1  **An optimized semi-empirical physical approach for satellite-based**
2  **PM$_{2.5}$ retrieval: embedding machine learning to simulate complex**
3  **physical parameters**
4
5  Caiyi Jin [a], Qiangqiang Yuan [a, c, d, *], Tongwen Li [b, *], Yuan Wang [a], Liangpei Zhang [c, e]
6
7  [a] School of Geodesy and Geomatics, Wuhan University, Wuhan 430079, China.
8  [b] School of Geospatial Engineering and Science, Sun Yat-Sen University, Zhuhai
9    519082, China
10  [c] The Collaborative Innovation Center of Geospatial Technology, Wuhan 430079, China.
11  [d] The Key Laboratory of Geospace Environment and Geodesy, Ministry of Education,
12    Wuhan University, Wuhan 430079, China.
13  [e] State Key Laboratory of Information Engineering in Surveying, Mapping and Remote
14    Sensing, Wuhan University, Wuhan 430079, China.
15
16  * Corresponding author.
17  E-mail address: yqiang86@gmail.com, litw8@mail.sysu.edu.cn
18

19  **ABSTRACT**

20    Satellite remote sensing of PM$_{2.5}$ mass concentration has become one of the most

21  popular atmospheric research aspects, resulting in the development of different models.

22  Among them, the semi-empirical physical approach constructs the transformation

23  relationship between the aerosol optical depth (AOD) and PM$_{2.5}$ based on the optical

24  properties of particles, which has strong physical significance.  Also, it performs the

25  PM$_{2.5}$ retrieval independently of the ground stations. However, due to the complex

26  physical relationship, the physical parameters in the semi-empirical approach are

27  difficult to calculate accurately, resulting in relatively limited accuracy. To achieve the

28  optimization effect, this study proposes a method of embedding machine learning into

29  a semi-physical empirical model (RF-PMRS). Specifically, based on the theory of the

30  physical PM$_{2.5}$ remote sensing approach (PMRS), the complex parameter (VE$_f$, a

31  columnar volume-to-extinction ratio of fine particles) is simulated by the random forest

32  model (RF). Also, a fine mode fraction product with higher quality is applied to make

33  up for the insufficient coverage of satellite products. Experiments in North China show

34  that the surface PM$_{2.5}$ concentration derived by RF-PMRS has an average annual value

35  of 57.92 μg/m³ versus the ground value of 60.23 μg/m³. Compared with the original

36    method, RMSE decreases by 39.95 μg/m³, and the relative deviation reduces by 44.87%.

37    Moreover, validation at two AERONET sites presents a trend closer to the true values,

38    with an R of about 0.80. This study is also a preliminary attempt to combine model-

39    driven and data-driven models, laying a foundation for further atmospheric research on

40    optimization methods.

41    **Keywords:** PM$_{2.5}$; Physical approach; Machine learning; Volume-to-extinction ratio;

42    Fine mode fraction

43

44    **1. Introduction**

45    Epidemiological studies have indicated that PM$_{2.5}$ (fine particulate matter with an

46    aerodynamic equivalent diameter no greater than 2.5 μm) can adversely affect human

47    health, such as increasing the risk of diabetes and respiratory diseases (Bowe et al.,

48    2018; Pope III et al., 2002; Xu et al., 2013), and accurate surface PM$_{2.5}$ concentration

49    is the basis of air pollution-health related research. Satellite remote sensing has the

50    advantages of high resolution and global coverage (Ma et al., 2014; Wu et al., 2020),

51    including variables strongly associated with PM$_{2.5}$ such as aerosol optical depth (AOD).

52    Therefore, it has become a mainstream method for fine particles estimation (Zhang et

53    al., 2021).

54    There are mainly three satellite-based ways of retrieving PM$_{2.5}$. 1) Chemical transport

55    models-based method. It calculates a scaling factor η between AOD and PM$_{2.5}$

56    simulated by atmospheric chemical transport models (CTM) (Lyu et al., 2022) and then

57    transfers the proportional relationship to satellite AOD data when calculating surface

58    PM$_{2.5}$ concentration (Geng et al., 2015; Van Donkelaar et al., 2006). However, the

59    assumption of a constant factor between simulated and observed values has large

60    spatiotemporal limitations. 2) Univariate/Multivariate regression. This kind of method

61    establishes a statistical model between AOD, auxiliary variables, and ground PM$_{2.5}$

62    observations. Machine learning is a common tool for such data-driven methods due to

63    its powerful nonlinear fitting ability between multiple variables (Irrgang et al., 2021).

64    But the regression is affected by the distribution and density of ground stations (Gupta

65    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. Raut and Chazette (2009) introduced a specific

74   extinction cross-section to simplify the expression of S and $PM_{2.5}$ concentration was

75   estimated based on 355nm-band radar observations. 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. It can detect complex nonlinear

88   relationships of multiple data and model their interaction (Yuan et al., 2020; Lee et

89   al.,2022), which provides an idea for improving the accuracy of physical parameter

90   acquisition, thereby estimating high-precision $PM_{2.5}$ through semi-physical empirical

91   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 illustrates the specific

102     derivation process of the proposed method. Section 3 describes the experimental

103     datasets and analyzes the evaluation results. Some supporting experiments are

104     discussed in section 4. And the final part provides the conclusion.

105

106     **2. Methods**

107     Based on the basic physical properties of atmospheric aerosols, the semi-physical

108     empirical approach starts from the integration of PM mass concentration and AOD.

109     Then it combines several key factors related to PM$_{2.5}$, to derive the in situ PM$_{2.5}$

110     concentration through multiple remote sensing variables (Koelemeijer et al., 2006). The

111     overall empirical relationship can be represented as:

112
$$PM_{2.5} = AOD \frac{\rho}{H \cdot f(RH)} S \tag{1}$$

113     where $\rho$ denotes the particle density and $H$ denotes the atmospheric boundary layer

114     height. $f(RH)$ represents the hygroscopic growth factor related to relative humidity

115     $(RH)$. $S$ is an optical characteristic parameter that should be simulated.

116

117     **2.1. PMRS method**

118     **2.1.1. The expression of VE$_f$**

119     To illustrate $S$ more precisely, PMRS defines the columnar volume-to-extinction

120     ratio of fine particles (i.e., $VE_f$), which can be regarded as the basis of our

121     optimization method. So equation (1) is transformed into:

122
$$PM_{2.5} = AOD \frac{\rho}{H \cdot f(RH)} VE_f \tag{2}$$

123     Related to particle size, aerosol extinction, and other properties, $VE_f$ can be

124    expressed as:

$$VE_f = \frac{V_{f,column}}{AOD_f} \qquad (3)$$

125

$$AOD_f = AOD \cdot FMF \qquad (4)$$

126

127    Here, $AOD_f$ is the fine particle AOD and $FMF$ is the fine mode fraction. $V_{f,column}$

128    can be expressed by the vertical integral of particle volume size distributions (PVSD)

129    within a certain aerodynamic diameter range:

$$V_{f,column} = \int_0^{D_{p,c}} V(D_p) dD_p \qquad (5)$$

130

131    $D_{p,c}$ represents the cutting diameter, and the empirical value of 2.0 μm is chosen based

132    on previous literature (Hand and Kreidenweis, 2002; Hänel and Thudium, 1977). And

133    $V(D_p)$ represents the PVSD corresponding to the geometric equivalent diameter ( $D_p$ ).

134

135    **2.1.2. Specific process and limitations**

136        The PMRS method is developed from equation (2). Based on satellite AOD, the near-

137    surface PM2.5 can be obtained through multi-step transformation. Fig. 1(a) shows its

138    specific process. Each arrow refers to a step, respectively: size cutting (output: $AOD_f$ ),

139    volume visualization (output: $V_{f,column}$ ), bottom isolation (output: $V_f$ , fine particle

140    volume near the ground), particle drying (output: $V_{f,dry}$ , dry $V_f$ ) and PM2.5 weighting.

141    The overall expression is as follows:

$$PM_{2.5} = AOD \frac{FMF \cdot VE_f \cdot \rho_{f,dry}}{PBLH \cdot f_0(RH)} \qquad (6)$$

142

$$f_0(RH) = \left(1 - \frac{RH}{100}\right)^{-1} \qquad (7)$$

143

144    where $FMF$ denotes the fine mode fraction, $\rho_{f,dry}$ denotes the dry mass density of

145    $PM_{2.5}$, and $PBLH$ represents the planet boundary layer height. $f_0(RH)$ represents

146    the approximation of $f(RH)$ in equation (2), as expressed as equation (7).

147    Considering the aerosol types in different regions, PMRS fits $VE_f$ to a quadratic

148    polynomial relation of $FMF$ :

$$VE_f = 0.2887FMF^2 - 0.4663FMF + 0.356 \quad (0.1 \le FMF \le 1.0) \qquad (8)$$

149

150    PMRS has strong physical significance, the calculation steps are well-defined and

151    site-independent. Zhang and Li (2015) tested the performance of PMRS on 15 stations,

152    and the validation results had an uncertainty of 34%. Compared with the ground value

153    of Jinhua city in China, a 31.3% relative error was generated in Li et al. (2016). Besides,

154    Zhang et al. (2020) applied it to the $PM_{2.5}$ change analysis and prediction experiments

155    in China over 20 years. However, there may be a more complex nonlinear relationship

156    between $VE_f$ with FMF, not just a simple quadratic formula. Since $VE_f$ is related to the

157    aerosol type, adding other spatiotemporal variables may optimize the fitting process.

158    Additionally, high-quality FMF data is the basic guarantee for the estimated $PM_{2.5}$

159    quality. In a word, to further improve the physical method, a better nonlinear model

160    between $VE_f$ and related variables from reliable datasets needs to be explored.

161



162
163    **Fig. 1.** Surface $PM_{2.5}$ estimation flow of RF-PMRS. a) The five steps of the PMRS method. Gray
164    boxes are the intermediate outputs, blue boxes are the input data, and orange ones denote the
165    variables to be optimized. b) The specific optimization of RF-PMRS: FMF dataset replacement and
166    $VE_f$ simulation by RF model.

167

## 2.2. Optimization method: RF-PMRS

169    Therefore, to overcome the above disadvantages, an optimized method called RF-

170    PMRS is proposed. Fig. 1(b) shows the process of our method, while optimizations for

171    FMF and VE$_f$ are described separately below.

**1) FMF dataset selection**

173    We introduce the Phy-DL FMF dataset into the PMRS method to improve the

174    accuracy of size-cutting results. In the comparison experiment against Aerosol Robotic

175    Network (AERONET) FMF, Phy-DL FMF shows a higher accuracy (R = 0.78, RMSE

176    = 0.100) than Moderate-resolution Imaging Spectroradiometer (MODIS) FMF (R =

177    0.37, RMSE = 0.282) (Yan et al., 2022). Also, it performs better spatiotemporal

178    continuity.

179



**Fig. 2.** Specific steps for simulating VE$_f$ based on ML in our RF-PMRS method. The map used in the 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.

185

**2) VE$_f$ simulation based on ML**

187     The main idea is to establish an ML model between the $VE_f$ truth obtained from

188     multiple AERONET sites and related variables, thus improving the subsequent $VE_f$-

189     simulation accuracy (Fig. 2).

190

191     **Step 1** $VE_f$ calculation

192     The $VE_f$ true values are calculated concerning equations (3)-(5). A total of 9

193     AERONET sites corresponding to four typical aerosol types participate in the training.

194     Table 1 shows the specific information.

195

196     **Table 1**. Data information of 9 AERONET sites classified by aerosol types. Location indicates the

197     latitude and longitude, where '-' means the south latitude and west longitude. Two sites in bold fonts

198     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 | / |

199

200     **Step 2** $VE_f$-related variables selection

201     According to the theory, FMF is selected as the most important modeling variable.

202     Previous studies have also shown that the FMF-$VE_f$ relationship has a good single-

203     value correspondence, which is not affected by AOD. Compared with $AOD_f$ and

204     $V_{f,column}$, FMF is a better indicator for estimation (Zhang and Li 2015). In addition,

205     considering the spatiotemporal heterogeneity of $VE_f$, the latitude, longitude (LAT,

206     LON), and data time (month, day) of each site are added to the training.

207

208 **Step 3** RF model establishment

209   From step 2, $VE_f$ can be expressed as:

210
$$VE_f = f(FMF, LAT, LON, month, day)$$
(9)

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

212   decision trees that can build high-accuracy models based on fewer variables (Yang et

213   al., 2020). This ensemble supervised learning method randomly samples the original

214   dataset into multiple sets and considers random subsets of features in node splitting,

215   which reduces correlation and the sensitivity to noise (Belgiu and Drăguţ, 2016). Note

216   that the station FMF values ($S\text{-}FMF$) are used when training.

217

218   **Step 4** Accuracy validation

219   The $VE_f$ estimation is also based on equation (9), where $f$ is the optimal relationship

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

221   model results from point to surface. 10-fold cross-validation (Rodriguez et al., 2009)

222   and isolated-validation are used to evaluate model performance (see Appendix A1).

223

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

225   Then, calculate PM$_{2.5}$ according to the corresponding process (equation (6)). The

226   statistical indicators used in the evaluation include correlation coefficient (R), mean

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

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

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

230   on these indicators.

231

232   **3. Experiment data and results**

233   **3.1. Data**

234   **3.1.1. MODIS AOD**

235   MCD19A2, the MODIS C6 Level-2 gridded (L2G) land AOD product, is selected in

236   this study. It is derived by the Multi-Angle Implementation of Atmospheric Correction

237    (MAIAC) algorithm, which can improve the accuracy in cloud detection and aerosol

238    retrieval (Lyapustin et al., 2011). Besides, this new advanced algorithm jointly

239    combines MODIS Terra and Aqua into a single sensor (Lyapustin et al., 2014). The

240    product is produced daily with a 1km resolution, including aerosol parameters such as

241    470nm/550nm AOD, quality assurance (QA), and uncertainty factors.

242    The processing of MCD19A2 data (HDF format) is mainly divided into five steps:

243    AOD/QA band extraction, best quality AOD selection, Terra/Aqua data synthesis,

244    missing information reconstruction, and mosaic. Finally, the daily AOD distribution in

245    GeoTiff format is obtained.

246

247    **3.1.2. Phy-DL FMF dataset**

248    To enhance the reliability of the global land FMF product, Yan et al. (2022) have

249    released a satellite-based dataset (daily scale) called Phy-DL FMF, which integrates

250    physical and deep learning methods. The product has a spatial resolution of 1° and

251    covers from 2001 to 2020. In terms of performance, it exhibits higher accuracy and

252    wider space-time coverage than satellite products (Yan, 2021).

253

254    **3.1.3. Meteorological data**

255    The values of PBLH and RH are obtained from the ERA5 dataset. As the fifth-

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

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

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

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

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

261
$$RH = 100 \times \frac{e_s(T_d)}{e_s(T)} \qquad (10)$$

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

263    (Simmons et al., 1999).

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

264                                                                                              (11)

265

### 3.1.4. AERONET data

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

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

269 25 years (Holben et al., 1998). Until now, the Version 3 dataset has been released (Giles

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

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

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

273 Distribution products with Level 2.0 (quality-assured) are applied to implement our

274 purpose.

275

### 3.1.5. Ground PM$_{2.5}$ measurements

277 The near-surface hourly PM$_{2.5}$ values are obtained from the China National

278 Environmental Monitoring Center (CNEMC). Nowadays, over 1600 ground-based

279 monitors are working continuously and a total of 232 stations (2017) in the North China

280 Region (NC) participate in this work. Fig. 3 displays the site distributions of our

281 validation area.

282

**Fig. 3.** The location of ground stations in the NC region (35°-45°N, 110°-120°E). The red points represent NC stations.

The above variables are spatially matched to ground sites at their respective resolutions. And based on UTC, the experiment is conducted on a daily scale in 2017. Note that we select the measured empirical value of $\rho_{f,dry}$ (i.e., 1.5 g/cm$^3$) for the NC region from Gao et al. (2007).

**3.2. Experimental results**

**3.2.1. RF model performance for training VE$_f$**

The simulation model of VE$_f$ is trained based on the data in Table 1 and see Appendix A3 for the adjustment of the model parameters. Table 2 shows that RF can capture the complex relationship between VE$_f$ and related variables well. R is as high as 0.974 (0.975), RMSE and MAE are both small, and RPE is around 30%, which suggests the desired estimation accuracy. Overall, the CV results represent the great performance of

12

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

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

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

302

303    **Table 2**. Performance statistics of the RF model for training $VE_f$. N represents the number of data,

304    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 |

305

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

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

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

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

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

311    with their respective nearest ground $PM_{2.5}$ stations (distances of 3.64 km and 3.91 km,

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

313    (Shi et al., 2018)).

314    Fig. 4 displays the $PM_{2.5}$ value trends of different models at two sites. The blue line

315    fits the red line better than the gray one, confirming that the $PM_{2.5}$ results of RF-PMRS

316    are closer to the true values. Within the range of the black circles at positions 1 and 2,

317    the variation trend of RF-PMRS results has better consistency with the ground truth,

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

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

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

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

322    of the AOD high values. It is worth noting that our method has well mitigated the

323    apparent overestimation of the original model (PMRS) in the case of above-normal

324    aerosol loadings. Furthermore, the average $PM_{2.5}$ values from ground stations, PMRS,

325    and RF-PMRS are compared. As for the two sites, the RF-PMRS results are satisfactory.

326    As depicted in Fig. 5, the RF-PMRS and station mean values are close, with a difference

327    of 4.82 μg/m³ (BJ) and 2.73 μg/m³ (BC), suggesting a good estimation. Nevertheless,

328    the PMRS results have deviations greater than 40 μg/m³, and overestimation basically

329    exists at both sites. It can be inferred that, in our proposed method, the optimization of

330    $VE_f$ can greatly improve the $PM_{2.5}$ estimation accuracy.

331



332

**Fig. 4.** Three $PM_{2.5}$ trends at the Beijing (BJ) and Beijing-CAMS (BC) sites under their respective
valid DOYs in 2017. 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.

336

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

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.

### 3.2.3. Generalization performance of RF-PMRS

Then, we estimate PM$_{2.5}$ based on PMRS and RF-PMRS within North China (Fig. 3 exhibits the distribution pattern of the validation stations). Table 3 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.

16

372    **Table 3**. 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 |

373

374    Meanwhile, $PM_{2.5}$ scatterplots are presented below. As depicted in Fig. 7, there are

375    sufficient estimated samples (28305) in the NC region, which guarantees the credibility

376    of our validation results. In general, the RF-PMRS $PM_{2.5}$ values are distributed around

377    the true values evenly, with a slightly higher R of 0.70 compared to that of the original

378    method. And the slope of the linear fitting relationship reaches 0.82, which indicates

379    that the proposed method greatly reduces the overestimation of PMRS with a linear

380    slope of 1.46. Although the overall performance of the RF-PMRS estimations maintains

381    an excellent level, defects do remain. To be specific, in areas with high $PM_{2.5}$

382    concentration (especially greater than 150 µg/m³), RF-PMRS results exist a slight

383    underestimation. It may be caused by the relatively small number of high-value points

384    (only 1319 out of 28305), which is difficult to adequately reflect the fitting effect of the
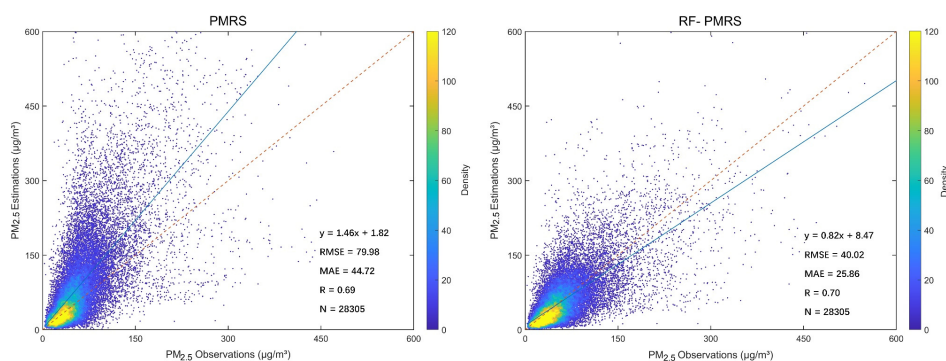
385    method.

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

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

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

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

390    small. With regard to the curve shape, it is high and narrow, manifesting that the bias

391    has a lower standard deviation (STD) and is more prone to appear around the mean.

392    However, PRMS shows a more discrete distribution pattern, and there are many outliers

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

394    the three boxes, within the bias range of ±20 µg/m³ and ±40 µg/m³, the data numbers of

395    RF-PMRS results increase by 8.32% and 12.81%, respectively. Outside the range of

396    ±100 µg/m³, the number decreases by 9.10%. Therefore, as far as the accuracy is

397    concerned, RF-PMRS results have lower bias and better stability.

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

399  polynomial relationship (equation (8)), the established RF model in RF-PMRS can

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

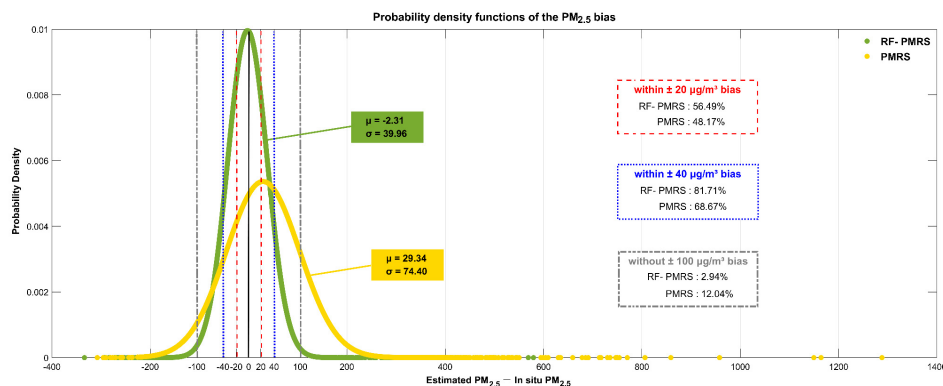401  improving the $PM_{2.5}$ estimation accuracy.

402



403
404  **Fig. 7.** Validation scatterplots of $PM_{2.5}$ results from PMRS (left) and RF-PMRS (right). Red dashed
405  lines are 1:1 reference lines, and blue solid lines stand for the linear fits. The right legends show the
406  point densities (frequency) represented by different colors.

407



408
409  **Fig. 8.** Probability density functions of PMRS (yellow) and RF-PMRS (green) $PM_{2.5}$ bias. The red,
410  blue and grey dotted lines indicate the bias boundaries of $\pm20$ µg/m³, $\pm40$ µg/m³, and $\pm100$ µg/m³,
411  respectively. µ and σ represent the mean value and standard deviation of each data.

412

413  **4. Discussion**

414  **4.1. Accuracy comparison of PMRS using MODIS/Phy-DL FMF**

415      To confirm the superiority of the Phy-FMF data adopted in our method framework,

416  taking the BJ and BC sites as examples, the experiment compares the $PM_{2.5}$ accuracy

417  and the number of effective days calculated by PMRS based on different FMF. Table 4

418  presents the overall day-level results. As can be seen, after the FMF replacement, the

419  valid DOY turns out to become more (an increase of 113 days), which illustrates that

420  the number of effective $PM_{2.5}$ concentration has gone up by about 5 times. Moreover,

421  the accuracy has been significantly enhanced, with R increased by about 0.30, RMSE

422  and MAE decreased by 26.14% and 16.47% accordingly. On the whole, Phy-DL FMF

423  contributes to the improvement of PMRS results, signifying the first step optimization

424  of the proposed RF-PMRS method is effective.

425

426  **Table 4**. Validation results of the PMRS method using different FMF data. The valid DOY refers to

427  the number of days that the AOD, FMF, and other data are not missing when calculating $PM_{2.5}$. Note

428  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 |

429

## 4.2. Performance compared with other ML models

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

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

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

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

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

436  models are presented in Table 5 and Table 6. By contrast, RF performs best in CV and

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

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

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

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

441

442 **Table 5**. Cross-validation results in comparison of the decision tree models for training $VE_f$. N
443 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 | |

444

445 **Table 6**. Isolated-validation results in comparison of the decision tree models for training $VE_f$. The
446 indicators are the same as those in Table 5.

| 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 | |

447

448 **5. Conclusion**

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

463

464 **Appendix A**

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

20

466    The sample-based 10-fold cross-validation method is applied to test the fitting and

467    predictive ability of our model. The original dataset is randomly divided into ten parts,

468    nine of which are used as the training set for model fitting, and the remaining one is

469    used for prediction, then the cross-validation process is repeated ten rounds until each

470    data has been used as the test set.

471    At the same time, when verifying the RF-based $VE_f$ model, the dataset in the time

472    period that did not participate in the training in Table 1 is used for isolated-validation.

473

474    **A2. Statistical indicators**

475
$$R = \frac{\sum_{i=1}^{m}(y_i - \bar{y})\sum_{i=1}^{m}(f_i - \bar{f})}{\sqrt{\sum_{i=1}^{m}(y_i - \bar{y})^2}\sqrt{\sum_{i=1}^{m}(f_i - \bar{f})^2}}$$

476
$$MB = \bar{y} - \bar{f}$$

477
$$RMB = \text{abs}\left(\frac{\bar{y} - \bar{f}}{\bar{y}}\right)$$

478
$$RMSE = \sqrt{\frac{1}{m}\sum_{i=1}^{m}(y_i - f_i)^2}$$

479
$$MAE = \frac{1}{m}\sum_{i=1}^{m}|y_i - f_i|$$

480
$$RPE = \frac{\sqrt{\frac{1}{m}\sum_{i=1}^{m}(y_i - f_i)^2}}{\bar{y}}$$

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

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

483    averages of all observations and estimates, respectively.

484

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

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

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

488 when splitting, (d) minimum number of split samples. Experiments shows that the

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

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

491 improving training efficiency.

492

493 **Code and data availability**

494 All relevant codes as well as the intermediate data of this work are archived at

495 https://doi.org/10.5281/zenodo.7183822 (Jin, 2022). The MCD19A2 data can be

496 downloaded on https://ladsweb.modaps.eosdis.nasa.gov (last access: 30-09-2022)

497 (Lyapustin and Wang, 2015). Detailed information about Phy-DL FMF dataset can be

498 found at https://doi.org/10.5281/zenodo.5105617 (Yan, 2021). Meteorological data

499 used in this work are obtained on

500 https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels (lass

501 access: 30-09-2022) (Hersbach et al., 2018). AERONET data was downloaded from

502 https://aeronet.gsfc.nasa.gov/ (last access: 30-09-2022) (Giles et al., 2019).

503

504 **Author contributions**

505 **Caiyi Jin:** Data curation, Methodology, Formal analysis, Writing - original draft.

506 **Qiangqiang Yuan:** Conceptualization, Supervision, Project administration, Writing -

507 review and editing. **Tongwen Li:** Resources, Methodology, Writing - review and

508 editing, Formal analysis. **Yuan Wang:** Methodology, Validation, Writing - review and

509 editing. **Liangpei Zhang:** Supervision, Writing - review and editing.

510

511 **Competing interests**

512 The contact author has declared that none of the authors has any competing interests.

513

514 **Acknowledgments**

515 We gratefully acknowledge the Atmosphere Archive and Distribution System

516 (LAADS), the ECMWF, the AERONET project, and the CNEMC for respectively

517  providing the MODIS products, the meteorological data, the ground aerosol data, and

518  the surface PM2.5 concentration. We also thank other institutions which provide related

519  data in this work.

520

525

526  **References**

527  Belgiu, M., and Drăguţ, L.: Random forest in remote sensing: A review of applications

528  and future directions, ISPRS J. Photogramm. Remote Sens., 114, 24-31,

529  https://doi.org/10.1016/j.isprsjprs.2016.01.011, 2016.

530  Bowe, B., Xie, Y., Li, T., Yan, Y., Xian, H., and Al-Aly, Z.: The 2016 global and

531  national burden of diabetes mellitus attributable to PM2.5 air pollution, Lancet Planet.

532  Health, 2, e301-e312, https://doi.org/10.1016/S2542-5196(18)30140-2, 2018.

533  Chen, X., de Leeuw, G., Arola, A., Liu, S., Liu, Y., Li, Z., and Zhang, K.: Joint retrieval

534  of the aerosol fine mode fraction and optical depth using MODIS spectral reflectance

535  over northern and eastern China: Artificial neural network method, Remote Sens

536  Environ, 249, 112006, https://doi.org/10.1016/j.rse.2020.112006, 2020.

537  Friedman, J.H.: Greedy function approximation: a gradient boosting machine, Ann Stat,

538  29(5), 1189–1232, http://www.jstor.org/stable/2699986, 2001.

539  Gao, J., Zhou, Y., Wang, J., Wang, T., and Wang, W.X.: Inter-comparison of WPSTM-

540  TEOMTM-MOUDITM and investigation on particle density, Huan Jing Ke Xue, 28,

541  1929-1934, https://doi.org/10.3321/j.issn:0250-3301.2007.09.005, 2007.

542  Gao, L., Li, J., Chen, L., Zhang, L., and Heidinger, A.K.: Retrieval and validation of

543  atmospheric aerosol optical depth from AVHRR over China, IEEE Trans Geosci

544  Remote Sens, 54, 6280-6291, https://doi.org/10.1109/TGRS.2016.2574756, 2016.

545  Geng, G., Zhang, Q., Martin, R.V., van Donkelaar, A., Huo, H., Che, H., Lin, J., and

546  He, K.: Estimating long-term PM2.5 concentrations in China using satellite-based

547  aerosol optical depth and a chemical transport model, Remote Sens Environ, 166, 262-

548  270, https://doi.org/10.1016/j.rse.2015.05.016, 2015.

549  Geurts, P., Ernst, D., and Wehenkel, L.: Extremely randomized trees, Mach Learn, 63,

550  3-42, https://doi.org/10.1007/s10994-006-6226-1, 2006.

551  Giles, D.M., Holben, B.N., Eck, T.F., Smirnov, A., Sinyuk, A., Schafer, J., Sorokin,

552  M.G., and Slutsker, I.: Aerosol robotic network (AERONET) version 3 aerosol optical

553  depth and inversion products, in: American Geophysical Union (AGU) 98th Fall

554  Meeting Abstracts, New Orleans, America, 11-15 December 2017, A11O-01, 2017.

555  Giles, D. M., Sinyuk, A., Sorokin, M. G., Schafer, J. S., Smirnov, A., Slutsker, I., Eck,

556  T. F., Holben, B. N., Lewis, J. R., Campbell, J. R., Welton, E. J., Korkin, S. V., and

557  Lyapustin, A. I.: Advancements in the Aerosol Robotic Network (AERONET) Version

558  3 database - automated near-real-time quality control algorithm with improved cloud

559  screening for Sun photometer aerosol optical depth (AOD) measurements, Atmos Meas

560  Tech, 12, 169–209, https://doi.org/10.5194/amt-12-169-2019, 2019.

561  Gupta, P., and Christopher, S.A.: Particulate matter air quality assessment using

562  integrated surface, satellite, and meteorological products: Multiple regression approach,

563  J. Geophys. Res. Atmos., 114, D14205, https://doi.org/10.1029/2008JD011496, 2009.

564  Hand, J.L., and Kreidenweis, S.M.: A new method for retrieving particle refractive

565  index and effective density from aerosol size distribution data, Aerosol Sci Technol, 36,

566  1012-1026, https://doi.org/10.1080/02786820290092276, 2002.

567  Hänel, G., and Thudium, J.: Mean bulk densities of samples of dry atmospheric aerosol

568  particles: A summary of measured data, Pure Appl. Geophys., 115, 799-803,

569  https://doi.org/10.1007/BF00881211, 1977.

570  Hersbach, H., Bell, B., Berrisford, P., Biavati, G., Horányi, A., Muñoz Sabater, J.,

571  Nicolas, J., Peubey, C., Radu, R., Rozum, I., Schepers, D., Simmons, A., Soci, C., Dee,

572  D., Thépaut, J-N.: ERA5 hourly data on single levels from 1979 to present, Copernicus

573  Climate Change Service (C3S) Climate Data Store (CDS) [data set], (Accessed on 30-

574  09-2022), https://doi.org/10.24381/cds.adbb2d47, 2018.

575  Holben, B.N., Eck, T.F., Slutsker, I., Tanré, D., Buis, J.P., Setzer, A., Vermote, E.,

576 Reagan, J.A., Kaufman, Y.J., Nakajima, T., Lavenu, F., Jankowiak, I., and Smirnov, A.:

577 AERONET — A federated instrument network and data archive for aerosol

578 characterization, Remote Sens Environ, 66, 1-16, https://doi.org/10.1016/S0034-

579 4257(98)00031-5, 1998.

580 Irrgang, C., Boers, N., Sonnewald, M., Barnes, E.A., Kadow, C., Staneva, J., and

581 Saynisch-Wagner, J.: Towards neural Earth system modelling by integrating artificial

582 intelligence in Earth system science, Nat. Mach. Intell., 3, 667-674,

583 https://doi.org/10.1038/s42256-021-00374-3, 2021.

584 Jin, C.: An optimized semi-empirical physical approach for satellite-based PM2.5

585 retrieval: using random forest model to simulate the complex parameter, Zenodo [code],

586 https://doi.org/10.5281/zenodo.7183822, 2022.

587 Koelemeijer, R.B.A., Homan, C.D., and Matthijsen, J.: Comparison of spatial and

588 temporal variations of aerosol optical thickness and particulate matter over Europe,

589 Atmospheric Environ., 40, 5304-5315, https://doi.org/10.1016/j.atmosenv.2006.04.044,

590 2006.

591 Kokhanovsky, A.A., Prikhach, A.S., Katsev, I.L., and Zege, E.P.: Determination of

592 particulate matter vertical columns using satellite observations, Atmos Meas Tech, 2,

593 327-335, https://doi.org/10.5194/amt-2-327-2009, 2009.

594 Lee, J.-B., Lee, J.-B., Koo, Y.-S., Kwon, H.-Y., Choi, M.-H., Park, H.-J., and Lee, D.-

595 G.: Development of a deep neural network for predicting 6 h average PM2.5

596 concentrations up to 2 subsequent days using various training data, Geosci. Model Dev.,

597 15, 3797–3813, https://doi.org/10.5194/gmd-15-3797-2022, 2022.

598 Li, T., Shen, H., Zeng, C., Yuan, Q., and Zhang, L.: Point-surface fusion of station

599 measurements and satellite observations for mapping PM2.5 distribution in China:

600 Methods and assessment, Atmospheric Environ., 152, 477-489,

601 https://doi.org/10.1016/j.atmosenv.2017.01.004, 2017.

602 Li, Z., Zhang, Y., Shao, J., Li, B., Hong, J., Liu, D., Li, D., Wei, P., Li, W., Li, L.,

603 Zhang, F., Guo, J., Deng, Q., Wang, B., Cui, C., Zhang, W., Wang, Z., Lv, Y., Xu, H.,

604 Chen, X., Li, L., and Qie, L.: Remote sensing of atmospheric particulate mass of dry

605 PM2.5 near the ground: Method validation using ground-based measurements, Remote

25

606  Sens Environ, 173, 59-68, https://doi.org/10.1016/j.rse.2015.11.019, 2016.

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

608  Reid, J.S.: Multiangle implementation of atmospheric correction (MAIAC): 2. Aerosol

609  algorithm, J. Geophys. Res. Atmos., 116, D03211,

610  https://doi.org/10.1029/2010JD014986, 2011.

611  Lyapustin, A., Wang, Y., Xiong, X., Meister, G., Platnick, S., Levy, R., Franz, B.,

612  Korkin, S., Hilker, T., Tucker, J., Hall, F., Sellers, P., Wu, A., and Angal, A.: Scientific

613  impact of MODIS C5 calibration degradation and C6+ improvements, Atmos Meas

614  Tech, 7, 4353-4365, https://doi.org/10.5194/amt-7-4353-2014, 2014.

615  Lyapustin, A., andWang, Y.: MCD19A2 MODIS/Terra+Aqua Aerosol Optical

616  Thickness Daily L2G Global 1km SIN Grid, NASA LP DAAC [data set], (Accessed

617  on 30-09-2022), http://doi.org/10.5067/MODIS/MCD19A2.006, 2015.

618  Lyu, B., Huang, R., Wang, X., Wang, W., and Hu, Y.: Deep-learning spatial principles

619  from deterministic chemical transport models for chemical reanalysis: an application in

620  China for PM2.5, Geosci. Model Dev., 15, 1583–1594, https://doi.org/10.5194/gmd-

621  15-1583-2022, 2022.

622  Ma, Z., Hu, X., Huang, L., Bi, J., and Liu, Y.: Estimating ground-Level PM2.5 in China

623  using satellite remote sensing, Environ. Sci. Technol., 48, 7436-7444,

624  https://doi.org/10.1021/es5009399, 2014.

625  Pope III, C.A., Burnett, R.T., Thun, M.J., Calle, E.E., Krewski, D., Ito, K., and Thurston,

626  G.D.: Lung cancer, cardiopulmonary mortality, and long-term exposure to fine

627  particulate air pollution, JAMA, 287, 1132-1141,

628  https://doi.org/10.1001/jama.287.9.1132, 2002.

629  Raut, J., and Chazette, P.: Assessment of vertically-resolved PM10 from mobile lidar

630  observations, Atmospheric Chem. Phys., 9, 8617-8638, https://doi.org/10.5194/acp-9-

631  8617-2009, 2009.

632  Rodriguez, J.D., Perez, A., and Lozano, J.A.: Sensitivity analysis of k-fold cross

633  validation in prediction error estimation, IEEE Trans. Pattern Anal. Mach. Intell., 32,

634  569-575, https://doi.org/10.1109/TPAMI.2009.187, 2009.

635  Shi, X., Zhao, C., Jiang, J.H., Wang, C., Yang, X., and Yung, Y.L.: Spatial

636 representativeness of PM2.5 concentrations obtained using observations from network

637 stations, J. Geophys. Res. Atmos., 123, 3145-3158,

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

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

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

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

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

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

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

645 111, D21201, https://doi.org/10.1029/2005JD006996, 2006.

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

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

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

649 280-301, https://doi.org/10.1016/j.atmosenv.2018.12.023, 2019.

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

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

652 1511, https://doi.org/10.5194/gmd-13-1499-2020, 2020.

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

654 2067, https://doi.org/10.1016/S0140-6736(13)62693-8, 2013.

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

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

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

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

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

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

661 2021.

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

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

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

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

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

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

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

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

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

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

672    Satellite remote sensing of atmospheric particulate matter mass concentration:

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

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

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

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

677    2025, Remote Sens., 12(16), 2518, https://doi.org/10.3390/rs12162518, 2020.

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

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

680    160, 252-262, https://doi.org/10.1016/j.rse.2015.02.005, 2015.