# An optimized semi-empirical physical approach for satellite-based PM$_{2.5}$ retrieval: embedding machine learning to simulate complex physical parameters

Caiyi Jin [a], Qiangqiang Yuan [a, c, d, *], Tongwen Li [b, *], Yuan Wang [a], Liangpei Zhang [c, e]

[a] School of Geodesy and Geomatics, Wuhan University, Wuhan 430079, China.

[b] School of Geospatial Engineering and Science, Sun Yat-Sen University, Zhuhai 519082, China

[c] The Collaborative Innovation Center of Geospatial Technology, Wuhan 430079, China.

[d] The Key Laboratory of Geospace Environment and Geodesy, Ministry of Education, Wuhan University, Wuhan 430079, China.

[e] State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University, Wuhan 430079, China.

* Corresponding author.

E-mail address: yqiang86@gmail.com, litw8@mail.sysu.edu.cn

## ABSTRACT

Satellite remote sensing of PM$_{2.5}$ mass concentration has become one of the most popular atmospheric research aspects, resulting in the development of different models. Among them, the semi-empirical physical approach constructs the transformation relationship between the aerosol optical depth (AOD) and PM$_{2.5}$ based on the optical properties of particles, which has strong physical significance. Also, it performs the PM$_{2.5}$ retrieval independently of the ground stations. However, due to the complex physical relationship, the physical parameters in the semi-empirical approach are difficult to calculate accurately, resulting in relatively limited accuracy. To achieve the optimization effect, this study proposes a method of embedding machine learning into a semi-physical empirical model (RF-PMRS). Specifically, based on the theory of the physical PM$_{2.5}$ remote sensing approach (PMRS), the complex parameter (VE$_f$, a columnar volume-to-extinction ratio of fine particles) is simulated by the random forest model (RF). Also, a fine mode fraction product with higher quality is applied to make up for the insufficient coverage of satellite products. Experiments in North China show that the surface PM$_{2.5}$ concentration derived by RF-PMRS has an average annual value of 57.92 μg/m³ versus the ground value of 60.23 μg/m³. Compared with the original

method, RMSE decreases by 39.95 μg/m³, and the relative deviation reduces by 44.87%. Moreover, validation at two AERONET sites presents a time series change closer to the true values, with an R of about 0.80. This study is also a preliminary attempt to combine model-driven and data-driven models, laying a foundation for further atmospheric research on optimization methods.

**Keywords:** $PM_{2.5}$; Physical approach; Machine learning; Volume-to-extinction ratio; Fine mode fraction

## 1. Introduction

Epidemiological studies have indicated that $PM_{2.5}$ (fine particulate matter with an aerodynamic equivalent diameter no greater than 2.5 μm) can adversely affect human health, such as increasing the risk of diabetes and respiratory diseases (Bowe et al., 2018; Pope III et al., 2002; Xu et al., 2013), and accurate surface $PM_{2.5}$ concentration is the basis of air pollution-health related research. Satellite remote sensing has the advantages of high resolution and global coverage (Ma et al., 2014; Wu et al., 2020; He et al., 2022), including variables strongly associated with $PM_{2.5}$ such as aerosol optical depth (AOD). Therefore, it has become a mainstream method for fine particle estimation (Zhang et al., 2021).

There are mainly three satellite-based ways of retrieving $PM_{2.5}$. 1) Chemical transport models-based method. It calculates a scaling factor η between AOD and $PM_{2.5}$ simulated by atmospheric chemical transport models (CTM) (Lyu et al., 2022) and then transfers the proportional relationship to satellite AOD data when calculating surface $PM_{2.5}$ concentration (Geng et al., 2015; Van Donkelaar et al., 2006). However, the assumption of a constant factor between simulated and observed values has large spatiotemporal limitations. 2) Univariate/Multivariate regression. This kind of method establishes a statistical model between AOD, auxiliary variables, and ground $PM_{2.5}$ observations. Machine learning is a common tool for such data-driven methods due to its powerful nonlinear fitting ability between multiple variables (Irrgang et al., 2021). But the regression is affected by the distribution and density of ground stations (Gupta and Christopher, 2009; Li et al., 2017). 3) Semi-empirical physical approach. Taking

66  the physical theory as the basis, surface $PM_{2.5}$ is derived through an empirical formula

67  constructed from AOD and some PM-related key parameters, including an important

68  empirical parameter related to the optical properties (S). The process steps are explicit

69  and independent of ground station observations. Meanwhile, this approach has stronger

70  physical interpretability than the previous two methods with a large space for

71  optimization.

72       Due to the complexity of the physical parameters, many studies have optimized the

73  semi-empirical physical approach. Based on 355nm-band radar observations, Raut and

74  Chazette (2009) introduced a specific extinction cross-section to simplify the

75  expression of S and $PM_{2.5}$ concentration was estimated. Kokhanovsky et al. (2009)

76  constructed a particle-effective radius model, which can obtain the particle

77  concentrations throughout the atmospheric column. Furthermore, Zhang and Li (2015)

78  proposed the physical $PM_{2.5}$ remote sensing method (PMRS). It replaced S by defining

79  a volume-to-extinction ratio of fine particles ($VE_f$) and used a quadratic polynomial of

80  fine mode fraction (FMF) to simulate $VE_f$, showing certain advantages (Li et al., 2016;

81  Zhang et al., 2020).

82       However, the above semi-physical empirical models have some shortcomings. Firstly,

83  the satellite data used in the models are blocked by clouds and fog in some areas, thus

84  high-coverage and high-precision products need to be excavated and applied; secondly,

85  there are still large uncertainties in estimating physical parameters (such as a simple

86  polynomial fit to S in the PMRS method) and their expressions need to be improved.

87  To date, machine learning (ML) has developed rapidly (He et al., 2021). It can detect

88  complex nonlinear relationships of multiple data and model their interaction (Yuan et

89  al., 2020; Lee et al.,2022). This provides an idea for improving the accuracy of physical

90  parameter acquisition, so as to estimate high-precision $PM_{2.5}$ through semi-physical

91  empirical models.

92       According to this idea, our study proposes an optimized semi-empirical physical

93  model (RF-PMRS) based on the PMRS theory, which attempts to explore the possibility

94  of combining physical models and ML. To be specific, we creatively embed ML (the

95  random forest model) into the PMRS method to simulate the physical parameter (i.e.,

96    VE$_f$) derived from FMF and related variables, thus optimizing the previous polynomial

97    expression. Besides, to further improve the PM$_{2.5}$ retrieval accuracy, the physical-deep

98    learning FMF (Phy-DL FMF) dataset generated by a hybrid retrieval algorithm of ML

99    and physical mechanisms is introduced. Ultimately, we comprehensively validate the

100   performance of the PM$_{2.5}$ obtained by our optimized approach.

101   The remained part of our article is as follows. Section 2 describes the experimental

102   datasets. Section 3 illustrates the specific derivation process of the proposed method.

103   Section 4 analyzes the evaluation results. Some supporting experiments are discussed

104   in section 5. And the final part provides the conclusion.

105

106   **2. Data**

107   **2.1. AERONET data**

108   The Aerosol Robotic Network (AERONET) is a federation of ground-based sun-sky

109   radiometer networks, providing worldwide remote sensing aerosol data for more than

110   25 years (Holben et al., 1998). Until now, the Version 3 dataset has been released (Giles

111   et al., 2017). Due to its high quality, the data from AERONET have been regarded as

112   theoretical true values to evaluate satellite-based products in related studies (Chen et

113   al., 2020; Gao et al., 2016; Wang et al., 2019). AOD, FMF, and Volume Size

114   Distribution products with Level 2.0 (quality-assured) are applied to calculate the true

115   values of the physical parameters, and then to implement our modeling purpose (not

116   involved in PM$_{2.5}$ calculations). A total of 9 AERONET sites corresponding to four

117   typical aerosol types participate in the training. Table 1 shows the specific information.

118

119   **Table 1**. Data information of 9 AERONET sites classified by aerosol types. Location indicates the
120   latitude and longitude, where '-' means the south latitude and west longitude. Two sites in bold fonts
121   participate in the PM$_{2.5}$ validation experiment.

| Aerosol Type | Site | Location (LAT, LON) | Training period | Isolated-validation period |
|---|---|---|---|---|
| **Urban–industrial** | **Beijing** | **39.98°, 116.38°** | 2001-2017 | 2018-2019 |
| | **Beijing-CAMS** | **39.93°, 116.32°** | 2012-2017 | 2018-2019 |
| | XiangHe | 39.75°, 116.96° | 2004-2017 | / |
| | Ascension Island | -7.98°, -14.41° | 2010-2017 | 2018-2019 |

| | Capo Verde | 16.73°, -22.94° | 2010-2017 | 2018 |
|---|---|---|---|---|
| **Biomass burning** | CUIABA MIRANDA | -15.73°, -56.07° | 2010-2017 | 2018-2019 |
| **Desert dust** | GSFC | 38.99°, -76.84° | 2010-2017 | 2018-2019 |
| | Mexico City | 19.33°, -99.18° | 2010-2017 | / |
| **Oceanic** | Solar Village | 24.91°, 46.40° | 2010-2013 | / |

## 2.2. MODIS AOD

MCD19A2, the Moderate-resolution Imaging Spectroradiometer (MODIS) C6 Level-2 gridded (L2G) land AOD product (Lyapustin and Wang, 2015), is selected in this study. It is derived from the Multi-Angle Implementation of the Atmospheric Correction (MAIAC) algorithm, which can improve the accuracy in cloud detection and aerosol retrieval (Lyapustin et al., 2011). Besides, this new advanced algorithm jointly combines MODIS Terra and Aqua into a single sensor (Lyapustin et al., 2014). The product is produced daily with a 1km resolution, including aerosol parameters such as 470nm/550nm AOD, quality assurance (QA), and uncertainty factors.

The processing of MCD19A2 data (HDF format) is mainly divided into five steps: AOD/QA band extraction, best quality AOD selection, Terra/Aqua data synthesis, missing information reconstruction, and mosaic. Finally, the daily AOD distribution in GeoTiff format is obtained.

## 2.3. Phy-DL FMF dataset

The original global land FMF products have poor data integrity and low accuracy. To enhance their reliability, Yan et al. (2022) have released a satellite-based dataset called Phy-DL FMF, which integrates physical and deep learning methods. Specifically, it selects the FMF data obtained by a physical method (i.e., Look-Up-Table-based Spectral Deconvolution Algorithm, LUT-SDA) as the optimization target (Yan et al., 2017). Then it combines the Phy-based FMF into a deep-learning model along with multiple auxiliary data such as satellite observations for the final Phy-DL results. Note that the process is trained with AERONET data as the ground truth. The product has a spatial resolution of 1° and covers from 2001 to 2020 (daily scale). In the comparison

147    experiment against the ground FMF, Phy-DL FMF shows a higher accuracy (R = 0.78,

148    RMSE = 0.100) than MODIS FMF (R = 0.37, RMSE = 0.282) (Yan et al., 2022).

149

150    **2.4. Meteorological data**

151    The meteorological data are obtained from the ERA5 dataset, including the values of

152    planetary boundary layer height (PBLH) and relative humidity (RH). As the fifth-

153    generation reanalysis product released by the European Center for Medium-Range

154    Weather Forecasts (ECMWF), ERA5 provides atmospheric data at 0.25° every hour

155    based on the data assimilation principle (Hersbach et al., 2018). It should be noted that

156    RH is not archived directly in ERA5, thus should be calculated by 2m temperature $T$

157    and dew point temperature $T_d$ (referred to ERA-Interim: documentation).

$$RH = 100 \times \frac{e_s(T_d)}{e_s(T)} \tag{1}$$

159    Here, $e_s(t)$ represents the saturation vapor pressure related to a Celsius temperature $t$

160    (Simmons et al., 1999).

$$e_s(t) = 6.112 \times \exp\left(\frac{17.67 \times t}{t + 243.5}\right) \tag{2}$$

162

163    **2.5. Ground PM$_{2.5}$ measurements**

164    The North China Region (NC) is chosen as the main experimental validation area for

165    the final PM$_{2.5}$ calculations. The near-surface hourly PM$_{2.5}$ values are obtained from the

166    China National Environmental Monitoring Center (CNEMC). Nowadays, over 1600

167    ground-based monitors are working continuously and a total of 232 stations (in 2017)

168    participate in this work. Fig. 1 displays the site distributions of the NC region.

169

**Fig. 1.** The location of ground stations in the NC region (35°-45°N, 110°-120°E). The red points represent NC stations.

## 3. Methods

Based on the basic physical properties of atmospheric aerosols, the semi-physical empirical approach starts from the integration of PM mass concentration and AOD. Then it combines several key factors related to $PM_{2.5}$, to derive the in situ $PM_{2.5}$ concentration through multiple remote sensing variables (Koelemeijer et al., 2006). The overall empirical relationship can be represented as:

$$PM_{2.5} = AOD \frac{\rho}{H \cdot f(RH)} S \tag{3}$$

where $\rho$ denotes the particle density and $H$ denotes the atmospheric boundary layer height. $f(RH)$ represents the hygroscopic growth factor related to relative humidity ($RH$). $S$ is an optical characteristic parameter that should be simulated.

### 3.1. PMRS method

### 3.1.1. The expression of VE$_f$

To illustrate S more precisely, PMRS defines the columnar volume-to-extinction ratio of fine particles (i.e., VE$_f$), which can be regarded as the basis of our optimization method. So equation (3) is transformed into:

$$PM_{2.5} = AOD \frac{\rho}{H \cdot f(RH)} VE_f \tag{4}$$

Related to particle size, aerosol extinction, and other properties, VE$_f$ can be expressed as:

$$VE_f = \frac{V_{f,column}}{AOD_f} \tag{5}$$

$$AOD_f = AOD \cdot FMF \tag{6}$$

Here, $AOD_f$ is the fine particle AOD and $FMF$ is the fine mode fraction. $V_{f,column}$ can be expressed by the vertical integral of particle volume size distributions (PVSD) within a certain aerodynamic diameter range:

$$V_{f,column} = \int_0^{D_{p,c}} V(D_p) dD_p \tag{7}$$

$D_{p,c}$ represents the cutting diameter, and the empirical value of 2.0 μm is chosen based on previous literature (Hand and Kreidenweis, 2002; Hänel and Thudium, 1977). And $V(D_p)$ represents the PVSD corresponding to the geometric equivalent diameter ($D_p$).


### 3.1.2. Specific process and limitations

The PMRS method is developed from equation (4). Based on satellite AOD, the near-surface PM$_{2.5}$ can be obtained through multi-step transformation. Fig. 2(a) shows its specific process. Each arrow refers to a step, respectively: size cutting (output: AOD$_f$), volume visualization (output: V$_{f,column}$), bottom isolation (output: V$_f$, fine particle volume near the ground), particle drying (output: V$_{f,dry}$, dry V$_f$) and PM$_{2.5}$ weighting. The overall expression is as follows:

$$PM_{2.5} = AOD \frac{FMF \cdot VE_f \cdot \rho_{f,dry}}{PBLH \cdot f_0(RH)} \tag{8}$$

$$f_0(RH) = \left(1 - \frac{RH}{100}\right)^{-1} \tag{9}$$

where $FMF$ denotes the fine mode fraction, $\rho_{f,dry}$ denotes the dry mass density of $PM_{2.5}$, and $PBLH$ represents the planet boundary layer height. $f_0(RH)$ represents the approximation of $f(RH)$ in equation (4), as expressed in equation (9). Considering the aerosol types in different regions, PMRS fits $VE_f$ to a quadratic polynomial relation of $FMF$ (Zhang and Li, 2015):

$$VE_f = 0.2887FMF^2 - 0.4663FMF + 0.356 \quad (0.1 \le FMF \le 1.0) \tag{10}$$



**Fig. 2.** Surface $PM_{2.5}$ estimation flow of RF-PMRS. a) The five steps of the PMRS method. Gray boxes are the intermediate outputs, blue boxes are the input data, and orange ones denote the variables to be optimized. b) The specific optimization of RF-PMRS: FMF dataset replacement and $VE_f$ simulation by RF model.

PMRS has strong physical significance, the calculation steps are well-defined and site-independent. Zhang and Li (2015) tested the performance of PMRS on 15 stations, and the validation results had an uncertainty of 34%. Compared with the ground value of Jinhua city in China, a 31.3% relative error was generated in Li et al. (2016). Besides, Zhang et al. (2020) applied it to the $PM_{2.5}$ change analysis and prediction experiments in China over 20 years. However, there may be a more complex nonlinear relationship

between $VE_f$ with FMF, not just a simple quadratic formula. Since $VE_f$ is related to the aerosol type, adding other spatiotemporal variables may optimize the fitting process. Additionally, high-quality FMF data is the basic guarantee for the estimated $PM_{2.5}$ quality. In a word, to further improve the physical method, a better nonlinear model between $VE_f$ and related variables from reliable datasets needs to be explored.

### 3.2. Optimization method: RF-PMRS

Therefore, to overcome the above disadvantages, an optimized method called RF-PMRS is proposed. Fig. 2(b) shows the process of our method, while optimizations for FMF and $VE_f$ are described separately below.

**1) FMF dataset selection**

We introduce the Phy-DL FMF dataset into the PMRS method to improve the accuracy of size-cutting results. In terms of performance, it exhibits higher accuracy and wider space-time coverage than satellite products (Yan, 2021). See the data section for details.

**2) $VE_f$ simulation based on ML**

The main idea is to establish an ML model between the $VE_f$ truth obtained from multiple AERONET sites and related variables, thus improving the subsequent $VE_f$-simulation accuracy (Fig. 3).

**Step 1** $VE_f$ calculation

The $VE_f$ true values are calculated concerning equations (5)-(7). Due to the spatiotemporal variability of different aerosol types, we calculate the $VE_f$ values at 9 AERONET stations around the world (Table 1) to train a universal model. The first step in Fig. 3 shows their distribution characteristics. Among them, Beijing and Beijing-CAMS sites are highlighted since they participate in the subsequent point validation experiment.

**Fig. 3.** Specific steps for simulating VE_f based on ML in our RF-PMRS method. The map used in step 1 is from NASA Visible Earth (https://visibleearth.nasa.gov/images/57752/blue-marble-land-surface-shallow-water-and-shaded-topography). The red points in step 1 represent the distribution of the 9 AERONET sites and the two yellow quadrangles in the zoom-in view highlight the Beijing (BJ) and Beijing-CAMS (BC) sites.

**Step 2** VE_f-related variables selection

According to the theory, FMF is selected as the most important modeling variable. Previous studies have also shown that the FMF-VE_f relationship has a good single-value correspondence, which is not affected by AOD. Compared with $AOD_f$ and $V_{f,column}$, FMF is a better indicator for estimation (Zhang and Li, 2015). In addition, considering the spatiotemporal heterogeneity of VE_f, the latitude, longitude (LAT, LON), and data time (month, day) of each site are added to the training.

**Step 3** RF model establishment

From step 2, VE_f can be expressed as:

$$VE_f = f(FMF, LAT, LON, month, day) \tag{11}$$

11

277     We optimize $VE_f$ expression based on random forest (RF). RF is made up of multiple

278 decision trees that can build high-accuracy models based on fewer variables (Svetnik

279 et al., 2003; Yang et al., 2020). This ensemble ML method randomly samples the

280 training dataset to form multiple subsets and random combinations of features are

281 selected in node splitting (Belgiu and Drăguţ, 2016). The specific process is to 1)

282 generate training subsets, 2) build an optimal model, and 3) calculate the result (Fig. 3

283 shows its flowchart). Note that the station FMF values (S-FMF) are used when training.

284

285 **Step 4** Accuracy validation

286     The $VE_f$ estimation is also based on equation (11), where *f* is the optimal relationship

287 after RF parameter adjustment, and Phy-DL FMF is applied to realize the extension of

288 model results from point to surface. 10-fold cross-validation (CV) (Rodriguez et al.,

289 2009) and isolated-validation (IV) are used to evaluate model performance (For details

290 of the validation methods, see Appendix A1).

291

292 **3) PM$_{2.5}$ value estimation and evaluation**

293     Then, calculate PM$_{2.5}$ according to the corresponding process (equation (8)). The

294 variables (in sections 2.2 to 2.4) are spatially matched to ground sites at their respective

295 resolutions. And based on UTC, the PM$_{2.5}$ validation is conducted on a daily scale in

296 2017. Because of the effective quantity of the AERONET public dataset and MODIS

297 data, we choose 2017 as the representative year. Note that we select the measured

298 empirical value of $\rho_{f,dry}$ (i.e., 1.5 g/cm$^3$) for the NC region from Gao et al. (2007).

299     The statistical indicators used in the evaluation include correlation coefficient (R),

300 mean bias (MB), relative mean bias (RMB), root mean square error (RMSE), and mean

301 absolute error (MAE). In addition, relative predictive error (RPE) is added to validate

302 the accuracy of the RF-based $VE_f$ model. See Appendix A2 for the specific information

303 on these indicators.

304

305 **4. Experiment results**

306     Three main experiments are conducted to verify the proposed RF-PMRS method,

307 and the specific information is shown in Table 2.

308 **Table 2**. A brief information summary of the experiments conducted in our study.

| Experiment | Object | Region | Period | Time scale |
|---|---|---|---|---|
| Model performance for training $VE_f$ | $VE_f$ | Global scale (Nine AERONET sites) | CV: Training period in Table 1 IV: Isolated-validation period in Table 1 (See Appendix A1) | Daily |
| Accuracy evaluation of PMRS/RF-PMRS | $PM_{2.5}$ | Two AERONET Sites: Beijing, Beijing-CAMS | 2017 | Daily |
| Generalization performance of RF-PMRS | $PM_{2.5}$ | North China region | 2017 | Daily |

309

## 4.1. RF model performance for training $VE_f$

311 The simulation model of $VE_f$ is trained based on the data in Table 1 and see Appendix

312 A3 for the adjustment of the model parameters. Table 3 shows that RF can capture the

313 complex relationship between $VE_f$ and related variables well. R is as high as 0.974

314 (0.975), RMSE and MAE are both small, and RPE is around 30%, which suggests the

315 desired estimation accuracy. Overall, the CV results represent the great performance of

316 the RF model for extracting information, that is, the relationship of multi-source data

317 to $VE_f$. In the meantime, the statistical results in CV and IV experiments are similar,

318 indicating that the RF model has no obvious overfitting phenomenon.

319

320 **Table 3**. Performance statistics of the RF model for training $VE_f$. N represents the number of data,
321 and $VE_f$ has no unit.

| | R | RMSE | RPE | MAE | N |
|---|---|---|---|---|---|
| **Cross-validation（CV）** | 0.974 | 0.076 | 32.9% | 0.034 | 6463 |
| **Isolated-validation（IV）** | 0.975 | 0.067 | 29.8% | 0.037 | 814 |

322

## 4.2. Accuracy evaluation of PMRS/RF-PMRS at AERONET stations

324 After applying the Phy-DL FMF data to the calculation process, the experiment

325 compares $PM_{2.5}$ results of PMRS and RF-PMRS at Beijing (BJ) and Beijing-CAMS

326 (BC) AERONET sites in 2017. Here, RF-PMRS simulates $VE_f$ based on RF, replacing

327 the polynomial of the PMRS method. Note that the results of the two sites are compared

328    with their respective nearest ground PM$_{2.5}$ stations (distances of 3.64 km and 3.91 km,

329    respectively, in line with the representative range of ground stations in previous studies

330    (Shi et al., 2018)).

331        Fig. 4 displays the time series of PM$_{2.5}$ values for different models at two sites. The

332    blue line fits the red line better than the gray one, confirming that the PM$_{2.5}$ results of

333    RF-PMRS are closer to the true values. Within the range of the black circles at positions

334    1 and 2, the variation of RF-PMRS results has better consistency with the ground truth,

335    while the PMRS results show dislocation and excessive growth. The overall

336    performance of the RF-PMRS estimations can signify the effectiveness of our proposed

337    method framework. As observed in the red boxes at positions 3 and 4, both models have

338    a certain degree of deviation, which is found to be consistent with the time regularity

339    of the AOD high values. Meanwhile, Fig. B1 (in Appendix B) plots the bias time series

340    between PMRS/RF-PMRS and in-situ values. As can be seen, the bias of the

341    optimization method (RF-PMRS) is stably distributed around zero, which greatly

342    reduces the numerical uncertainty. And it is worth noting that our method has well

343    mitigated the apparent overestimation of the original model (PMRS) in the case of

344    above-normal aerosol loadings. Furthermore, the average PM$_{2.5}$ values from ground

345    stations, PMRS, and RF-PMRS are compared. As for the two sites, the RF-PMRS

346    results are satisfactory. As depicted in Fig. 5, the RF-PMRS and station mean values

347    are close, with a difference of 4.82 μg/m³ (BJ) and 2.73 μg/m³ (BC), suggesting a good

348    estimation. Nevertheless, the PMRS results have deviations greater than 40 μg/m³, and

349    overestimation exists at both sites. It can be inferred that, in our proposed method, the

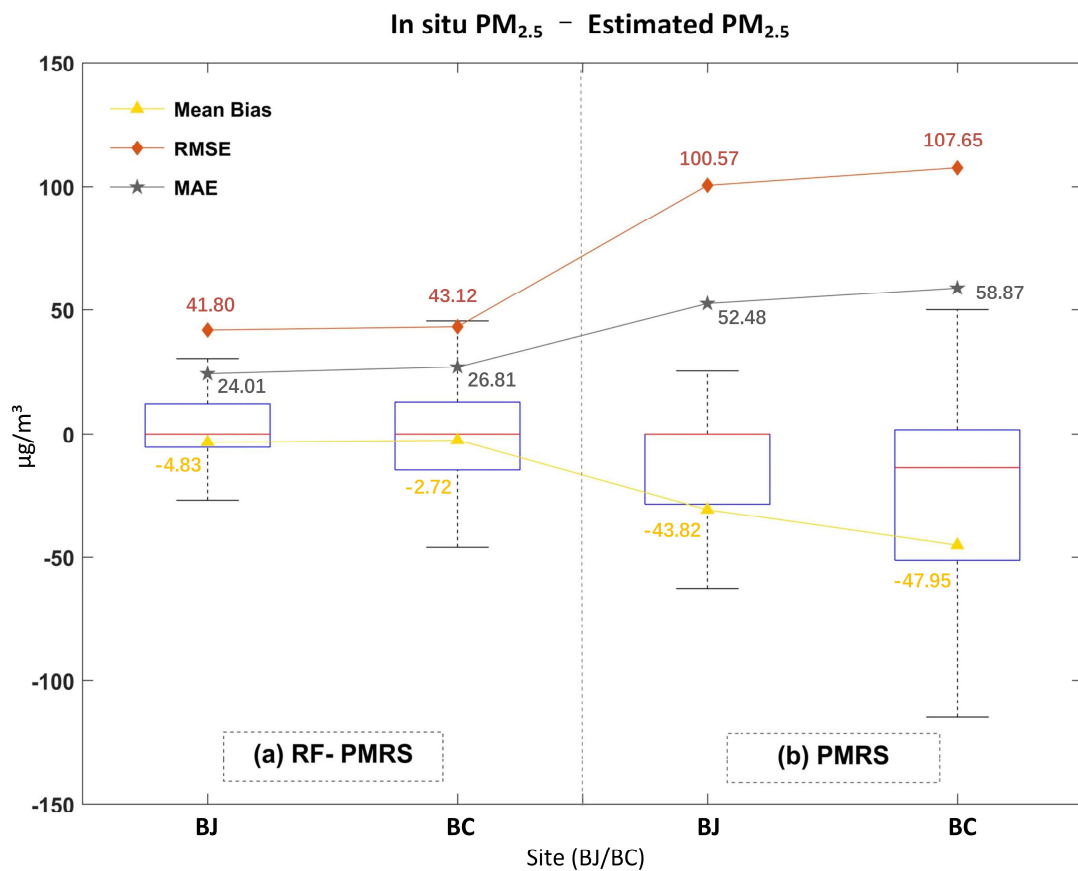350    optimization of VE$_f$ can greatly improve the PM$_{2.5}$ estimation accuracy.

**Fig. 4.** Three PM$_{2.5}$ time series at the Beijing (BJ) and Beijing-CAMS (BC) sites under their respective DOYs in 2017. Here, DOY (valid) means the day of the year with valid AOD, FMF, and other PM$_{2.5}$-related data. Grey, blue, and red lines represent PM$_{2.5}$ values of PMRS, RF-PMRS, and stations (STA), respectively. The red boxes and black circles select a specific period for analysis.



**Fig. 5.** Annual average PM$_{2.5}$ values from stations (left), RF-PMRS (middle), and PMRS model (right) at the BJ and BC sites.

15

Aiming at visually comparing the optimization effect, Fig. 6 plots the $PM_{2.5}$ bias distribution patterns for two methods. From the boxplot, the average $PM_{2.5}$ bias of RF-PMRS is close to zero (less than 5 μg/m³), which is greatly lower than that of PMRS. Besides, PMRS $PM_{2.5}$ has a larger deviation range, which manifests in two aspects. One is the maximum bias, specifically, it has exceeded 100 μg/m³ at the BC site. The other is the overall distribution of the data bias, the BJ site ones are mostly distributed below zero, indicating an obvious overestimation. As for RF-PMRS, the above circumstances are not obviously reflected in it. In addition, as can be seen from the indicators, RMSE and MAE of RF-PMRS $PM_{2.5}$ decrease by about half in comparison with PMRS. And the experiment has confirmed that the RF-PMRS $PM_{2.5}$ values have a strong linear relationship with the ground truth at both sites, with R around 0.8 (0.82 at BJ and 0.78 at BC). Such a large optimization effect is attributed to the $VE_f$ expression replacement to the fitted RF model.



**Fig. 6.** Boxplots of RF-PMRS (a) and PMRS (b) $PM_{2.5}$ bias at the BJ and BC sites. The upper (lower)

black line of each box represents the largest (smallest) value, the blue upper (lower) border represents the upper (lower) quartile, and the red line denotes the median. Besides, the yellow, orange, and gray symbols are the MB, RMSE, and MAE of the corresponding $PM_{2.5}$ concentration.

## 4.3. Generalization performance of RF-PMRS

Then, we estimate $PM_{2.5}$ based on PMRS and RF-PMRS within North China in 2017 (Fig. 1 exhibits the distribution pattern of the validation stations). Table 4 shows the accuracy statistics. It can be seen that RF-PMRS greatly reduces the bias (about 44.87%), with MB of about 2.31 µg/m³. Similar to the results at the sites, the RF-PMRS method can derive $PM_{2.5}$ concentration with practically no overestimation (underestimation). Although there is not much difference in R values of the two models (R of RF-PMRS is only improved by 0.01), RMSE and MAE of which decrease by about 39.96 µg/m³ and 18.86 µg/m³, respectively. As a result, the optimized method deserves to be considered excellent.

**Table 4**. Validation results of PMRS and RF-PMRS $PM_{2.5}$ in North China.

| Method | R | MB (µg/m³) | RMB (%) | RMSE (µg/m³) | MAE (µg/m³) |
|---|---|---|---|---|---|
| **PMRS** | 0.69 | -29.34 | 48.71% | 79.98 | 44.72 |
| **RF-PMRS** | 0.70 | 2.31 | 3.84% | 40.02 | 25.86 |

Meanwhile, $PM_{2.5}$ scatterplots are presented below. As depicted in Fig. 7, there are sufficient estimated samples (28305) in the NC region, which guarantees the credibility of our validation results. In general, the RF-PMRS $PM_{2.5}$ values are distributed around the 1:1 reference line evenly, with a slightly higher R of 0.70 compared to that of the original method. And the slope of the linear fitting relationship reaches 0.82, which indicates that the proposed method greatly reduces the overestimation of PMRS with a linear slope of 1.46. Although the overall performance of the RF-PMRS estimations maintains an excellent level, defects do remain. To be specific, in areas with high $PM_{2.5}$ concentration (especially greater than 150 µg/m³), RF-PMRS results exist a slight underestimation. It may be caused by the relatively small number of high-value $PM_{2.5}$ points (only 1319 out of 28305), which is difficult to adequately reflect the fitting effect
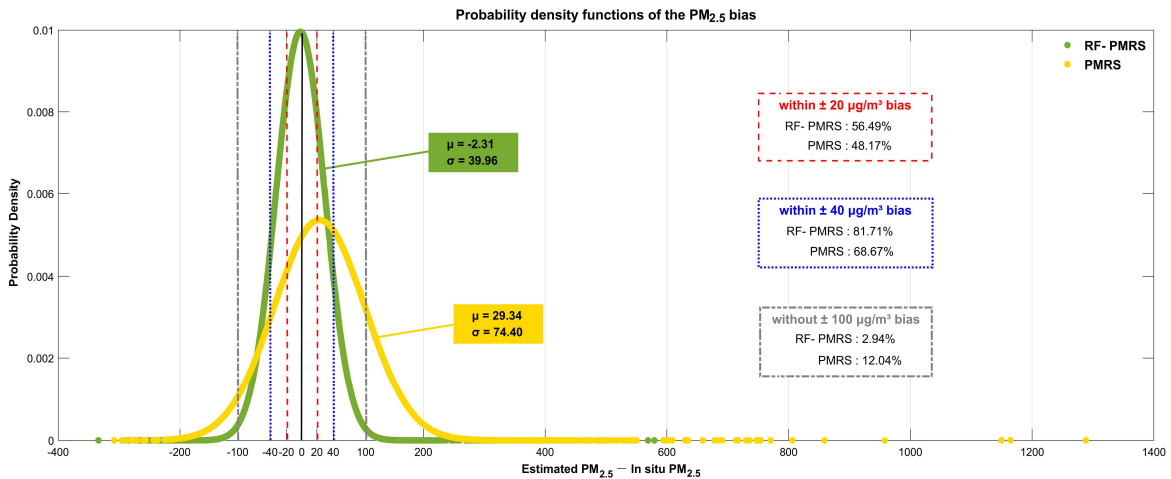
404    of the method.

405

410

411    As for RF-PMRS, the deviation is reduced to a large extent, so the probability density

412    function maps based on the bias of PMRS and RF-PMRS are further drawn. Fig. 8

413    visualizes the probability densities within different bias ranges. In terms of distribution

414    characteristics, the overall bias of RF-PMRS from the zero value (black solid line) is

415    small. About the curve shape, it is high and narrow, manifesting that the bias has a lower

416    standard deviation (STD) and is more prone to appear around the mean. However,

417    PRMS shows a more discrete distribution pattern, and there are many outliers outside

418    the range of greater than 600 μg/m³. Simultaneously, as can be concluded from the three

419    boxes, within the bias range of ±20 μg/m³ and ±40 μg/m³, the data numbers of RF-

420    PMRS results increase by 8.32% and 12.81%, respectively. Outside the range of ±100

421    μg/m³, the number decreases by 9.10%. Therefore, as far as the accuracy is concerned,

422    RF-PMRS results have lower bias and better stability.

423

424    In addition to the above general performance comparison in Section 4.3, Fig. 9

425    presents the annual average RMSE spatial distribution of PMRS and RF-PMRS $PM_{2.5}$

426    at NC stations. The two methods show a large deviation in the middle and southeast,

427    and the RMSE map of PMRS has more red points. However, RF-PMRS can weaken

428  this phenomenon very well since its RMSE representative colors are generally light. In

429  particular, the proportion of dark red sites (RMSE greater than 60 μg/m³) decreases

430  from 65.44% (PMRS) to 4.15% (RF-PMRS). In the areas where the ground stations are

431  clustered, the deviation also reduces significantly.

432



433

434  **Fig. 8.** Probability density functions of PMRS (yellow) and RF-PMRS (green) PM$_{2.5}$ bias. The red,

435  blue and grey dotted lines indicate the bias boundaries of ±20 μg/m³, ±40 μg/m³, and ±100 μg/m³,

436  respectively. μ and σ represent the mean value and standard deviation of each data.

437



438

439  **Fig. 9.** RMSE of the year-average PM$_{2.5}$ concentration values between different models and ground

440  stations (left: PMRS PM$_{2.5}$, right: RF-PMRS PM$_{2.5}$). Note that the top red of the RMSE legend

441  indicates RMSE values equal to or greater than 60 μg/m³.

442

443    In a word, the above analysis demonstrates that compared with the simple quadratic

444    polynomial relationship (equation (10)), the established RF model in RF-PMRS can

445    more accurately capture the relationship between $VE_f$ and multiple variables, thereby

446    improving the $PM_{2.5}$ estimation accuracy.

447

448    **5. Discussion**

449    **5.1. Accuracy comparison of PMRS using MODIS/Phy-DL FMF**

450    To confirm the superiority of the Phy-DL FMF data adopted in our method

451    framework, taking the BJ and BC sites as examples (in 2017), the experiment compares

452    the $PM_{2.5}$ accuracy and the number of effective days calculated by PMRS based on

453    different FMF. Table 5 presents the overall day-level results. Here, 'DOY' means the

454    day of the year and 'valid' means that all variables related to the $PM_{2.5}$ calculation are

455    valid. As can be seen, after the FMF replacement, the valid DOY turns out to become

456    more (an increase of 113 days), which illustrates that the number of effective $PM_{2.5}$

457    concentrations has gone up by about 5 times. Moreover, the accuracy has been

458    significantly enhanced, with R increased by about 0.30, RMSE and MAE decreased by

459    26.14% and 16.47% accordingly. On the whole, Phy-DL FMF contributes to the

460    improvement of PMRS results, signifying the first step optimization of the proposed

461    RF-PMRS method is effective.

462

463    **Table 5**. Validation results of the PMRS method using different FMF data. The valid DOY refers to
464    the number of days that the AOD, FMF, and other data are not missing when calculating $PM_{2.5}$. Note
465    that since the valid days of the two schemes are different, the MB and RMB are not compared.

|  | Valid DOY | R | RMSE ($\mu g/m^3$) | MAE ($\mu g/m^3$) |
|---|---|---|---|---|
| **PMRS with MODIS FMF** | 30 | 0.38 | 63.01 | 35.64 |
| **PMRS with Phy-DL FMF** | 143 | 0.68 | 46.54 | 29.77 |

466

467    **5.2. Performance compared with other ML models**

468    Different machine learning models are suitable for diverse research data, and

469    decision tree (DT) models can better fit experiments with fewer variables, such as this

470    study. For comparison, except for RF, the Extremely Randomized Tree (ERT) (Geurts

471 et al., 2006) and Gradient Boosting Decision Tree (GBDT) (Friedman, 2001) models

472 have also been established. The results of training $VE_f$ based on the above three DT

473 models are presented in Table 6 and Table 7. By contrast, RF performs best in CV and

474 IV experiments, as indicated by the multiple accuracy indicators. Although ERT and

475 GBDT models are comparable to RF in some indicators, there exists a certain degree

476 of overfitting in the above two models, which is manifested in that their IV results are

477 clearly worse than their respective CV ones. Thus, the RF model is applied to our study.

478

479 **Table 6**. Cross-validation results in comparison of the decision tree models for training $VE_f$. N
480 represents the number of data, and $VE_f$ has no unit.

| CV results | | | | | |
|---|---|---|---|---|---|
| | **R** | **RMSE** | **RPE** | **MAE** | **N** |
| **RF** | 0.974 | 0.076 | 0.330 | 0.034 | |
| **ERT** | 0.972 | 0.079 | 0.343 | 0.035 | 6463 |
| **GBDT** | 0.973 | 0.078 | 0.339 | 0.036 | |

481

482 **Table 7**. Isolated-validation results in comparison of the decision tree models for training $VE_f$. The
483 indicators are the same as those in Table 6.

| IV results | | | | | |
|---|---|---|---|---|---|
| | **R** | **RMSE** | **RPE** | **MAE** | **N** |
| **RF** | 0.975 | 0.067 | 0.299 | 0.037 | |
| **ERT** | 0.967 | 0.076 | 0.340 | 0.042 | 814 |
| **GBDT** | 0.969 | 0.074 | 0.331 | 0.040 | |

484

485 **5.3. Feature importance of the embedded RF model**

486     Additionally, the feature importance of RF is calculated to evaluate the contribution

487 of model predictors to $VE_f$ simulation. Fig. B2 (in Appendix B) shows the results by

488 normalization (taking 100 as the total). Without a doubt, FMF accounts for the largest

489 proportion, about 76.4%, which is consistent with the analysis when selecting the $VE_f$-

490 related variables (see Section 3.2). The contribution of spatiotemporal variables is about

491 1/3 of FMF, which indirectly affirms the credibility of RF feature learning. Also, it

492 provides a basis for further uncertainty optimization of $VE_f$ and $PM_{2.5}$ accuracy.

493

**6. Conclusion**

Among various satellite remote sensing methods for PM$_{2.5}$ retrieval, the semi-empirical physical approach has strong physical significance and clear calculation steps and derives the PM$_{2.5}$ mass concentration independently of in situ observations. However, the parameters with the meaning of optical properties are difficult to express, which need to be optimized. Hence, the study proposes a method (RF-PMRS) that embeds machine learning in a physical model to obtain surface PM$_{2.5}$: 1) Based on the PMRS method and select the Phy-DL FMF product with a combined mechanism; 2) Use the RF model to fit the parameter VE$_f$, rather than a simple quadratic polynomial. In the point-to-surface validation, RF-PMRS shows great optimized performance. Experiments at two AERONET sites show that R reaches up to 0.8. And in North China, RMSE decreases by 39.95 μg/m³ with a 44.87% reduction in relative deviation. In the future, we will further explore the combination of atmospheric mechanism and machine learning, then research the PM$_{2.5}$ retrieval methods with physical meaning and higher accuracy.

**Appendix A: Supplementary description**

**A1. 10-fold cross-validation and isolated-validation**

The sample-based 10-fold cross-validation method is applied to test the fitting and predictive ability of our model. The original dataset is randomly divided into ten parts, nine of which are used as the training set for model fitting, and the remaining one is used for prediction, then the cross-validation process is repeated ten rounds until each data has been used as the test set.

At the same time, when verifying the RF-based VE$_f$ model, the dataset in the period that did not participate in the training in Table 1 is used for isolated-validation.

**A2. Statistical indicators**

$$R = \frac{\sum_{i=1}^{m}\left(y_i - \overline{y}\right)\sum_{i=1}^{m}\left(f_i - \overline{f}\right)}{\sqrt{\sum_{i=1}^{m}\left(y_i - \overline{y}\right)^2}\sqrt{\sum_{i=1}^{m}\left(f_i - \overline{f}\right)^2}}$$

522
$$MB = \bar{y} - \bar{f}$$

523
$$RMB = \text{abs}\left(\frac{\bar{y} - \bar{f}}{\bar{y}}\right)$$

524
$$RMSE = \sqrt{\frac{1}{m}\sum_{i=1}^{m}\left(y_i - f_i\right)^2}$$

525
$$MAE = \frac{1}{m}\sum_{i=1}^{m}\left|y_i - f_i\right|$$

526
$$RPE = \frac{\sqrt{\frac{1}{m}\sum_{i=1}^{m}\left(y_i - f_i\right)^2}}{\bar{y}}$$

527 where $m$ is the total number of observations, $i$ is the number of measurements, $y_i$ is the

528 i-th observation, $f_i$ is the corresponding estimation result. And $\bar{y}$ and $\bar{f}$ are the

529 averages of all observations and estimates, respectively.

530

531 **A3. Parameter adjustments of the RF model**

532 The four parameters of RF are adjusted, that is the correlation coefficient r changes

533 with (a) the number of trees, (b) maximum depth, (c) maximum number of features

534 when splitting, (d) minimum number of split samples. Experiments show that the

535 maximum depth varies greatly in a small range. To prevent overfitting, the four

536 parameters of RF are adjusted to 60, 10, 2, and 8. It can ensure high accuracy while

537 improving training efficiency.

538

**Appendix B: Figures**

**Fig. B1.** The time series of PMRS/RF-PMRS $PM_{2.5}$ bias at the Beijing and Beijing-CAMS sites under their respective DOYs in 2017. The orange line represents the bias between the $PM_{2.5}$ values of PMRS and stations, while the blue one indicates the $PM_{2.5}$ difference between RF-PMRS and stations.



**Fig. B2.** The predictor importance results (normalized) of the RF model for training $VE_f$.

**Code and data availability**

All relevant codes as well as the intermediate data of this work are archived at https://doi.org/10.5281/zenodo.7183822 (Jin, 2022). The MCD19A2 data can be

downloaded on https://ladsweb.modaps.eosdis.nasa.gov (last access: 30-09-2022) (Lyapustin and Wang, 2015). Detailed information about the Phy-DL FMF dataset can be found at https://doi.org/10.5281/zenodo.5105617 (Yan, 2021). Meteorological data used in this work are obtained at https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels (last access: 30-09-2022) (Hersbach et al., 2018). AERONET data was downloaded from https://aeronet.gsfc.nasa.gov/ (last access: 30-09-2022) (Giles et al., 2019).

**Author contributions**

**Caiyi Jin:** Data curation, Methodology, Formal analysis, Writing - original draft. **Qiangqiang Yuan:** Conceptualization, Supervision, Project administration, Writing - review and editing. **Tongwen Li:** Resources, Methodology, Writing - review and editing, Formal analysis. **Yuan Wang:** Methodology, Validation, Writing - review and editing. **Liangpei Zhang:** Supervision, Writing - review and editing.

582

**References**

Belgiu, M., and Drăguţ, L.: Random forest in remote sensing: A review of applications and future directions, ISPRS J. Photogramm. Remote Sens., 114, 24-31, https://doi.org/10.1016/j.isprsjprs.2016.01.011, 2016.

Bowe, B., Xie, Y., Li, T., Yan, Y., Xian, H., and Al-Aly, Z.: The 2016 global and national burden of diabetes mellitus attributable to PM2.5 air pollution, Lancet Planet. Health, 2, e301-e312, https://doi.org/10.1016/S2542-5196(18)30140-2, 2018.

Chen, X., de Leeuw, G., Arola, A., Liu, S., Liu, Y., Li, Z., and Zhang, K.: Joint retrieval of the aerosol fine mode fraction and optical depth using MODIS spectral reflectance over northern and eastern China: Artificial neural network method, Remote Sens Environ, 249, 112006, https://doi.org/10.1016/j.rse.2020.112006, 2020.

Friedman, J.H.: Greedy function approximation: a gradient boosting machine, Ann Stat, 29(5), 1189–1232, http://www.jstor.org/stable/2699986, 2001.

Gao, J., Zhou, Y., Wang, J., Wang, T., and Wang, W.X.: Inter-comparison of WPSTM-TEOMTM-MOUDITM and investigation on particle density, Huan Jing Ke Xue, 28, 1929-1934, https://doi.org/10.3321/j.issn:0250-3301.2007.09.005, 2007.

Gao, L., Li, J., Chen, L., Zhang, L., and Heidinger, A.K.: Retrieval and validation of atmospheric aerosol optical depth from AVHRR over China, IEEE Trans Geosci Remote Sens, 54, 6280-6291, https://doi.org/10.1109/TGRS.2016.2574756, 2016.

Geng, G., Zhang, Q., Martin, R.V., van Donkelaar, A., Huo, H., Che, H., Lin, J., and He, K.: Estimating long-term PM2.5 concentrations in China using satellite-based aerosol optical depth and a chemical transport model, Remote Sens Environ, 166, 262-270, https://doi.org/10.1016/j.rse.2015.05.016, 2015.

Geurts, P., Ernst, D., and Wehenkel, L.: Extremely randomized trees, Mach Learn, 63, 3-42, https://doi.org/10.1007/s10994-006-6226-1, 2006.

Giles, D.M., Holben, B.N., Eck, T.F., Smirnov, A., Sinyuk, A., Schafer, J., Sorokin, M.G., and Slutsker, I.: Aerosol robotic network (AERONET) version 3 aerosol optical depth and inversion products, in: American Geophysical Union (AGU) 98th Fall Meeting Abstracts, New Orleans, America, 11-15 December 2017, A11O-01, 2017.

Giles, D. M., Sinyuk, A., Sorokin, M. G., Schafer, J. S., Smirnov, A., Slutsker, I, Eck, T. F., Holben, B. N., Lewis, J. R., Campbell, J. R., Welton, E. J., Korkin, S. V., and Lyapustin, A. I.: Advancements in the Aerosol Robotic Network (AERONET) Version 3 database - automated near-real-time quality control algorithm with improved cloud screening for Sun photometer aerosol optical depth (AOD) measurements, Atmos Meas Tech, 12, 169–209, https://doi.org/10.5194/amt-12-169-2019, 2019.

Gupta, P., and Christopher, S.A.: Particulate matter air quality assessment using integrated surface, satellite, and meteorological products: Multiple regression approach, J. Geophys. Res. Atmos., 114, D14205, https://doi.org/10.1029/2008JD011496, 2009.

Hand, J.L., and Kreidenweis, S.M.: A new method for retrieving particle refractive index and effective density from aerosol size distribution data, Aerosol Sci Technol, 36, 1012-1026, https://doi.org/10.1080/02786820290092276, 2002.

Hänel, G., and Thudium, J.: Mean bulk densities of samples of dry atmospheric aerosol particles: A summary of measured data, Pure Appl. Geophys., 115, 799-803, https://doi.org/10.1007/BF00881211, 1977.

He, J., Yuan, Q., Li, J., and Zhang, L.: PoNet: A universal physical optimization-based spectral super-resolution network for arbitrary multispectral images, Inform Fusion, 80, 205-225, https://doi.org/10.1016/j.inffus.2021.10.016, 2022.

He, J., Li, J., Yuan, Q., Shen, H., and Zhang, L.: Spectral Response Function-Guided Deep Optimization-Driven Network for Spectral Super-Resolution, IEEE Trans Neural Netw. Learn. Syst., PP(99), 1-15, https://doi.org/10.1109/TNNLS.2021.3056181, 2021.

Hersbach, H., Bell, B., Berrisford, P., Biavati, G., Horányi, A., Muñoz Sabater, J., Nicolas, J., Peubey, C., Radu, R., Rozum, I., Schepers, D., Simmons, A., Soci, C., Dee, D., Thépaut, J-N.: ERA5 hourly data on single levels from 1979 to present, Copernicus Climate Change Service (C3S) Climate Data Store (CDS) [data set], (Accessed on 30-09-2022), https://doi.org/10.24381/cds.adbb2d47, 2018.

Holben, B.N., Eck, T.F., Slutsker, I., Tanré, D., Buis, J.P., Setzer, A., Vermote, E., Reagan, J.A., Kaufman, Y.J., Nakajima, T., Lavenu, F., Jankowiak, I., and Smirnov, A.: AERONET — A federated instrument network and data archive for aerosol

characterization, Remote Sens Environ, 66, 1-16, https://doi.org/10.1016/S0034-4257(98)00031-5, 1998.

Irrgang, C., Boers, N., Sonnewald, M., Barnes, E.A., Kadow, C., Staneva, J., and Saynisch-Wagner, J.: Towards neural Earth system modelling by integrating artificial intelligence in Earth system science, Nat. Mach. Intell., 3, 667-674, https://doi.org/10.1038/s42256-021-00374-3, 2021.

Jin, C.: An optimized semi-empirical physical approach for satellite-based PM2.5 retrieval: using random forest model to simulate the complex parameter, Zenodo [code], https://doi.org/10.5281/zenodo.7183822, 2022.

Koelemeijer, R.B.A., Homan, C.D., and Matthijsen, J.: Comparison of spatial and temporal variations of aerosol optical thickness and particulate matter over Europe, Atmospheric Environ., 40, 5304-5315, https://doi.org/10.1016/j.atmosenv.2006.04.044, 2006.

Kokhanovsky, A.A., Prikhach, A.S., Katsev, I.L., and Zege, E.P.: Determination of particulate matter vertical columns using satellite observations, Atmos Meas Tech, 2, 327-335, https://doi.org/10.5194/amt-2-327-2009, 2009.

Lee, J.-B., Lee, J.-B., Koo, Y.-S., Kwon, H.-Y., Choi, M.-H., Park, H.-J., and Lee, D.-G.: Development of a deep neural network for predicting 6 h average PM2.5 concentrations up to 2 subsequent days using various training data, Geosci. Model Dev., 15, 3797–3813, https://doi.org/10.5194/gmd-15-3797-2022, 2022.

Li, T., Shen, H., Zeng, C., Yuan, Q., and Zhang, L.: Point-surface fusion of station measurements and satellite observations for mapping PM2.5 distribution in China: Methods and assessment, Atmospheric Environ., 152, 477-489, https://doi.org/10.1016/j.atmosenv.2017.01.004, 2017.

Li, Z., Zhang, Y., Shao, J., Li, B., Hong, J., Liu, D., Li, D., Wei, P., Li, W., Li, L., Zhang, F., Guo, J., Deng, Q., Wang, B., Cui, C., Zhang, W., Wang, Z., Lv, Y., Xu, H., Chen, X., Li, L., and Qie, L.: Remote sensing of atmospheric particulate mass of dry PM2.5 near the ground: Method validation using ground-based measurements, Remote Sens Environ, 173, 59-68, https://doi.org/10.1016/j.rse.2015.11.019, 2016.

Lyapustin, A., Wang, Y., Laszlo, I., Kahn, R., Korkin, S., Remer, L., Levy, R., and

Reid, J.S.: Multiangle implementation of atmospheric correction (MAIAC): 2. Aerosol algorithm, J. Geophys. Res. Atmos., 116, D03211, https://doi.org/10.1029/2010JD014986, 2011.

Lyapustin, A., Wang, Y., Xiong, X., Meister, G., Platnick, S., Levy, R., Franz, B., Korkin, S., Hilker, T., Tucker, J., Hall, F., Sellers, P., Wu, A., and Angal, A.: Scientific impact of MODIS C5 calibration degradation and C6+ improvements, Atmos Meas Tech, 7, 4353-4365, https://doi.org/10.5194/amt-7-4353-2014, 2014.

Lyapustin, A., andWang, Y.: MCD19A2 MODIS/Terra+Aqua Aerosol Optical Thickness Daily L2G Global 1km SIN Grid, NASA LP DAAC [data set], (Accessed on 30-09-2022), http://doi.org/10.5067/MODIS/MCD19A2.006, 2015.

Lyu, B., Huang, R., Wang, X., Wang, W., and Hu, Y.: Deep-learning spatial principles from deterministic chemical transport models for chemical reanalysis: an application in China for PM2.5, Geosci. Model Dev., 15, 1583–1594, https://doi.org/10.5194/gmd-15-1583-2022, 2022.

Ma, Z., Hu, X., Huang, L., Bi, J., and Liu, Y.: Estimating ground-Level PM2.5 in China using satellite remote sensing, Environ. Sci. Technol., 48, 7436-7444, https://doi.org/10.1021/es5009399, 2014.

Pope III, C.A., Burnett, R.T., Thun, M.J., Calle, E.E., Krewski, D., Ito, K., and Thurston, G.D.: Lung cancer, cardiopulmonary mortality, and long-term exposure to fine particulate air pollution, JAMA, 287, 1132-1141, https://doi.org/10.1001/jama.287.9.1132, 2002.

Raut, J., and Chazette, P.: Assessment of vertically-resolved PM10 from mobile lidar observations, Atmospheric Chem. Phys., 9, 8617-8638, https://doi.org/10.5194/acp-9-8617-2009, 2009.

Rodriguez, J.D., Perez, A., and Lozano, J.A.: Sensitivity analysis of k-fold cross validation in prediction error estimation, IEEE Trans. Pattern Anal. Mach. Intell., 32, 569-575, https://doi.org/10.1109/TPAMI.2009.187, 2009.

Shi, X., Zhao, C., Jiang, J.H., Wang, C., Yang, X., and Yung, Y.L.: Spatial representativeness of PM2.5 concentrations obtained using observations from network stations, J. Geophys. Res. Atmos., 123, 3145-3158,

701 https://doi.org/10.1002/2017JD027913, 2018.

702 Simmons, A.J., Untch, A., Jakob, C., Kållberg, P., and Undén, P.: Stratospheric water

703 vapour and tropical tropopause temperatures in ECMWF analyses and multi-year

704 simulations, Q J R Meteorol Soc, 125, 353-386,

705 https://doi.org/10.1002/qj.49712555318, 1999.

706 Svetnik, V., Liaw, A., Tong, C., Culberson J.C., Sheridan R.P., and Feuston, B.P.:

707 Random Forest: A Classification and Regression Tool for Compound Classification

708 and QSAR Modeling, J. Chem. Inf. Comput. Sci., 43, 6, 1947-1958,

709 https://doi.org/10.1021/ci034160g, 2003.

710 Van Donkelaar, A., Martin, R.V., and Park, R.J.: Estimating ground-level PM2. 5 using

711 aerosol optical depth determined from satellite remote sensing, J. Geophys. Res. Atmos.,

712 111, D21201, https://doi.org/10.1029/2005JD006996, 2006.

713 Wang, Y., Yuan, Q., Li, T., Shen, H., Zheng, L., and Zhang, L.: Evaluation and

714 comparison of MODIS Collection 6.1 aerosol optical depth against AERONET over

715 regions in China with multifarious underlying surfaces, Atmospheric Environ., 200,

716 280-301, https://doi.org/10.1016/j.atmosenv.2018.12.023, 2019.

717 Wu, X., Wang, Y., He, S., and Wu, Z.: PM2.5／PM10 ratio prediction based on a long

718 short-term memory neural network in Wuhan, China, Geosci. Model Dev., 13, 1499–

719 1511, https://doi.org/10.5194/gmd-13-1499-2020, 2020.

720 Xu, P., Chen, Y., and Ye, X.: Haze, air pollution, and health in China, Lancet, 382,

721 2067, https://doi.org/10.1016/S0140-6736(13)62693-8, 2013.

722 Yan, X., Zang, Z., Li, Z., Luo, N., Zuo, C., Jiang, Y., Li, D., Guo, Y., Zhao, W., Shi,

723 W., and Cribb, M.: A global land aerosol fine-mode fraction dataset (2001--2020)

724 retrieved from MODIS using hybrid physical and deep learning approaches, Earth Syst.

725 Sci. Data, 14, 1193-1213, https://doi.org/10.5194/essd-14-1193-2022, 2022.

726 Yan, X., Li, Z., Shi, W., Luo, N., Wu, T., and Zhao, W.: An improved algorithm for

727 retrieving the fine-mode fraction of aerosol optical thickness, part 1: Algorithm

728 development, Remote Sensing of Environment, 192, 87-97,

729 https://doi.org/10.1016/j.rse.2017.02.005, 2017.

730 Yan, X.: Physical and deep learning retrieved fine mode fraction (Phy-DL FMF),

731 Zenodo [data set], (Accessed on 30-09-2022), https://doi.org/10.5281/zenodo.5105617,

732 2021.

733 Yang, Q., Yuan, Q., Li, T., and Yue, L.: Mapping PM2.5 concentration at high

734 resolution using a cascade random forest based downscaling model: Evaluation and

735 application, J. Clean. Prod., 277, 123887,

736 https://doi.org/10.1016/j.jclepro.2020.123887, 2020.

737 Yuan, Q., Shen, H., Li, T., Li, Z., Li, S., Jiang, Y., Xu, H., Tan, W., Yang, Q., Wang,

738 J., Gao, J., and Zhang, L.: Deep learning in environmental remote sensing:

739 Achievements and challenges, Remote Sens Environ, 241, 111716,

740 https://doi.org/10.1016/j.rse.2020.111716, 2020.

741 Zhang, Y., Li, Z., Bai, K., Wei, Y., Xie, Y., Zhang, Y., Ou, Y., Cohen, J., Zhang, Y.,

742 Peng, Z., Zhang, X., Chen, C., Hong, J., Xu, H., Guang, J., Lv, Y., Li, K., and Li, D.:

743 Satellite remote sensing of atmospheric particulate matter mass concentration:

744 Advances, challenges, and perspectives, Fundamental Research, 1, 240-258,

745 https://doi.org/10.1016/j.fmre.2021.04.007, 2021.

746 Zhang, Y., Li, Z., Chang, W., Zhang, Y., de Leeuw, G., and Schauer, J.J.: Satellite

747 observations of PM2.5 changes and driving factors based forecasting over China 2000−

748 2025, Remote Sens., 12(16), 2518, https://doi.org/10.3390/rs12162518, 2020.

749 Zhang, Y., and Li, Z.: Remote sensing of atmospheric fine particulate matter (PM2.5)

750 mass concentration near the ground from satellite observation, Remote Sens Environ,

751 160, 252-262, https://doi.org/10.1016/j.rse.2015.02.005, 2015.