



Calibration of ERA5 daily precipitation using MLP, D-Tree, and KNN algorithms in Razavi Khorasan province

Majid Rajabi Jaghargh¹, Seyed Mohammad Mousavi Baygi^{*}, Seyed Alireza Araghi³, Hadi Jabari Noghabi⁴

1. Ph.D. Student, Department of Water Science and Engineering, Faculty of Agriculture, Ferdowsi University of Mashhad, Mashhad, Iran. Email: magidrajabijaghargh@mail.um.ac.ir
2. Professor, Department of Water Science and Engineering, Faculty of Agriculture, Ferdowsi University of Mashhad, Mashhad, Iran. Email: mousavib@um.ac.ir
3. Assistant Professor, Department of Water Science and Engineering, Faculty of Agriculture, Ferdowsi University of Mashhad, Mashhad, Iran. Email: a.araghi@um.ac.ir
4. Associate Professor, Department of Statistics, Faculty of Mathematical Sciences, Ferdowsi University of Mashhad, Mashhad, Iran. Email: jabbarimh@um.ac.ir

ARTICLE INFO	ABSTRACT
Article type: Research Paper	Satellite products are the only available data source with adequate spatial coverage, however, their data do not match the observed values and have biases, although this discrepancy cannot be fixed precisely, however, a solution to reduce the bias is data recalibration. Currently, machine learning techniques are used to improve the accuracy of forecasting various types of weather phenomena, so regression solving such problems through methods based on machine learning and deep learning is very efficient. The daily precipitation of 19 rain gauge stations of the Ministry of Energy between 2010 and 2021 was extracted and compared to the average values of their corresponding daily precipitation pixels in the ERA5 database. To measure the data, three algorithms D-Tree, KNN, and MLP were used. The range of changes of correlation coefficient in MLP, D-Tree, and KNN is equal to [0.87, 0.98], [0.75, 0.97], and [0.4, 0.87], respectively. In addition, the range of changes for RMSE in MLP varies from 0.7 to 2.4 mm per day, and these changes for D-Tree and KNN are calculated between 0.8 to 2.2 and 1.2 to 2.5, respectively. In 75% of stations, RMSE in MLP, D-Tree, and KNN algorithms is less than 1.5, 1.9, and 2.2 mm per day, respectively. The range of bias changes in MLP is [0.18, -0.6 mm per day] and this range of changes for D-Tree and KNN is respectively [0.16, 0.5 mm per day] and [0.6, -0.8 mm per day] have been calculated. The bias of corrected data and observed values in MLP, D-Tree, and KNN algorithms for the middle of the stations is -0.09, -0.11, and -0.16 mm per day, respectively. The evaluation of the performance of three machine learning algorithms (MLP, D-Tree, and KNN) in correcting the daily precipitation of the ERA5 database and the comparison of CC, RMSE, and bias statistical indices for the reproduced data compared to ground values showed that in all three statistical indices, the MLP algorithm works better than the others and has good accuracy for correcting the daily precipitation.
Article history	
Received: 27 November 2023	
Revised: 11 January 2024	
Accepted: 14 January 2024	
Published online: 06 January 2024	
Keywords: Calibration, database, statistical indicators, machine learning	

Citation: Rajabi Jaghargh, M., Mousavi Baygi, S. M., Araghi, S.A., & Jabari Noghabi, H. (2024). Calibration of ERA5 daily precipitation using MLP, D-Tree, and KNN algorithms in Razavi Khorasan province. *Iranian Journal of Rainwater Catchment Systems*, 12(1), 129-147.

DOR: 10.1001.1.24235970.1403.12.1.8.1

Publisher: Iranian Rainwater Catchment Systems Association

© Author(s)



*Corresponding author: Seyed Mohammad Mousavi Baygi

Address: Department of Water Science and Engineering, Ferdowsi University of Mashhad, Mashhad, Iran.

Tel: +989153167311

Email: mousavib@um.ac.ir



Calibration of ERA5 daily precipitation using MLP, D-Tree, and KNN algorithms in Razavi Khorasan province

Majid Rajabi Jaghargh¹, Seyed Mohammad Mousavi Baygi^{*2}, Seyed Alireza Araghi³, Hadi Jabari Noghabi⁴

1. Ph.D. Student, Department of Water Science and Engineering, Faculty of Agriculture, Ferdowsi University of Mashhad, Mashhad, Iran. Email: magidrajabijaghargh@mail.um.ac.ir
2. Professor, Department of Water Science and Engineering, Faculty of Agriculture, Ferdowsi University of Mashhad, Mashhad, Iran. Email: mousavib@um.ac.ir
3. Assistant Professor, Department of Water Science and Engineering, Faculty of Agriculture, Ferdowsi University of Mashhad, Mashhad, Iran. Email: a.araghi@um.ac.ir
4. Associate Professor, Department of Statistics, Faculty of Mathematical Sciences, Ferdowsi University of Mashhad, Mashhad, Iran. Email: jabbarimh@um.ac.ir

EXTENDED ABSTRACT

Introduction: Spatial and temporal changes in precipitation significantly affect access to water for various human needs and environmental conservation. Although ground-based rain gauge systems are the most reliable tools for monitoring changes in precipitation at a point scale, various factors such as topography, remote locations, and budget constraints limit their geographical use. Therefore, ground-based rain gauge networks in many parts of the world are often spatially sparse and have weaknesses in spatial coverage. Satellite products are the only available data source with global coverage; however, there is a mismatch between them and ground observations. The values provided by sensors and data platforms cannot accurately estimate precipitation data due to multiple issues. These problems are not easily solvable, however, one solution to reduce ambiguities is to calibrate the estimated data. Currently, machine learning techniques are being employed to improve the accuracy of predicting various weather phenomena. Therefore, solving regression problems of this nature through machine learning and deep learning methods is not only possible but also very efficient. To this end, this study aimed to calibrate daily precipitation values in the ERA5 database by fitting these values against daily precipitation from 19 rain gauge stations of the Ministry of Energy in the Khorasan Razavi province, using three algorithms: D-Tree, K-NN, and MLP.

Methodology: Khorasan-Razavi province is located in the northeast of Iran with an area of about 117,000 km². This province contains arid and semi-arid areas with a wide range of temperature changes and precipitation patterns. Half of its area is made up of mountainous areas and the other half is plain and low-altitude areas. The climate of this province is cold dry and its average annual precipitation is about 254 mm. The southern regions of this province have less than 150 mm of precipitation per year, and the central and northwestern regions receive more than 320 mm. In this study, we used three popular and widely used algorithms, decision trees, K-nearest neighbors, and artificial neural networks for calibrating precipitation values. Daily precipitation values from 19 rain gauge stations of the Ministry of Energy have been extracted since 2010 and compared to the average daily precipitation pixel values from the ERA5 database. Using the above algorithms, terrestrial precipitation values were estimated based on ERA5 data, and correlation coefficient, RMSE, and bias indices were extracted and compared for each of the three algorithms.

Results and Discussion: Performance evaluation of machine learning algorithms (MLP, D-Tree, KNN) in correcting daily precipitation from ERA5 database and comparing them with statistical indices CC, RMSE, and bias for the reproduced data relative to ground values showed that in all three statistical indices, the MLP algorithm performed better than the other two algorithms and had suitable accuracy for correcting daily precipitation from ERA5. The range of changes in correlation coefficient in MLP, D-Tree, and KNN is [0.87, 0.98], [0.75, 0.97], and [0.4, 0.87], respectively, and these changes for RMSE in MLP are between 0.7 and 2.4 mm per days. These

***Corresponding author:** Seyed Mohammad Mousavi Baygi

Address: Department of Water Science and Engineering, Ferdowsi University of Mashhad, Mashhad, Iran.

Tel: +989153167311

Email: mousavib@um.ac.ir

changes for D-Tree and KNN are calculated between 0.8 to 2.2 and 1.2 to 2.5, respectively, in 75 stations, the RMSE in MLP, D-Tree, and KNN algorithms is less than 5.5, 1.9, and 2.2 mm per day. The range of bias changes in MLP [0.18, -0.6 mm d⁻¹] and this range of changes for D-Tree and KNN is respectively calculated as [0.16, 0.5 mm d⁻¹] and [0.6, -0.8 mm d⁻¹]. The corrected and observed bias in MLP, D-Tree, and KNN algorithms for the middle of the stations are equal to -0.09, -0.11, and -0.16 mm per day, respectively.

Conclusion: With the advent of advanced technologies and the presence of satellites and associated databases in the field of meteorological science, the time series of satellite precipitation values have become long enough to analyze their applicability for water resources management. Although these data are cheap and easily accessible, they do not have enough accuracy and they need to be calibrated with ground values. This study was shown that MLP performed better in calibrating daily precipitation values in ERA5 compared to the other two algorithms, as it effectively increased the correlation coefficient and desirably reduced RMSE and bias. For further informations, it is suggested to use deep learning techniques for calibrating monthly and annual precipitation values, as well as to investigate and analyze the accuracy of microscaling and correction of satellite precipitation data in the Google Earth Engine system.

Ethical Considerations

Data availability statement: The datasets are available upon a reasonable request to the corresponding author.

Funding: This research has been extracted from a doctoral thesis under the support of Ferdowsi University of Mashhad.

Authors' contribution: Majid Rajabi Jaghargh prepared the first draft and it was edited by Seyed Mohammad Mousavi Baygi, Seyed Alireza Araghi, and Hadi Jabari Noghabi.

Conflicts of interest: The authors of this article declared no conflict of interest regarding the authorship or publication of this article.

Acknowledgment: Khorasan Razavi Regional Water Company is thanked for providing the required statistics and information.



واسنجی بارش روزانه ERA5 با استفاده از الگوریتم‌های MLP، D-Tree و KNN در استان خراسان رضوی

مجید رجبی جاغرق^۱، سید محمد موسوی بایگی^{۲*}، سید علیرضا عراقی^۳، هادی جباری نوقابی^۴

۱. دانشجوی دکتری، گروه علوم و مهندسی آب، دانشکده کشاورزی، دانشگاه فردوسی مشهد، مشهد، ایران.

magidrajabijaghargh@mail.um.ac.ir

۲. استاد، گروه علوم و مهندسی آب، دانشکده کشاورزی، دانشگاه فردوسی مشهد، مشهد، ایران.

mousavib@um.ac.ir

۳. استادیار، گروه علوم و مهندسی آب، دانشکده کشاورزی، دانشگاه فردوسی مشهد، مشهد، ایران.

a.araghi@um.ac.ir

۴. دانشیار، گروه آمار، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد، مشهد، ایران.

jabbarimh@um.ac.ir

مشخصات مقاله	چکیده
نوع مقاله: پژوهشی	محصولات ماهواره‌ای تنها منبع داده موجود با پوشش فضایی مناسب است، با این وجود، داده‌های آن‌ها بر مقادیر مشاهداتی متنطبق نموده و دارای انحراف است، هرچند این عدم تطابق به طور دقیق قبل رفع نیست، با این حال، یک راه حل کاهش سوگیری، واسنجی داده‌های است. در حال حاضر، تکنیک‌های یادگیری ماشین برای بهبود دقت پیش‌بینی انواع مختلف پدیده‌های آب و هوایی به کار گرفته‌می‌شوند، لذا حل رگرسیونی مسائلی از این قبیل از طریق روش‌های مبتنی بر یادگیری ماشین و آموزش عمیق بسیار کارآمد است. بارش روزانه ۱۹ ایستگاه باران سنج ثبات وزارت نیرو بین سال‌های ۲۰۲۱ تا ۲۰۲۱ میلادی استخراج شد و در مقابل مقادیر متوسط پیکسل‌های بارش روزانه متناظر آن‌ها در پایگاه داده ERA5 قرار گرفت. به منظور واسنجی داده‌ها، از سه الگوریتم MLP، KNN و D-Tree و KNN و D-Tree به ترتیب برابر [۰.۹۸، ۰.۹۷، ۰.۹۷] و [۰.۷۵، ۰.۹۷] و [۰.۴، ۰.۸۷] است. هم‌چنین این دامنه تغییرات برای RMSE بین ۰/۷ تا ۲/۴ میلی‌متر در روز متغیر بوده و این تغییرات برای KNN و D-Tree به ترتیب بین ۰/۸ تا ۲/۵ و ۱/۲ تا ۲/۲ محسوبه شده‌اند. در ۷۵ درصد ایستگاه‌ها در الگوریتم‌های MLP، KNN و D-Tree به ترتیب کمتر از ۱/۵، ۱/۶ و ۱/۹ و ۲/۲ میلی‌متر در روز است. دامنه تغییرات سوگیری در MLP، KNN و D-Tree به ترتیب [۰/۰، ۰/۱۸]، [۰/۰، ۰/۶] و [۰/۰، ۰/۸] میلی‌متر در روز بوده و این دامنه تغییرات برای سوگیری داده‌های اصلاحی و مقادیر مشاهده شده، در الگوریتم‌های MLP، KNN و D-Tree برای ایستگاه‌ها به ترتیب برابر [۰/۰، ۰/۵]، [۰/۰، ۰/۵] و [۰/۰، ۰/۱۶] میلی‌متر در روز است. ارزیابی عملکرد سه الگوریتم یادگیری ماشین (MLP، KNN و D-Tree) در تصحیح بارش روزانه پایگاه داده ERA5 و مقابله شاخص‌های آماری CC، RMSE و سوگیری برای داده‌های بازتولید شده نسبت به مقادیر زمینی نشان داد که در هر سه شاخص آماری الگوریتم MLP نسبت به دوالگوی دیگر بهتر عمل نموده و از دقت مناسبی برای تصحیح بارش روزانه برخوردار است.
دریافت: ۶ آذر ۱۴۰۲ بازنگری: ۲۱ دی ۱۴۰۲ پذیرش: ۲۴ دی ۱۴۰۲ انتشار برخط: ۱۶ خرداد ۱۴۰۳	
واژه‌های کلیدی: پایگاه داده، شاخص‌های آماری، یادگیری ماشین، واسنجی	
استناد: مجید رجبی جاغرق، سید محمد موسوی بایگی، سید علیرضا عراقی، هادی جباری نوقابی، داده ERA5 با استفاده از الگوریتم‌های MLP، D-Tree و KNN در استان خراسان رضوی. سامانه‌های سطوح آبگیر باران، (۱۲)، (۱)، ۱۴۷-۱۲۹.	
DOR: 20.1001.1.24235970.1403.12.1.8.1	ناشر: انجمن علمی سیستم‌های سطوح آبگیر باران ایران



© نویسنده‌گان

* نویسنده مسئول: سید محمد موسوی بایگی

نشانی: گروه علوم و مهندسی آب، دانشکده کشاورزی، دانشگاه فردوسی مشهد، مشهد، ایران
تلفن: ۰۹۱۵۳۱۶۷۳۱۱

پست الکترونیکی: mousavib@um.ac.ir

مقدمه

بارش یکی از متغیرهای کلیدی در چرخه آب و اصلی‌ترین داده در مدل‌سازی‌های هیدرولوژیکی محسوب می‌شود (Yuan et al., 2018, Beck et al., 2017, Liu et al., 2017). بنابراین داده‌های دقیق و قابل اعتماد بارش نه تنها برای بهبود دقت شبیه‌سازی‌های هیدرولوژیکی بسیار مهم هستند (Sun et al., 2018, Prakash et al., 2013, Pena et al., 2013). امروزه مقادیر بارش‌های تخمینی مبتنی بر اطلاعات ماهواره‌ای، به عنوان فرآیندهای آب و هواشناسی ضروری هستند (Ma et al., 2015). اما از آن‌جا که جایگزینی مناسب در مطالعات مرتبط با علوم آب و محیط زیست، به طور فزاینده‌ای در دسترس عموم قرار گرفته است (Li et al., 2018)، اگرچه این محصولات دارای یک پوشش فضایی وسیع و واضح مکانی و زمانی مناسب هستند (Zubieta et al., 2015)، اما از آن‌جا که الگوریتم‌های تخمین بارش ماهواره‌ای، به صورت جهانی تنظیم شده‌اند و امکان نادیده گرفتن عوامل تأثیرگذار محلی، در تخمین بارش بسیار محتمل است، لذا عدم تطابق بین داده‌های ماهواره‌ای و زمینی دور از انتظار نیست، بنابراین ضروری است تا نسبت به اصلاح و واستنجی این مقادیر با استفاده از روش‌های موجود اقدام شود. پژوهشگران برای اصلاح داده‌های تخمینی، بسته به میزان دقت مورد نیاز خود از روش‌های مختلفی استفاده می‌کنند. برای نمونه، (Chen et al., 2013) روش‌های تصحیح سوگیری را به دو دسته طبقه‌بندی نمودند: ۱) رویکردهای متکی بر میانگین، که معمولاً از عوامل تصحیح ماهانه بر اساس نسبت بین بارش و شبکه و مقادیر مشاهده شده استفاده می‌کنند، ۲) رویکردهای مبتنی بر توزیع، که معمولاً در آن‌ها از توابع احتمالی استفاده می‌شود، مانند تابع توزیع تجمعی (CDF) (Guo et al., Hamill et al., 2018, Vrac et al., 2016, Amengual et al., 2012, Yang et al., 2010, Piani et al., 2010) (2018). رویکردهای متکی بر میانگین، روشی به نسبت ساده و آسان برای پیاده‌سازی در مقایسه با روش تابع توزیع تجمعی است، لذا در مطالعاتی از این قبیل با اقبال همراه بوده است (Chen et al., 2013, Themeßl et al., 2011).

وجود خطاهای تصادفی، برای تخمین بارش ماهواره‌ای در مقیاس روزانه، باعث سوگیری این داده‌ها نسبت به مقادیر بارش ایستگاه‌های زمینی می‌شود، لذا برای تصحیح این ناسازگاری فضایی، محققین اقدام به ایجاد رابطه بین بارش و متغیرهای مکانی نموده‌اند. ژانگ (Zhang et al., 2019a), کل منطقه حوضه آبریز رودخانه لنکانگ-مکونگ را با توجه به نتایج تجزیه و تحلیل REOF برای بارش زمینی و ماهواره‌ای به هفت زیرمنطقه تقسیم کرد و دریافت، کاهش مقیاس، باعث ارائه نتایج بهتر می‌شود. (Zhang et al., 2019b) عملکرد و کاربرد هیدرولوژیکی^۱ TMPA را ارزیابی نمودند. آن‌ها منطقه مورد مطالعه را به چهار زیرحوضه تقسیم کردند و مدل تنظیم جغرافیایی را برای تصحیح داده‌های TMPA بر اساس چند ضلعی‌های تیسن^۲ در هر زیرحوضه به کار بردند و نتایج رضایت‌بخشی را به دست آوردند. (Heredia et al., 2018) برای انتباق فضایی بارش، رویکرد تابع توزیع تجمعی^۳ (CDF) را برای تصحیح مقیاس نقطه‌ای و اصلاح شبکه، با تقسیم ناحیه مورد نظر، با استفاده از چند ضلعی‌های تیسن پیشنهاد کرد. (Guo et al., 2018) معتقد است، بارش تخمینی توسط ماهواره، به شدت بارش وابسته بوده و برای اثبات ادعای خود، شش حوضه آبریز درکشور چین را که از نظرلویژگی‌هایی مانند اقلیم، عرض جغرافیایی، طول جغرافیایی و ارتفاع متفاوت بودند را انتخاب نمود و نسبت به اعتبارسنجی بارش در آن‌ها اقدام کرد، او دریافت که روش CDF را باید در یک منطقه خاص به کار بست تا بتواند به طور موثری ویژگی‌های بارش را در رژیمهای مختلف اقلیمی معنکس نماید. دانستن مقدار بارش در یک شبکه متراکم فضایی و برای یک دوره زمانی طولانی می‌تواند در حل انواع مسائل مهندسی، مشمر ثمر باشد (Blöschl et al., 2019)، اگرچه، داده‌های شبکه بارانسنج زمینی از دقت بالایی برخوردارند، اما نگهداری از چنین شبکه‌ای با تراکم فضایی بالا و برای مدت زمان طولانی بسیار هزینه‌بر است (Mega et al., 2019, Tang et al., 2022). لذا باید به دنبال روش‌های جایگزین گشت که قادر باشد هم از نظر زمانی و هم به لحاظ مکانی مقادیر بارش را با دقت مطلوبی ارائه نماید.

با ظهور فناوری‌های پیشرفته و حضور ماهواره‌ها و پایگاه‌های داده مرتبط در حوزه علم هواشناسی، سری‌های زمانی مقادیر بارش ماهواره‌ای به اندازه کافی طولانی شده‌اند تا قابلیت استفاده از آن‌ها را برای مدیریت منابع آب، تجزیه و تحلیل نمود، هرچند این داده‌های تولیدی ارزان و سهل‌الوصول هستند اما از دقت کافی برخوردار نیستند و برای استفاده از آن‌ها نیاز است با مقادیر زمینی واستنجی شوند. (Mega et al., 2019, Tang et al., 2022)

¹ Rotated Empirical Orthogonal Function

² TRMM Multi-Satellite Precipitation Analysis

³ Thiessen

⁴ Cumulative Distribution Function

با ادغام محصولات بارش ماهواره‌ای شبکه‌بندی شده و اندازه‌گیری‌های زمینی، این امکان وجود دارد به داده‌های با دقت و کیفیت بالاتر دست یافت، که نه تنها دقیق‌تر از داده‌های ماهواره‌ای است بلکه به طور همزمان، فضای را با تراکم بسیار بالاتر در مقایسه با اندازه‌گیری‌های زمینی پوشش می‌دهد. این ادغام، عملاً حل رگرسیونی مسئله در یک محیط مکانی است، این روش‌ها معمولاً تحت عنوان کاهش مقیاس نیز معرفی می‌شوند و نوع خاصی از درون‌یابی فضایی هستند. مسئله اخیر در زمینه‌های مختلفی قابل مشاهده است و می‌توان انواع روش‌های مرتبط با کاهش مقیاس بارش را در (Hu et al. 2019) Abdollahipour et al. (2022) به طور جامع‌تری یافته. درون‌یابی فضایی بارش با ادغام محصولات بارش ماهواره‌ای و اندازه‌گیری‌های زمینی در مقیاس‌های زمانی و مکانی مختلف و با استفاده از انواع الگوریتم‌های رگرسیونی از جمله الگوریتم‌های یادگیری ماشین توسط محققان مختلفی انجام شده است. خلاصه اطلاعات روش‌شنختی برخی از این مطالعات در جدول (۱) خلاصه شده است.

جدول ۱- خلاصه اطلاعات روش‌شنختی الگوریتم‌های یادگیری ماشین

Table 1- Summary of methodological information of machine learning algorithms

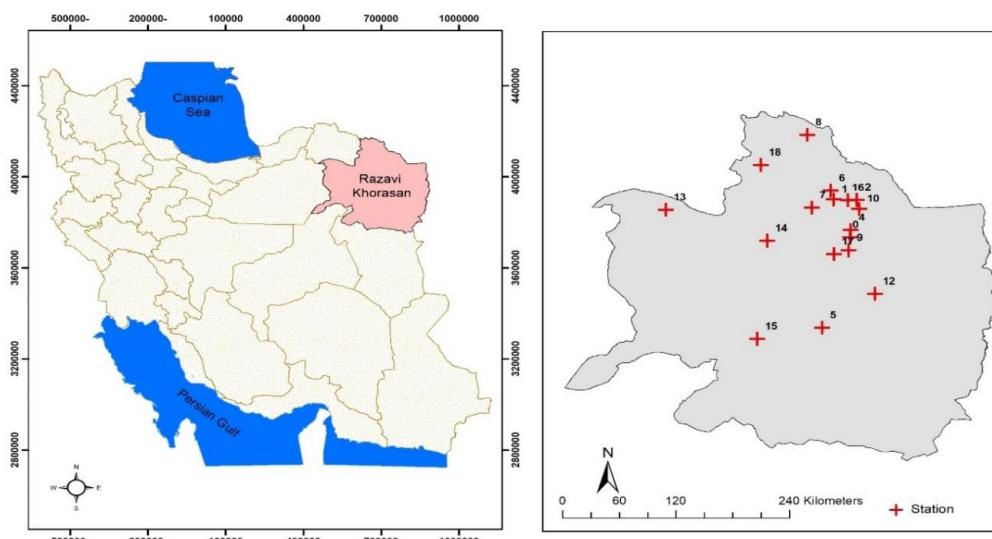
Study	Time Scale	Spatial Scale	Algorithms
He et al. 2016	Hourly	South-western, central, north-eastern and south-eastern United States	Random forests
Meyer et al. 2016	Daily	Germany	Random forests, artificial neural networks, support vector regression
Tao et al. 2016	Daily	Central United States	Deep learning
Yang et al. 2016	Daily	Chile	Quantile mapping
Baez-Villanueva et al. 2020	Daily	Chile	Random forests
Chen et al. 2020	Daily	Dallas–Fort Worth in the United States	Deep learning
Chen et al. 2020	Daily	Xijiang basin in China	Geographically weighted ridge regression
Rata et al. 2020	Annual	Chélieff watershed in Algeria	Kriging
Chen et al. 2021	Monthly	Sichuan Province in China	Artificial neural networks, geographically weighted regression, kriging, random forests
Nguyen et al. 2021	Daily	South Korea	Random forests
Shen and Yong 2021	Annual	China	Gradient boosting decision trees, random forests, support vector regression
Zhang et al. 2021	Daily	China	Artificial neural networks, extreme learning machines, random forests, support vector regression
Chen et al. 2021	Daily	Coastal mountain region in the western United States	Deep learning
Fernandez Palomino et al. 2022	Daily	Ecuador and Peru	Random forests
Lin et al. 2022	Daily	Three Gorges Reservoir area in China	Adaptive boosting decision trees, decision trees, random forests
Yang et al. 2022	Daily	Kelantan river basin in Malaysia	Deep learning
Zandi et al. 2022	Monthly	Alborz and Zagros mountain ranges in Iran	Artificial neural networks, locally weighted linear regression, random forests, stacked generalization, support vector regression
Militino et al. 2023	Daily	Navarre in Spain	K-nearest neighbors, random forests, artificial neural networks

یادگیری ماشین روشی است که در سال‌های اخیر به طور گسترده مورد استفاده قرار گرفته است، با پیشرفت فناوری، این روش به یک روش عمومی برای طبقه‌بندی، و همچنین برای مشکلات تشخیص تغییرات و ناهنجاری در علوم زمین تبدیل شده است (Gómez- Chova et al., 2015, Muhlbauer et al., 2014). برای هواشناسی، که شاخه‌ای از علوم زمین است، روش‌های یادگیری ماشین نیز با موفقیت به کار گرفته شده، به عنوان مثال، شبکه‌های عصبی مصنوعی (ANN) برای بهبود تخمین بارش (Pan et al., 2019)، شناسایی مرحله ENSO و پیش‌بینی تأثیر آن از جمله این موارد هستند (Toms et al., 2020).

انتخاب استان خراسان‌رضوی به عنوان منطقه هدف در این تحلیل، به واسطه دسترسی به گراف‌های باران‌سنج‌های ثبات وزارت نیرو برای تحلیل‌های آماری بود. با توجه به این که مقادیر بارش سنجنده‌ها و پایگاه‌های داده می‌تواند در آینده‌ای نزدیک جایگزین شبکه ایستگاه‌های زمینی شود، بنابراین استفاده از محصولات آن‌ها به سرعت در حال افزایش است، اما مقادیر بارش تخمینی توسعه ماهواره‌ها و پایگاه‌های داده بر مقادیر مشاهداتی منطبق نبوده لذا واسنجی داده‌های تخمینی را ضروری می‌نماید، لذا باید تا نسبت به ارزیابی دقت داده‌های ماهواره‌ای و واسنجی آن‌ها با استفاده از روش‌های نوین از جمله یادگیری ماشین اقدام شود تا تطابق داده‌ها حداکثر و میزان خطاهای به حداقل ممکن برسد. آن‌چه نقطه قوت و وجه تمایز این پژوهش است شامل استفاده از سه الگوریتم یادگیری ماشین (درخت تصمیم، K نزدیک‌ترین همسایگان، پرسپترون چندلایه) و انتخاب بهترین برآوردگر، برای واسنجی و کاهش سوگیری مقادیر بارش روزانه پایگاه ERA5 در استان خراسان‌رضوی است.

مواد و روش تحقیق محدوده مورد مطالعه

استان خراسان‌رضوی در شمال شرقی ایران واقع شده و دارای مساحتی در حدود ۱۱۷ هزار کیلومترمربع است، این استان، مناطق خشک و نیمه خشک با طیف گسترده‌ای از تغییرات دمایی و الگوهای بارش را در خود جای داده است. نیمه از وسعت آن را مناطق کوهستانی و نیمه دیگر را مناطق دشتی و کم ارتفاع تشکیل داده‌اند، مرتفع‌ترین نقطه آن در قله بینالود با ارتفاع ۳۳۳۹ متر و پست‌ترین نقطه آن در دشت سرخس با ارتفاع ۲۹۹ متر از سطح دریا واقع شده‌اند، اقلیم این استان براساس معیار طبقه‌بندی آمبرژه، اقلیم خشک سرد است و متوسط بارش سالانه آن حدود ۲۵۴ میلی‌متر در سال است. مناطق جنوبی این استان دارای بارشی کمتر از ۱۵۰ میلی‌متر در سال بوده و نواحی مرکزی و شمال غربی آن بیش از ۳۲۰ میلی‌متر بارش در سال را دریافت می‌نمایند (رجی و همکاران، ۱۴۰۲). شکل (۱) نقشه استان خراسان‌رضوی به همراه توزیع مکانی ایستگاه‌های باران‌سنج ثبات وزارت نیرو را نمایش می‌دهد.



شکل ۱- نقشه استان خراسان‌رضوی و پراکندگی مکانی ایستگاه‌های باران‌سنج ثبات وزارت نیرو

Figure 1- Map of Razavi Khorasan Province and spatial distribution of automatic rain gauge stations of the Ministry of Energy

داده‌های مورد استفاده

در این تحقیق از داده‌های بارش روزانه ۱۹ ایستگاه باران‌سنج ثبات وزارت نیرو که مشخصات، موقعیت و توزیع مکانی آن‌ها در جدول (۲) و شکل (۱) ارائه شده است، به همراه مقادیر بارش روزانه پایگاه داده ERA5 در محل ایستگاه‌های زمینی، برای سال‌های ۲۰۲۱ تا ۲۰۱۰ استفاده شده است.

داده‌های ایستگاه‌های باران سنجی

برای دستیابی به بارش روزانه ایستگاه‌های باران سنج ثبات، اقدام به استخراج مقادیر بارش‌های رخداده با گام زمانی نیم ساعته با استفاده از نرم‌افزار RRDgitizer ver 2.1.1.1 از گراف آن‌ها شد. با توجه به اختلاف زمان موجود بین زمان تنظیمی باران سنج‌ها در منطقه مورد مطالعه با زمان مرجع بین‌المللی (UTC)^۱، زمان و میزان بارش‌ها با زمان مرجع منطبق شد و بارش روزانه با جمع مقادیر ساعتی محاسبه شد. این عمل برای کاهش خطای میزان تطابق بارش از دو مرجع متفاوت (ایستگاه‌های زمینی و داده‌های ERA5) و در دو سیستم زمانی مختلف صورت گرفت.

جدول ۲- مشخصات ایستگاه‌های باران سنجی ثبات وزارت نیرو مورد استفاده در پژوهش

Table 2 - Specifications of the automatic rain gauge stations of the Ministry of Energy used in the research

مشخصات جغرافیایی				مشخصات جغرافیایی				ایستگاه	شناخت
(m) ارتفاع	UTM_Y	UTM_X	ایستگاه	ID	(m) ارتفاع	UTM_Y	UTM_X		
1265	4056093	740695	سد کاردنه	11	970	4021779	731051	اداره مشهد	1
273	4044721	334933	سرخس	12	1320	4067524	713648	ارداک بند ساروج	2
1415	3953876	757541	فریمان	13	1452	4066787	738162	آل	3
1138	4054892	536166	کارخانه قند چوبین	14	938	3900564	281783	باغسنگان تربت‌جام	4
1199	4017696	643474	کارخانه قند نیشابور	15	1920	4030506	731838	بلغور	5
1054	3899535	632645	کاشمر	16	1468	3913151	701468	تربت‌حیدریه	6
1561	4066333	728919	گوش بالا	17	1540	4078183	710544	تلغور	7
1770	4001739	714070	منان	18	1170	4057482	690538	چنان	8
1333	4108912	636561	هی‌هی قوچان	19	429	4145323	685794	درگز	9
					1240	4006132	729885	سد طرق	10

پایگاه داده ERA5

نسل پنجم تجزیه و باز تحلیل مقادیر پدیده‌های جوی است که داده‌های خروجی آن برگرفته از باز تحلیل اطلاعات استخراج شده از مدل ECMWF^۲ است که تجزیه و باز تحلیل اطلاعات در آن با منظور نمودن تأثیر مقادیر شبکه ایستگاه‌های زمینی، انجام می‌شود، متغیرهای جوی در این پایگاه، دارای وضوح زمانی حداقل یک ساعته و دقت مکانی 0.25×0.25 درجه هستند و از سال ۱۹۷۵ میلادی با تأخیر اندکی نسبت به زمان واقعی در دسترس قرار دارند (رجی و همکاران، ۱۴۰۲). مقادیر بارش روزانه پایگاه داده ERA5 به لحاظ کمی و کیفی دارای داده‌هایی مناسب دراستان خراسان‌رضوی هستند (رجی و همکاران، ۱۴۰۲)، از این رو از مقادیر بارش روزانه این پایگاه داده در فرایند تجزیه، تحلیل و واسنجی استفاده شد. مقادیر بارش روزانه پایگاه داده ERA5 از سامانه متن باز و تحت وب GEE^۳ به آدرس ERA5_LAND/HOURLY استخراج شد.

روش تحقیق

برای ارزیابی کارایی الگوریتم‌های رگرسیونی یادگیری ماشین، برای محصول بارش روزانه ERA5 در محل ایستگاه‌های باران سنج ثبات واقع در محدوده مورد مطالعه تحلیل‌های آماری انجام شد. مقادیر بارش روزانه ارائه شده توسط پایگاه داده ERA5، طی سال‌های ۲۰۱۰ تا ۲۰۲۱ (منطبق بر مهر ۱۳۸۹ تا پایان شهریور ۱۴۰۰) به مدت ۱۱ سال آبی، در پیکسل‌های متناظر با محل ایستگاه باران سنج ثبات واقع در محدوده مورد مطالعه استخراج شد و سپس با استفاده از الگوریتم‌های MLP، KNN، D-Tree، اقدام به محاسبه مقادیر تصحیح شده بارش روزانه این پایگاه داده شد، پارامترهای آماری ضریب همبستگی (CC)، ریشه دوم میانگین مربع خطای (RMSE) و سوگیری (Bias) در محل هر ایستگاه، محاسبه و بهترین برآوردگر، جهت واسنجی وافزایش دقت، تصحیح سوگیری و کاهش خطای داده‌های بارش روزانه شناسایی شد. مقادیر بارش ایستگاه اداره مشهد به دلیل نامطلوب بودن داده‌های آن از روند تجزیه و تحلیل‌ها کنار گذاشته شد.

¹ Coordinated Universal Time

² European Centre for Medium-Range Weather

³ Google Earth Engine

واستنجی داده‌های بارش حاصل از مدل با ایستگاه‌های شاهد، به وسیله الگوریتم‌های یادگیری ماشین به طور گسترده‌ای در حال استفاده و توسعه است، لذا تشخیص الگوی مناسب، نیازمند تحلیل خروجی داده‌های آزمایشی است، بنابراین تحلیل‌های آماری مقادیر آزمایشی ۱۸ ایستگاه زمینی و برای سه الگوریتم یاد شده بهمنظور بهینه‌سازی و واستنجی داده‌های بارش روزانه ERA5 صورت پذیرفت تا مشخص شود، کدام الگوریتم شاخص‌های خطأ (Bias, RMSE) را کمینه و ضریب همبستگی (CC) را بیشینه می‌نماید و در کل، عملکرد بهتری از خود به نمایش می‌گذارد.

پردازش داده‌ها

مقادیر بارش روزانه ایستگاه‌های باران سنج ثبات و داده‌های متناظر با آن‌ها در پایگاه داده ERA5 برای ۷۲،۲۷۰ روز بین سال‌های ۲۰۱۰ تا ۲۰۲۱ استخراج شد، با توجه به اختلاف زمان ثبت اطلاعات، بین باران سنج وزارت نیرو و این پایگاه داده (اختلاف زمان محلی و بین‌المللی)، مقادیر بارش ایستگاه‌های زمینی از نظر زمانی بر داده‌های ERA5 منطبق شد و سپس با جمع مقادیر بارش کوتاه‌مدت، بارش روزانه برای هر دو منبع محاسبه شد. برای شناسایی و حذف داده‌های دورافتاده روش‌های مختلفی وجود دارد از جمله رسم نمودار جعبه‌ای با نمایه‌های چارکی، رسم نمودار پراکنده‌ی داده‌ها، رسم نمودار Q-Q و همچنین استفاده از مقاهم آماری مانند تعریف یک حد آستانه، دراین پژوهش برای شناسایی داده‌های دورافتاده از حد آستانه استفاده شد، داده‌های بزرگ‌تر و کوچک‌تر از سه برابر انحراف معیار شناسایی و پس از بررسی و تطبیق آن با مقادیر مشابه در ایستگاه‌های مجاور در ارتباط با حذف و یا استفاده از آن‌ها تصمیم‌گیری شد. ۲۵ درصد داده‌ها به عنوان داده‌های آزمایشی و ۷۵ درصد باقی‌مانده برای آموزش مدل انتخاب شدند.

الگوریتم‌های یادگیری ماشین

در حوزه هوش مصنوعی و یادگیری ماشین، روش‌ها و الگوهای گوناگونی برای حل مسائل رگرسیونی وجود دارد که بسته به شرایط حاکم بر پژوهش و نظر محقق از آن‌ها استفاده می‌شود. در این پژوهش از سه الگوریتم KNN، D-Tree، KNN و MLP بهمنظور واستنجی مقادیر بارش روزانه پایگاه داده ERA5 استفاده شد.

رگرسیون K نزدیک‌ترین همسایگان (K-Nearest Neighbors)

الگوریتم K نزدیک‌ترین همسایگان یکی از محبوب‌ترین رویکردهای داده کاوی است و فرض اصلی آن است که اگر همه K نزدیک‌ترین همسایه‌های یک نقطه در مجموعه آموزشی به یک دسته تعلق داشته باشند، فرض می‌شود که نقطه دارای ویژگی‌ها و کیفیت‌های یکسانی است. با استفاده از ویژگی‌های داده، KNN می‌تواند از روش حل معادله دشوار اجتناب نموده و بر همبستگی تمرکز نماید، الگوریتم KNN از شباهت ویژگی برای پیش‌بینی مقادیر هر نقطه داده جدید استفاده می‌کند. این بدان معنی است که نقطه جدید بر اساس شباهت بسیار نزدیک به نقاط مجموعه آموزشی است و بر همین اساس مقداری به آن اختصاص داده می‌شود. (فتاحی و جیری‌بابی، ۱۴۰۰؛ Huang et al., 2017). اما آن‌چه در KNN اتفاق می‌افتد، به صورت گام به گام در زیر به آن اشاره شده است:

- ابتدا فاصله بین نقطه جدید و هر نقطه در داده‌های آموزشی محاسبه می‌شود. -۲- نزدیک‌ترین k نقطه داده بر اساس فاصله انتخاب خواهند شد. -۳- میانگین این نقاط داده، پیش‌بینی نهایی برای مقداردهی به نقطه جدید است.
- اولین مرحله، محاسبه فاصله بین نقطه جدید و هر نقطه در داده‌های آموزشی است. روش‌های مختلفی برای محاسبه این فاصله وجود دارد که رایج‌ترین روش‌ها عبارتند از: اقلیدسی، منهتن و مینکوفسکی (برای داده‌های پیوسته)، فاصله مینکوفسکی شکل تعمیم یافته فاصله اقلیدسی و منهتن است.

- فاصله اقلیدسی: این فاصله به عنوان جذر مجدور تفاوت‌های یک نقطه جدید (x) و یک نقطه موجود (y) محاسبه می‌شود.

$$Euclidean = \sqrt{\sum_{i=1}^k (x_i - y_i)^2} \quad (1)$$

- فاصله منهتن: این فاصله بین بردارهای واقعی و برابر مجموع اختلاف مطلق آن‌ها است.

$$Manhattan = \sum_{i=1}^k |x_i - y_i| \quad (2)$$

مقدار K، تعداد همسایه‌هایی را مشخص می‌کند که برای برآورد یک مشاهده جدید، به میانگین این تعداد همسایه نزدیک به آن مشاهده رجوع می‌شود. با توجه به تغییر در مقدار K، نتیجه نهایی برآورد مقدار جدید، تمایل به تغییر پیدا می‌کند، لذا برای بهدست آوردن مقدار بهینه K باید براساس محاسبه خطای داده‌های آموزشی و اعتبارسنجی آن تصمیم گیری شود. مقدار بهینه این ضریب، عددی است که در آن خطای در مقادیر آموزشی و آزمایشی به حداقل خود می‌رسد، زیرا به حداقل رسیدن عدم قطعیت هدف نهایی در حل مسائلی از این دست است.

رگرسیون درخت تصمیم (Decision Tree)

رگرسیون درخت تصمیم، یک الگوریتم یادگیری نظارت شده است که برای پیش‌بینی مقادیر پیوسته استفاده می‌شود. این الگوریتم با تقسیم فضای ویژگی، به مجموعه‌ای از مستطیل‌ها، و سپس اختصاص دادن مقدار ثابت به هر مستطیل عمل می‌نماید. بنابراین پیش‌بینی یک مشاهده جدید، برابر مقدار ثابت است که در آن قرار می‌گیرد. در واقع، الگوریتم رگرسیون درخت تصمیم با تقسیم بازگشتی فضای ویژگی به مناطق کوچک‌تر و برازش یک مدل ساده (عموماً یک مقدار ثابت) برای هر منطقه کار می‌کند. این فرآیند توسط یکتابع ضرر Chaudhary & Dhanya، 2008 می‌شود که خطای بین مقادیر پیش‌بینی شده و مقادیر واقعی در داده‌های آموزشی را اندازه‌گیری می‌کند (Sandri & Zuccolotto، 2020).

- معیار تقسیم: در هر گره درخت تصمیم، این الگوریتم، یک ویژگی و یک مقدار آستانه را برای تقسیم داده‌ها به دو زیر مجموعه انتخاب می‌کند. این معیار تقسیم برای به حداقل رساندن تابع ضرراست، که با انتخاب یکی از شاخص‌های خطای مربعات میانگین مرتعات خطای میانگین خطای مطلق، صورت می‌گیرد.
- پارتیشن‌بندی بازگشتی: الگوریتم درخت تصمیم، به تقسیم داده‌ها به زیر مجموعه‌های کوچک‌تر در هر گره، بر اساس معیار تقسیم انتخابی ادامه می‌دهد. این فرآیند تا زمانی ادامه می‌یابد که یک معیار توقف، مانند رسیدن به حداقل عمق درخت یا حداقل داده در یک گره، تامین شود.

- برازش مدل: هنگامی که داده‌ها به مناطق کوچک‌تر تقسیم می‌شوند، یک مدل ساده (یا یک مقدار ثابت) برای هر منطقه برازش داده می‌شود. این مدل پیش‌بینی نقاط داده را در آن منطقه نشان می‌دهد.
- پیش‌بینی: برای پیش‌بینی یک نقطه داده جدید، الگوریتم، مسیر درخت تصمیم را طی می‌کند و مقدار ثابت ناحیه‌ای را که مشاهده جدید در آن قرار می‌گیرد را به عنوان مقدار پیش‌بینی شده اختصاص می‌دهد.

- پرسپترون‌های چند لایه (MLP)

پرسپترون چند لایه، یک شبکه عصبی مصنوعی با ساختار رو به جلو است که دارای لایه‌های مختلفی از گره‌ها است و هر کدام به لایه بعدی توسط همین گره‌ها متصل می‌شود. هر گره علاوه بر گره‌های ورودی از یک نورون با تابع فعال ساز غیرخطی تشکیل شده است. MLPها اغلب با استفاده از الگوریتم‌های پس انتشار آموزش داده می‌شوند و می‌توانند برای حل رگرسیون و طبقه‌بندی مسائل استفاده شوند. در ساختار یک MLP، گره‌های (نورون‌های) زیادی وجود دارد و در مجموع شامل سه لایه است. لایه ورودی، اولین لایه بوده و شامل متغیرهای مستقل برای پردازش است، لایه دوم در ساختار MLP، لایه‌های پنهان است که باید در میان لایه‌های خروجی و ورودی متمرکز شوند. راهبردهای های پیش‌خور، مقادیر را از یک لایه به لایه دیگر منتقل می‌کنند و بسته به نوع مسئله از یک یا چند لایه پنهان، بین ورودی و خروجی عبور نموده و از توابع انتقال برای تشخیص الگو موجود بین ورودی‌ها و خروجی‌ها استفاده می‌نمایند. MLP می‌تواند از انواع توابع انتقال استفاده کند که رایج‌ترین آن‌ها "Sigmoid" و "Tansig" است. لایه سوم، که به لایه خروجی موسوم است، در برگیرنده تخمین و یا پیش‌بینی مقادیر هدف است و در بیش‌تر موارد، "Pureline" به عنوان تابع انتقال لایه خروجی مورد استفاده قرار می‌گیرد. عملکرد یک شبکه تحت تأثیر تعداد لایه‌ها و نورون‌های پنهان، در آن لایه‌ها است. آزمون و خطای یک روش معمول و آسان برای بهینه‌سازی این مقادیر است. گام بعدی پس از ترسیم ساختار MLP، انتخاب وزن‌ها و بایاس‌های متعدد برای هر نورون در بخش آموزش است، تا کارایی شبکه بهبود یابد. شبکه‌های MLP با استفاده از تکنیک‌های مختلفی آموزش داده می‌شود، از جمله منظم‌سازی بیزی (BR)، لونبرگ-مارکوارت (LM) و گرادیان مزدوج مقیاس شده (SCG) که این الگوریتم‌ها وظیفه محاسبه بایاس‌ها و وزن‌ها را بر عهده دارند. بر حسب نیاز می‌توان از شبکه‌های عمیق و یا ساده استفاده نمود و این الگوریتم توانایی بسیار بالایی در یادگیری و آموزش دارد (قبائی سوق و همکاران، ۱۳۸۹؛ Khalili et al., 2016؛ Esteves et al., 2018).

توابع هزینه که می‌تواند در MLP رگرسیونی استفاده شوند عبارتند از میانگین مربعات خطای (MSE) و میانگین خطای مطلق (MAE)، هرچه مقدار خطای کمتر باشد مدل بهترین عملکرد را در آن محدوده از خود به نمایش گذاشته است. این مدل رگرسیونی مربع

خطاهای را از طریق توابع Stochastic gradient descent و ya بهینه می‌نماید، آنچه در مدل Rگرسیون پرسپترون چند لایه باید مدنظر قرار گیرد، تنظیم فرآپارامترهای اجرای آن از جمله، تعداد نورون‌ها، تابع فعال‌سازی برای لایه پنهان و تابع بهینه‌ساز، تعداد لایه پنهان و تعداد تکرار است که سعی و خطا در کمینه نمودن انحراف و به حداقل رساندن عدم قطعیت، ساده‌ترین روش برای بدست آوردن مقادیر بهینه است.

متغیرهای ارزیابی

به طور کلی برای مقایسه بین چندین منبع تولید اطلاعات یکسان و تشخیص بهترین منبع از شاخص‌های آماری برای تعیین دقت محصولات استفاده می‌شود. در این پژوهش نیز از معیارهای آماری مختلفی برای ارزیابی و مقایسه داده‌های تخمین بارش ارائه شده توسط الگوریتم‌های یادگیری ماشین استفاده شد. توابع آماری و شاخص‌های مورد استفاده در جدول (۳) ارائه شده است. از ضریب همبستگی پیرسون (CC)، برای نشان دادن میزان همبستگی بین مقادیر ثبت شده توسط باران‌سنج‌ها و مقادیر تخمینی توسط الگوریتم‌های یادگیری ماشین استفاده شد. همچنان، از میانگین ریشه دوم خطاهای (RMSE)^۲ برای ارزیابی مقدار خطای مجموعه داده‌های تخمینی استفاده شد. از شاخص آماری سوگیری، برای ارزیابی تمایل الگوریتم‌ها، به بیش برآورد یا کم برآورد مقادیر بارش استفاده شد (Sorooshian et al., 1993).

جدول ۳- شاخص‌های آماری مورد استفاده در ارزیابی عملکرد الگوریتم‌های یادگیری ماشین در تخمین بارش روزانه

Table 3- Statistical indicators used in evaluating the performance of machine algorithms in estimating daily precipitation

رابطه	معادله	واحد	دامنه تغییرات	مقدار بهینه
1	$e_i = S_i - G_i$	متغیر	$(-\infty, +\infty)$	0
2	$CC = \frac{\frac{1}{n} \sum_{i=1}^n (S_i - \bar{S})(G_i - \bar{G})}{\sigma_S \sigma_G}$	بدون واحد	$[-1, 1]$	+1, -1
3	$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n e_i^2}$	متغیر	$(0, +\infty)$	0
4	$Bias = \frac{\sum_{i=1}^n e_i}{\sum_{i=1}^n G_i}$	متغیر	$(-\infty, +\infty)$	0

در روابط مندرج در جدول فوق، n تعداد نمونه، G مقادیر ثبت شده توسط باران‌سنج، S مقادیر شبیه‌سازی شده توسط الگوریتم، \bar{S} و \bar{G} به ترتیب میانگین مقادیر تخمین زده شده بارش توسط الگوریتم و باران‌سنج است. همچنان تغییرات بارش توسط دو پارامتر σ_G و σ_S برآورد شد که نشان‌دهنده انحراف استاندارد مقادیر ثبت شده باران توسط باران‌سنج و الگوریتم است.

نتایج

ارزیابی شاخص‌های آماری

به منظور شناسایی بهترین الگوریتم در بهینه نمودن و کاهش انحراف مقادیر بارش روزانه پایگاه داده ERA5، اقدام به مقایسه، ارزیابی و تحلیل پارامترهای آماری سه الگوریتم یادگیری ماشینی مورد استفاده در این پژوهش شد.

ضریب همبستگی (CC)

ضریب همبستگی معیاری برای تشخیص درجه وابستگی داده‌های مدل با مقادیر مرجع است و دامنه تغییرات آن از منفی یک تا مثبت یک متغیر است و هرچه به کران‌های بالا و پایین نزدیک‌تر باشد بیانگر شدت وابستگی است و مقدار صفراین شاخص حاکی از عدم وجود ارتباط معنی‌دار بین دو دسته داده است، برای بررسی وضعیت ارتباط بین مقادیر بارش اصلاح شده ERA5 توسط الگوریتم‌های MLP و KNN و D-

¹ Correlation Coefficient

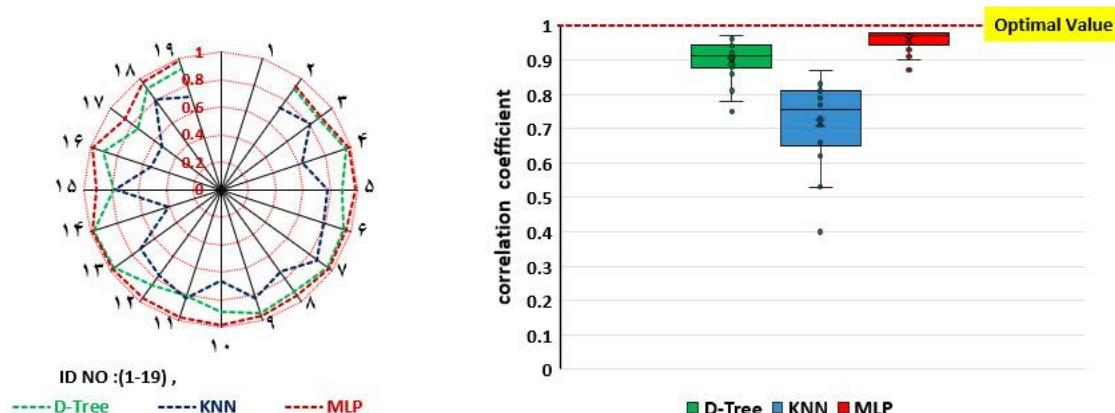
² Root Mean Square Error

Tree با بارش‌های ثبت شده در ایستگاه‌های زمینی، اقدام به محاسبه این ضریب برای هر ایستگاه و به تفکیک هر الگوریتم شد که نتایج آن در جدول (۴) و شکل (۲) ارائه شده است. دامنه تغییرات ضریب همبستگی در MLP بین ۰/۸۷ تا ۰/۹۸ متفاوت بود و این تغییرات برای D-Tree و KNN به ترتیب بین ۰/۷۵ تا ۰/۸۷ و ۰/۹۷ تا ۰/۹۰ متفاوت بودند، در ۷۵ درصد ایستگاه‌ها ضریب همبستگی بین داده‌های اصلاحی و مقادیر مشاهده شده، در الگوریتم‌های MLP و KNN به ترتیب بیشتر از ۶۵ و ۸۸، ۹۵ درصد است. لذا، ضریب همبستگی در کلیه ایستگاه‌های مرجع با خروجی الگوریتم MLP ارتباط قوی‌تری داشته و از این حیث دارای عملکرد مناسب‌تری است.

جدول ۴- ضریب همبستگی مقادیرداده‌های آزمایشی بارش روزانه الگوریتم‌های مورد استفاده با ایستگاه‌های شاهد

Table 4- Correlation coefficient of experimental daily precipitation data values of used algorithms with control stations

ردیف	ایستگاه	ضریب همبستگی			ردیف	ایستگاه	ضریب همبستگی		
		D-Tree	KNN	MLP			D-Tree	KNN	MLP
1	اردک	0.91	0.73	0.93	10	سدکارده	0.81	0.83	0.91
2	آل	0.88	0.81	0.9	11	سرخس	0.86	0.77	0.97
3	باغسنگان	0.96	0.62	0.98	12	فریمان	0.96	0.73	0.98
4	بلغور	0.89	0.78	0.98	13	جوین	0.97	0.40	0.98
5	تربیت‌حیدریه	0.94	0.79	0.96	14	نیشاپور	0.78	0.77	0.91
6	تلغور	0.96	0.87	0.97	15	کاشمر	0.9	0.53	0.98
7	چنانان	0.91	0.74	0.95	16	گوش	0.75	0.53	0.87
8	درگز	0.94	0.83	0.96	17	مغان	0.91	0.81	0.96
9	سدطرق	0.84	0.66	0.96	18	هی‌هی	0.92	0.71	0.96



شکل ۲- مقایسه شاخص آماری ضریب همبستگی بارش تخمینی داده‌های آزمایشی و مقادیر ایستگاهی با استفاده از الگوریتم‌های مورد استفاده

Figure 2- Comparison of the statistical index of the correlation coefficient of the estimated precipitation of experimental data and station values using the used algorithms

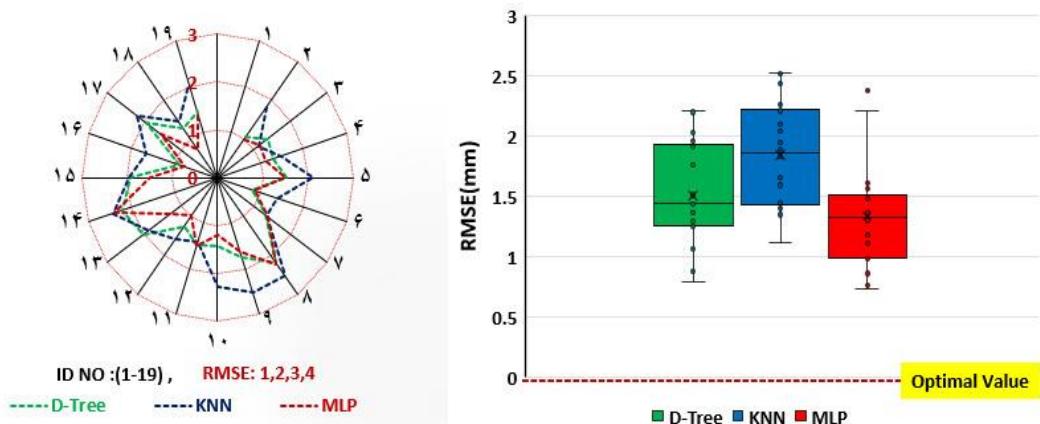
ریشه میانگین مربع خطاهای (RMSE)

ریشه میانگین مربع خطاء، معیاری برای تشخیص عملکرد مدل در تولید داده است، دامنه تغییرات این شاخص بین صفر تا بی‌نهایت بوده و هرچه به صفر نزدیک‌تر باشد، نشان از عملکرد مناسب مدل در تخمین مقادیر بارش است. این شاخص، دارای مقیاس و واحد مشابه با داده‌های مورد تحلیل است و همواره مثبت و بدون جهت است. بهمنظور بررسی وضعیت تغییرات این پارامتر و عملکرد الگوریتم‌های مورد استفاده، اقدام به محاسبه، مقایسه و تحلیل این شاخص برای هر ایستگاه و به تفکیک هر الگوریتم (MLP, KNN, D-Tree) شد، نتایج محاسبه این پارامتر آماری در جدول (۵) و شکل (۳) ارائه شده است. دامنه تغییرات در MLP بین ۰/۷ تا ۰/۹۸ میلی‌متر در روز متغیر بوده و این تغییرات برای D-Tree و KNN به ترتیب بین ۰/۸ تا ۰/۲ و ۰/۲ تا ۰/۵ میلی‌متر در ۷۵ درصد ایستگاه‌ها داده‌های اصلاحی و مقادیر مشاهده شده، در الگوریتم‌های MLP و D-Tree به ترتیب کمتر از ۰/۹ و ۰/۹۵ میلی‌متر در روز است. RMSE در اغلب ایستگاه‌های مرجع، در الگوریتم MLP به مقدار بهینه، نزدیک‌تر از KNN و D-Tree بوده لذا دارای عملکرد مناسب‌تری در مقایسه با دو الگوریتم دیگر است.

جدول ۵- مقادیر RMSE داده‌های آزمایشی بارش روزانه الگوریتم‌های مورد استفاده با ایستگاه‌های شاهد

Table 5- RMSE of daily test data of algorithms used with reference stations

ردیف	ایستگاه	ریشه میانگین مربع خطأ-میلی‌متر			ردیف	ایستگاه	ریشه میانگین مربع خطأ-میلی‌متر		
		D-Tree	KNN	MLP			D-Tree	KNN	MLP
1	ارداک	1.1	1.8	1.0	10	سدکارده	1.5	1.4	1.4
2	آل	1.4	1.1	1.2	11	سرخس	1.2	1.6	1.6
3	باغسنگان	1.3	1.5	1.1	12	فریمان	2.0	1.9	1.9
4	بلغور	1.5	2.1	1.5	13	جوین	2.2	2.4	2.4
5	تربت‌حیدریه	0.8	1.4	0.9	14	نیشابور	1.9	2.0	2.0
6	تلغور	1.3	1.3	1.4	15	کاشمر	0.9	1.7	1.7
7	چناران	2.2	2.5	2.2	16	گوش	2.0	2.2	2.2
8	درگز	1.8	2.5	1.6	17	مغان	1.3	1.4	1.4
9	سدطرق	1.4	2.3	1.2	18	هی‌هی	1.4	2.0	2.0



شکل ۳- مقایسه شاخص آماری RMSE بارش تخمینی داده‌های آزمایشی و مقادیر ایستگاهی با استفاده از الگوریتم‌های مورد استفاده

Figure 3- Comparison of RMSE statistical index of estimated precipitation of experimental data and station values using the used algorithms

سوگیری

این شاخص معیاری برای تشخیص سوگیری داده‌های تولید شده توسط مدل است که دارای مقدار و جهت است، دامنه تغییرات این شاخص محدودیتی نداشته و مقدار بهینه آن صفر بوده و دیمانسیون این پارامتر آماری با داده‌های مورد تحلیل یکسان است. مقدار عددی سوگیری میزان انحراف داده‌های اصلاحی از مقادیر حقیقی را ارائه می‌دهد و علامت آن تمایل الگوریتم به بیش برآورد و یا کم برآورد نسبت به مقادیر مشاهداتی را نشان می‌دهد. بهمنظور بررسی وضعیت تغییرات این پارامتر، جهت‌گیری آن و عملکرد الگوریتم‌های مورد استفاده، اقدام به محاسبه، مقایسه و تحلیل این شاخص برای کلیه ایستگاه‌های مرجع و به تفکیک الگوریتم‌های مورد استفاده، شد، نتایج محاسبه این پارامتر آماری در جدول (۶) و شکل (۴) ارائه شده است. دامنه تغییرات سوگیری در MLP، [-0.6, 0.18] میلی‌متر در روز بوده و این دامنه تغییرات برای D-Tree و KNN به ترتیب [۰/۱۶, ۰/۰۵] و [۰/۸, ۰/۰۶] میلی‌متر در روز بوده و این دامنه تغییرات برای MLP برابر اصلاح شده در این شاخص میانه ایستگاه‌ها به ترتیب برابر [-۰/۰۹, ۰/۱۱] و [-۰/۱۶, ۰/۰۹] میلی‌متر در روز است. سوگیری مقادیر مشاهده شده، در الگوریتم‌های D-Tree، KNN و MLP برای این شاخص میانه ایستگاه‌ها منفی بوده و نشان از تمایل الگوریتم به کم برآورد داده‌های تخمینی نسبت به مقادیر واقعی دارد. با توجه به اکثر ایستگاه‌ها دارای جهت‌گیری منفی بوده و نشان از تمایل الگوریتم به کم برآورد داده‌های تخمینی نسبت به مقادیر واقعی دارد. با توجه به مقادیر سوگیری درسه الگوریتم و تراکم این مقادیر حول مقدار بهینه در MLP، لذا پرسپترون چندلایه دارای عملکرد مناسب‌تری در مقایسه با دو الگوریتم دیگر در این زمینه است.

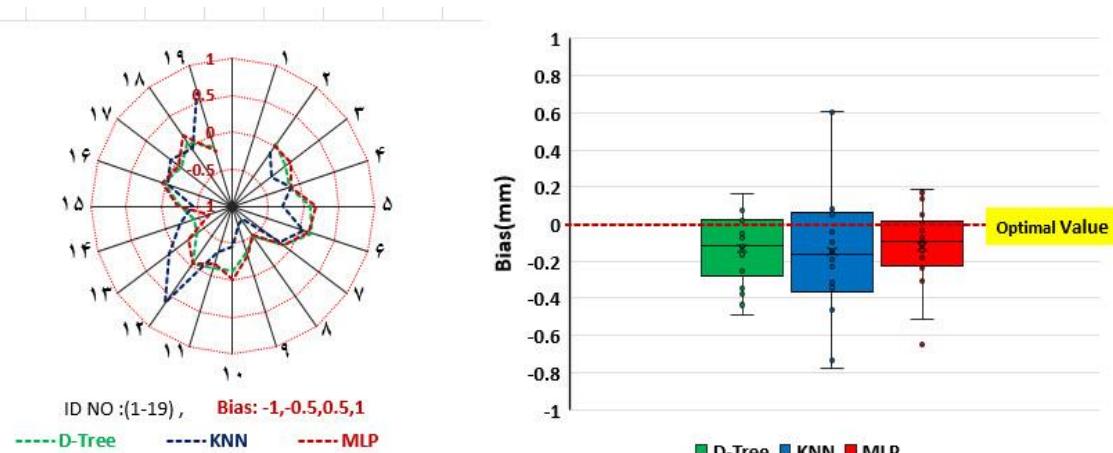
عملکرد الگوریتم منتخب

با توجه به تحلیل‌های صورت گرفته روی داده‌های آزمایشی الگوریتم‌های مورد استفاده و مقادیر بارش روزانه ایستگاه‌های زمینی، مدل MLP دارای عملکرد بهتری نسبت به دو الگوریتم دیگر است، بنابراین به عنوان الگوی منتخب در واسنجی داده‌های بارش روزانه ERA5 در محدوده استان خراسان‌رضوی مناسب تشخیص داده شد. اما مقادیر بارش روزانه پایگاه داده ERA5 نسبت به مقادیر ایستگاه‌های زمینی دارای ضریب همبستگی، ریشه میانگین مربع خطأ و سوگیری است، لذا باید توانایی الگوریتم منتخب در بهبود وضعیت داده‌ها و ارتقاء شاخص‌های آماری مورد بررسی قرار گیرد. برای این منظور نسبت به ارزیابی، مقایسه و تحلیل شاخص‌های CC، Bias و RMSE مقادیر بارش روزانه پایگاه داده و مقادیر واسنجی شده بارش توسط الگوریتم MLP برای ایستگاه‌های بارانسنجی اقدام شد.

جدول ۶- مقادیر سوگیری داده‌های آزمایشی بارش روزانه الگوریتم‌های مورد استفاده با ایستگاه‌های شاهد

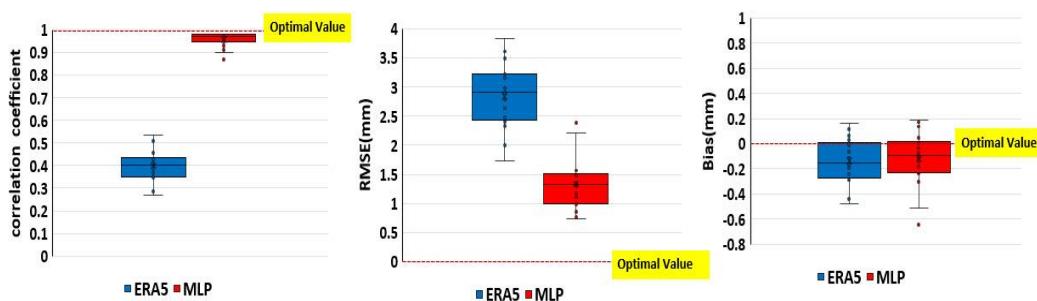
Table 6- Bias values of experimental daily precipitation data of algorithms used with reference stations

ردیف	ایستگاه	سوگیری- میلی‌متر در روز			ردیف	ایستگاه	سوگیری- میلی‌متر در روز		
		D-Tree	KNN	MLP			D-Tree	KNN	MLP
1	اردک	-0.037	-0.09	0.008	10	سدکارده	-0.123	-0.343	-0.18
2	آل	-0.078	-0.31	0.003	11	سرخس	-0.051	0.061	-0.03
3	باغستان	-0.147	-0.14	-0.14	12	فریمان	-0.378	0.099	-0.23
4	بلقور	-0.091	-0.28	0.171	13	جوین	-0.435	-0.232	-0.64
5	تریت‌حیدریه	0.162	0.054	0.137	14	نیشاپور	-0.254	-0.445	-0.21
6	تلغور	-0.109	-0.18	-0.11	15	کاشمر	-0.021	-0.036	-0.05
7	چناران	-0.488	-0.78	-0.52	16	گوش	-0.104	0.083	-0.07
8	درگز	-0.343	-0.73	-0.31	17	معان	0.075	-0.041	0.18
9	سدطرق	-0.122	-0.46	-0.01	18	هی‌هی	-0.163	0.603	-0.22



شکل ۴- مقایسه شاخص آماری سوگیری بارش تخمینی داده‌های آزمایشی و مقادیر ایستگاهی با استفاده از الگوریتم‌های مورد استفاده
Figure 4- Comparison of the statistical bias index of the estimated precipitation of experimental data and station values using the used algorithms

با توجه به شاخص‌های آماری به دست آمده، دامنه تغییرات ضریب همبستگی مقادیر واسنجی شده بین ۰/۹۸ تا ۰/۸۷ متغیر بوده در صورتی که تغییرات این شاخص برای داده‌های ERA5 بین ۰/۵۳ تا ۰/۲۷ است که نشان‌دهنده بهبود داده‌های واسنجی شده توسط الگوریتم منتخب است. همچنین دامنه تغییرات ریشه میانگین مربع خطأ مربوط به داده‌های واسنجی شده بارش روزانه بین ۰/۷۵ تا ۰/۳۸ میلی‌متر در روز است در صورتی که این دامنه برای مقادیر ERA5 بین ۰/۱۷ تا ۰/۳۸ است که نشان‌دهنده بهبود وضعیت در مقادیر واسنجی شده است. مقدار سوگیری بارش روزانه واسنجی شده توسط الگوریتم منتخب برای میانه ایستگاه‌ها ۰/۰۹- میلی‌متر بوده و همین شاخص، برای مقادیر متناظر ERA5 ۰/۱۶- است که نشان‌دهنده عملکرد مناسب الگوریتم منتخب در واسنجی داده‌ها است.



شکل ۵- مقایسه شاخص‌های آماری بارش روزانه ERA5 واستنجی شده برای ایستگاه‌های شاهد

Figure 5- Comparison of ERA5 daily precipitation statistical indices calibrated for control stations

بحث

مقادیر ارائه شده توسط سنجنده‌ها و پایگاه‌های داده، با کاهش مقیاس زمانی از داده‌های مشاهداتی فاصله بیشتری می‌گیرند، لذا برای استفاده از آن‌ها ضروریست تا این مقادیر با استفاده از روش‌ها و الگوهای پرکاربرد اصلاح شوند. این پژوهش نشان داد که در مقیاس روزانه، الگوریتم‌های یادگیری ماشین می‌توانند کمیت داده‌های بارش ERA5 را بهبود بخشد و شاخص‌های آماری خطرا کاهش دهنده و بارش را تا حد زیادی به مقادیر واقعی نزدیک نمایند. هر چند این الگوریتم‌ها قادر به بهبود شرایط هستند، اما عملکرد آن‌ها یکسان نبوده و نتایج متفاوتی را ارائه می‌نمایند، لذا ایجاد می‌نماید تا بهترین الگوریتم درواستنجی داده‌ها شناسایی شود. ارزیابی و تحلیل‌های آماری صورت گرفته می‌بین این است که الگوریتم MLP نسبت به D-Tree و KNN از عملکرد بهتری برخوردار بوده و می‌تواند شاخص‌های CC، RMSE و Bias را به طور موثری بهبود بخشد. همچنین ارزیابی و تحلیل پارامترهای آماری برای مقادیر بارش روزانه اصلاح شده توسط این الگوریتم در مقایسه با داده‌های اولیه ERA5 نشان داد که MLP در ارتقاء ضریب همبستگی و کاهش RMSE و سوگیری بسیار خوب عمل نموده و توانایی بهبود داده‌های ERA5 را دارد.

با توسعه هوش مصنوعی و جنبه‌های کاربردی آن در حوزه‌ها و علوم مختلف، امروزه شاهد استفاده از یادگیری ماشین در علوم آب و هوافضای هستیم، لذا در سطح دنیا مطالعاتی در خصوص قابلیت استفاده از الگوریتم‌ها و الگوهای مختلف برای اصلاح سوگیری و کاهش خطرا ارائه شده است، هرچند روش‌های قدیمی قادر به اصلاح داده‌ها هستند، اما فرایند انجام و تحلیل، توسط آن‌ها پیچیده و زمان بر است، تکنیک‌های یادگیری ماشین با توجه به سرعت، دقت، سادگی در قابلیت تکرار و ارائه نتایج خیره‌کننده، مورد اقبال قرار گرفته به طوری که استفاده از آن‌ها با سرعت در حال فرآیند شدن است. نتایج این پژوهش، همانند اغلب تحقیقات صورت گرفته استفاده از الگوریتم‌های یادگیری ماشین را برای بهبود عملکرد و کاهش خطرا و استنجی داده‌های بارش توصیه می‌نماید.

برای کاهش سوگیری داده‌های بارش سنجنده Chaudhary & GPM (۲۰۲۰) طی مطالعه‌ای اقدام به واستنجی مقادیر بارش ماهانه این سنجنده نسبت به داده‌های شبکه‌ای ایستگاه‌های زمینی در شبکه قاره هند نمود، اوبا تحلیل داده‌های بارش ماهانه بین سال‌های ۲۰۰۱ تا ۲۰۱۶ میلادی و استفاده از روش درخت تصمیم، سوگیری آن‌ها را تصحیح نمود و نشان داد این الگوریتم قادر است به طور متوسط ضریب همبستگی را به ۷۷/۰ رسانده و شاخص RMSE را بین ۸/۷ تا ۳۷/۳ و MAB را بین ۲۹/۲ تا ۶/۳ بهبود بخشد. این پژوهش نیز نشان داد که الگوریتم درخت تصمیم در افزایش ضریب همبستگی و کاهش شاخص‌های خطرا موثر بوده و می‌تواند در بهبود شرایط نقش بسزایی ایفا نماید (Chaudhary & Dhanyav, 2020)، همچنین مینگ مینگ (۲۰۱۷)، با استفاده از الگوریتم KNN و توسعه آن نشان داد این روش از عملکرد مناسبی در اصلاح و پیش‌بینی بارش در حوضه‌های آبریز کشور چین برخوردار است (Huang, 2017)، نتایج به دست آمده از این تحقیق نیز نشان داد که الگوریتم KNN قادر است شاخص‌های آماری خطرا را کاهش داده و همبستگی را بهبود بخشد. خلیلی و همکاران (۲۰۱۶) از الگوریتم ANN برای پیش‌بینی بارش ماهانه استفاده نمود، او از آمار ماهانه ایستگاه سینوپتیک مشهد طی سال‌های ۱۹۵۳ تا ۲۰۰۳ و یک شبکه پرسپترون سه لایه با ریتم الگوی انتشار برگشتی بهره جست و با تنظیم فرایامترهای مدل برای بارش ماهانه نشان داد که ضریب همبستگی ۹۳/۰ و مقادیر خطای RMSE و MAE به ترتیب برابر ۹۹/۰ و ۰/۰۲ است که با مقایسه با روش سنتی و رایج در سازمان هوافضایی (پیش‌بینی از روی نقشه‌های سینوپتیک) دارای دقت قابل قبولی در پیش‌بینی بارش ماهانه است (Khalili, et al., 2016)، خروجی‌های این مطالعه نیز نشان داد که الگوریتم MLP توانایی اصلاح و بهبود داده‌های اولیه را به طور مطلوبی دارد و می‌توان از آن در اصلاح سوگیری استفاده نمود. در مطالعه‌ای که به منظور پیش‌بینی بارش توسط الگوریتم پرسپترون چند لایه روی ده ناحیه متفاوت در کشور بزرگی انجام شد، بارش و دمای ماهانه یک دوره ۶۰ ساله را به عنوان ورودی مدل جمع‌آوری شد و استفاده از الگوریتم MLP

نشان داد که دقت مدل در فصل‌های مختلف متفاوت بوده و رابطه مستقیمی با ارتفاع و حجم نرمال بارش دارد (Esteves et al., 2018). Meyer و همکاران (۲۰۱۶) اقدام به اصلاح مقادیر بارش روزانه در کشور آلمان نمود و از میان الگوریتم‌های مورد استفاده (RF, ANN, MLP, SVR) جنگل تصادفی به عنوان بهترین الگوریتم تصحیح بارش تشخیص داده شد (Meyer et al., 2016). در این پژوهش به عنوان بهترین الگوریتم شناسایی شد و نتایج آن با RF و SVR مقایسه نشده است و در مطالعات تکمیلی باید به ارزیابی آن‌ها پرداخته شود. TAO و همکاران (۲۰۱۶) نیز کار مشابهی در مرکز آمریکا انجام و الگوهای یادگیری عمیق را مورد بررسی قرار دادند. آن‌ها پیشنهاد نمود که از شبکه عصبی پیچشی برای تصحیح بارش می‌توان به نحو مناسبی استفاده نمود. Baez و همکاران (۲۰۲۰) با پژوهشی در کشور شیلی، الگوریتم جنگل تصادفی (RF) را برای تصحیح سوگیری مناسب تشخیص دادند. Chen و همکاران (۲۰۲۰) با مطالعه‌ای در منطقه دالاس کشور آمریکا، الگوریتم‌های یادگیری عمیق را برای تصحیح بارش مناسب ارزیابی نمودند. Nguyen (۲۰۲۱) با مطالعه‌ای در کشور کره الگوریتم جنگل تصادفی را برای تصحیح بارش ماهواره‌ای پیشنهاد نمودند. این تحقیق نیز نتایج بدست آمده از تحقیقات سایر محققان در خصوص قابلیت الگوریتم‌های یادگیری ماشین در کاهش سوگیری را تایید می‌نماید.

نتیجه‌گیری

الگوریتم‌های یادگیری ماشین، قابلیت تصحیح مقادیر بارش از منابع دیگر از جمله ماهواره‌ها و پایگاه‌های داده را به خوبی دارا بوده و می‌توانند در مطالعات مورد استفاده قرار گیرند. نتایج به دست آمده از این پژوهش مبنی آن است که استفاده از الگوریتم MLP برای تصحیح بارش‌های کوتاه‌مدت ماهواره‌ای در استان خراسان‌رضوی بسیار کارآمد بوده و از نتایج آن می‌توان در مطالعات مرتبط با علوم آب و هوافضای استفاده نمود. تکنیک‌های یادگیری ماشین در تصحیح Bias، تنها نقطه ضعف سنجنده‌ها که عدم تطابق مقادیر آن‌ها با داده‌های مشاهداتی است را تا حد زیادی مرتفع نموده و نقطه عطفی در پردازش داده‌های ماهواره‌ای و استفاده گسترده از آن‌ها گشوده است.

در این پژوهش از داده‌های کوتاه‌مدت پایگاه ERA5 برای شناسایی بهترین الگوریتم از میان MLP، KNN، D-Tree به‌منظور تصحیح بارش روزانه استفاده شد، لذا پیشنهاد می‌شود تا در تحقیقات دیگر:

- بارش‌های ماهانه و سالانه مورد تجزیه و تحلیل قرار گیرد.
- از سایر سنجنده‌ها و پایگاه‌های داده در مقیاس زمانی کوتاه‌مدت و بلندمدت استفاده شود.
- از سایر الگوریتم‌های رگرسیونی مبتنی بر یادگیری ماشین برای تصحیح بارش استفاده به عمل آید.
- دقت تکنیک‌های یادگیری عمیق در تصحیح بارش، مورد آزمون قرار گیرد.
- الگوریتم‌های یادگیری ماشین در ریزمقیاس‌نمایی و تصحیح مقادیر بارش ماهواره‌ای و پایگاه داده بر سامانه GEE مورد بررسی و واسنجی قرار گیرد.

ملاحظات اخلاقی

دسترسی به داده‌ها: دسترسی به داده‌ها و نتایج استفاده شده در این پژوهش از طریق مکاتبه با نویسنده مسئول امکان‌پذیر است.

حمایت مالی: این پژوهش مستخرج از رساله دکتری بوده و تحت حمایت دانشگاه فردوسی مشهد به انجام رسیده است.

مشارکت نویسنده‌گان: این پژوهش توسط مجید رجبی جاغرق بصورت طرح اولیه تهیه شد. سپس توسط آقایان دکتر سید محمد موسوی بایگی، دکتر سید علیرضا عراقی و دکترهادی جباری نویسندگان این پژوهش شدند.

تضاد منافع نویسنده‌گان: نویسنده‌گان این مقاله اعلام می‌دارند که هیچ‌گونه تضاد منافعی در خصوص نگارش این پژوهش ندارند.

سپاس‌گزاری: از شرکت آب منطقه‌ای خراسان‌رضوی به‌دلیل همکاری صمیمانه در ارائه آمار و اطلاعات مورد نیاز تشکر و قدردانی می‌شود.

منابع

۱. رجبی جاغرق، مجید، موسوی بایگی، سید محمد، عراقی، سید علی رضا و جباری، نویسندگان، هادی (۱۴۰۲). ارزیابی دقت مقادیر بارش روزانه TRMM، ERA5، GPM و PERSIANN در استان خراسان‌رضوی. سامانه‌های سطوح آبگیر باران، ۱۱(۲)، ۷۹-۱۰۱.
۲. فتاحی، هادی و جبری‌ایی، فاطمه (۱۴۰۰). ارزیابی پتانسیل روانگرایی خاک در اثر وقوع زمین لرزه بالاستفاده از چند الگوریتم طبقه‌بندی <http://jircsa.ir/article-1-505-fa.html> هوشمند در نرم‌افزار Orange، عمران‌فردوسي، ۳۴(۳)، ۳۹-۵۲.

۳. قبائی سوق، محمد، مساعدی، ابوالفضل، حسام، موسی و هزارجریبی، ابوطالب (۱۳۸۹). ارزیابی تأثیر پیش‌پردازش پارامترهای ورودی به شبکه عصبی مصنوعی (ANNs) با استفاده از روش‌های رگرسیون گام به گام و گاما تست به منظور تخمین سریع‌تر تبخیر و تعرق روزانه، آب و خاک، ۲۴(۳)، ۶۱۰-۶۲۴.

doi.org/10.22067/jsw.v0i0.3631

References

- Abdollahipour, A., Ahmadi, H., & Aminnejad, B. (2022). A review of downscaling methods of satellite-based precipitation estimates. *Earth Science Informatics*, 15(1), 1-20. doi.org/10.1007/s12145-021-00669-4
- Amengual, A., Homar, V., Romero, R., Alonso, S., & Ramis, C. (2012). A statistical adjustment of regional climate model outputs to local scales: application to Platja de Palma, Spain. *Journal of Climate*, 25(3), 939-957. doi.org/10.1175/jcli-d-10-05024.1
- Baez-Villanueva, O. M., Zambrano-Bigiarini, M., Beck, H. E., McNamara, I., Ribbe, L., Nauditt, A., & Thinh, N. X. (2020). RF-MEP: A novel Random Forest method for merging gridded precipitation products and ground-based measurements. *Remote Sensing of Environment*, 239, 111606. doi.org/10.1016/j.rse.2019.111606
- Beck, H. E., Vergopolan, N., Pan, M., Levizzani, V., Van Dijk, A. I., Weedon, G. P., & Wood, E. F. (2017). Global-scale evaluation of 22 precipitation datasets using gauge observations and hydrological modeling. *Hydrology and Earth System Sciences*, 21(12), 6201-6217. doi.org/10.5194/hess-21-6201-2017
- Blöschl, G., Bierkens, M. F., Chambel, A., Cudennec, C., Destouni, G., Fiori, A., & Renner, M. (2019). Twenty-three unsolved problems in hydrology (UPH)—a community perspective. *Hydrological Sciences Journal*, 64(10), 1141-1158. doi.org/10.1080/02626667.2019.1620507
- Chaudhary, S., & Dhanya, C. T. (2020). Decision tree-based reduction of bias in monthly IMERG satellite precipitation dataset over India. *h2oj*, 3(1), 236-255. doi.org/10.2166/h2oj.2020.124
- Chen, C., Hu, B., & Li, Y. (2021). Easy-to-use spatial Random Forest-based downscaling-calibration method for producing high resolution and accurate precipitation data. *Hydrology and Earth System Sciences Discussions*, 2021, 1-50. 25, doi.org/10.5194/hess-25-5667-2021
- Chen, H., Chandrasekar, V., Cifelli, R., & Xie, P. (2019). A machine learning system for precipitation estimation using satellite and ground radar network observations. *IEEE Transactions on Geoscience and Remote Sensing*, 58(2), 982-994. doi.org/10.1109/tgrs.2019.2942280
- Chen, H., Sun, L., Cifelli, R., & Xie, P. (2021). Deep learning for bias correction of satellite retrievals of orographic precipitation. *IEEE Transactions on Geoscience and Remote Sensing*, 60, 1-11. doi.org/10.1109/tgrs.2021.3105438
- Chen, J., Brissette, F. P., Chaumont, D., & Braun, M. (2013). Finding appropriate bias correction methods in downscaling precipitation for hydrologic impact studies over North America. *Water Resources Research*, 49(7), 4187-4205. doi.org/10.1002/wrcr.20331
- Chen, S., Xiong, L., Ma, Q., Kim, J. S., Chen, J., & Xu, C. Y. (2020). Improving daily spatial precipitation estimates by merging gauge observation with multiple satellite-based precipitation products based on the geographically weighted ridge regression method. *Journal of Hydrology*, 589, 125156. doi.org/10.1016/j.jhydrol.2020.125156
- Curtis, S., Crawford, T. W., & Lecce, S. A. (2007). A comparison of TRMM to other basin-scale estimates of rainfall during the 1999 Hurricane Floyd flood. *Natural Hazards*, 43, 187-198. http://dx.doi.org/10.1007/s11069-006-9093-y
- Esteves, J. T., de Souza Rolim, G., & Ferrando, A. S. (2019). Rainfall prediction methodology with binary multilayer perceptron neural networks. *Climate Dynamics*, 52, 2319-2331. doi.org/10.1007/s00382-018-4252-x
- Fattahی, H., & Jiryae, F. (2021). Evaluation of Soil Liquefaction Potential due to Earthquake using Intelligent Classification Algorithm in Orange Software. *Journal of Ferdowsi Civil Engineering*, 34(3), 39-59. https://dorl.net/dor/20.1001.1.27832805.1400.34.3.3.9 [in Persian].
- Fernandez-Palomino, C. A., Hattermann, F. F., Krysanova, V.; Lobanova, A., Vega-Jácome, F., Lavado, W., ... & Bronstert, A. (2022). A novel high-resolution gridded precipitation dataset for Peruvian and Ecuadorian watersheds: Development and hydrological evaluation. *Journal of Hydrometeorology*, 23(3), 309-336. doi.org/10.2139/ssrn.4602668
- Ghabaei Sough, M., Mosaedi, A., Hesam, M. O. U. S. A., & Hezarjaribi, A. (2010). Evaluation effect of input parameters preprocessing in artificial neural networks (Anns) by using stepwise regression and gamma test techniques for fast estimation of daily evapotranspiration. *Water and Soil*, 24(3), 610-624. doi.org/10.22067/jsw.v0i0.3631 [in Persian].
- Gómez-Chova, L., Tuia, D., Moser, G., & Camps-Valls, G. (2015). Multimodal classification of remote sensing images: A review and future directions. *Proceedings of the IEEE*, 103(9), 1560-1584. doi.org/10.1109/jproc.2015.2449668
- Guo RuiFang, G. R., Liu YuanBo, L. Y., Zhou Han, Z. H., & Zhu YaQiao, Z. Y. (2018). Precipitation downscaling using a probability-matching approach and geostationary infrared data: an evaluation over six climate regions. *Hydrology and Earth System Sciences*, 22(7), 3685-3699. doi.org/10.5194/hess-22-3685-2018

19. Hamill, T. M., & Schaeferer, M. (2018). Probabilistic precipitation forecast postprocessing using quantile mapping and rank-weighted best-member dressing. *Monthly Weather Review*, 146(12), 4079-4098. doi.org/10.1175/mwr-d-18-0147.1
20. He, X., Chaney, N. W., Schleiss, M., & Sheffield, J. (2016). Spatial downscaling of precipitation using adaptable random forests. *Water Resources Research*, 52(10), 8217-8237. doi.org/10.1002/2016wr019034
21. Heredia, M. B., Junquas, C., Prieur, C., & Condom, T. (2018). New statistical methods for precipitation bias correction applied to WRF model simulations in the Antisana region, Ecuador. *Journal of Hydrometeorology*, 19(12), 2021-2040. doi.org/10.1175/jhm-d-18-0032.1
22. Hu, Q., Li, Z., Wang, L., Huang, Y., Wang, Y., & Li, L. (2019). Rainfall spatial estimations: A review from spatial interpolation to multi-source data merging. *Water*, 11(3), 579. doi.org/10.3390/w11030579
23. Huang, M., Lin, R., Huang, S., & Xing, T. (2017). A novel approach for precipitation forecast via improved K-nearest neighbor algorithm. *Advanced Engineering Informatics*, 33, 89-95. doi.org/10.1016/j.aei.2017.05.003
24. Jakob Themeßl, M., Gobiet, A., & Leuprecht, A. (2011). Empirical-statistical downscaling and error correction of daily precipitation from regional climate models. *International Journal of Climatology*, 31(10), 1530-1544. doi.org/10.1002/joc.2168
25. Khalili, N., Khodashenas, S. R., Davary, K., Baygi, M. M., & Karimaldini, F. (2016). Prediction of rainfall using artificial neural networks for synoptic station of Mashhad: a case study. *Arabian Journal of Geosciences*, 9, 1-9. doi.org/10.1007/s12517-016-2633-1
26. Li, C. Y. Climate Dynamics, 2nd ed.; Chapter 1; Meteorological Press: Beijing, China, 2000; pp. 2-3.
27. Li, H., Haugen, J. E., & Xu, C. Y. (2018). Precipitation pattern in the Western Himalayas revealed by four datasets. *Hydrology and Earth System Sciences*, 22(10), 5097-5110. doi.org/10.5194/hess-22-5097-2018.
28. Lin, Q., Peng, T., Wu, Z., Guo, J., Chang, W., & Xu, Z. (2022). Performance evaluation, error decomposition and Tree-based Machine Learning error correction of GPM IMERG and TRMM 3B42 products in the Three Gorges Reservoir Area. *Atmospheric Research*, 268, 105988. doi.org/10.1016/j.atmosres.2021.105988
29. Liu, X., Yang, T., Hsu, K., Liu, C., & Sorooshian, S. (2017). Evaluating the streamflow simulation capability of PERSIANN-CDR daily rainfall products in two river basins on the Tibetan Plateau. *Hydrology and Earth System Sciences*, 21(1), 169-181. doi.org/10.5194/hess-21-169-2017
30. Ma, Y., Zhang, Y., Yang, D., & Farhan, S. B. (2015). Precipitation bias variability versus various gauges under different climatic conditions over the Third Pole Environment (TPE) region. *International Journal of Climatology*, 35(7), 1201-1211. doi.org/10.1002/joc.4045
31. Mega, T., Ushio, T., Takahiro, M., Kubota, T., Kachi, M., & Oki, R. (2018). Gauge-adjusted global satellite mapping of precipitation. *IEEE Transactions on Geoscience and Remote Sensing*, 57(4), 1928-1935. doi.org/10.1109/tgrs.2018.2870199
32. Militino, A. F., Ugarte, M. D., & Pérez-Goya, U. (2022). Machine learning procedures for daily interpolation of rainfall in Navarre (Spain). *Trends in Mathematical, Information and Data Sciences: A Tribute to Leandro Pardo*, 399-413. doi.org/10.1007/978-3-031-04137-2_34
33. Muhlbauer, A., McCoy, I. L., & Wood, R. (2014). Climatology of stratocumulus cloud morphologies: microphysical properties and radiative effects. *Atmospheric Chemistry and Physics*, 14(13), 6695-6716. doi.org/10.5194/acp-14-6695-2014
34. Nguyen, G. V., Le, X. H., Van, L. N., Jung, S., Yeon, M., & Lee, G. (2021). Application of random forest algorithm for merging multiple satellite precipitation products across South Korea. *Remote Sensing*, 13(20), 4033. doi.org/10.3390/rs13204033
35. Pan, B., Hsu, K., AghaKouchak, A., & Sorooshian, S. (2019). Improving precipitation estimation using convolutional neural network. *Water Resources Research*, 55(3), 2301-2321. doi.org/10.1029/2018wr024090
36. Peña-Arancibia, J. L., Van Dijk, A. I., Renzullo, L. J., & Mulligan, M. (2013). Evaluation of precipitation estimation accuracy in reanalyses, satellite products, and an ensemble method for regions in Australia and South and East Asia. *Journal of Hydrometeorology*, 14(4), 1323-1333. doi.org/10.1175/jhm-d-12-0132.1
37. Piani, C., Weedon, G. P., Best, M., Gomes, S. M., Viterbo, P., Hagemann, S., & Haerter, J. O. (2010). Statistical bias correction of global simulated daily precipitation and temperature for the application of hydrological models. *Journal of Hydrology*, 395(3-4), 199-215. doi.org/10.1016/j.jhydrol.2010.10.024
38. Prakash, S., Mitra, A. K., AghaKouchak, A., & Pai, D. S. (2015). Error characterization of TRMM Multisatellite Precipitation Analysis (TMPA-3B42) products over India for different seasons. *Journal of Hydrology*, 529, 1302-1312. doi.org/10.1016/j.jhydrol.2015.08.062
39. Rajabi Jaghargh, M., Mousavi Baygi, S. M., Araghi, S. A., & Jabari Noghabi, H. (2023). Evaluation of accuracy of daily rainfall values TRMM, GPM, ERA5, and PERSIANN in Razavi Khorasan Province. *Journal of Rainwater Catchment Systems*, 11 (2), 79-101. <http://jircsa.ir/article-1-505-fa.html> [in Persian].
40. Rata, M., Douaoui, A., Larid, M., & Douaik, A. (2020). Comparison of geostatistical interpolation methods to map annual rainfall in the Chélieff watershed, Algeria. *Theoretical and Applied Climatology*, 141, 1009-1024. doi.org/10.1007/s00704-020-03218-z
41. Sandri, M., & Zuccolotto, P. (2008). A bias correction algorithm for the Gini variable importance measure in classification trees. *Journal of Computational and Graphical Statistics*, 17(3), 611-628. doi.org/10.1198/106186008x344522

42. Scheel, M. L. M., Rohrer, M., Huggel, C., Santos Villar, D., Silvestre, E., & Huffman, G. J. (2011). Evaluation of TRMM Multi-satellite Precipitation Analysis (TMPA) performance in the Central Andes region and its dependency on spatial and temporal resolution. *Hydrology and Earth System Sciences*, 15(8), 2649-2663. doi.org/10.5194/hess-15-2649-2011
43. Shen, Z., & Yong, B. (2021). Downscaling the GPM-based satellite precipitation retrievals using gradient boosting decision tree approach over Mainland China. *Journal of Hydrology*, 602, 126803. doi.org/10.1016/j.jhydrol.2021.126803
44. Sorooshian, S., Duan, Q., & Gupta, V. K. (1993). Calibration of rainfall-runoff models: Application of global optimization to the Sacramento Soil Moisture Accounting Model. *Water Resources Research*, 29(4), 1185-1194. doi.org/10.1029/92wr02617
45. Sun, Q., Miao, C., Duan, Q., Ashouri, H., Sorooshian, S., & Hsu, K. L. (2018). A review of global precipitation data sets: Data sources, estimation, and intercomparisons. *Reviews of Geophysics*, 56(1), 79-107. doi.org/10.1002/2017rg000574
46. Tang, T., Chen, T., & Gui, G. (2022). A comparative evaluation of gauge-satellite-based merging products over multiregional complex terrain basin. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 15, 5275-5287. doi.org/10.1109/jstars.2022.3187983
47. Tao, Y., Gao, X., Hsu, K., Sorooshian, S., & Ihler, A. (2016). A deep neural network modeling framework to reduce bias in satellite precipitation products. *Journal of Hydrometeorology*, 17(3), 931-945. doi.org/10.1175/jhm-d-15-0075.1
48. Toms, B. A., Barnes, E. A., & Ebert-Uphoff, I. (2020). Physically interpretable neural networks for the geosciences: Applications to earth system variability. *Journal of Advances in Modeling Earth Systems*, 12(9), e2019MS002002. doi.org/10.1029/2019ms002002
49. Villarini, G., & Krajewski, W. F. (2007). Evaluation of the research version TMPA three-hourly 0.25×0.25 rainfall estimates over Oklahoma. *Geophysical Research Letters*, 34(5). doi.org/10.1029/2006gl029147
50. Vrac, M., Noël, T., & Vautard, R. (2016). Bias correction of precipitation through Singularity Stochastic Removal: Because occurrences matter. *Journal of Geophysical Research: Atmospheres*, 121(10), 5237-5258. doi.org/10.1002/2015jd024511
51. Yang, W., Andréasson, J., Phil Graham, L., Olsson, J., Rosberg, J., & Wetterhall, F. (2010). Distribution-based scaling to improve usability of regional climate model projections for hydrological climate change impacts studies. *Hydrology Research*, 41(3-4), 211-229. doi.org/10.2166/nh.2010.004
52. Yang, X., Yang, S., Tan, M. L., Pan, H., Zhang, H., Wang, G., ... & Wang, Z. (2022). Correcting the bias of daily satellite precipitation estimates in tropical regions using deep neural network. *Journal of Hydrology*, 608, 127656. doi.org/10.1016/j.jhydrol.2022.127656
53. Yang, Z., Hsu, K., Sorooshian, S., Xu, X., Braithwaite, D., & Verbist, K. M. (2016). Bias adjustment of satellite-based precipitation estimation using gauge observations: A case study in Chile. *Journal of Geophysical Research: Atmospheres*, 121(8), 3790-3806. doi.org/10.1002/2015jd024540
54. Yuan, F., Wang, B., Shi, C., Cui, W., Zhao, C., Liu, Y., ... & Yang, X. (2018). Evaluation of hydrological utility of IMERG Final run V05 and TMPA 3B42V7 satellite precipitation products in the Yellow River source region, China. *Journal of Hydrology*, 567, 696-711. doi.org/10.1016/j.jhydrol.2018.06.045
55. Zandi, O., Zahraie, B., Nasseri, M., & Behrangi, A. (2022). Stacking machine learning models versus a locally weighted linear model to generate high-resolution monthly precipitation over a topographically complex area. *Atmospheric Research*, 272, 106159. doi.org/10.1016/j.atmosres.2022.106159
56. Zhang, J., Fan, H., He, D., & Chen, J. (2019). Integrating precipitation zoning with random forest regression for the spatial downscaling of satellite-based precipitation: A case study of the Lancang-Mekong River Basin. *International Journal of Climatology*, 39(10), 3947-3961. doi.org/10.1002/joc.6050
57. Zhang, L., Li, X., Zheng, D., Zhang, K., Ma, Q., Zhao, Y., & Ge, Y. (2021). Merging multiple satellite-based precipitation products and gauge observations using a novel double machine learning approach. *Journal of Hydrology*, 594, 125969. doi.org/10.1016/j.jhydrol.2021.125969
58. Zhang, Y., Zheng, H., Herron, N., Liu, X., Wang, Z., Chiew, F. H., & Parajka, J. (2019). A framework estimating cumulative impact of damming on downstream water availability. *Journal of Hydrology*, 575, 612-627. doi.org/10.1016/j.jhydrol.2019.05.061
59. Zubíeta, R., Getirana, A., Espinoza, J. C., & Lavado, W. (2015). Impacts of satellite-based precipitation datasets on rainfall-runoff modeling of the Western Amazon basin of Peru and Ecuador. *Journal of Hydrology*, 528, 599-612. doi.org/10.1016/j.jhydrol.2015.06.064