

1 Estimation of PM_{2.5} Concentration in China Using 2 Linear Hybrid Machine Learning Model

3 Zhihao Song¹, Bin Chen¹, Yue Huang¹, Li Dong¹, Tingting Yang²

4 ¹Atmospheric Science College of Lanzhou University, Lanzhou 730000, China

5 ²Gansu Seed General Station, Lanzhou 730030, China

6 *Correspondence to:* Bin Chen (chenbin@lzu.edu.cn)

7 **Abstract.** The satellite remote-sensing aerosol optical depth (AOD) and meteorological elements
8 were employed to invert PM_{2.5} (The fine particulate matter with a diameter below 2.5μm) in order to
9 control air pollution more effectively. This paper proposes a restricted gradient-descent linear hybrid
10 machine learning model (RGD-LHMLM) by integrating a random forest (RF), a gradient boosting
11 regression tree (GBRT), and a deep neural network (DNN) to estimate the concentration of PM_{2.5} in
12 China in 2019. The research data included Himawari-8 AOD with high spatiotemporal resolution,
13 ERA-5 meteorological data, and geographic information. The results showed that, in the hybrid model
14 developed by linear fitting, the DNN accounted for the largest proportion, whereas the weight
15 coefficient was 0.62. The R² values of RF, GBRT, and DNN were reported 0.79, 0.81, and 0.8,
16 respectively. Preferably, the generalization ability of the mixed model was better than that of each
17 sub-model, and R² (determination coefficient) reached 0.84, whereas RMSE (root mean square error)
18 and MAE (mean Absolute Error) were reported 12.92 μg/m³ and 8.01 μg/m³, respectively. For the
19 RGD-LHMLM, R² was above 0.7 in more than 70% of the sites, whereas RMSE and MAE were
20 below 20 μg/m³ and 15 μg/m³, respectively, in more than 70% of the sites due to the correlation
21 coefficient having seasonal difference between the meteorological factor and PM_{2.5}. Furthermore, the
22 hybrid model performed best in winter (mean R² was 0.84) and worst in summer (mean R² was 0.71).
23 The spatiotemporal distribution characteristics of PM_{2.5} in China were then estimated and analyzed.
24 According to the results, there was severe pollution in winter with an average concentration of PM_{2.5}
25 being reported 62.10 μg/m³. However, there was slight pollution in summer with an average
26 concentration of PM_{2.5} being reported 47.39 μg/m³. The period from 10:00 to 15:00 every day is the
27 best time for model inversion, also at this time the pollution is high. The findings also indicate that
28 North China and East China are more polluted than other areas and that their average annual
29 concentration of PM_{2.5} was reported 82.68 μg/m³. Moreover, there was relatively low pollution in

1 Inner Mongolia, Qinghai, and Tibet, for their average PM_{2.5} concentrations were reported below 40
2 µg/m³.

3 **1 Background**

4 In recent years, pollutants have been discharged increasingly in China where air pollution is
5 becoming worse than ever before due to rapid urbanization and industrialization (Wang et al., 2019a).
6 The fine particulate matter (PM_{2.5}) with a diameter below 2.5µm is the main component of air pollutants
7 having considerable impacts on human health, atmospheric visibility, and climate change (Gao et al.,
8 2015; Pan et al., 2018; Pun et al., 2017; Qin et al., 2017). The global concern about PM_{2.5} has increased
9 significantly since it was listed as a top carcinogen (Apte et al., 2015; Lim et al., 2020). Currently, ground
10 monitoring is the most efficient method of measuring PM_{2.5} (Yang et al., 2018). However, monitoring
11 stations are not evenly distributed due to terrain and construction costs; therefore, it is difficult to obtain
12 a wide range of accurate PM_{2.5} concentration data (Han et al., 2015). To solve the problem, the method
13 of estimating PM_{2.5} with satellite remote-sensing was developed. Satellite remote-sensing is
14 characterized by a wide coverage and high resolution (Hoff and Christopher, 2009; Xu et al., 2021).
15 There is also a high correlation between AOD, obtained from satellite remote sensing inversion, and
16 PM_{2.5}; therefore, AOD is a very effective method of monitoring the spatiotemporal concentration
17 characteristics of PM_{2.5}.

18 After Engel-Cox et al. (2004) proposed using satellite AOD to estimate PM_{2.5} concentration, several
19 studies are reported in the literature to address this theory. Based on the regression model, Liu et al. (2005)
20 introduced AOD, boundary layer height, relative humidity, and geographical parameters as the main
21 controlling factors to estimate PM_{2.5} in the eastern part of the United States, and the verification
22 coefficient R² obtained was 0.46. Tian and Chen (2010) used AOD, PM_{2.5}, and meteorological parameters
23 in Southern Ontario, Canada, to establish a semi-empirical model to predict PM_{2.5} concentration per hour,
24 and the verification coefficient R² obtained in rural and urban areas was 0.7 and 0.64, respectively. Hu et
25 al. (2013) proposed a geography weighted regression model to estimate the surface PM_{2.5} concentration
26 in southeastern America by combining AOD, meteorological parameters, and land use information. Their
27 model average R² was 0.6. Lee et al. (2012) believed that the satellite remote sensing AOD data would
28 be interfered by clouds and snow and ice, and the reliability of the data was questionable. They proposed

1 a mixed model based on AOD calibration to predict the ground $PM_{2.5}$ concentration in New England,
2 USA, and achieved good results ($R^2 = 0.83$). Li et al. (2016) used PMRS method to remote sensing
3 ground $PM_{2.5}$. Combined with MODIS (Moderate-resolution Imaging Spectroradiometer) AOD and
4 ground observation data, Lv et al. (2017) estimated the daily surface $PM_{2.5}$ concentration in the Beijing-
5 Tianjin-Hebei region and improved the data resolution to 4 km. Using interpretable self-adaptive deep
6 neural network, Chen et al. (2021) estimated daily spatially-continuous $PM_{2.5}$ concentrations across
7 China, and analyzed the contribution of various characteristics to the $PM_{2.5}$ model. The data used in
8 these early studies are AOD products obtained from polar-orbit satellite sensors. The daily observation
9 frequency is limited. Due to the influence of cloud and ground reflection, the dynamic change
10 information of $PM_{2.5}$ cannot be obtained. As a result, geostationary satellite observations can be used to
11 overcome the problem of low temporal resolution for estimating surface $PM_{2.5}$ (Emili et al., 2010).

12 The Himawari-8 satellite commonly used in the Asia-Pacific region is a geostationary satellite
13 launched by the Japan Meteorological Agency in 2014. The observation frequency is 10 minutes, and the
14 observation results can characterize the aerosol and provide AOD data with a resolution of 5 km (Bessho
15 et al., 2016; Yumimoto et al., 2016). Due to its excellent performance, Wei et al. (2021a) use Himawari-
16 8 data to estimate ground $PM_{2.5}$, result shows that the $CV-R^2$ (cross-validation coefficient of
17 determination) is 0.85, with a root-mean-square error (RMSE) and mean absolute error (MAE) of 13.62
18 and $8.49 \mu g/m^3$, respectively. Wang et al. (2017) proposed an improved linear model, introduced AOD,
19 meteorological parameters, geographic information to estimate $PM_{2.5}$ in the Beijing-Tianjin-Hebei region,
20 and the verification coefficient R^2 was 0.86. Zhang et al. (2019b) used Himawari-8 hourly AOD product
21 to estimate ground $PM_{2.5}$ in China's four major urban agglomerations. The results showed significant
22 diurnal, seasonal, and spatial changes and improved the temporal resolution of estimating $PM_{2.5}$
23 concentration to the hourly level. Yin et al. (2021) used Himawari-8 hourly TOAR (top-of-the-
24 atmosphere reflectance) data to estimate ground $PM_{2.5}$ in China, improved data coverage area.

25 As research into ground-based $PM_{2.5}$ estimation deepens, traditional linear or nonlinear models
26 cannot meet the requirements of large-scale estimation and are gradually being replaced by machine
27 learning algorithms with strong nonlinear fitting ability (Guo et al., 2021; Mao et al., 2021). Liu et al.
28 (2018) combined Kriging interpolation and random forest algorithm to obtain the concentration of high-
29 resolution ground $PM_{2.5}$ in the United States. To demonstrate the accuracy and superiority of the proposed
30 method, the results were compared with the $PM_{2.5}$ concentration in ground measurement stations. Chen

1 et al. (2019) stacked and predicted $PM_{2.5}$ concentration based on a variety of machine learning algorithms,
2 discussed the influence of meteorological factors on $PM_{2.5}$ and achieved an $R^2 = 0.85$. Li et al. (2017a)
3 established a GRNN (Generalized regression neural networks) model for the whole of China to estimate
4 $PM_{2.5}$ concentration, and the results demonstrated that the performance of the deep learning model was
5 better than that of the traditional linear model. In addition, there are some novel algorithms such as space-
6 time extra-trees (STET) (Wei et al., 2021b) and space-time random forest (STRF) (Wei et al., 2019a) that
7 are also used for $PM_{2.5}$ inversion research.

8 A large number of existing studies in the broader literature have examined the estimation of ground
9 $PM_{2.5}$ concentrations using satellite remote sensing AOD. However, the performance of $PM_{2.5}$ estimation
10 models established in the existing studies varies greatly and the performance of the models is not stable
11 in different seasons and regions. To overcome this limitation, in this paper, a linear hybrid machine
12 learning model (RGD-LHMLM) based on random forest (RF), gradient lifting regression tree (GBRT),
13 and deep neural network (DNN) is proposed to estimate ground $PM_{2.5}$ concentration. The model
14 performance is evaluated from time and space to analyze its causes. Finally, spatiotemporal distribution
15 of $PM_{2.5}$ concentration in China in 2019 is obtained.

16 **2 Data**

17 **2.1 Ground $PM_{2.5}$ Monitoring Data**

18 $PM_{2.5}$ concentration data for 2019 used in this study are available from the China Environmental
19 Monitoring Center's Air Quality Real-Time Publication System. The $PM_{2.5}$ datasets are calibrated and
20 quality-controlled according to national standards GB 3095-2012 (China's National Ambient air quality
21 standards)(China, 2012).The system extracts hourly mean $PM_{2.5}$ data. By the end of 2019, China had
22 1641 monitoring stations built and in operation. Figure 1 shows the spatial distribution of monitoring
23 stations in China.

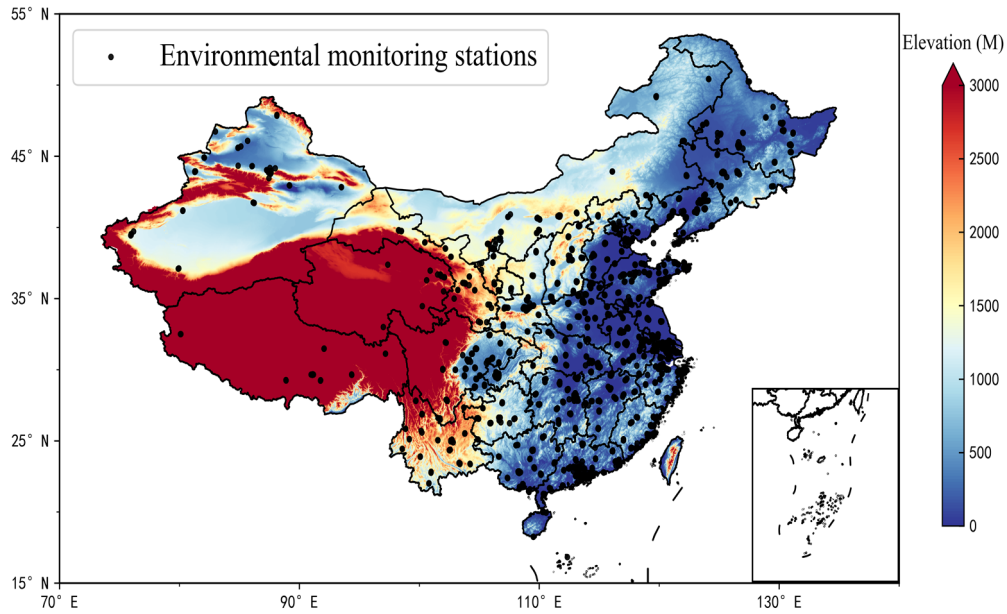


Figure 1 Distribution diagram of Environmental monitoring stations in China

2.2 Satellite AOD Data

The Advance Himawari Imager (AHI) on the Himawari-8 satellite launched by the Japan Meteorological Agency is a highly improved multi-wavelength imager. It adopts the whole disk observation method and has 16 visible and infrared channels. It has the characteristics of fast imaging speed, flexible observation area, and time. Himawari-8 AOD is obtained by an aerosol retrieval algorithm based on Lambertian-surface-assumed developed by Yoshida et al. (2018). The Level-3-hour AOD product, released by the Japan Aerospace Space Agency (JAXA), provides 500 nm AOD data with a spatial resolution of 5km during the day. In previous studies (Zang et al., 2018), Himawari-8 AOD was compared with the AOD data of AERONET (Aerosol Robotic Network) in China and achieved good performance (Zhang et al., 2019c), so that the results show that they are consistent ($R^2=0.75$), RMSE and MAE were achieved 0.39 and 0.21, respectively(Wei et al., 2019b). The AOD data used in this study is the Himawari-8 Level 3-hour AOD data in 2019 obtained from the Himawari Monitor website of the Japan Meteorological Agency. In the study, we selected AOD with strict cloud screening, that is, AOD data with low uncertainty.

2.3 Meteorological Data

ERA-5 reanalysis data is an hourly collection of atmospheric and land-surface meteorological elements since 1979 that the European Centre (ECMWF) has used its prediction model and data

1 assimilation system to "Reanalyse" archived observations(Jiang et al., 2021). Data used in this paper
2 include surface relative humidity (RH, expressed as a percentage), air temperature at a height of 2 m
3 (TM, expressed as K), Wind speed (U10, V10, in m/s), surface pressure (SP, in Pa), boundary layer height
4 (BLH, in m) and cumulative precipitation (RAIN, in m) at 10 m above the ground. A series of studies
5 has indicated that these parameters can affect the concentration of PM_{2.5} (Fang et al., 2016; Guo et al.,
6 2017; Li et al., 2017b; Wang et al., 2019b; Zheng et al., 2017; Gui et al., 2019). Uncertainty estimation
7 of ERA5 data has described in detail in the following website:

8 <https://confluence.ecmwf.int/display/CKB/ERA5%3A+uncertainty+estimation>.

9 **2.4 Auxiliary Data**

10 The auxiliary data used in this study include high and low vegetation index (LH, LL), ground
11 elevation data (DEM), and population density data (PD). The high and low vegetation index is derived
12 from ERA5 reanalysis data, which respectively represent half of the total green leaf area per unit level
13 ground area of high and low vegetation type. The ground elevation data are derived from SRTM-3
14 measurements jointly conducted by NASA and the Defense Department's National Mapping Agency
15 (NIMA), with a spatial resolution of 90 m. The population data come from the 2015 United Nations
16 Adjust Population Density data provided by NASA's Center for Socio-Economic Data and Applications
17 (SEDAC), which is based on national censuses and adjusted for relative spatial distribution.

18 **3 Method**

19 **3.1 Random Forest**

20 Random Forest (RF) is built based on the combination of the Bagging algorithm and decision
21 tree(Breiman, 2001), which is an extended variant of the parallel ensemble learning method (Stafoggia
22 et al., 2019). To construct a large number of decision trees, the random forest model takes multiple
23 samples of the sample data. In the decision tree, the nodes are divided into sub-nodes by using the
24 randomly selected optimal features until all the training samples of the node belong to the same class.
25 Finally, all the decision trees are merged to form the random forest. This method has proved to be
26 effective in regression and classification problems and is one of the most well-known Machine learning
27 algorithms used in many different fields (Yesilkanat, 2020).

1 3.2 Gradient Boosted Regression Trees

2 Different from the random forest, Gradient Boosting Regression Tree (GBRT) is based on Boosting
3 algorithm and decision tree(Friedman, 2001). The basic principle of GBRT is to construct M different
4 basic learners through multiple iterations, and constantly add the weight of the learners with a small error
5 probability, to eventually generate a strong learner (Johnson et al., 2018). The core of this method is that
6 after each iteration, a learner will be built in the direction of residual reduction (gradient direction) to
7 make the residual decrease in the gradient direction (Schonlau, 2005). The basic learner of GBRT is the
8 regression tree in the decision tree. During the prediction, a predicted value is calculated according to the
9 model obtained. The minimum square root error is used to select the optimal feature to split the dataset,
10 and the average value of the child node is then taken as the predicted value.

11 3.3 Deep Neural Networks

12 Deep Neural Networks (DNN) is a supervised learning technique that uses a backpropagation
13 algorithm to minimize the loss function. It adjusts the parameters through an optimizer, and has high
14 computational power, making it ideal for solving classification and regression problems (Wang and Sun,
15 2019). The structure of DNN includes an input layer, an output layer, and several hidden layers. Each
16 layer takes the output of all nodes of the previous layer as the input, and this process requires activation
17 functions. Compared with other activation functions, the linear rectifying function (ReLU) has the
18 advantages of simple derivation, faster convergence, and higher efficiency. At the same time, among the
19 adaptive learning rate optimizers, the Adamx optimizer performs the best. It not only has the advantages
20 of Adam in determining the learning rate range and having stable parameters in each iteration but also
21 simplifies the method of defining the upper limit range of the learning rate and improves the iteration
22 efficiency (Diederik and Jimmy, 2015). Therefore, in this paper, we selected the Adamx optimizer and
23 ReLU activation function to train the DNN.

24 3.4 Model Establishment and Verification

25 After data processing, RF, GBRT, and DNN are used for modeling.

$$\begin{aligned} 26 \quad PM_{2.5i,j} = & AOD_{i,j} + BLH_{i,j} + RH_{i,j} + TM_{i,j} + LL_{i,j} + LH_{i,j} + SP_{i,j} \\ 27 \quad & + RAIN_{i,j} + U_{10i,j} + V_{10i,j} + PD_{i,j} + HEIGHT_{i,j} + LON_{i,j} + LAT_{i,j} \\ 28 \quad & + MONTH_{i,j} + HOUR_{i,j} \end{aligned} \quad (1)$$

1 Formula (1) is applicable to RF, GBRT and DNN, where $PM_{2.5i,j}$ is the $PM_{2.5}$ at time i on station j .

2 To prevent model parameters from being controlled by large or small range data and speed up the
3 convergence rate of the model, the data must be normalized before starting the training process. Finally,
4 the three optimal sub-models are linear combined to achieve the final mixed model. To verify the model
5 performance, this paper uses the "10-fold cross-validation" method (Adams et al., 2020). In this method,
6 the data is split into 10 copies, 9 copies for training and 1 copy for verification; this process is repeated
7 10 times, and then the average of the 10 predictions is computed as the final result. Finally, the predicted
8 value and the measured value are fitted linearly. At the same time, several indicators are used to evaluate
9 the model, including the mean absolute error (MAE, when the predicted value and the true value are
10 exactly equal to 0, that is, perfect model; The larger the error, the greater the value), the root mean square
11 error (RMSE, when the predicted value and the real value are completely consistent is equal to 0, that is,
12 the perfect model; The larger the error, the greater the value), the slope of the fitting equation and the
13 determination coefficient R^2 (the greater the value, the better the model fitting effect), the bias (Bias, is
14 the difference between the predicted values and the true values, so that models with larger bias performed
15 worse), and the GME (generalization error of the bias, It is generally believed that bias should be
16 expressed as a square when using generalization error). The calculation formula of each indicator is
17 shown as follows:

$$18 \quad R^2 = 1 - \frac{SS_{res}}{SS_{tot}} \quad (2)$$

$$19 \quad MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (3)$$

$$20 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (4)$$

$$21 \quad Bias = \frac{\sum_{i=1}^N \hat{y}_i - y_i}{N} \quad (5)$$

$$22 \quad GEB = \frac{\sum_{i=1}^N (\hat{y}_i - y_i)^2}{N} \quad (6)$$

23 Where \hat{y}_i represents the predicted value, y_i shows the true value, SS_{res} denotes the error between
24 the regression data and the mean value, SS_{tot} represents the error between the real data and the mean
25 value, and the mean value is the mean value of the true value.

26 The research process is illustrated in Figure 2:

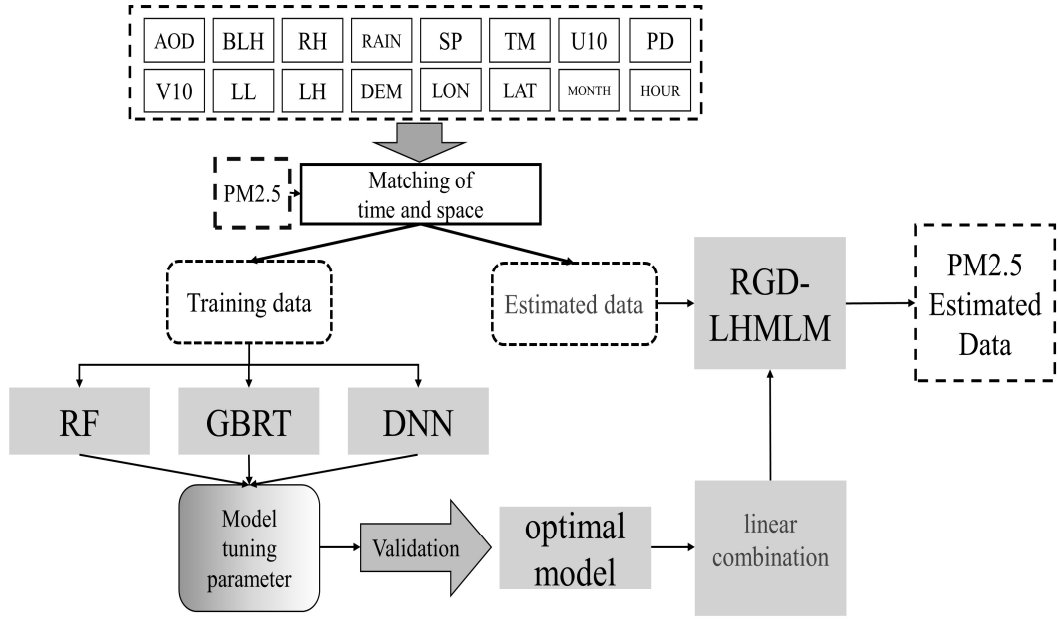


Figure 2 Schematic diagram of model

4 Results and Discussion

4.1 Modeling Results

According to the above steps, the mixed model RGD-LHMLM is obtained through modeling verification, and is compared with RF, GBRT, and DNN. The fitting and verification accuracy results of each model are shown in Table 1.

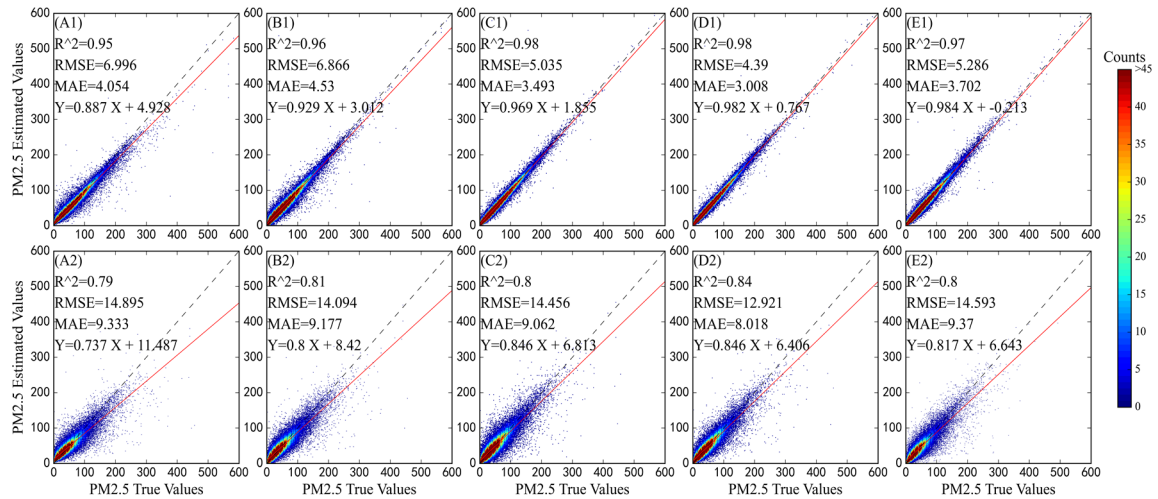
Table 1 Comparison of model accuracy

Model	Fitting				Validation			
	R ²	RMSE	MAE	GEB	R ²	RMSE	MAE	GEB
RF	0.95	6.99	4.05	114.19	0.79	14.89	9.33	208.97
GBRT	0.96	6.87	4.52	110.00	0.81	14.09	9.18	198.65
DNN	0.97	5.03	3.49	59.16	0.80	14.45	9.06	221.86
RGD-LHMLM	0.98	4.39	3.00	44.97	0.84	12.92	8.01	166.95

The PM_{2.5} inversion results of a single machine learning model show that DNN has the best inversion performance, followed by GBRT, and RF has the worst performance. The expression of the mixing model obtained after linear mixing is as follows:

$$PM_{2.5RGD-LHMLM} = 0.25PM_{2.5RF} + 0.17PM_{2.5GBRT} + 0.62PM_{2.5DNN} - 2.13 \quad (7)$$

1 The weight coefficient of DNN in the mixed model was the largest (0.62). The R^2 of RGD-LHMLM in
 2 the training set was 0.98, and the RMSE was only $4.39 \mu\text{g}/\text{m}^3$, indicating that the model had an excellent
 3 data fitting effect. Meanwhile, the generalization ability of the mixed model is also good, with R^2 of 0.84
 4 and RMSE of $12.92 \mu\text{g}/\text{m}^3$ on the validation data set. Among all the models, the deviation generalization
 5 error of the linear mixed model is also the lowest, indicating that the difference between the results
 6 obtained by this model and the real value is the least. Compared with RF, GBRT, and DNN, the inversion
 7 performance of RGD-LHMLM is improved. In other words, the combination of multiple models can
 8 improve the robustness and generalization ability of the model (Wolpert, 1992). The linear fitting
 9 equation coefficients between the predicted and measured values in the training set and the verification
 10 set were 0.98 and 0.84, respectively, indicating that the prediction accuracy of the model reached a high
 11 level. The fitting curve between the model predicted value and the real value is shown in Figure 3. The
 12 RGD-LHMLM model has the smallest degree of data dispersion, and the slope of the fitting line reaches
 13 0.84, indicating that 84% of the prediction results are accurate, higher than the three sub-models. The
 14 accuracy of the model decreased in the site-based validation, in which the R^2 and RMSE values are 0.8
 15 and $14.59 \mu\text{g}/\text{m}^3$, respectively.

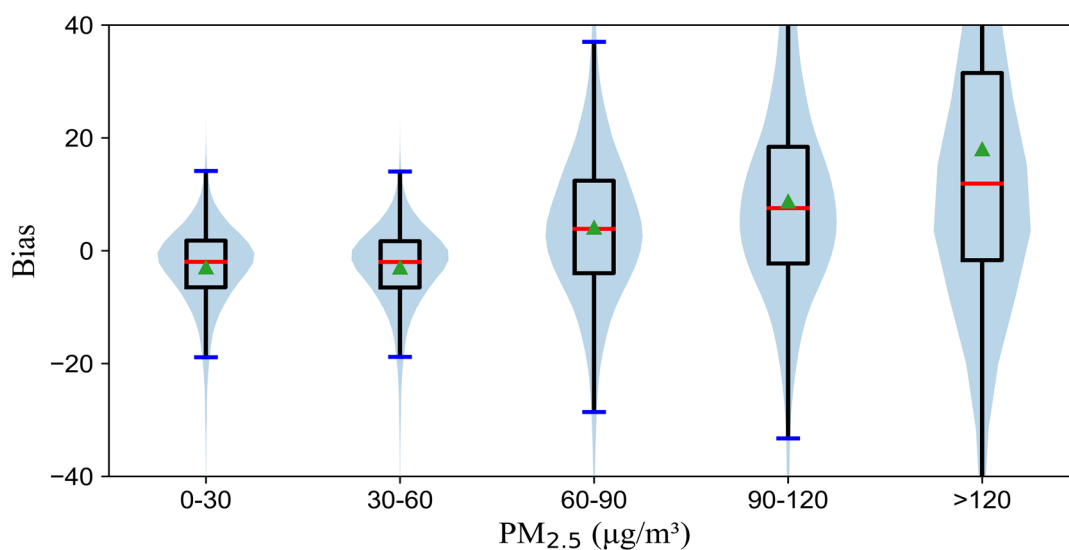


16
 17 **Figure 3 Accuracy of model Fitting (The first line) and Validation (The second line) (A: RF, B: GBRT, C:**
 18 **DNN, D: RGD-LHMLM (Based on sample), E: RGD-LHMLM (Based on site)). R^2 represents**
 19 **determination coefficient, RMSE represents root mean square error, MAE represents mean Absolute Error,**
 20 **N represents the number of samples. The equation Y and X represent the fitting relationship between the**
 21 **actual and estimated $\text{PM}_{2.5}$ values. Black dashed line represents 1:1 line, and red line represents best-fit line**
 22 **from linear regression.**

1 4.2 Model Performance Analysis

2 4.2.1 Bias analysis of Model

3 The average bias of the mixed model in different $PM_{2.5}$ concentration ranges was analyzed, and the
4 result is shown in figure 4. When the $PM_{2.5}$ concentration is less than $60 \mu\text{g}/\text{m}^3$, the average bias of the
5 model is less than 0. As the $PM_{2.5}$ concentration increases, the model deviation gradually increases. In
6 other words, when the $PM_{2.5}$ concentration is small, the predicted value of the model will generally
7 overestimate $PM_{2.5}$, and when the $PM_{2.5}$ further increases, it will underestimate the $PM_{2.5}$ concentration.

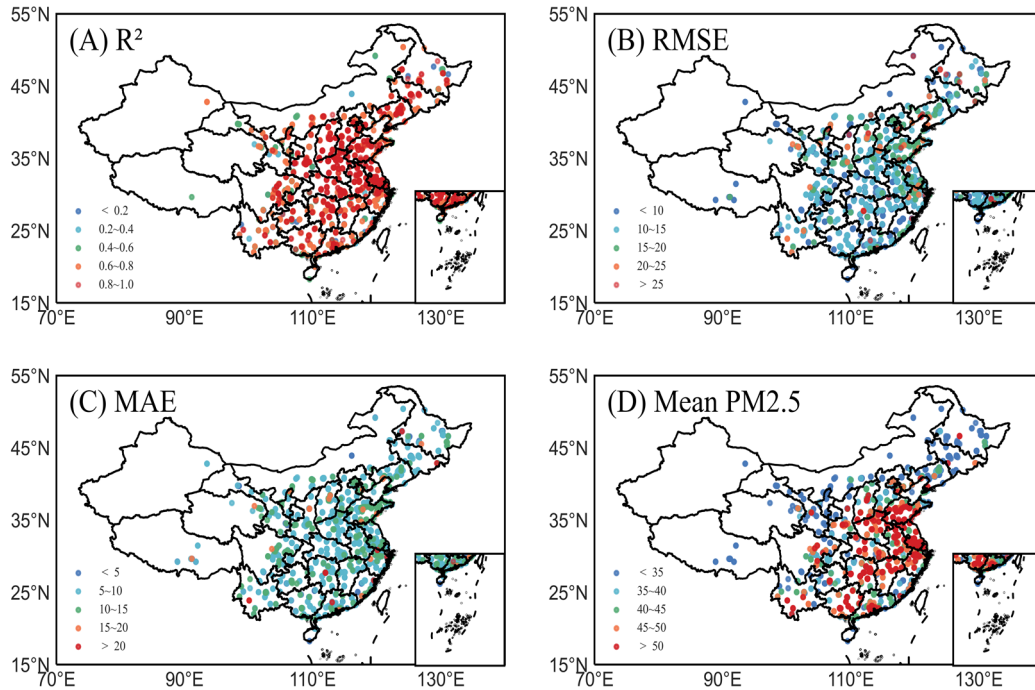


8

9 **Figure 4** Boxplots of resulting bias (y-axis) for different $PM_{2.5}$ concentration ranges in $\mu\text{g}/\text{m}^3$ (x-axis) (The
10 green arrow symbol, dark blue and red marks represent the average Bias, the median of Bias and the
11 extremum of Bias, respectively. Data density is represented by the light blue shading.)

12 4.2.2 Performance Analysis of Monitoring Station Model

13 The spatial performance of the model was analyzed by measuring R^2 , RMSE, and MAE at the
14 monitoring stations. According to Figure 5, there are regional differences in the inversion performance
15 of RGD-LHMLM. At all monitoring stations, the average R^2 was reported 0.74, and R^2 was above 0.7 at
16 more than 70% of the stations, especially in the densely populated and industrially developed areas. The
17 model prediction accuracy was reported low ($R^2 < 0.6$) in Xinjiang, Tibet, Qinghai, Western Sichuan, and
18 a few other areas of Northeast China. The mean values of RMSE and MAE were reported $11.4 \mu\text{g}/\text{m}^3$
19 and $8.01 \mu\text{g}/\text{m}^3$, respectively. In fact, the mean values of RMSE and MAE were below $20 \mu\text{g}/\text{m}^3$ and 15
20 $\mu\text{g}/\text{m}^3$ in more than 95% of stations, something showed a low estimation error.



1

2

Figure 5 Spatial distributions of model precision in terms of (A) determination coefficient (R^2), (B) root mean square error (RMSE), (C) mean Absolute Error (MAE) and (D) mean $PM_{2.5}$ concentration at each site in China. Color circles represent different value ranges of shown statistical parameters.

3

4

5

6

7

8

9

10

11

12

13

14

15

Based on the analysis of spatial differences in the RGD-LHMLM inversion performance, the following deductions can be made. First, the environmental monitoring stations in the central and eastern regions with better inversion performance were distributed densely, and there are large data available; therefore, the model had a satisfactory training effect. Moreover, data matching was lower in the western region than in other regions, something which resulted in model over-fitting and reduced accuracy (Zhang et al., 2018). Second, some areas of western and northeastern China are covered by snow and the Gobi Desert with high surface albedo. This reduces the accuracy of AOD obtained by satellite observation and brings errors to model training. Finally, the Himawari-8 scanning range is limited, and the satellite observation data obtained in Western China are limited in terms of quantity and accuracy. In general, the RGD-LHMLM has a satisfactory spatial performance, especially in areas with high annual average concentration of $PM_{2.5}$; therefore, it can leave a good inversion effect.

16

4.2.3 Time-Scale Model Performance Analysis

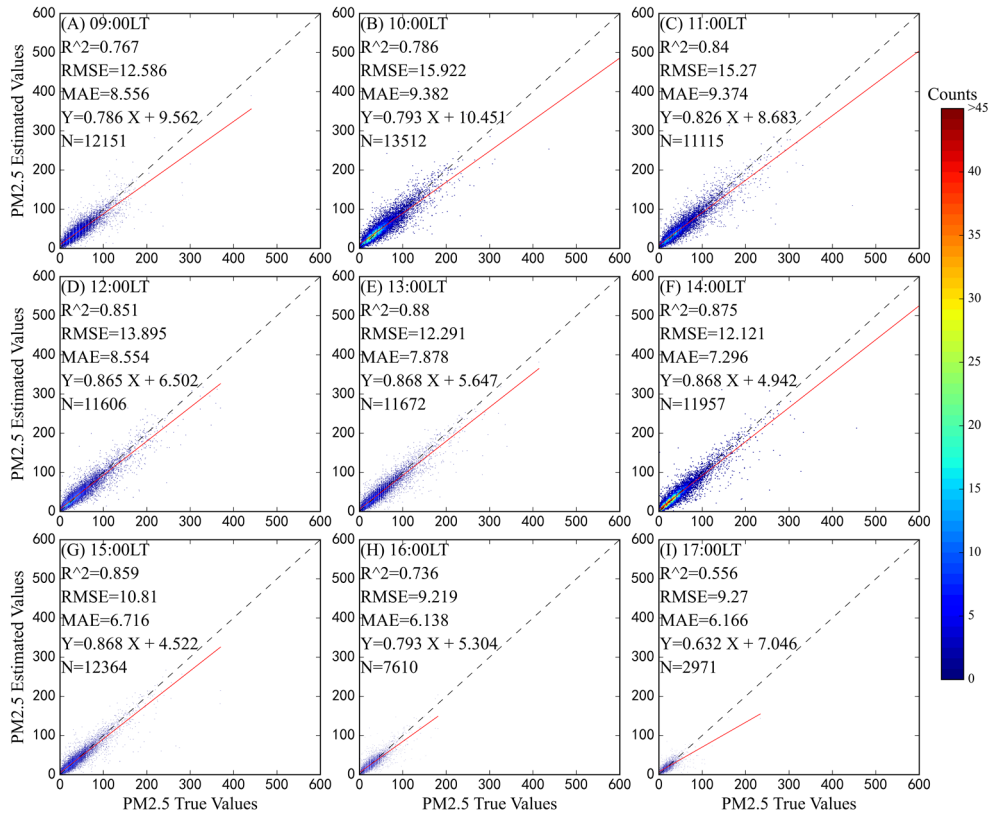
17

18

19

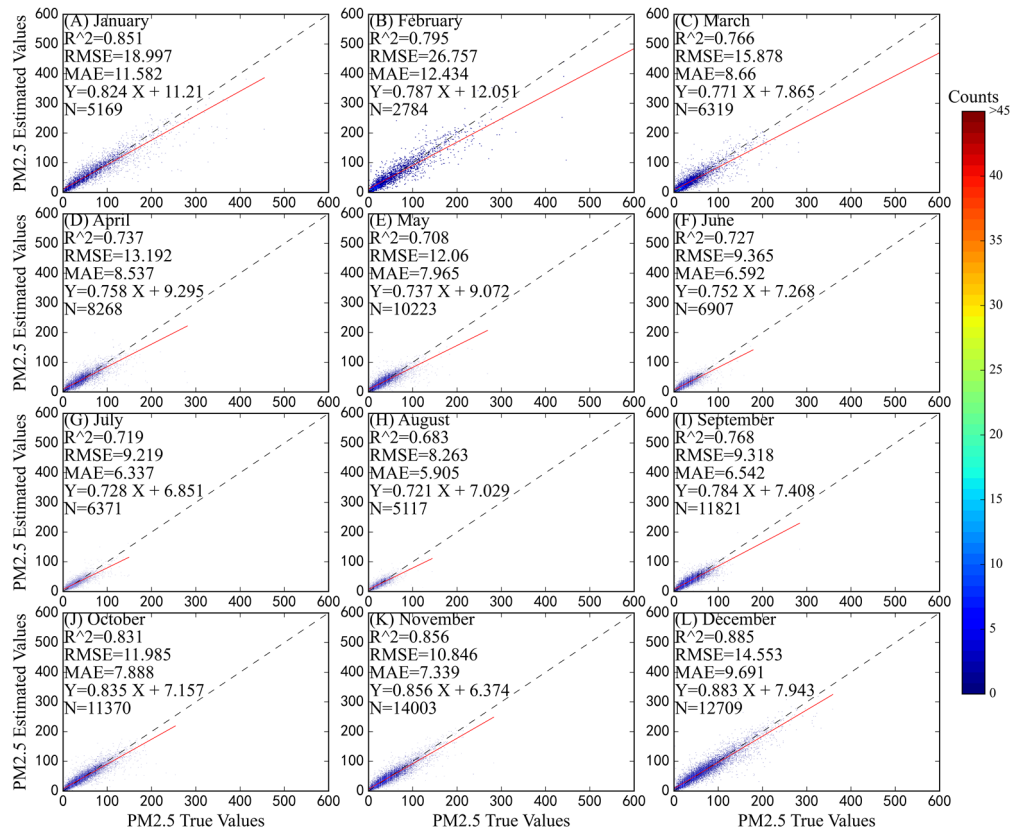
Figure 6 shows the scatterplot fitted with the inversion results of the mixed model from 9:00-17:00 local Time. The model R^2 ranged from 0.556 to 0.88 at different times. Except for 17:00 when the model had the worst performance, the model R^2 exceeded 0.7 at other times, indicating that the model had a

1 good performance. The optimal performance time is 13:00, and R^2 is 0.88. According to the results, the
 2 hourly differences in model performance were significant.



3
 4 **Figure 6 Density scatterplot of Actual hourly PM_{2.5} values (x-axis) model estimated values (y-axis) in hourly**
 5 **PM_{2.5} estimates in China from (A) 09:00 LT to (I) 17:00 LT. R^2 represents determination coefficient, RMSE**
 6 **represents root mean square error, MAE represents mean Absolute Error, N represents the number of**
 7 **samples. The equation Y and X represent the fitting relationship between the actual and estimated PM_{2.5}**
 8 **values. Black dashed line represents 1:1 line, and red line represents best-fit line from linear regression.**

9 Figure 7 shows the inversion performance results of the hybrid model collected from January to
 10 December 2019. The model performed the worst in summer months because R^2 was reported 0.73, 0.72,
 11 and 0.68, respectively; however, RMSE and MAE were only 9.37, 9.22, 8.26 $\mu\text{g}/\text{m}^3$ and 6.59, 6.34, and
 12 5.91 $\mu\text{g}/\text{m}^3$, respectively, due to the lower average concentration of PM_{2.5} in summer. Winter and autumn
 13 models gained better performance results with an average R^2 over 0.8. However, in contrast to summer,
 14 the estimation errors of these two seasons were relatively large, with average RMSE of 20.10 $\mu\text{g}/\text{m}^3$ and
 15 10.72 $\mu\text{g}/\text{m}^3$ and average MAE of 11.20 $\mu\text{g}/\text{m}^3$ and 7.25 $\mu\text{g}/\text{m}^3$, respectively. The mean R^2 was 0.74,
 16 whereas the mean RMSE and MAE were 13.71 $\mu\text{g}/\text{m}^3$ and 8.39 $\mu\text{g}/\text{m}^3$, respectively.



1
2

Figure 7 Same as Figure 6, but for monthly PM_{2.5} estimates.

3 4.2.4 Feature importance analysis

4 The model performance differences were also analyzed to extract and rank the model features of
 5 RF and GBRT based on the feature importance. The higher the feature importance, the greater the
 6 contribution of factors to the model. Figure 8 shows that AOD, boundary layer height, 2 m surface
 7 temperature, and relative humidity had the greatest effect on the mixed model performance out of all
 8 variable characteristic parameters. Accordingly, AOD is greatly affected by the fine particulate matter
 9 and is the main factor in the inversion of PM_{2.5}. Changes of the boundary layer height can affect the
 10 diffusion ability of the atmosphere. If the boundary layer height is low, the accumulation of pollutants
 11 will be caused. At the same time, the 2 m surface temperature has a great impact on the boundary layer
 12 height (Miao et al., 2018). Finally, higher rates of atmospheric humidity can improve the fine particulate
 13 matter accumulation.

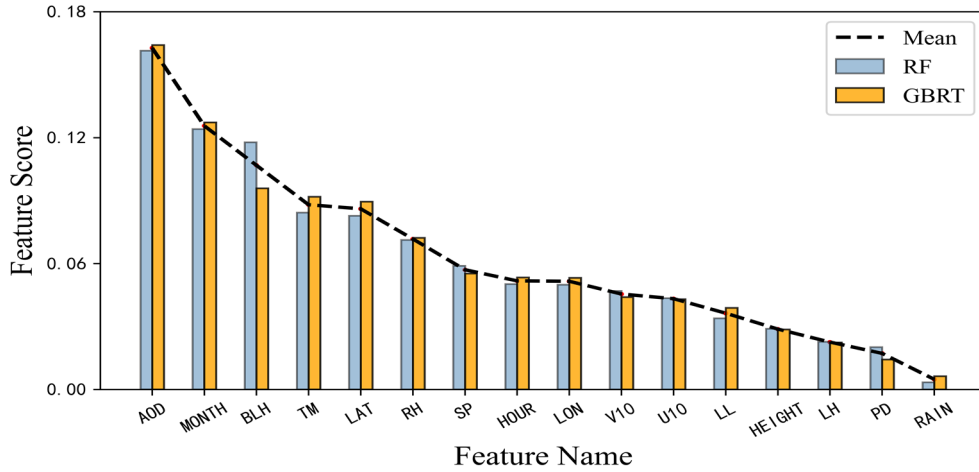
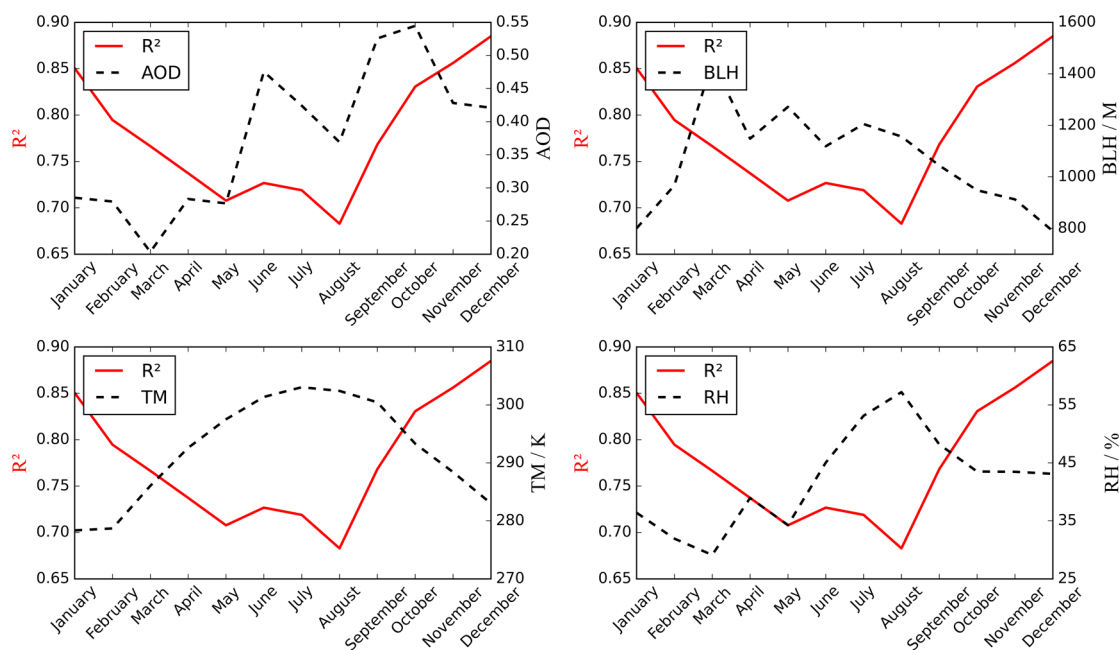


Figure 8 Score (y-axis) for each model contributing feature factor (y-axis) for the RF (blue) and GBRT (orange). Dashed line represents the mean values.

The correlation coefficients between the monthly mean values of important meteorological parameters (AOD, BLH, TM and RH) and R^2 were also analyzed. According to the results, the correlation coefficients between the meteorological parameters and $PM_{2.5}$ were lower in summer. Furthermore, there are many rainy days and large cloud coverage, which is not conducive to satellite observation and decreases the accuracy of AOD data in summer. Therefore, the summer model performance is poor. There was a strong correlation between meteorological parameters and $PM_{2.5}$ in autumn. There were also similar correlations between spring and winter; however, the winter model performed was better. The reasons can be interpreted as below. The winter temperature and boundary layer height are low, whereas the atmosphere is stable but not conducive to the diffusion of pollutants. Moreover, during the heating period in winter, pollutant emissions soar greatly and result in a sharp rise in the concentration of $PM_{2.5}$. The increased pollution in winter ensures the quality and quantity of data, thereby improving the model performance effectively.

Table 2 Correlation coefficient between meteorological parameters with $PM_{2.5}$

Season	AOD	BLH	TM	RH
Spring	0.47	-0.33	0.12	0.36
Summer	0.42	-0.21	0.06	0.19
Autumn	0.38	-0.29	0.24	0.41
Winter	0.44	-0.33	0.12	0.35

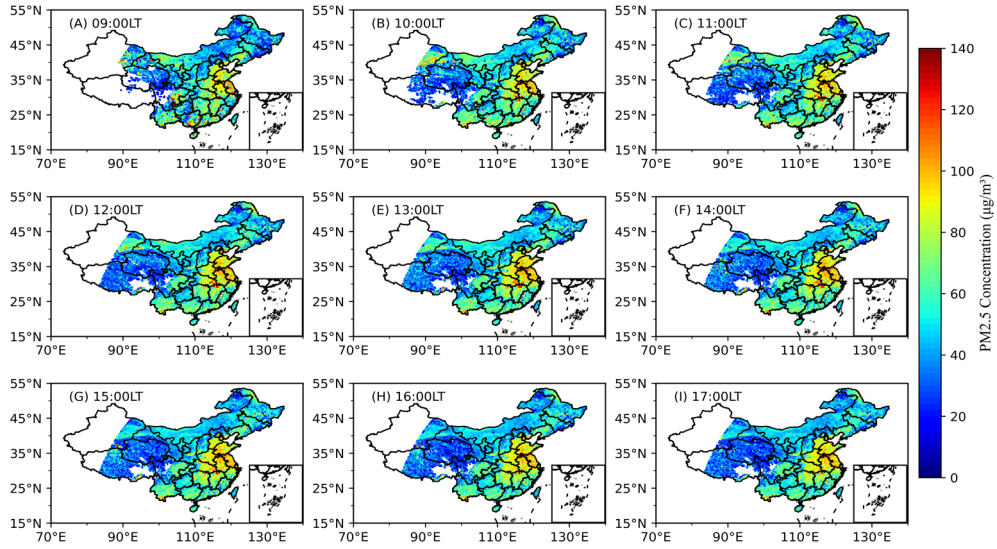


1
2 **Figure 9 Annual variability (x-axis) of monthly average of meteorological parameters AOD, BLH (m), TM**
3 **(K), RH (%) (right y-axis) and R^2 (left y-axis)**

4 **4.3 Temporal and Spatial Distribution Characteristics of $PM_{2.5}$ Concentration in China**

5 In terms of spatial distribution, Shandong, Henan, Jiangsu, Anhui, as well as parts of Hubei and
6 Hebei were the most polluted areas in China in 2019, with an annual average $PM_{2.5}$ concentration of
7 $82.86 \mu\text{g}/\text{m}^3$. On the one hand, these areas are economically developed and densely populated, resulting
8 in a large amount of pollutant emissions. On the other hand, the barrier of the peripheral mountains
9 (Taihang Mountains, Qinling Mountains and the Southern Hills) leads to the accumulation of pollutants
10 that are difficult to diffuse. Sichuan Basin is a rare area with a high $PM_{2.5}$ value due to its unique
11 topography (Zhang et al., 2019a), with an annual average $PM_{2.5}$ concentration of $64.69 \mu\text{g}/\text{m}^3$. In addition,
12 Inner Mongolia, Qinghai, Tibet and other places, the pollution level is low, the average annual $PM_{2.5}$
13 concentration is less than $40 \mu\text{g}/\text{m}^3$.

14 The temporal distribution of $PM_{2.5}$ is shown in Figure 10, The $PM_{2.5}$ concentration began to rise
15 from 9:00, and peaked at $55.65 \mu\text{g}/\text{m}^3$ between 10:00 and 11:00 every day. After that, it maintained a high
16 concentration until 15:00; and began to decrease. In the most polluted areas of China, the peak
17 concentration of $PM_{2.5}$ can reach $85.05 \mu\text{g}/\text{m}^3$, while the peak in the less polluted areas is only about
18 $40 \mu\text{g}/\text{m}^3$. On a national scale, daily $PM_{2.5}$ concentrations fluctuates slightly.

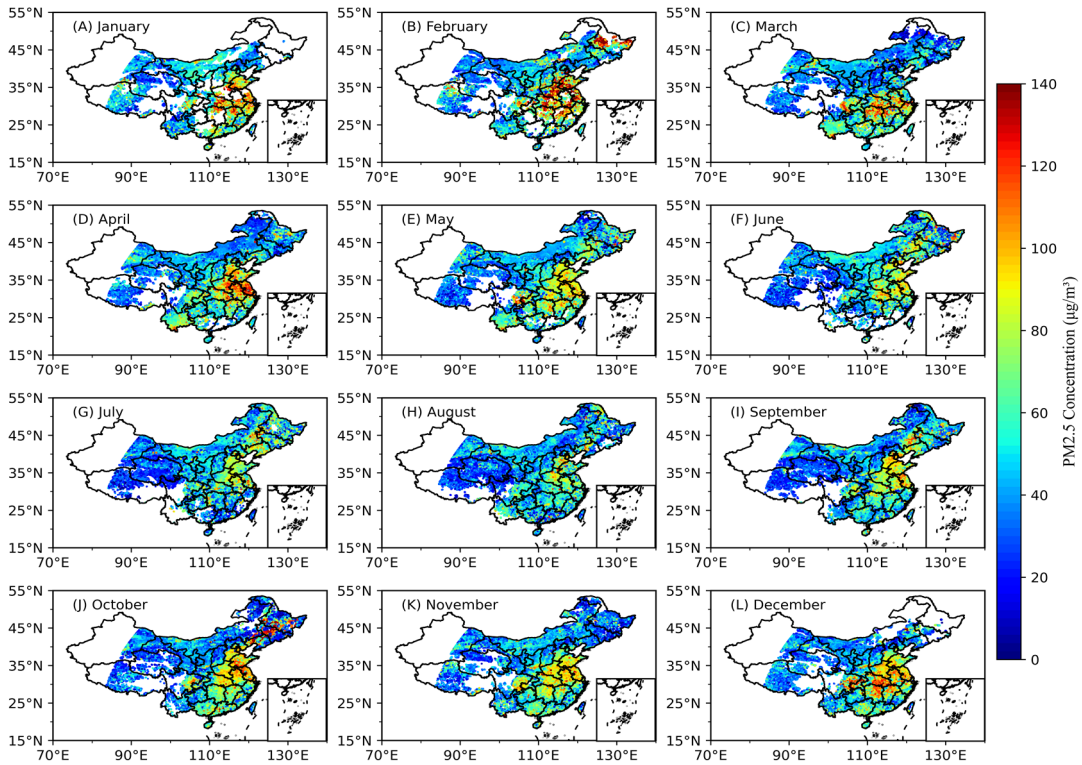


1
2
3

Figure 10 Hourly Spatial distribution of PM_{2.5} concentration in China at different local times from (A) 09:00 LT to (I) 17:00 LT.

4
5
6
7
8

PM_{2.5} concentration in China varies significantly with the seasons. As shown in Figure 11, PM_{2.5} concentration in winter is the highest, with an average value of 62.10µg/m³. January 2019 was the most polluted month in China, with the average PM_{2.5} concentration reaching 63.58µg/m³. The average PM_{2.5} concentration was 47.39 µg/m³ in summer. The average concentration of PM_{2.5} in spring and autumn was 54.21µg/m³ and 52.26 µg/m³, respectively, indicating similar levels of pollution.



9
10

Figure 11 Same as Fig. 10, but for monthly spatial distribution

1 **5 Conclusion**

2 It is essential to collect the spatiotemporal evolution characteristics regarding the concentration of
3 $PM_{2.5}$ for air pollution prevention and containment. Based on the linear hybrid machine learning model,
4 this paper used the AOD data of Himawari-8 to invert the concentration of $PM_{2.5}$ in China and obtain its
5 distribution characteristics. The model performance and inversion results are analyzed and summarized
6 below:

7 (1) In the RGD-LHMLM obtained from linear fitting, the DNN accounted for the largest proportion
8 with a weight coefficient of 0.62. The R^2 of RGD-LHMLM was 0.84, whereas its generalization ability
9 was significantly better than that of a single model (DNN: 0.80; GBRT: 0.81; RF: 0.79). Moreover,
10 RMSE and MAE were $12.92 \mu\text{g}/\text{m}^3$ and $8.01 \mu\text{g}/\text{m}^3$, respectively.

11 (2) The RGD-LHMLM was spatially stable, with $R^2 > 0.7$ in more than 70% of sites as well as
12 $\text{RMSE} < 20 \mu\text{g}/\text{m}^3$ and $\text{MAE} < 15 \mu\text{g}/\text{m}^3$ in more than 95% of sites. These sites are mainly located in densely
13 populated and industrially developed areas. The correlation difference between the inversion factor and
14 $PM_{2.5}$ in various seasons would lead to seasonal variations in the model performance. In addition, the
15 performance was the worst in summer with an average R^2 of 0.71; however, winter showed the best
16 performance with an average R^2 of 0.84. The diurnal variation of the model inversion effect is also
17 obvious, and the 11:00-14:00 model usually has better performance.

18 (3) Changes in the spatiotemporal characteristics were obvious in the concentration of $PM_{2.5}$ in
19 China. In other words, North China and East China had the highest concentration of $PM_{2.5}$ with an
20 average annual concentration of $82.86 \mu\text{g}/\text{m}^3$, whereas Inner Mongolia, Qinghai, Tibet, and other regions
21 had low pollution levels with an average annual concentration of $PM_{2.5}$ below $40 \mu\text{g}/\text{m}^3$. In winter, the
22 concentration of $PM_{2.5}$ was higher with an average of $62.10 \mu\text{g}/\text{m}^3$, whereas the pollution was lighter in
23 summer with an average concentration of $PM_{2.5}$ being reported $47.39 \mu\text{g}/\text{m}^3$. In the most polluted areas,
24 the peak concentration of $PM_{2.5}$ can reach $85.05 \mu\text{g}/\text{m}^3$, but the daily $PM_{2.5}$ concentration fluctuates
25 slightly.

26 In conclusion, the RGD-LHMLM can accurately measure the concentration of $PM_{2.5}$ and perform
27 the seasonal evolution of pollutants. These results can help control the local pollution. This study also
28 indicated that integrating multiple Machine learning models improved the accuracy of fitting results
29 effectively. For more accurate pollutant data, such models can be employed to fit the $PM_{2.5}$ in the future

1 with more parameters closely related to PM_{2.5}. However, there are some vacant values in the results of
2 this study. There are also no data for some areas. Thus, other satellite data can be used in future studies
3 to solve this problem.

4 **Code/Data availability**

5 Datasets and Code related to this paper can be requested from the corresponding author
6 (chenbin@lzu.edu.cn). The PM_{2.5} data download address is : <http://106.37.208.233:20035/>; Himawari-8
7 AOD data provided by the Japan Meteorological Agency, download from:
8 <http://www.eorc.jaxa.jp/ptree/index.html>; ERA-5 meteorological data can be downloaded from the
9 European Centre for Medium-Range Weather Forecasts (ECMWF) at : <https://cds.climate.copernicus.eu>;
10 Ground elevation SRTM3 data download address is: <http://srtm.csi.cgiar.org/index.asp>; NASA's social
11 and economic data and the population density of data center, download address is:
12 <http://sedac.ciesin.columbia.edu/data/collection/gpw-v4/documentation>.

13 **Author contributions**

14 Chen proposed the content of the study. Song performed data processing, model building, result analysis,
15 and article writing. Huang, Dong and Yang checked the content of the article.

16 **Competing interests**

17 The authors declare that they have no conflict of interest.

18 **Acknowledgments**

19 We thank China National Environmental Monitoring Center, Japan Meteorological Agency, European
20 Centre for Medium-Range Weather Forecasts, NASA, and the National Mapping Service of the
21 Department of Defense.

22 **Financial support**

23 The National Key Research and Development Program of China (Grant number 2019YFA0606800), the

1 National Natural Science Foundation of China (Grant 41775021), The Fundamental Research Funds for
2 the Central Universities (Grant lzujbky-2019-43).

3 **References**

- 4 Adams, M. D., Massey, F., Chastko, K., and Cupini, C.: Spatial modelling of particulate matter air
5 pollution sensor measurements collected by community scientists while cycling, land use regression with
6 spatial cross-validation, and applications of machine learning for data correction, *Atmos Environ*, 230,
7 <https://doi.org/10.1016/j.atmosenv.2020.117479>, 2020.
- 8 Apte, J. S., Marshall, J. D., Cohen, A. J., and Brauer, M.: Addressing Global Mortality from Ambient
9 PM_{2.5}, *Environ Sci Technol*, 49, 8057-8066, <https://doi.org/10.1021/acs.est.5b01236>, 2015.
- 10 Bessho, K., Date, K., Hayashi, M., Ikeda, A., Imai, T., Inoue, H., Kumagai, Y., Miyakawa, T., Murata,
11 H., Ohno, T., Okuyama, A., Oyama, R., Sasaki, Y., Shimazu, Y., Shimoji, K., Sumida, Y., Suzuki, M.,
12 Taniguchi, H., Tsuchiyama, H., Uesawa, D., Yokota, H., and Yoshida, R.: An Introduction to Himawari-
13 8/9-Japan's New-Generation Geostationary Meteorological Satellites, *J Meteorol Soc Jpn*, 94, 151-183,
14 <https://doi.org/10.2151/jmsj.2016-009>, 2016.
- 15 Breiman, L.: Random forests, *Mach Learn*, 45, 5-32, <https://doi.org/10.1023/A:1010933404324>, 2001.
- 16 Chen, B. J., You, S. X., Ye, Y., Fu, Y. Y., Ye, Z. R., Deng, J. S., Wang, K., and Hong, Y.: An interpretable
17 self-adaptive deep neural network for estimating daily spatially-continuous PM_{2.5} concentrations across
18 China, *Sci Total Environ*, 768, <https://doi.org/10.1016/j.scitotenv.2020.144724>, 2021.
- 19 Chen, J. P., Yin, J. H., Zang, L., Zhang, T. X., and Zhao, M. D.: Stacking machine learning model for
20 estimating hourly PM_{2.5} in China based on Himawari 8 aerosol optical depth data, *Sci Total Environ*,
21 697, <https://doi.org/10.1016/j.scitotenv.2019.134021>, 2019.
- 22 China: Ambient air quality standards. GB 3095-2012., 2012.
- 23 Diederik, P. K. and Jimmy, B.: Adam: A Method for Stochastic Optimization, ICLR2015.
- 24 Emili, E., Popp, C., Petitta, M., Riffler, M., Wunderle, S., and Zebisch, M.: PM₁₀ remote sensing from
25 geostationary SEVIRI and polar-orbiting MODIS sensors over the complex terrain of the European
26 Alpine region, *Remote Sens Environ*, 114, 2485-2499, <https://doi.org/10.1016/j.rse.2010.05.024>, 2010.
- 27 Engel-Cox, J. A., Holloman, C. H., Coutant, B. W., and Hoff, R. M.: Qualitative and quantitative
28 evaluation of MODIS satellite sensor data for regional and urban scale air quality, *Atmos Environ*, 38,
29 2495-2509, <https://doi.org/10.1016/j.atmosenv.2004.01.039>, 2004.
- 30 Fang, X., Zou, B., Liu, X. P., Sternberg, T., and Zhai, L.: Satellite-based ground PM_{2.5} estimation using
31 timely structure adaptive modeling, *Remote Sens Environ*, 186, 152-163,
32 <https://doi.org/10.1016/j.rse.2016.08.027>, 2016.
- 33 Friedman, J. H.: Greedy function approximation: A gradient boosting machine, *Ann Stat*, 29, 1189-1232,
34 <https://doi.org/10.1214/aos/1013203451>, 2001.
- 35 Gao, M., Guttikunda, S. K., Carmichael, G. R., Wang, Y. S., Liu, Z. R., Stanier, C. O., Saide, P. E., and
36 Yu, M.: Health impacts and economic losses assessment of the 2013 severe haze event in Beijing area,
37 *Sci Total Environ*, 511, 553-561, <https://doi.org/10.1016/j.scitotenv.2015.01.005>, 2015.
- 38 Gui, K., Che, H. Z., Wang, Y. Q., Wang, H., Zhang, L., Zhao, H. J., Zheng, Y., Sun, T. Z., and Zhang, X.
39 Y.: Satellite-derived PM_{2.5} concentration trends over Eastern China from 1998 to 2016: Relationships
40 to emissions and meteorological parameters, *Environ Pollut*, 247, 1125-1133,
41 <https://doi.org/10.1016/j.envpol.2019.01.056>, 2019.

1 Guo, B., Zhang, D. M., Pei, L., Su, Y., Wang, X. X., Bian, Y., Zhang, D. H., Yao, W. Q., Zhou, Z. X., and
2 Guo, L. Y.: Estimating PM_{2.5} concentrations via random forest method using satellite, auxiliary, and
3 ground-level station dataset at multiple temporal scales across China in 2017, *Sci Total Environ*, 778,
4 <https://doi.org/10.1016/j.scitotenv.2021.146288>, 2021.

5 Guo, J. P., Xia, F., Zhang, Y., Liu, H., Li, J., Lou, M. Y., He, J., Yan, Y., Wang, F., Min, M., and Zhai, P.
6 M.: Impact of diurnal variability and meteorological factors on the PM_{2.5} - AOD relationship:
7 Implications for PM_{2.5} remote sensing, *Environ Pollut*, 221, 94-104,
8 <https://doi.org/10.1016/j.envpol.2016.11.043>, 2017.

9 Han, Y., Wu, Y. H., Wang, T. J., Zhuang, B. L., Li, S., and Zhao, K.: Impacts of elevated-aerosol-layer
10 and aerosol type on the correlation of AOD and particulate matter with ground-based and satellite
11 measurements in Nanjing, southeast China, *Sci Total Environ*, 532, 195-207,
12 <https://doi.org/10.1016/j.scitotenv.2015.05.136>, 2015.

13 Hoff, R. M. and Christopher, S. A.: Remote Sensing of Particulate Pollution from Space: Have We
14 Reached the Promised Land?, *J Air Waste Manage*, 59, 645-675, <https://doi.org/10.3155/1047-3289.59.6.645>, 2009.

16 Hu, X. F., Waller, L. A., Al-Hamdan, M. Z., Crosson, W. L., Estes, M. G., Estes, S. M., Quattrochi, D.
17 A., Sarnat, J. A., and Liu, Y.: Estimating ground-level PM_{2.5} concentrations in the southeastern US using
18 geographically weighted regression, *Environ Res*, 121, 1-10,
19 <https://doi.org/10.1016/j.envres.2012.11.003>, 2013.

20 Jiang, Y., Yang, K., Shao, C., Zhou, X., Zhao, L., Chen, Y., and Wu, H.: A downscaling approach for
21 constructing high-resolution precipitation dataset over the Tibetan Plateau from ERA5 reanalysis, *Atmos*
22 *Res*, 256, 105574, <https://doi.org/10.1016/j.atmosres.2021.105574>, 2021.

23 Johnson, N. E., Bonczak, B., and Kontokosta, C. E.: Using a gradient boosting model to improve the
24 performance of low-cost aerosol monitors in a dense, heterogeneous urban environment, *Atmos Environ*,
25 184, 9-16, <https://doi.org/10.1016/j.atmosenv.2018.04.019>, 2018.

26 Lee, H. J., Coull, B. A., Bell, M. L., and Koutrakis, P.: Use of satellite-based aerosol optical depth and
27 spatial clustering to predict ambient PM_{2.5} concentrations, *Environ Res*, 118, 8-15,
28 <https://doi.org/10.1016/j.envres.2012.06.011>, 2012.

29 Li, T. W., Shen, H. F., Yuan, Q. Q., Zhang, X. C., and Zhang, L. P.: Estimating Ground-Level PM_{2.5} by
30 Fusing Satellite and Station Observations: A Geo-Intelligent Deep Learning Approach, *Geophys Res Lett*,
31 44, 11985-11993, <https://doi.org/10.1002/2017gl075710>, 2017b.

32 Li, T. W., Shen, H. F., Zeng, C., Yuan, Q. Q., and Zhang, L. P.: Point-surface fusion of station
33 measurements and satellite observations for mapping PM_{2.5} distribution in China: Methods and
34 assessment, *Atmos Environ*, 152, 477-489, <https://doi.org/10.1016/j.atmosenv.2017.01.004>, 2017a.

35 Li, Z. Q., Zhang, Y., Shao, J., Li, B. S., Hong, J., Liu, D., Li, D. H., Wei, P., Li, W., Li, L., Zhang, F. X.,
36 Guo, J., Deng, Q., Wang, B. X., Cui, C. L., Zhang, W. C., Wang, Z. Z., Lv, Y., Xu, H., Chen, X. F., Li,
37 L., and Qie, L. L.: Remote sensing of atmospheric particulate mass of dry PM_{2.5} near the ground: Method
38 validation using ground-based measurements, *Remote Sens Environ*, 173, 59-68,
39 <https://doi.org/10.1016/j.rse.2015.11.019>, 2016.

40 Lim, C. H., Ryu, J., Choi, Y., Jeon, S. W., and Lee, W. K.: Understanding global PM_{2.5} concentrations
41 and their drivers in recent decades (1998-2016), *Environ Int*, 144,
42 <https://doi.org/10.1016/j.envint.2020.106011>, 2020.

43 Liu, Y., Cao, G. F., Zhao, N. Z., Mulligan, K., and Ye, X. Y.: Improve ground-level PM_{2.5} concentration
44 mapping using a random forests-based geostatistical approach, *Environ Pollut*, 235, 272-282,

1 <https://doi.org/10.1016/j.envpol.2017.12.070>, 2018.

2 Liu, Y., Sarnat, J. A., Kilaru, A., Jacob, D. J., and Koutrakis, P.: Estimating ground-level PM_{2.5} in the
3 eastern united states using satellite remote sensing, *Environ Sci Technol*, 39, 3269-3278,
4 <https://doi.org/10.1021/es049352m>, 2005.

5 Lv, B. L., Hu, Y. T., Chang, H. H., Russell, A. G., Cai, J., Xu, B., and Bai, Y. Q.: Daily estimation of
6 ground-level PM_{2.5} concentrations at 4 km resolution over Beijing-Tianjin-Hebei by fusing MODIS
7 AOD and ground observations, *Sci Total Environ*, 580, 235-244,
8 <https://doi.org/10.1016/j.scitotenv.2016.12.049>, 2017.

9 Mao, F. Y., Hong, J., Min, Q. L., Gong, W., Zang, L., and Yin, J. H.: Estimating hourly full-coverage
10 PM_{2.5} over China based on TOA reflectance data from the Fengyun-4A satellite, *Environ Pollut*, 270,
11 <https://doi.org/10.1016/j.envpol.2020.116119>, 2021.

12 Miao, Y. C., Liu, S. H., Guo, J. P., Huang, S. X., Yan, Y., and Lou, M. Y.: Unraveling the relationships
13 between boundary layer height and PM_{2.5} pollution in China based on four-year radiosonde
14 measurements, *Environ Pollut*, 243, 1186-1195, <https://doi.org/10.1016/j.envpol.2018.09.070>, 2018.

15 Pan, Z. X., Mao, F. Y., Wang, W., Zhu, B., Lu, X., and Gong, W.: Impacts of 3D Aerosol, Cloud, and
16 Water Vapor Variations on the Recent Brightening during the South Asian Monsoon Season, *Remote
17 Sens-Basel*, 10, <https://doi.org/10.3390/rs10040651>, 2018.

18 Pun, V. C., Kazemiparkouhi, F., Manjourides, J., and Suh, H. H.: Long-Term PM_{2.5} Exposure and
19 Respiratory, Cancer, and Cardiovascular Mortality in Older US Adults, *Am J Epidemiol*, 186, 961-969,
20 <https://doi.org/10.1093/aje/kwx166>, 2017.

21 Qin, K., Wang, L. Y., Wu, L. X., Xu, J., Rao, L. L., Letu, H., Shi, T. W., and Wang, R. F.: A campaign for
22 investigating aerosol optical properties during winter hazes over Shijiazhuang, China, *Atmos Res*, 198,
23 113-122, <https://doi.org/10.1016/j.atmosres.2017.08.018>, 2017.

24 Schonlau, M.: Boosted regression (boosting): An introductory tutorial and a Stata plugin, *Stata J*, 5, 330-
25 354, <https://doi.org/10.1177/1536867x0500500304>, 2005.

26 Stafoggia, M., Bellander, T., Bucci, S., Davoli, M., de Hoogh, K., de'Donato, F., Gariazzo, C., Lyapustin,
27 A., Michelozzi, P., Renzi, M., Scortichini, M., Shtein, A., Viegi, G., Kloog, I., and Schwartz, J.:
28 Estimation of daily PM₁₀ and PM_{2.5} concentrations in Italy, 2013-2015, using a spatiotemporal land-
29 use random-forest model, *Environ Int*, 124, 170-179, <https://doi.org/10.1016/j.envint.2019.01.016>, 2019.

30 Tian, J. and Chen, D. M.: A semi-empirical model for predicting hourly ground-level fine particulate
31 matter (PM_{2.5}) concentration in southern Ontario from satellite remote sensing and ground-based
32 meteorological measurements, *Remote Sens Environ*, 114, 221-229,
33 <https://doi.org/10.1016/j.rse.2009.09.011>, 2010.

34 Wang, W., Mao, F. Y., Du, L., Pan, Z. X., Gong, W., and Fang, S. H.: Deriving Hourly PM_{2.5}
35 Concentrations from Himawari-8 AODs over Beijing-Tianjin-Hebei in China, *Remote Sens-Basel*, 9,
36 <https://doi.org/10.3390/rs9080858>, 2017.

37 Wang, X. H., Zhong, S. Y., Bian, X. D., and Yu, L. J.: Impact of 2015-2016 El Nino and 2017-2018 La
38 Nina on PM_{2.5} concentrations across China, *Atmos Environ*, 208, 61-73,
39 <https://doi.org/10.1016/j.atmosenv.2019.03.035>, 2019a.

40 Wang, X. P. and Sun, W. B.: Meteorological parameters and gaseous pollutant concentrations as
41 predictors of daily continuous PM_{2.5} concentrations using deep neural network in Beijing-Tianjin-Hebei,
42 China, *Atmos Environ*, 211, 128-137, <https://doi.org/10.1016/j.atmosenv.2019.05.004>, 2019.

43 Wang, X. Q., Wei, W., Cheng, S. Y., Yao, S., Zhang, H. Y., and Zhang, C.: Characteristics of PM_{2.5} and
44 SNA components and meteorological factors impact on air pollution through 2013-2017 in Beijing,

1 China, Atmospheric Pollution Research, 10, 1976-1984, <https://doi.org/10.1016/j.apr.2019.09.004>,
2 2019b.

3 Wei, J., Huang, W., Li, Z. Q., Xue, W. H., Peng, Y. R., Sun, L., and Cribb, M.: Estimating 1-km-resolution
4 PM2.5 concentrations across China using the space-time random forest approach, Remote Sens Environ,
5 231, <https://doi.org/10.1016/j.rse.2019.111221>, 2019a.

6 Wei, J., Li, Z., Pinker, R. T., Wang, J., Sun, L., Xue, W., Li, R., and Cribb, M.: Himawari-8-derived
7 diurnal variations in ground-level PM2.5 pollution across China using the fast space-time Light Gradient
8 Boosting Machine (LightGBM), Atmos. Chem. Phys., 21, 7863-7880, <https://doi.org/10.5194/acp-21-7863-2021>, 2021a.

10 Wei, J., Li, Z. Q., Lyapustin, A., Sun, L., Peng, Y. R., Xue, W. H., Su, T. N., and Cribb, M.: Reconstructing
11 1-km-resolution high-quality PM2.5 data records from 2000 to 2018 in China: spatiotemporal variations
12 and policy implications, Remote Sens Environ, 252, <https://doi.org/10.1016/j.rse.2020.112136>, 2021b.

13 Wei, J., Li, Z., Sun, L., Peng, Y., Zhang, Z., Li, Z., Su, T., Feng, L., Cai, Z., and Wu, H.: Evaluation and
14 uncertainty estimate of next-generation geostationary meteorological Himawari-8/AHI aerosol products,
15 Sci Total Environ, 692, 879-891, <https://doi.org/10.1016/j.scitotenv.2019.07.326>, 2019b.

16 Wolpert, D. H.: Stacked Generalization, Neural Networks, 5, 241-259, [https://doi.org/10.1016/S0893-6080\(05\)80023-1](https://doi.org/10.1016/S0893-6080(05)80023-1), 1992.

18 Xu, J. H., Lindqvist, H., Liu, Q. F., Wang, K., and Wang, L.: Estimating the spatial and temporal
19 variability of the ground-level NO2 concentration in China during 2005–2019 based on satellite remote
20 sensing,, Atmospheric Pollution Research, 12, 57-67,
21 <https://doi.org/https://doi.org/10.1016/j.apr.2020.10.008>, 2021.

22 Yang, X. C., Jiang, L., Zhao, W. J., Xiong, Q. L., Zhao, W. H., and Yan, X.: Comparison of Ground-
23 Based PM2.5 and PM10 Concentrations in China, India, and the US, Int J Env Res Pub He, 15,
24 <https://doi.org/10.3390/ijerph15071382>, 2018.

25 Yesilkanat, C. M.: Spatio-temporal estimation of the daily cases of COVID-19 in worldwide using
26 random forest machine learning algorithm, Chaos Soliton Fract, 140, <https://doi.org/ARTN 110210>
27 10.1016/j.chaos.2020.110210, 2020.

28 Yin, J. H., Mao, F. Y., Zang, L., Chen, J. P., Lu, X., and Hong, J.: Retrieving PM2.5 with high spatio-
29 temporal coverage by TOA reflectance of Himawari-8, Atmospheric Pollution Research, 12, 14-20,
30 <https://doi.org/10.1016/j.apr.2021.02.007>, 2021.

31 Yoshida, M., Kikuchi, M., Nagao, T. M., Murakami, H., Nomaki, T., and Higurashi, A.: Common
32 Retrieval of Aerosol Properties for Imaging Satellite Sensors, Journal of the Meteorological Society of
33 Japan. Ser. II, 96B, 193-209, <https://doi.org/10.2151/jmsj.2018-039>, 2018.

34 Yumimoto, K., Nagao, T. M., Kikuchi, M., Sekiyama, T. T., Murakami, H., Tanaka, T. Y., Ogi, A., Irie,
35 H., Khatri, P., Okumura, H., Arai, K., Morino, I., Uchino, O., and Maki, T.: Aerosol data assimilation
36 using data from Himawari-8, a next-generation geostationary meteorological satellite, Geophys Res Lett,
37 43, 5886-5894, <https://doi.org/10.1002/2016gl069298>, 2016.

38 Zang, L., Mao, F. Y., Guo, J. P., Gong, W., Wang, W., and Pan, Z. X.: Estimating hourly PM1
39 concentrations from Himawari-8 aerosol optical depth in China, Environ Pollut, 241, 654-663,
40 <https://doi.org/10.1016/j.envpol.2018.05.100>, 2018.

41 Zhang, L., Guo, X. M., Zhao, T. L., Gong, S. L., Xu, X. D., Li, Y. Q., Luo, L., Gui, K., Wang, H. L.,
42 Zheng, Y., and Yin, X. F.: A modelling study of the terrain effects on haze pollution in the Sichuan Basin,
43 Atmos Environ, 196, 77-85, <https://doi.org/10.1016/j.atmosenv.2018.10.007>, 2019a.

44 Zhang, T. H., Zhu, Z. M., Gong, W., Zhu, Z. R., Sun, K., Wang, L. C., Huang, Y. S., Mao, F. Y., Shen, H.

1 F., Li, Z. W., and Xu, K.: Estimation of ultrahigh resolution PM2.5 concentrations in urban areas using
2 160 m Gaofen-1 AOD retrievals, *Remote Sens Environ*, 216, 91-104,
3 <https://doi.org/10.1016/j.rse.2018.06.030>, 2018.

4 Zhang, T. X., Zang, L., Wan, Y. C., Wang, W., and Zhang, Y.: Ground-level PM2.5 estimation over urban
5 agglomerations in China with high spatiotemporal resolution based on Himawari-8, *Sci Total Environ*,
6 676, 535-544, <https://doi.org/10.1016/j.scitotenv.2019.04.299>, 2019b.

7 Zhang, Z., Wu, W., Fan, M., Tao, M., Wei, J., Jin, J., Tan, Y., and Wang, Q.: Validation of Himawari-8
8 aerosol optical depth retrievals over China, *Atmos Environ*, 199, 32-44,
9 <https://doi.org/10.1016/j.atmosenv.2018.11.024>, 2019c.

10 Zheng, C. W., Zhao, C. F., Zhu, Y. N., Wang, Y., Shi, X. Q., Wu, X. L., Chen, T. M., Wu, F., and Qiu, Y.
11 M.: Analysis of influential factors for the relationship between PM2.5 and AOD in Beijing, *Atmos Chem*
12 *Phys*, 17, 13473-13489, <https://doi.org/10.5194/acp-17-13473-2017>, 2017.

13