

## Response to Reviewer #2:

The study developed a machine learning model to predict UV radiation and highlighted the model performance. This research topic is very important given the rise in the UV radiation in recent years.

Response: Thank you for the positive comments and constructive suggestions to help improve our manuscript. We have fully responded to the comments below point-to-point and revised the manuscript accordingly. The line numbers referred to in this response document corresponded to those in the revised manuscript with tracked changes.

Overall comments: Technically the manuscript seems strong however the writing can be improved.

Highly suggest the authors to go through the language and make changes wherever necessary throughout the manuscript.

Response: Thanks for your suggestion. We carefully reviewed and revised the language throughout the entire manuscript. Also, we have had the manuscript language-edited by Elsevier. The certificate is provided below.



## Certificate of Elsevier Language Editing Services

The following article was edited by Elsevier Language Editing Services:

A 10 km daily-level ultraviolet radiation predicting dataset  
based on machine learning models in China from  
2005 to 2020

Ordered by:

Yichen Jiang

Estimated Delivery date:

2024-07-30

Order reference:

ASLESTD1069282



Comments:

Line 14: Seems grammatically incorrect. Reword it to "but limited studies have implemented it for UV radiation"

Response: Thank you for the suggestion. The sentence has been revised as "Machine learning algorithms have been widely used to predict environmental factors with high accuracy, but limited studies have implemented it for UV radiation." in lines 13-14 in the revision as suggested.

Line 14-15: The language can be improved. Reword these lines to "The main aim of this study is to develop UV radiation prediction model using the random forest approach and predict the UV radiation at daily and 10km resolution in mainland China from 2005 to 2020".

Response: Thank you for the suggestion. The sentence has been revised as "The main aim of this study is to develop UV radiation prediction model using the random forest approach and predict the UV radiation at daily and 10 km resolution in mainland China from 2005 to 2020." in lines 14-16 in the revision as suggested.

Line 16: It is already mentioned above that random forest model was employed to predict UV radiation. Reword this line.

Response: Thank you for the suggestion. The sentence has been revised as "The model was developed with multiple predictors such as UV radiation data from satellites as independent variables and ground UV radiation measurements from monitoring stations as the dependent variable." in lines 16-17 in the revision as suggested.

Line 21: OMI EDD is used for the first time, write the full form of EDD before introducing the acronym

Response: Thank you for the suggestion. The sentence has been revised as "The model that incorporated erythemal daily dose (EDD) retrieved from the Ozone Monitoring Instrument (OMI) had a higher prediction accuracy than that without it." in lines 21-22 in the revision as suggested.

Line 26: Consider rewording this line as it is not flowing well. May be change it to something like this: "Using machine learning this study generated gridded UV radiation dataset with extensive spatiotemporal coverage which can be utilized for future health-related research".

Response: Thank you for the suggestion. The sentence has been revised as "Using machine learning algorithm, this study generated gridded UV radiation dataset with extensive spatiotemporal coverage, which can be utilized for future health-related research." in lines 27-28 in the revision as suggested.

Line 35 - 36: Please consider rewording these lines.

Response: Thank you for the suggestion. The sentence has been revised as "Further studies are required to ascertain the effects of UV radiation on human health; however, the lack of high-accurate exposure data of UV radiation hinders such health-related investigations." in lines 35-37 in the revision as suggested.

Line 43: remove stands, change it to "despite being"

Response: Thank you for the suggestion. The sentence has been revised as "For example, erythemal UV irradiance from the Total Ozone Mapping Spectrometer (TOMS), despite being one of the initial instruments for evaluating the UV radiation backscattered by the Earth's atmospheric layers, it exhibits lower spatial resolution of 50 km×50 km, and it has limited accuracy." in lines 42-44 in the revision as suggested.

Line 70: Remove "What's more, missingness of satellite-based", change it to "The missing satellite-based"

Response: Thank you for the suggestion. The sentence has been revised as "The missing satellite-based UV radiation were filled to improve the spatial coverage of the final UV radiation predictions." in lines 71-72 in the revision as suggested.

Line 108: Why did the author use O3 concentrations predicted from random forest and not use

directly the monitoring data? Clarify this and explain it in the text clearly.

Response: Thank you for the comments and suggestions. As suggested, we further clarified this issue in the revision in lines 110-118 as “This study used gridded O<sub>3</sub> data instead of O<sub>3</sub> monitoring data from station sites, primarily due to considerations of data coverage in both temporal and spatial dimensions. Regarding the temporal coverage, the air quality monitoring network in China has not established until 2013, which could not fully cover the study period of 2005-2020 in this study. For the spatial coverage, the density of air quality monitoring stations is relatively low, with the majority of them are located in urban areas and eastern China, which could not capture the spatial variability within city and reflect the O<sub>3</sub> pollution level in rural areas and western regions (Geyh et al., 2000). While the gridded O<sub>3</sub> predictions used in this study are available from 2005-2020, have full spatial coverage in mainland China and achieved relatively high accuracy comparing with ground measurements with cross-validation (CV) R<sup>2</sup> and root mean square error of 0.80 and 20.93 ug/m<sup>3</sup>, respectively (Meng et al., 2022).”

#### References:

- Geyh, A. S., Xue, J., Ozkaynak, H., and Spengler, J. D.: The Harvard Southern California Chronic Ozone Exposure Study: Assessing Ozone Exposure of Grade-School-Age Children in Two Southern California Communities, *Environmental Health Perspectives*, 108, 265–270, <https://doi.org/10.1289/ehp.00108265>, 2000.
- Meng, X., Wang, W., Shi, S., Zhu, S., Wang, P., Chen, R., Xiao, Q., Xue, T., Geng, G., Zhang, Q., Kan, H., and Zhang, H.: Evaluating the spatiotemporal ozone characteristics with high-resolution predictions in mainland China, 2013-2019, *Environ Pollut*, 299, 118865, <https://doi.org/10.1016/j.envpol.2022.118865>, 2022.

Line 139: provide reference for 10-fold cross validation if it was used in previous studies and explain the cross-validation process details and the differences between the various (temporal, spatial and year) 10-fold CV?

Response: Thanks for the suggestion. We have added references, refined the details of the CV process, and explained the differences among various CVs in lines 155-177 in the revision as “CV is commonly utilized to assess model performance in regard of overfitting and predicting accuracy, especially in studies of model development for UV radiation (Wu et al., 2022), particulate matter (Chen et al., 2018; Park et al., 2022; Wongnakae et al., 2023), O<sub>3</sub> (Hsu et al., 2019; Wu et al., 2021), and nitrogen dioxide (Lu et al., 2021a). In this study, model performance was tested through overall 10-fold CV, temporal 10-fold CV, spatial 10-fold CV, and by-year

temporal CV, which is a stricter temporal CV. Overall 10-fold CV is the most commonly used form of CV, offering a dependable evaluation of overall model performance and assessing model overfitting (Wu et al., 2022; Wongnakae et al., 2023; Hsu et al., 2019). Temporal 10-fold CV can evaluate the models' capacity of temporal extrapolation for predicting UV radiation levels on days without measurements (He et al., 2023b; Lu et al., 2021b; Bi et al., 2020; Zhu et al., 2022). Spatial 10-fold CV is able to evaluate the models' capacity of spatial extrapolation in locations without monitoring stations (Wang et al., 2018; Zhu et al., 2022; Bi et al., 2020). By-year temporal CV can be used to evaluate the predicting accuracy of our models in years out of the study period of model development (Meng et al., 2021; He et al., 2023a; He et al., 2021).

The overall 10-fold CV was conducted by randomly dividing the dataset into ten parts, with nine parts used as a training dataset to train a random forest model and one part used as a test dataset for predictions. This process was repeated ten times and all measurements were compared with the corresponding predictions. Temporal 10-fold CV was done by randomly dividing the dataset into ten parts based on days, in which data on 90% of the days were used to develop a training model to predict UV radiation on the remaining 10% days each time, and this process was repeated ten times. Similarly, spatial 10-fold CV involved randomly dividing the dataset into ten parts based on the locations of monitoring stations, with data from 90% of the sites were used to develop a training model to predict the UV radiation for the remaining 10% of the sites each time and this process was repeated ten times. In order to further validate the predicting accuracy of our models beyond 2005-2020, this study performed another stricter temporal CV, by-year temporal CV, which left an entire year of data as the testing dataset each time, while data from the remaining years are used as the training dataset. Regression  $R^2$  and root mean square error (RMSE; the square root of the average of the squared differences between the predictions and measurements) between the UV radiation measurements and predictions from model development and CVs were calculated to indicate the model performance.”

References:

- Bi, J., Wildani, A., Chang, H. H., and Liu, Y.: Incorporating Low-Cost Sensor Measurements into High-Resolution PM(2.5) Modeling at a Large Spatial Scale, *Environ Sci Technol*, 54, 2152-2162, <https://doi.org/10.1021/acs.est.9b06046>, 2020.
- Chen, G., Knibbs, L. D., Zhang, W., Li, S., Cao, W., Guo, J., Ren, H., Wang, B., Wang, H., Williams, G., Hamm, N. A. S., and Guo, Y.: Estimating spatiotemporal distribution of PM(1) concentrations in China with satellite remote sensing, meteorology, and land use information, *Environ Pollut*, 233, 1086-1094, <https://doi.org/10.1016/j.envpol.2017.10.011>, 2018.
- He, Q., Ye, T., Zhang, M., and Yuan, Y.: Enhancing the reliability of hindcast modeling for air pollution using history-informed machine learning and satellite remote sensing in China, *Atmospheric Environment*, 312, <https://doi.org/10.1016/j.atmosenv.2023.119994>, 2023a.
- He, Q., Gao, K., Zhang, L., Song, Y., and Zhang, M.: Satellite-derived 1-km estimates and long-term trends of PM<sub>2.5</sub> concentrations in China from 2000 to 2018, *Environment International*, 156, <https://doi.org/10.1016/j.envint.2021.106726>, 2021.
- He, Q., Ye, T., Chen, X., Dong, H., Wang, W., Liang, Y., and Li, Y.: Full-coverage mapping high-resolution atmospheric CO<sub>2</sub> concentrations in China from 2015 to 2020: Spatiotemporal variations and coupled trends with particulate pollution, *Journal of Cleaner Production*, 428, <https://doi.org/10.1016/j.jclepro.2023.139290>, 2023b.
- Hsu, C. Y., Wu, J. Y., Chen, Y. C., Chen, N. T., Chen, M. J., Pan, W. C., Lung, S. C., Guo, Y. L., and Wu, C. D.: Asian Culturally Specific Predictors in a Large-Scale Land Use Regression Model to Predict Spatial-Temporal Variability of Ozone Concentration, *Int J Environ Res Public Health*, 16, <https://doi.org/10.3390/ijerph16071300>, 2019.
- Lu, T., Marshall, J. D., Zhang, W., Hystad, P., Kim, S. Y., Bechle, M. J., Demuzere, M., and Hankey, S.: National Empirical Models of Air Pollution Using Microscale Measures of the Urban Environment, *Environ Sci Technol*, 55, 15519-15530, <https://doi.org/10.1021/acs.est.1c04047>, 2021a.
- Lu, Y., Giuliano, G., and Habre, R.: Estimating hourly PM(2.5) concentrations at the neighborhood scale using a low-cost air sensor network: A Los Angeles case study, *Environ Res*, 195, 110653, <https://doi.org/10.1016/j.envres.2020.110653>, 2021b.
- Meng, X., Liu, C., Zhang, L., Wang, W., Stowell, J., Kan, H., and Liu, Y.: Estimating PM(2.5) concentrations in Northeastern China with full spatiotemporal coverage, 2005-2016, *Remote Sens Environ*, 253, <https://doi.org/10.1016/j.rse.2020.112203>, 2021.
- Park, S., Im, J., Kim, J., and Kim, S. M.: Geostationary satellite-derived ground-level particulate matter concentrations using real-time machine learning in Northeast Asia, *Environ Pollut*, 306, 119425, <https://doi.org/10.1016/j.envpol.2022.119425>, 2022.
- Wang, Y., Hu, X., Chang, H. H., Waller, L. A., Belle, J. H., and Liu, Y.: A Bayesian Downscaler Model to Estimate Daily PM<sub>2.5</sub> Levels in the Conterminous US, *International Journal of Environmental Research and Public Health*, 15, <https://doi.org/10.3390/ijerph15091999>, 2018.
- Wongnakae, P., Chitchum, P., Sripramong, R., and Phosri, A.: Application of satellite remote sensing data and random forest approach to estimate ground-level PM(2.5) concentration in Northern region of Thailand, *Environ Sci Pollut Res Int*, 30, 88905-88917, <https://doi.org/10.1007/s11356-023-28698-0>, 2023.
- Wu, J., Qin, W., Wang, L., Hu, B., Song, Y., and Zhang, M.: Mapping clear-sky surface solar ultraviolet radiation in China at 1 km spatial resolution using Machine Learning technique and Google Earth Engine, *Atmospheric Environment*, 286, <https://doi.org/10.1016/j.atmosenv.2022.119219>, 2022.
- Wu, J., Wang, Y., Liang, J., and Yao, F.: Exploring common factors influencing PM(2.5) and O(3)

concentrations in the Pearl River Delta: Tradeoffs and synergies, *Environ Pollut*, 285, 117138, <https://doi.org/10.1016/j.envpol.2021.117138>, 2021.

Zhu, Q., Bi, J., Liu, X., Li, S., Wang, W., Zhao, Y., and Liu, Y.: Satellite-Based Long-Term Spatiotemporal Patterns of Surface Ozone Concentrations in China: 2005-2019, *Environ Health Perspect*, 130, 27004, <https://doi.org/10.1289/EHP9406>, 2022.

Line 123: Why did the authors use random forest compared to the other machine learning algorithm? Include the necessary information that supports the argument.

**Response:** Thanks for the constructive comment.

Overall, random forest method is a widely used machine learning algorithm for predicting multiple environmental factors (Araki et al., 2018; Guo et al., 2021; Huang et al., 2018; Liu et al., 2020), with several advantages comparing with other machine learning methods. First, random forest exhibits high flexibility in processing various types of data and strong tolerance to multicollinearity among predictors (Breiman, 2001; Fox et al., 2017; Strobl et al., 2008; Bamrah et al., 2020). Second, comparing to some other black-box machine learning models, random forest method is able to provide feature importance rankings and facilitate a deeper understanding in contribution of all predictors in predictions, which makes the models easier to be understood and explained (Hu et al., 2017; Wei et al., 2019). Third, the predicting errors in random forest models are generally lower, due to the reduction in variance achieved by aggregating multiple trees (Ameer et al., 2019; Ding and Qie, 2022). Forth, random forest is user-friendly with relatively small number of parameter settings and a relatively fast processing speed (Ameer et al., 2019; Hu et al., 2017). Due to the above advantages, many previous studies found that random forest method could achieve the higher or at least comparable predicting accuracy over other machine learning models. A study in Taiwan, China, predicting the air quality index showed that compared to methods of adaptive boosting, artificial neural networks, stacking ensemble, and support vector machines, the random forest model performed better (Liang et al., 2020). In a study predicting CO, NO, PM<sub>2.5</sub>, and NO<sub>2</sub> in Spain, random forest outperformed other machine learning models (decision tree for regression, support vector machines, and neural networks) in predicting almost all pollutants (Ochando et al., 2015). A study compared multiple models of decision tree, random forest, gradient boosting, and artificial neural network multi-layer perceptron in predicting PM<sub>2.5</sub> in multiple cities in China

and found that random forest model performed the best (Ameer et al., 2019). In another study conducted in Valencia, Spain, comparing a decision tree for regression and random forest in predicting NO, NO<sub>2</sub>, SO<sub>2</sub>, and O<sub>3</sub>, the random forest model produced better results (Contreras and Ferri, 2016). A study in India predicting the Air Quality Index compared decision tree, support vector regression, and random forest, with random forest having the highest accuracy (Bamrah et al., 2020).

In the revision, we also added an extra analysis by developing an eXtreme Gradient Boosting (XGBoost) model, another machine learning model based on decision trees with relatively high predicting accuracy (Zamani Joharestani et al., 2019; Nasabpour Molaei et al., 2023; Dai et al., 2023; Wu et al., 2022). Based on data in this study, the XGBoost model yielded an R<sup>2</sup> (RMSE) of 0.81 (39.25 W m<sup>-2</sup>) in predicting UV radiation levels with the same predictors in the random forest model, while the random forest model achieved better performance with slightly higher R<sup>2</sup> (0.83) and lower RMSE (37.44 W m<sup>-2</sup>). Therefore, this study employs random forest to construct the models.

The explanations were also summarized in the lines 289-306 in the Discussion section of the revision as “In this study we employed random forest method to develop the models as it is a widely used machine learning algorithm with several advantages for predicting multiple environmental factors (Araki et al., 2018; Guo et al., 2021; Huang et al., 2018; Liu et al., 2020). First, random forest exhibits high flexibility in processing various types of data and strong tolerance to multicollinearity among predictors (Breiman, 2001; Fox et al., 2017; Strobl et al., 2008; Bamrah et al., 2020). Second, comparing to some other black-box machine learning models, random forest method is able to provide feature importance rankings and facilitate a deeper understanding in contribution of all predictors in predictions, which makes the models easier to be understood and explained (Hu et al., 2017; Wei et al., 2019). Third, the predicting errors in random forest models are generally lower, due to the reduction in variance achieved by aggregating multiple trees (Ameer et al., 2019; Ding and Qie, 2022). Forth, random forest is user-friendly with relatively small number of parameter settings and a relatively fast processing speed (Ameer et al., 2019; Hu et al., 2017). Due to the above advantages, many previous studies



found that random forest method could achieve higher or at least comparable predicting accuracy over other machine learning models in predicting environmental factors (Liang et al., 2020; Julián et al., 2015; Contreras and Ferri, 2016; Ameer et al., 2019). In this study, we also compared results from random forest model and eXtreme Gradient Boosting (XGBoost) model, which is another machine learning model based on decision trees with relatively high predicting accuracy (Zamani Joharestani et al., 2019; Nasabpour Molaei et al., 2023; Dai et al., 2023; Wu et al., 2022). The results indicated that the predicting accuracy from XGBoost method was comparable but slightly lower than those from random forest method with lower  $R^2$  (XGBoost: 0.81 v.s. random forest: 0.83) and higher RMSE (XGBoost: 39.25  $W m^{-2}$  v.s. random forest: 37.44  $W m^{-2}$ )”.

#### References:

- Ameer, S., Shah, M. A., Khan, A., Song, H., Maple, C., Islam, S. U., and Asghar, M. N.: Comparative Analysis of Machine Learning Techniques for Predicting Air Quality in Smart Cities, *IEEE Access*, 7, 128325-128338, <https://doi.org/10.1109/access.2019.2925082>, 2019.
- Araki, S., Shima, M., and Yamamoto, K.: Spatiotemporal land use random forest model for estimating metropolitan NO<sub>2</sub> exposure in Japan, *Science of The Total Environment*, 634, 1269-1277, 2018.
- Bamrah, S. K., Saiharshith, K., and Gayathri, K.: Application of random forests for air quality estimation in india by adopting terrain features, 2020 4th International Conference on Computer, Communication and Signal Processing (ICCCSP), 1-6.
- Breiman, L.: Random Forests, *Machine Learning*, 45, 5-32, <https://doi.org/10.1023/A:1010933404324>, 2001.
- Contreras, L. and Ferri, C.: Wind-sensitive interpolation of urban air pollution forecasts, *Procedia Computer Science*, 80, 313-323, 2016.
- Dai, H., Huang, G., Wang, J., and Zeng, H.: VAR-tree model based spatio-temporal characterization and prediction of O<sub>3</sub> concentration in China, *Ecotoxicol Environ Saf*, 257, 114960, <https://doi.org/10.1016/j.ecoenv.2023.114960>, 2023.
- Ding, W. and Qie, X.: Prediction of Air Pollutant Concentrations via RANDOM Forest Regressor Coupled with Uncertainty Analysis — A Case Study in Ningxia, *Atmosphere*, 13, <https://doi.org/10.3390/atmos13060960>, 2022.
- Fox, E. W., Hill, R. A., Leibowitz, S. G., Olsen, A. R., Thornbrugh, D. J., and Weber, M. H.: Assessing the accuracy and stability of variable selection methods for random forest modeling in ecology, *Environ Monit Assess*, 189, 316, <https://doi.org/10.1007/s10661-017-6025-0>, 2017.
- Guo, B., Zhang, D., Pei, L., Su, Y., Wang, X., Bian, Y., Zhang, D., Yao, W., Zhou, Z., and Guo, L.: Estimating PM<sub>2.5</sub> concentrations via random forest method using satellite, auxiliary, and ground-level station dataset at multiple temporal scales across China in 2017, *Sci Total Environ*, 778, 146288, <https://doi.org/10.1016/j.scitotenv.2021.146288>, 2021.
- Hu, X., Belle, J. H., Meng, X., Wildani, A., Waller, L. A., Strickland, M. J., and Liu, Y.: Estimating PM<sub>2.5</sub>

- Concentrations in the Conterminous United States Using the Random Forest Approach, *Environmental Science & Technology*, 51, 6936-6944, <https://doi.org/10.1021/acs.est.7b01210>, 2017.
- Huang, K., Xiao, Q., Meng, X., Geng, G., Wang, Y., Lyapustin, A., Gu, D., and Liu, Y.: Predicting monthly high-resolution PM<sub>2.5</sub> concentrations with random forest model in the North China Plain, *Environmental Pollution*, 242, 675-683, <https://doi.org/10.1016/j.envpol.2018.07.016>, 2018.
- Julián, C. I. F., ES, U., and Ferri, C.: Airvlc: An application for real-time forecasting urban air pollution, *Proceedings of the 2nd international workshop on mining urban data*, Lille, France.
- Liang, Y.-C., Maimury, Y., Chen, A. H.-L., and Juarez, J. R. C.: Machine Learning-Based Prediction of Air Quality, *Applied Sciences*, 10, <https://doi.org/10.3390/app10249151>, 2020.
- Liu, H., Liu, J., Liu, Y., Ouyang, B., Xiang, S., Yi, K., and Tao, S.: Analysis of wintertime O(3) variability using a random forest model and high-frequency observations in Zhangjiakou-an area with background pollution level of the North China Plain, *Environ Pollut*, 262, 114191, <https://doi.org/10.1016/j.envpol.2020.114191>, 2020.
- Nasabpour Molaei, S., Salajegheh, A., Khosravi, H., Nasiri, A., and Ranjbar Saadat Abadi, A.: Prediction of hourly PM<sub>10</sub> concentration through a hybrid deep learning-based method, *Earth Science Informatics*, 17, 37-49, <https://doi.org/10.1007/s12145-023-01146-w>, 2023.
- Ochando, L. C., Julian, C. I. F., Ochando, F. C., and Ferri, C.: Airvlc: An application for real-time forecasting urban air pollution, *Proceedings of the 2nd international workshop on mining urban data*, Lille, France, 2015.
- Strobl, C., Boulesteix, A. L., Kneib, T., Augustin, T., and Zeileis, A.: Conditional variable importance for random forests, *BMC Bioinformatics*, 9, 307, <https://doi.org/10.1186/1471-2105-9-307>, 2008.
- Wei, J., Huang, W., Li, Z., Xue, W., Peng, Y., Sun, L., and Cribb, M.: Estimating 1-km-resolution PM<sub>2.5</sub> concentrations across China using the space-time random forest approach, *Remote Sensing of Environment*, 231, <https://doi.org/10.1016/j.rse.2019.111221>, 2019.
- Wu, J., Qin, W., Wang, L., Hu, B., Song, Y., and Zhang, M.: Mapping clear-sky surface solar ultraviolet radiation in China at 1 km spatial resolution using Machine Learning technique and Google Earth Engine, *Atmospheric Environment*, 286, <https://doi.org/10.1016/j.atmosenv.2022.119219>, 2022.
- Zamani Joharestani, M., Cao, C., Ni, X., Bashir, B., and Talebiesfandarani, S.: PM<sub>2.5</sub> Prediction Based on Random Forest, XGBoost, and Deep Learning Using Multisource Remote Sensing Data, *Atmosphere*, 10, <https://doi.org/10.3390/atmos10070373>, 2019.

Line 218: Fix the typo. It is Figure 5 not 3.

Response: Thanks for pointing out the issue. It has been corrected.

Line 269: reword the line to "there is no atmospheric UV standards"

Response: Thank you for the suggestion. The sentence has been revised as “The threshold for the health effects of UV radiation on the population is still unclear, and there are no atmospheric UV radiation standards so far, which requires support from further epidemiological studies.” in lines 332-334 in the revision as suggested.