. در مقابل، در مناطقی با اقلیم‌های خشک، ناهمگن یا دارای داده‌های مرجع محدود، دقت مدل در نتایج اعتبارسنجی مستقل به طور محسوسی کاهش یافت. این تفاوت‌های منطقه‌ای نشان می‌دهد اعتبارسنجی ریزمقیاس‌نمایی به شدت به کیفیت، تعداد و پراکنش داده‌های رطوبت خاک مشاهده‌ای و پیچیدگی ویژگی‌های منطقه‌ای وابسته است. تمایل مدل به بیش‌برآورد رطوبت خاک می‌تواند تا حدی ناشی از رفتار ذاتی الگوریتم‌های یادگیری ماشین در همگرایی به مقادیر میانگین و کاهش حساسیت به مقادیر حدی باشد؛ پدیده‌ای که در مطالعات ریزمقیاس‌نمایی رطوبت خاک به طور گسترده گزارش شده است. بر این اساس، اگرچه استفاده از خروجی‌های مدل برای کاربردهایی که نیازمند دقت نقطه‌ای بالا هستند، مانند مدیریت آبیاری در مقیاس مزرعه، توصیه نمی‌شود، اما نتایج مدل می‌تواند برای تحلیل‌های منطقه‌ای، پایش تغییرات رطوبت خاک، ارزیابی خشکسالی و پشتیبانی از تصمیم‌گیری‌های مدیریت منابع آب در مقیاس حوضه یا استان قابل اتکا باشد. لازم به تأکید است که هدف اصلی این مطالعه، به‌کارگیری الگوریتم‌های یادگیری ماشین برای ریزمقیاس‌نمایی داده‌های

رطوبت خاک حاصل از محصول GLDAS و ارتقای تفکیک مکانی این داده‌ها بوده است؛ از این رو، روش پیشنهادی بیش از آنکه یک ابزار پیش‌بینی دقیق نقطه‌ای مبتنی بر داده‌های محلی باشد، چارچوبی کارآمد برای بازتولید الگوهای مکانی - زمانی رطوبت خاک GLDAS و تحلیل‌های مقیاس منطقه‌ای و میانی تا کلان محسوب می‌شود.

پیشنهادها

بر اساس نتایج و تجربیات حاصل از این پژوهش، راهکارهای زیر برای ارتقای دقت مدل‌های ریزمقیاس‌نمایی رطوبت خاک در تحقیقات آینده پیشنهاد می‌شود:

تلفیق متغیرهای فیزیکی و فیزیوگرافی: در کنار متغیرهای اقلیمی، دخالت دادن پارامترهای پویای زمین مانند شاخص‌های پوشش گیاهی (مانند NDVI)، کاربری اراضی و ویژگی‌های ثابت نظیر بافت خاک و توپوگرافی (شیب و جهت) می‌تواند توانایی مدل را در تفکیک نوسانات محلی رطوبت خاک به‌ویژه در مناطق با ناهمگنی بالا افزایش دهد.

بهینه‌سازی هوشمند هایپرپارامترها: به منظور جلوگیری از بیش‌برازش (به‌ویژه در الگوریتم‌های تقویت گرادیانی) و افزایش قدرت تعمیم‌پذیری مدل در مناطق فاقد ایستگاه، پیشنهاد می‌شود از روش‌های پیشرفته بهینه‌سازی نظیر جست‌وجوی بیزی (Bayesian Optimization) برای تنظیم دقیق پارامترهای مدل CatBoost استفاده شود.

پیاده‌سازی روش‌های اصلاح بایاس (Bias Correction): به منظور کاهش شکاف ناشی از اختلاف مقیاس (Mismatch Scale) بین پیکسل‌های GLDAS و مشاهدات نقطه‌ای، استفاده از تکنیک‌های پس‌پردازشی نظیر نگاشت چندکی (Quantile Mapping) یا تصحیح خطی توصیه می‌شود. این روش‌ها می‌توانند بدون تغییر در ساختار یادگیری ماشین، دقت کاربردی داده‌های ریزمقیاس‌شده را برای مصارف حساس نظیر مدیریت آبیاری در مقیاس محلی بهبود بخشند.

- 1 Abdeh kolahchi, A. , Miri, M. , Zand, M. and Porhemmat, J. (2023). Comparative Evaluation of GLDAS, ESA CCI SM and SMAP Soil Moisture with in situ Measurements (Case Study: Lorestan Province). *Environment and Water Engineering*, 9(4), 548-562. doi: 10.22034/ewe.2023.367471.1819. (in Persian)
- 2 Ali, Z., Hamed, M. M., Nashwan, M. S., & Shahid, S. (2023). Spatiotemporal analysis of groundwater resources sustainability in South Asia and China using GLDAS data sets. *Environmental Earth Sciences*, 82(24), 586.
- 3 Amini, A., Moghadam, M. K., Kolahchi, A. A., Raheli-Namin, M., & Ahmed, K. O. (2023). Evaluation of GLDAS soil moisture product over Kermanshah province, Iran. *H2Open Journal*, 6(3), 373–386.
- 4 Arah, O. A., Chiba, Y., & Greenland, S. (2008). Bias formulas for external adjustment and sensitivity analysis of unmeasured confounders. *Annals of Epidemiology*, 18(8), 637–646.
- 5 Azur, M. J., Stuart, E. A., Frangakis, C., & Leaf, P. J. (2011). Multiple imputation by chained equations: what is it and how does it work? *International Journal of Methods in Psychiatric Research*, 20(1), 40–49.
- 6 Batchu, V., Nearing, G., & Gulshan, V. (2023). A deep learning data fusion model using sentinel-1/2, SoilGrids, SMAP, and GLDAS for soil moisture retrieval. *Journal of Hydrometeorology*, 24(10), 1789–1823.
- 7 Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
- 8 Cai, J., Zhang, Y., Li, Y., Liang, X. S., & Jiang, T. (2017). Analyzing the characteristics of soil moisture using GLDAS data: A case study in eastern China. *Applied Sciences*, 7(6), 566.
- 9 Chahine, M. T. (1992). The hydrological cycle and its influence on climate. *Nature*, 359(6394), 373–380.
- 10 Chen, L., He, Q., Liu, K., Li, J., & Jing, C. (2019). Downscaling of GRACE-derived groundwater storage based on the random forest model. *Remote Sensing*, 11(24), 2979.
- 11 Chen, S., She, D., Zhang, L., Guo, M., & Liu, X. (2019). Spatial downscaling methods of soil moisture based on multisource remote sensing data and its application. *Water*, 11(7), 1401.
- 12 Chen, Y., Yang, K., Qin, J., Zhao, L., Tang, W., & Han, M. (2013). Evaluation of AMSR-E retrievals and GLDAS simulations against observations of a soil moisture network on the central Tibetan Plateau. *Journal of Geophysical Research: Atmospheres*, 118(10), 4466–4475.
- 13 Chicco, D., Warrens, M. J., & Jurman, G. (2021). The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis evaluation. *PeerJ Computer Science*, 7, e623.
- 14 Dormann, C. F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carré, G., Marquéz, J. R. G., Gruber, B., Lafourcade, B., & Leitão, P. J. (2013). Collinearity: a review of methods to deal with it and a simulation study evaluating their performance. *Ecography*, 36(1), 27–46.
- 15 Evans, S., Williams, G. P., Jones, N. L., Ames, D. P., & Nelson, E. J. (2020). Exploiting earth observation data to impute groundwater level measurements with an extreme learning machine. *Remote Sensing*, 12(12), 2044.
- 16 Famiglietti, J. S., Lo, M., Ho, S. L., Bethune, J., Anderson, K. J., Syed, T. H., Swenson, S. C., de Linage, C. R., & Rodell, M. (2011). Satellites measure recent rates of groundwater depletion in California's Central Valley. *Geophysical Research Letters*, 38(3).
- 17 Fatolazadeh, F., Eshagh, M., & Goïta, K. (2020). A new approach for generating optimal GLDAS hydrological products and uncertainties. *Science of the Total Environment*, 730, 138932.
- 18 Gaona, J., Benito-Verdugo, P., Martínez-Fernández, J., González-Zamora, Á., Almendra-Martín, L., & Herrero-Jiménez, C. M. (2023). Predictive value of soil moisture and concurrent variables in the multivariate modelling of cereal yields in water-limited environments. *Agricultural Water Management*, 282, 108280.
- 19 Gedara, S. M., Wasantha, P. L. P., Teodosio, B., Yaghoubi, E., van Staden, R., & Guerrieri, M. (2025). Investigation of seasonal soil moisture and temperature variations underneath a waffle raft foundation built on reactive soil. *Scientific Reports*, 15(1), 34499.
- 20 GIOVANNI. (2024). *The Bridge Between Data and Science*. NASA.

- <https://giovanni.gsfc.nasa.gov/giovanni/>
- 21 Hasan, F., Medley, P., Drake, J., & Chen, G. (2024). Advancing hydrology through machine learning: insights, challenges, and future directions using the CAMELS, caravan, GRDC, CHIRPS, PERSIANN, NLDAS, GLDAS, and GRACE datasets. *Water*, 16(13), 1904.
 - 22 Jiménez, C., Prigent, C., Mueller, B., Seneviratne, S. I., McCabe, M. F., Wood, E. F., Rossow, W. B., Balsamo, G., Betts, A. K., & Dirmeyer, P. A. (2011). Global intercomparison of 12 land surface heat flux estimates. *Journal of Geophysical Research: Atmospheres*, 116(D2).
 - 23 Koster, R. D., & Suarez, M. J. (1992). Modeling the land surface boundary in climate models as a composite of independent vegetation stands. *Journal of Geophysical Research: Atmospheres*, 97(D3), 2697–2715.
 - 24 Liang, X. (1994). A simple hydrologically based model of land surface water and energy fluxes for general circulation models. *J. Geophys. Res.*, 99(7), 14–415.
 - 25 Little, R. J. A., & Rubin, D. B. (2019). *Statistical analysis with missing data*. John Wiley & Sons.
 - 26 Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30.
 - 27 Mehmood, K., Anees, S. A., Muhammad, S., Shahzad, F., Liu, Q., Khan, W. R., Shrahili, M., Ansari, M. J., & Dube, T. (2025). Machine learning and spatio temporal analysis for assessing ecological impacts of the billion tree afforestation project. *Ecology and Evolution*, 15(2), e70736.
 - 28 Mitchell, K. E., Lohmann, D., Houser, P. R., Wood, E. F., Schaake, J. C., Robock, A., Cosgrove, B. A., Sheffield, J., Duan, Q., & Luo, L. (2004). The multi-institution North American Land Data Assimilation System (NLDAS): Utilizing multiple GCIP products and partners in a continental distributed hydrological modeling system. *Journal of Geophysical Research: Atmospheres*, 109(D7).
 - 29 Mousavimehr, S. M., & Kavianpour, M. R. (2025). Estimating Groundwater Levels in Tehran Province Using Ensemble Learning Algorithms. *Contributions of Science and Technology for Engineering*, 2(1), 51–63.
 - 30 Nabavi, S. N., Alizadeh, A. and Faridhosseini, A. (2020). Evaluation of groundwater resources using GRACE Satellite Gravimetric Data (Case study: Khorasan Razavi). *Iranian Journal of Irrigation & Drainage*, 14(3), 855-866. (in Persian)
 - 31 Nouraki, A., Golabi, M., Albaji, M., Naseri, A. and Homayouni, S. (2023). Spatial-temporal modeling of soil moisture using optical and thermal remote sensing data and machine learning algorithms. *Iranian Journal of Soil and Water Research*, 54(4), 637-653. doi: 10.22059/ijswr.2023.356707.669469. (in Persian)
 - 32 Park, S., Park, S., Im, J., Rhee, J., Shin, J., & Park, J. D. (2017). Downscaling GLDAS soil moisture data in East Asia through fusion of multi-sensors by optimizing modified regression trees. *Water*, 9(5), 332.
 - 33 Razinei, T. (2022). Climate of Iran according to Köppen-Geiger, Feddema, and UNEP climate classifications. *Theoretical & Applied Climatology*, 148.
 - 34 Rodell, M., Chen, J., Kato, H., Famiglietti, J. S., Nigro, J., & Wilson, C. R. (2007). Estimating groundwater storage changes in the Mississippi River basin (USA) using GRACE. *Hydrogeology Journal*, 15, 159–166.
 - 35 Rodell, M., Houser, P. R., Jambor, U. E. A., Gottschalck, J., Mitchell, K., Meng, C.-J., Arsenault, K., Cosgrove, B., Radakovich, J., & Bosilovich, M. (2004). The global land data assimilation system. *Bulletin of the American Meteorological Society*, 85(3), 381–394.
 - 36 Rodell, M., Velicogna, I., & Famiglietti, J. S. (2009). Satellite-based estimates of groundwater depletion in India. *Nature*, 460(7258), 999–1002.
 - 37 San Liang, X., & Zhang, Y. (2018). *Coastal Environment, Disaster, and Infrastructure: A Case Study of China's Coastline*. BoD–Books on Demand.
 - 38 Senanayake, I. P., Pathira Arachchilage, K. R. L., Yeo, I.-Y., Khaki, M., Han, S.-C., & Dahlhaus, P. G. (2024). Spatial downscaling of satellite-based soil moisture products using machine learning techniques: A review. *Remote Sensing*, 16(12), 2067.
 - 39 Shang, K. Z., Wang, S. G., Ma, Y. X., Zhou, Z. J., Wang, J. Y., Liu, H. L., & Wang, Y. Q. (2007). A scheme for calculating soil moisture content by using routine weather data. *Atmospheric Chemistry and Physics*, 7(19), 5197–5206.

- 40 Strassberg, G., Scanlon, B. R., & Rodell, M. (2007). Comparison of seasonal terrestrial water storage variations from GRACE with groundwater-level measurements from the High Plains Aquifer (USA). *Geophysical Research Letters*, 34(14).
- 41 Ting, Y.-S. (2024). Why Machine Learning Models Systematically Underestimate Extreme Values. *ArXiv Preprint ArXiv:2412.05806*.
- 42 Tiwari, V. M., Wahr, J., & Swenson, S. (2009). Dwindling groundwater resources in northern India, from satellite gravity observations. *Geophysical Research Letters*, 36(18).
- 43 Wang, L., & Gao, Y. (2025). Estimating and downscaling ESA-CCI soil moisture using Multi-Source remote sensing images and Stacking-Based ensemble learning algorithms in the Shandian River Basin, China. *Remote Sensing*, 17(4), 716.
- 44 Wu, Q., Si, B., He, H., & Wu, P. (2019). Determining regional-scale groundwater recharge with GRACE and GLDAS. *Remote Sensing*, 11(2), 154.
- 45 Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., Luo, L., Alonge, C., Wei, H., & Meng, J. (2012). Continental-scale water and energy flux analysis and validation for the North American Land Data Assimilation System project phase 2 (NLDAS-2): 1. Intercomparison and application of model products. *Journal of Geophysical Research: Atmospheres*, 117(D3).
- 46 Xu, J., Su, Q., Li, X., Ma, J., Song, W., Zhang, L., & Su, X. (2024). A Spatial Downscaling Framework for SMAP Soil Moisture Based on Stacking Strategy. *Remote Sensing*, 16(1), 200.
- 47 Xu, Z., Sun, H., Gao, J., Wang, Y., Wu, D., Zhang, T., & Xu, H. (2024). PhySoilNet: A deep learning downscaling model for microwave satellite soil moisture with physical rule constraint. *International Journal of Applied Earth Observation and Geoinformation*, 135, 104290.
- 48 Yin, W., Hu, L., Zhang, M., Wang, J., & Han, S. (2018). Statistical downscaling of GRACE-derived groundwater storage using ET data in the North China plain. *Journal of Geophysical Research: Atmospheres*, 123(11), 5973–5987.
- 49 Zaitchik, B. F., Rodell, M., & Olivera, F. (2010). Evaluation of the Global Land Data Assimilation System using global river discharge data and a source-to-sink routing scheme. *Water Resources Research*, 46(6).
- 50 Zuo, J., Xu, J., Li, W., & Yang, D. (2019). Understanding shallow soil moisture variation in the data-scarce area and its relationship with climate change by GLDAS data. *Plos One*, 14(5), e0217020.