

# Predicting PM<sub>2.5</sub> Concentrations Across USA Using Machine Learning

P. Preetham Vignesh<sup>1</sup>, Jonathan H Jiang<sup>2</sup>, and Pangaluru Kishore<sup>3</sup>

<sup>1</sup>University High School

<sup>2</sup>Jet Propulsion Laboratory, California Institute of Technology

<sup>3</sup>University of California, Irvine

February 13, 2023

## Abstract

Fine particulate matter with a size less than 2.5  $\mu\text{m}$  (PM<sub>2.5</sub>) is increasing due to economic growth, air pollution, and forest fires in some states in the United States. Although previous studies have attempted to retrieve the spatial and temporal behavior of PM<sub>2.5</sub> using aerosol remote sensing and geostatistical estimation methods the coarse resolution and accuracy limit these methods. In this paper the performance of machine learning models on predicting PM<sub>2.5</sub> is assessed with Linear Regression (LR), Decision Tree (DT), Gradient Boosting Regression (GBR), AdaBoost Regression (ABR), XG Boost (XGB), k-nearest neighbors (KNN), Long Short-Term Memory (LSTM), Random Forest (RF), and support vector machine (SVM) using PM<sub>2.5</sub> station data from 2017-2021. To compare the accuracy of all the nine machine learning models the coefficient of determination (R<sup>2</sup>), root mean square error (RMSE), Nash-Sutcliffe efficiency (NSE), root mean square error ratio (RSR), and percent bias (PBIAS) were evaluated. Among all nine models the RF and SVM models were the best for predicting PM<sub>2.5</sub> concentrations. Comparison of the PM<sub>2.5</sub> performance metrics displayed that the models had better predictive behavior in the western United States than that in the eastern United States.

# Predicting PM<sub>2.5</sub> Concentrations Across USA Using Machine Learning

P. Preetham Vignesh<sup>1</sup>, Jonathan H. Jiang<sup>2</sup>, P. Kishore

<sup>1</sup> University of California, Los Angeles, USA

<sup>2</sup> Jet Propulsion Laboratory, California Institute of Technology, Pasadena, California, USA.

<sup>3</sup> Retired, University of California, Irvine, USA

Copyright ©2022, All Rights Reserved.

Correspondence: Jonathan.H.Jiang@jpl.nasa.gov

Keywords: Surface Temperature, Climate Model, Global Warming Projection

## Abstract:

Fine particulate matter with a size less than 2.5  $\mu\text{m}$  (PM<sub>2.5</sub>) is increasing due to economic growth, air pollution, and forest fires in some states in the United States. Although previous studies have attempted to retrieve the spatial and temporal behavior of PM<sub>2.5</sub> using aerosol remote sensing and geostatistical estimation methods the coarse resolution and accuracy limit these methods. In this paper the performance of machine learning models on predicting PM<sub>2.5</sub> is assessed with Linear Regression (LR), Decision Tree (DT), Gradient Boosting Regression (GBR), AdaBoost Regression (ABR), XG Boost (XGB), k-nearest neighbors (KNN), Long Short-Term Memory (LSTM), Random Forest (RF), and support vector machine (SVM) using PM<sub>2.5</sub> station data from 2017-2021. To compare the accuracy of all the nine machine learning models the coefficient of determination ( $R^2$ ), root mean square error (RMSE), Nash-Sutcliffe efficiency (NSE), root mean square error ratio (RSR), and percent bias (PBIAS) were evaluated. Among all nine models the RF and SVM models were the best for predicting PM<sub>2.5</sub> concentrations. Comparison of the PM<sub>2.5</sub> performance metrics displayed that the models had better predictive behavior in the western United States than that in the eastern United States.

42

### 43 **1. Introduction:**

44

45 Air pollution has had negative effects on human health and has interfered with social functions;  
46 particles with diameters less than  $2.5 \mu\text{m}$  ( $\text{PM}_{2.5}$ ) have especially been the primary pollutants in  
47 many cities in the USA. Among air pollutants,  $\text{PM}_{2.5}$  is among the most harmful and can easily cross  
48 the human defense barrier, enter the lungs, and cause human disease and even death because of its  
49 small particle size and potential for long-term exposure (Wu et al., 2018; Chen et al., 2019c; Wei et  
50 al., 2019). The  $\text{PM}_{2.5}$  observations were from environmental monitoring stations, however, the  
51 quantity of available  $\text{PM}_{2.5}$  data presented regional differences due to the uneven station distribution.  
52 He et al. (2016) conducted research that indicates the  $\text{PM}_{2.5}$  pollution index was positively correlated  
53 with the emergency admission rate of female acute myocardial infarction and with the increased  
54 incidence of diabetes and hypertension. According to the latest urban air quality database, 98% of  
55 low and middle income countries with more than 100,000 inhabitants do not meet the World Health  
56 Organization (WHO) air quality guidelines [2].

57 Several researchers have used satellite remote sensing data for spatial monitoring coverage in  
58 their studies to estimate  $\text{PM}_{2.5}$  concentrations (Fang et al., 2016; Hu et al., 2017; Park et al., 2019).  
59 One way of using remote sensing satellites for estimating  $\text{PM}_{2.5}$  levels is through the aerosol optical  
60 depth (AOD) parameter, which refers to the solar radiation attenuation due to the scattering and  
61 absorption characteristics of aerosols within the atmosphere (Hutschison et., 2005; Van Donkelaar et  
62 al., 2010; Soni et al., 2018). Wang and Christopher (2003) was the first estimated  $\text{PM}_{2.5}$  using AOD  
63 measurements from Moderate Resolution Imaging Spectrometer (MODIS). Several researchers noted  
64 that satellite AOD as well as monitoring sources and transport of aerosols are key variables in  
65 estimating  $\text{PM}_{2.5}$  and air quality (Gupta and Christopher, 2009). Most have used linear regression  
66 models to correlate AOD and  $\text{PM}_{2.5}$  (Gupta and Christopher, 2009). Grahremanloo et al., 2021

67 examined seasonal behavior of PM<sub>2.5</sub> over Texas using the Random Forest model. Liu et al. (2005)  
68 studied PM<sub>2.5</sub> levels in three different areas such as urban, suburban, and county in the Eastern United  
69 States using multiple linear regression (MLR). They concluded that the model performance may  
70 decrease since the satellite images have a relatively coarse spatial resolution since each pixel  
71 represents a large area on the ground.

72 The design of a model for time series prediction focuses on the application of algorithms to predict  
73 future events based on past trends. The model captures the variables with certain assumptions and  
74 represents the existing dynamic relations, summarizing them to better understand the process that  
75 produced the past data to better predict the future. Most of the above studies have used linear and  
76 non-linear regressions to correlate various parameters with PM<sub>2.5</sub> concentrations over a particular  
77 region. In our study we focused on the entire United States and predicted PM<sub>2.5</sub> concentrations over  
78 various regions using different machine learning models.

79 Recently, due to an increase in the application of machine learning models to various fields  
80 in order to increase the accuracy of predictions, machine learning has also been used to predict particle  
81 concentrations (Kuremoto et al., 2014; Ong et al., 2016; Gui et al., 2020). However, the data mining  
82 does not only differ from one study to another but also in terms of classification algorithms and used  
83 features. The regression, boosting models, and deep learning-based methods display remarkable  
84 performance in time-series data processing to make predictions (Hochreiter and Schmidhuber, 1997).  
85 The estimation using traditional statistical methods requires a large amount of historical data to  
86 construct the relationship between explanatory variables and target variables (Breiman, 2001b). Since  
87 machine learning is a very promising tool to forecast pollution, we proposed applying this approach  
88 to predict PM<sub>2.5</sub> concentrations in the USA. The model predictions based on ML algorithms were  
89 checked by cross-validation and evaluated using appropriate metrics such as root mean square  
90 (RMSE) and mean absolute error (MAE).

91 Earlier studies used a limited number of statistical models, but in our study, we used nearly six  
92 machine learning models to find the best accuracy of predictions. In addition to this, our research  
93 paper took a novel approach in PM<sub>2.5</sub> concentration research by exploring concentrations over USA  
94 as opposed to China where many existing PM<sub>2.5</sub> studies have already been conducted. The purpose of  
95 this paper is to present the predictions of PM<sub>2.5</sub> over different states over the USA. The data collection  
96 and different machine learning techniques applied in the context of time series predictions are adopted  
97 for the present study as described in Section 2. Results and discussion are given in Section 3 and  
98 finally the overall conclusions are drawn from the present study presented in Section 4.

## 99 **2. Datasets:**

### 100 **2.1 Ground PM<sub>2.5</sub> Measurements:**

101 Daily PM<sub>2.5</sub> observational data was collected from January 2015 to December 2021 from the  
102 openaq air quality database (<https://openaq.org/>). These datasets are available from nearly 1081  
103 stations around the USA. The PM<sub>2.5</sub> concentrations of ground sites were taken as the dependent  
104 variable of the model. In this paper, the daily PM<sub>2.5</sub> concentration data of 1081 ground monitoring  
105 stations were sorted in to monthly and seasonal data from January 2015 to December 2021, and the  
106 data integrity exceeded 97%. The datasets were calibrated and quality-controlled according to  
107 national standards. Figure 1 shows the ground-level monitoring site coverage over the United States;  
108 these sites collected 7 years of daily continuous observations. From this figure, we can see that PM<sub>2.5</sub>  
109 monitoring sites are greater in number in the eastern part than in the western part of USA. We  
110 observed small data gaps and therefore applied linear interpolation for filling the gaps of PM<sub>2.5</sub>  
111 datasets. However, stations are sparsely located, therefore ground level PM<sub>2.5</sub> monitoring sites face  
112 difficulties in meeting the data requirements (Lin et al., 2015). As expected, the PM<sub>2.5</sub> concentrations  
113 were much lower at remote sites compared to urban areas, mainly due to the absence of anthropogenic  
114 sources.

115 This study aims to achieve the best statistical comparison of nine machine learning models: Linear  
116 Regression, K-Nearest Neighbors Regressor, Logistic Regression, Gradient Boosting Regressor, Ada  
117 Boost Regressor, Decision Tree Regressor, XG Boost, Support Vector Regressor, Random Forest,  
118 Support Vector Machine, and LSTM for estimating the PM<sub>2.5</sub> concentrations over the specified  
119 period. The datasets are split into 80% and 20% as training and testing datasets, respectively. The  
120 training datasets are used to build the model, and the testing dataset is used to verify the model  
121 performance of the trained model.

## 122 **2.2 K Nearest Neighbors (K-NN):**

123 The K-NN model is one of the earliest ML models (reference). The K-NN model categorizes each  
124 unknown instance in the training set by choosing the majority class label among its k nearest  
125 neighbors. Its performance is also crucially dependent on the Euclidean distance metric used to define  
126 the most immediate neighbors. After determining the Euclidean distance between the data, the  
127 database samples are sorted in ascending order from the least distance (maximal similarity) to  
128 maximum distance (minimal similarity) [Wu et al. 2008]. The k nearest distances are looked at, and  
129 the highest occurring class label of these k nearest points to the instance is decided to be the class  
130 label of the previously unknown instance in the training set. Selecting an optimal value of k becomes  
131 challenging since too low of a value for k can result in overfitting while a larger value of k can cause  
132 the opposite to occur.

## 133 **2.3 Random Forest (RF):**

134 RF is a machine learning algorithm and was proposed by Breiman (2001); it integrates multiple  
135 trees through the idea of ensemble learning, utilizes classification and regression tree (CART) as  
136 learning algorithms of decision trees. The RF is a set of decision trees, where the structure of each  
137 one, and the space of the variables is divided into smaller subspaces so that the data in each region is  
138 as uniform as possible [Hastie et al., 2005 and Breiman, 2001]. It uses the bootstrap resampling

139 technique to randomly extract  $k$  samples (with replacement) from the original training set to generate  
140 new training samples. RF uses multiple base classifiers to obtain higher accuracy classification results  
141 by voting or averaging. RF excels because of its ability to leverage several different independent  
142 decision trees in order to classify better, thereby reducing the error from using a single decision tree  
143 because oftentimes viewing classification in independent directions can lead to lower error than a  
144 single decision tree's direction.

#### 145 **2.4 XGBoost:**

146 This is a highly efficient and optimized distributed gradient boosting algorithm. XGBoost  
147 supports a range of different predictive modeling problems such as classification and regression. It is  
148 trained by minimizing the loss of an objective function against a dataset, and the loss function is a  
149 critical hyperparameter which is tied directly to the type of problem being solved. Regular gradient  
150 boosting, stochastic gradient boosting, and regularized gradient boosting are the three main forms of  
151 gradient boosting. For efficiency, the system features include parallelization, distributed computing,  
152 out-of-core computing, cache optimization, and optimization of data structures to achieve the best  
153 global minimum and run time.

#### 154 **2.5 Long Short-Term Memory (LSTM):**

155 LSTM is well suited for prediction based on time-series data, with better performance, to learn  
156 long-term dependency, and it deals with exploding and vanishing gradient problems [Alahi et al.,  
157 2016, Kong et al., 2017]. LSTM is superior to traditional ML methods in processing large input data  
158 and is a type of Recurrent Neural Network (RNN) [Rumelhart et al., 1986], that has been proposed  
159 to predict future outputs using past inputs. LSTM is great at processing time-series data because the  
160  $PM_{2.5}$  concentrations are time-dependent, and it can better predict future air pollution concentrations  
161 by learning features contained in past air pollution concentration time-series data.

#### 162 **2.6 Decision Tree (DT):**

163 Decision Trees are one of the most commonly used machine learning models in classification and  
164 regression problems. To split a node into two or more sub-nodes DT uses mean squared error (MSE).  
165 It is a tree structure with three types of nodes. The root node is the initial node, which may get split  
166 into further nodes of the branched tree that finally leads to a terminal node (leaf node) that represents  
167 the prediction or final outcome of the model. The interior nodes and branches represent features of  
168 a data set and decision rules respectively. The final prediction is the average of the value of the  
169 dependent variable in that particular leaf node.

### 170 **2.7 Gradient Boosting Regression (GBR):**

171 The type of boosting that combines simple models called weak learners into a single composite  
172 model. Gradient boosting involves optimizing the loss function and a weak learner which makes  
173 predictions. Generally, the gradient descent procedure is used to minimize a set of parameters, such  
174 as coefficients in a regression equation or weights in a neural network. After estimating loss or error,  
175 the weights are updated to minimize that error. Gradient Boosting algorithms minimize the bias error  
176 of the model. The Gradient Boosting algorithm predicts the target variable using a regressor and Mean  
177 Square Error (MSE) as the cost function (for regression problems) or predicts the target variable with  
178 a classifier using a Log Loss cost function (for classification problems).

### 179 **2.8 Support Vector Regression (SVR):**

180 The SVR model is widely applied to time series prediction problems. It is a novel forecasting  
181 approach, which is trained independently based on the same training data with different targets. The  
182 SVR can be used with functions that are linear or non-linear (called kernel functions). The linear  
183 function is used for the linear regression model and evaluates results with metrics such as Root Mean  
184 Square Error (RMSE) and Mean Absolute Error (MAE) to estimate the performance of the model.

### 185 **2.9 AdaBoost Regressor (ABR):**

186 AdaBoost (Adaptive Boosting) is a popular technique, as it combines multiple weak classifiers to  
187 build one strong classifier. The boosting approach is a class of ensembles of ML algorithms and is  
188 described by Schapire (1990). Generally, the boosting approach requires a large amount of training  
189 data which is not possible for many cases, and one way of mitigating this issue is by using AdaBoost  
190 (Freund and Schapire, 1997). The main difference of AdaBoosting from most of the other boosting  
191 approaches is in computing loss functions using relative error rather than absolute error. AdaBoost  
192 regressor fits the data set and adjusts the weights according to the error rate of the current prediction,  
193 and reduces the bias as well as the variance for supervised learning.

#### 194 **2.10 Linear Regression:**

195 Linear Regression is a great statistical tool that achieves to model and predict variables by fitting  
196 the predicted values to the observed values with a straight line or surface. This fitting process is  
197 implemented by reducing the average perpendicular distance from the straight line/surface (which  
198 are the predictions) to the observed values which oftentimes are scattered. The lower this  
199 perpendicular distance, the better the line of best fit; based on this line of best fit's equation future  
200 values can be predicted. In this case, the line of best fit's equation uses the  $PM_{2.5}$  values as the  
201 dependent and output variable whereas time is the independent variable.

#### 202 **3.0 Results and Discussion:**

203 Before proceeding to apply machine learning models on the  $PM_{2.5}$  data we will first discuss the  
204  $PM_{2.5}$  concentrations monthly mean structures, a common method of data exploration to better  
205 understand the data and potentially adjust hyperparameters of the models. Figure 2 shows the USA  
206 monthly anomalies and quantiles for four years using daily  $PM_{2.5}$  values. The monthly anomalies are  
207 in percent form, so we subtracted 100 to set the average value to zero. In addition, we estimated the  
208 anomaly to be positive or negative. Using anomalies we estimated the minimum, maximum values,  
209 the 25%, 75% quantiles, and the interquartile ranges for each month of the entire time period, and the

210 resultant plot is shown in Figure 2. During 2018, in USA, the highest levels of PM<sub>2.5</sub> were observed  
211 in the inland locations and they declined nearly 20% in the year 2019. In the inland areas, PM<sub>2.5</sub>  
212 concentrations are primarily influenced by the secondary particles' formation resulting from the  
213 oxidation of gaseous precursors (NO<sub>x</sub>, SO<sub>x</sub>, and NH<sub>3</sub>) (South Coast Air Quality Management  
214 District, 2017). PM<sub>2.5</sub> concentrations show a drastic change before and during pandemic years. Before  
215 pandemic years the PM<sub>2.5</sub> concentrations are higher in the spring and summer months especially  
216 towards the end of summer (August) and early fall (September) during summer years.

217 The monthly PM<sub>2.5</sub> concentrations are greatest in 2018 when compared to other years. The  
218 positive anomalies are observed on a higher frequency in August 2018 whereas negative anomalies  
219 are observed more in September 2018. This indicates that before COVID-19 the PM<sub>2.5</sub> concentrations  
220 were a little higher than in other years throughout the USA. PM<sub>2.5</sub> values were also higher in the  
221 Eastern USA than in Western USA (Figure not shown). The decrease was moderate (in absolute and  
222 relative terms) in urban areas and progressively became lower from the urban to the rural sites. From  
223 our review of recent sources, primary traffic emissions are highest at traffic sites in absolute and  
224 relative terms (Masiol et al., 2015; Khan et al., 2016, Pietrogrande et al., 2016). Before proceeding  
225 with applying machine learning models to the data, a preliminary statistical analysis was performed  
226 for each state's PM<sub>2.5</sub> values and all time series values were freed of trend and outliers. This was done  
227 because otherwise the time-series data values would give rise to several issues during training like  
228 overfitting or significantly decreasing the performance of the model. The seasonal and annual  
229 variations were removed from all states' time series data points from the entire time period. This  
230 ensured stationarity in the time series data, which is a preprocessing prerequisite before applying  
231 different machine learning algorithms. This is because it is better to observe statistical properties of  
232 a time series which do not change over time, since statistical properties would have to be averaged  
233 for the entire time period, which is not as accurate.

### 234 **3.1 Evaluation Parameters:**

235 For model evaluation, the errors between the estimated and true values were evaluated using  
236 several evaluation indices (Chadalawada & Babovic 2017; Shahid et al., 2018; Yi et al., 2019). The  
237 statistical metrics selected for comparing the performance of the models and error-values between  
238 computed and observed data are evaluated by Root Mean Square Error (RMSE): square root of the  
239 mean squared differences between observed and predicted, and suggests the dispersion of the sample.  
240 Smaller RMSE indicates better performance, and as performance decreases, the RMSE increases.  
241 The coefficient of determination ( $R^2$ ) indicates the collinearity (relationship) between the observed  
242 and predicted data. The  $R^2$  value ranges from 0 to 1 (Santhi et al., 2001 and Van Liew et al., 2003).  
243 Mean absolute error (MAE): average of the absolute differences between the observed and predicted  
244 values where a small value of MAE indicates better performance. Mean absolute percentage error  
245 (MAPE): this index indicates the ratio between errors and observations, the lower the MAPE the  
246 higher the accuracy (Chen et al., 2018). Root mean square error ratio (RSR): the ratio of the RMSE  
247 to the standard deviation of measured data (Stajkowski et al., 2020). RSR is classified into four  
248 intervals: very good ( $0.0 \leq RSR \leq 0.50$ ), good ( $0.50 < RSR \leq 0.60$ ), acceptable ( $0.60 < RSR \leq 0.70$ ),  
249 and unacceptable ( $RSR > 0.70$ ), respectively (Khosravi et al., 2018). Nash-Sutcliffe efficiency (NSE):  
250 is a normalized statistical metric to determine the relative magnitude of the residual variance relative  
251 to the variance or noise (Nash and Sutcliffe 1970). NSE performance ratings are very good ( $0.75 <$   
252  $NSE \leq 1.0$ ), good ( $0.65 < NSE \leq 0.75$ ), satisfactory ( $0.50 < NSE \leq 0.65$ ), and unsatisfactory ( $NSE \leq$   
253  $0.50$ ). Percent bias (PBIAS): it measures the average percent of the predicted value that is smaller or  
254 larger than the observed value (Malik et al., 2018; Nury et al., 2017). The PBIAS is classified into  
255 four ranges, very good ( $PBIAS < \pm 10$ ), good ( $\pm 10 \leq PBIAS < \pm 15$ ), satisfactory ( $\pm 15 \leq PBIAS <$   
256  $\pm 25$ ), and unsatisfactory ( $PBIAS \geq \pm 25$ ).

257 
$$MSE = \frac{\sum_{i=1}^n (x_{oi} - x_{pi})^2}{N}$$

258 
$$MAE = \frac{1}{N} \sum_{i=1}^n |x_{oi} - x_{pi}|$$

259  
260 
$$R^2 = 1 - \frac{\sum_{i=1}^n (x_{oi} - x_{pi})^2}{\sum_{i=1}^n (x_{oi} - x_{mean})^2}$$

261  
262  
263 
$$RSR = \frac{RMSE}{STDEV_{obj}} = \frac{\sqrt{\sum_{i=1}^n (x_{oi} - x_{pi})^2}}{\sqrt{\sum_{i=1}^n (x_{oi} - x_{mean})^2}}$$

264  
265  
266 
$$PBIAS = \left| \frac{\sum_{i=1}^n (x_{oi} - x_{pi})}{\sum_{i=1}^n x_{oi}} \right| * 100$$

267  
268  
269 
$$NORM = \sqrt{\sum_{i=1}^n (x_{oi} - x_{pi})^2}$$

270  
271  
272 
$$MAPE = \frac{\sum_{i=1}^n \frac{|x_{oi} - x_{pi}|}{x_{oi}}}{N} * 100\%$$

273  
274 
$$NSE = 1 - \left[ \frac{\sum_{i=1}^n (x_{oi} - x_{pi})^2}{\sum_{i=1}^n (x_{oi} - x_{mean})^2} \right]$$

275  
276 where N refers to the number of data points,  $x_{oi}$ ,  $x_{pi}$  are the observed and predicted daily  $PM_{2.5}$   
277 concentrations, respectively.

278 The nine machine learning models can describe daily variations of observed and estimated values  
279 of  $PM_{2.5}$  concentrations as shown in Figure 3 and Figure 4, in which the blue curve represents the  
280 observed  $PM_{2.5}$  concentrations, while the red curve represents the estimated  $PM_{2.5}$  concentrations.  
281 We generated time series plots for all states but we showed one state from the western side of the  
282 USA: California (Figure 3) and another state from east USA: New York (Figure 4). All nine machine  
283 learning models show that the seasonal variability of  $PM_{2.5}$  concentration is lower in the spring and

284 summer and higher in autumn and winter, maybe due to atmospheric circulation of autumn and  
285 winter. The  $PM_{2.5}$  concentrations in the autumn and winter are less accurate because air pollution is  
286 more severe than that in spring and summer. The SVM and RF models give better agreement with  
287 observed  $PM_{2.5}$  concentrations. However, the California  $PM_{2.5}$  estimations are less accurate than  
288 those of the New York because pollution is more severe due to forest fires in the summer. Sulfate  
289 concentrations may reflect regional influences of  $PM_{2.5}$ ; these concentrations decreased from east to  
290 west but with higher amounts in California (Meng et al. 2018).

291         Figures 5 and 6 display California and New York's scatter plots of the observed vs estimated  
292 daily  $PM_{2.5}$  concentrations during the period of observations using different machine learning models  
293 respectively. The scatter plot of the two variables suggests a positive linear relationship between  
294 them. All points on the scatter plot lie on a straight line; this indicates the differences are zero and  
295 suggest a strong correlation between the observed and estimated  $PM_{2.5}$  concentrations. Tables 1 and  
296 2 indicate the performance and statistical metrics as estimated for New York and California. The  
297 metrics of all models in Table 1 are for New York: Random Forest with  $R^2 = 0.899$ ,  $MAE = 2.122$ ,  
298 and  $RMSE = 3.121$  has less error than the other models. The next model with the lowest error is  
299 Support Vector Machine with  $R^2 = 0.857$ ,  $MAE = 2.145$ , and  $RMSE = 3.125$ .

300         The performance of the models at different states are good at most sites, as 73% of them show an  
301  $R^2 > 0.62$  and 10% show an  $R^2$  less than 0.3. Moreover, an average  $RMSE$  less than  $4.5 \text{ Mg/m}^3$  in  
302 70% of the states and more than  $5 \text{ Mg/m}^3$  in rest of the states demonstrates good performance.  $PM_{2.5}$   
303 estimations are lower and higher than observations with high and low  $PM_{2.5}$  concentration scenarios  
304 respectively, indicating that estimation accuracy will decline in extreme cases in both states. Zhan et  
305 al. (2017) also found similar behavior using  $PM_{2.5}$  concentration in some parts of China. This may be  
306 due to the model's lack of performance caused by a smaller amount of training data, especially  
307 during extreme  $PM_{2.5}$  concentrations. Ghahremanloo et al. 2021 observed  $PM_{2.5}$  levels in Texas are

308 maximal in the summer and are attributed to higher temperatures and humidity that accelerate the  
309 formation of nitrate and sulfate from NO<sub>2</sub> and SO<sub>2</sub> (Lin et al., 2019). Overall, the performance of RF  
310 is reasonable, with California's R<sup>2</sup>, RMSE, and MAE values of 0.77, 3.051 mg/m<sup>3</sup>, and 2.233 mg/m<sup>3</sup>,  
311 respectively. New York's R<sup>2</sup>, RMSE, and MAE values were 0.899, 3.121 mg/m<sup>3</sup>, and 2.12 mg/m<sup>3</sup>,  
312 respectively. Comparing California's to New York's results, we observe that the California PM<sub>2.5</sub>  
313 concentration values and biases were slightly higher. Overall, the average error values are slightly  
314 lower in the Eastern states than in the Western states. Each state's R<sup>2</sup>, RMSE, MAE, and bias values  
315 are estimated for each model and we observed RF and SVM models produce better estimates than  
316 the other models. On average, the R<sup>2</sup> of the SVM model is 5% higher than that of the RF model. The  
317 biases are 15% lower in the Eastern states than in the Western states of the USA. The high sulfate  
318 concentrations around Los Angeles and Long Beach may be due to the ship emissions, since these  
319 two areas combined have one-fourth of all container cargo traffic in the United States  
320 (<http://www.dot.ca.gov>) (Vutukuru and Dabdub, 2008). However, the PM<sub>2.5</sub> estimations in the  
321 autumn and winter are less accurate because air pollution is more severe than that present in the spring  
322 and summer. Among the nine machine learning models, only the SVM and RF models give desirable  
323 results in the mildest air pollution cases. The LSTM model performs the outperformed among all  
324 models, which can neither reflect the variations of PM<sub>2.5</sub> concentrations significantly nor estimate the  
325 PM<sub>2.5</sub> concentrations accurately.

326 A Taylor diagram can display multiple metrics in a single plot and can be used to summarize  
327 the relative skill with several states' PM<sub>2.5</sub> model outputs. The Taylor diagram characterizes the  
328 statistical relationship between two fields (Taylor, 2001). In this paper, observed is representing the  
329 values based on observations, and predicted indicates that the values were simulated by a machine  
330 learning model. Figures 7 and 8 illustrate the Random Forest and Support Vector Machine of standard  
331 deviation and correlation of all states of USA. Metrics of RF and SVM were computed at each state,

332 and a number was assigned to each state considered. The position of each number appearing on the  
333 plot quantifies how closely model  $PM_{2.5}$  values matches with different states. Consider state 50, for  
334 example and its correlation is about 0.78. The centered standard deviation difference between the  
335 observed and predicted patterns is proportional to the point on the x-axis identified as observed. The  
336 dotted line contours indicate the normalized standard deviation values, and it can be seen that in the  
337 case of state 50 it is centered at about 1.65. Predicted patterns that agree well with observed test data  
338 will lie nearest to the observed marked point. The state values lie near or on the observed dotted line,  
339 and it indicates a small predicted pattern difference. Some of the state values are slightly further from  
340 the observed value, it also shows that the predicted values are larger than the observed.

#### 341 **4. Conclusion:**

342 In this paper, we present the prediction of  $PM_{2.5}$  concentrations over USA using various machine  
343 learning algorithms with the goal of improving our understanding of the differences among them.  
344 Machine learning algorithms are new approaches for analyzing large datasets due to the  
345 computational speed and easy implementation for massive data. In this paper we studied and  
346 examined nine machine learning models (Linear Regression, Decision Tree, Gradient Boost, Ada  
347 Boost, XG Boost, K-Nearest Neighbors, LSTM, Random Forest, and SVM) and their performance  
348 in predicting  $PM_{2.5}$  concentrations.

349 The obtained machine learning-based methods' accuracies vary in all of USA's states, but the  
350 performance of RF (California:  $R^2=0.77$ ,  $NSE = 0.817$ ,  $PBIAS=7.022$ , and  $RSR=0.355$ ; New York:  
351  $R^2=0.899$ ,  $NSE=0.811$ ,  $PBIAS=2.989$ , and  $RSR=0.331$ ) and SVM (California:  $R^2=0.71$ ,  $NSE=0.897$ ,  
352  $PBIAS=7.027$ , and  $RSR=0.424$ ; New York:  $R^2= 0.857$ ,  $NSE=0.280$ ,  $PBIAS=3.011$ , and  $RSR=0.338$ )  
353 were better than the other examined methods. Moreover, it should be noted that the accuracy and  
354 performance of these machine learning methods are not constant in different climates and regions.

355 Both RF and SVM models'  $R^2$  scores were between 0.71 and 0.899, RMSE scores ranged between  
356 3.05 to 3.714, NSE values ranged between 0.811 to 0.899, PBIAS ranged between 2.989-7.027, and  
357 RSR scores ranged between 0.331-0.424 for California and New York states. These metrics revealed  
358 high model reliability and performed well for both RF and SVM and larger datasets produced better  
359 prediction results.

360 Our study can also contribute to limiting human health exposure risks and helping future  
361 epidemiological studies of air pollution. With the improved computational efficiency, machine  
362 learning models improved prediction performance and served as a better scientific tool for decision-  
363 makers to make sound  $PM_{2.5}$  control policies. Real-time measurements of the chemical composition  
364 of  $PM_{2.5}$  taken as regulatory air quality measurements are needed in the future.

365 Several parameters affect  $PM_{2.5}$  concentrations; in the future, it is possible to improve the  
366 performance of our machine learning models with GDP per capita, urbanization data, and other  
367 atmospheric parameters which would be investigated for model development. In the United States  
368 more extensive ground monitoring is needed, as the total number of stations is 1000, suggesting the  
369 network of stations is too sparse for a large nation (See Figure 1). This becomes much more apparent  
370 in some states as also displayed in Figure 1. However, understanding the spatial and temporal  
371 distribution of each region over the United States is helpful, especially over rural areas. Considering  
372 these areas, a larger amount of data for these locations and other ground-based locations would  
373 enhance predicting  $PM_{2.5}$  concentrations. Furthermore, the machine learning models can always be  
374 updated to yield better results as new data becomes available, therefore, the expansion of sources of  
375 data becomes even more important as models can be updated.

376 **Acknowledgements:** The first author (PPV) acknowledges the Jet Propulsion Laboratory (JPL) for  
377 providing him the opportunity with their summer internship program. Author JHJ conducted research

378 at the Jet Propulsion Laboratory and California Institute of Technology under contract by NASA. We  
379 sincerely acknowledge the open air quality group for providing PM<sub>2.5</sub> station data used in this study.  
380 **Data availability:** All PM<sub>2.5</sub> data used for this study can be downloaded from the public website  
381 <https://openaq.org>. For additional questions regarding the data sharing, please contact the  
382 corresponding author at [Jonathan.H.Jiang@jpl.nasa.gov](mailto:Jonathan.H.Jiang@jpl.nasa.gov).

### 383 **References:**

- 384 Alahi, A., K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, “Social LSTM:  
385 Human trajectory prediction in crowded spaces”, in Proc. IEEE Conf. Comput.Vis. Pattern  
386 Recognit., Jun, 2016, pp.961-971.
- 387 Breiman, L. Random Forests, Mach. Learn. 2001, 45, 5-32. [https://doi.org/10.1023/A:](https://doi.org/10.1023/A:1010933404324)  
388 1010933404324.
- 389 Breiman, L., 2001b, Statistical modeling: the two cultures. Stat. Sci., 16 (3), 199-215,  
390 <https://doi.org/10.1214/ss/1009213726>.
- 391 Chadalawada, J., and Babovic, V., 2017. Review and comparison of performance indices for  
392 automatic model induction. J. of Hydroinformatics, 21, 13-31,  
393 <https://doi.org/10.2166/hydro.2017.078>.
- 394 Chen, S., Li, D.C., Zhang, H.Y., Yu, D.K., Chen, R., Zhang, B., Tan, Y.F. et al., 2019c. The  
395 development of a cell-based model for the assessment of carcinogenic potential upon long-term  
396 PM<sub>2.5</sub> exposure. Environ. Int. 131, <https://doi.org/10.1016/j.envint.2019.104943>.
- 397 Fang, X., Zou, B., Liu, X., Sternberg, T., Zhai, I., 2016. Satellite-based ground PM<sub>2.5</sub>  
398 estimation using timely structure adaptive modeling. Rem. Sens. Environ, 186, 152-163.
- 399 Freund, Y., Schapire, R.E., 1997. A decision-theoretic generalization of on-line learning and an  
400 application to boosting, J. computer and System Sciences, 55 (1), 119-139.

401 Ghahremanloo, M., Lops, Y., Choi, Y., Mousavinezhad, S., 2021. Impact of the COVID-19  
402 outbreak on air pollution levels in East Asia. *Sci. Total Environ.* 142226.

403 Gui, K., Che, H., Zeng, Z., Wang, Y., Zhai, S., Wang, Z., Luo, M., Zhang, L., Liao, T., Zhao,  
404 H., Li, L., Zheng, Y., Zhang, X., 2020. Construction of a virtual PM<sub>2.5</sub> observation network in  
405 China based on high-density surface meteorological observations using the extreme gradient  
406 boosting model. *Environ. Int.* 141, 105801. <https://doi.org/10.1016/j.envint.2020.105801>.

407 Gupta, P., Christopher, S.A., 2009. Particulate matter air quality assessment using integrated  
408 surface, satellite, and meteorological products: 2. A neural network approach. *J. Geophys. Res.*  
409 *Atmosphere* 114 (D20).

410 Hastie, T.; Tibshirani, R.; Friedman, J.; Franklin, J. The elements of statistical learning: Data mining,  
411 inference and prediction. *Math. Intell.* 2005, 27, 83-85.

412 He, X. N., Chen, P., Zhang, C., Chen, J.Y. Study on the correlation between PM<sub>2.5</sub> and onset of acute  
413 myocardial infarction among female patients. *Child Care China* 31, 22, 4626-4629, 2016.

414 Hu, X., Belle, J.H., Meng, X., Wildani, A., Waller, L.A., et al. Estimating PM<sub>2.5</sub> concentrations in  
415 the conterminous United States using the Random Forest approach. *Environ, Sci., Technol.* 2017,  
416 51, 6936-6944. <https://doi.org.10.1021/acs.est.7b01210>.

417 Hochreiter, S., and Schmidhuber, J., 1997. Long short-term memory, *Neural Comput.*, 9, 8, 1735-  
418 1780.

419 Hutschison, K.D., Smith, S., Faruqui, S.J., 2005. Correlating MODIS aerosol optical thickness data  
420 with ground-based PM<sub>2.5</sub> observations across Texas for use in a real time air-quality prediction  
421 system. *Atmos. Environ.* 39 (37), 7190-7203.

422 Lin, C., Li, Y., Yuan, Z., Lau, A.K.H., Li, C., Fung, J.C.H., 2015. Using satellite remote sensing  
423 data to estimate the high-resolution distribution of ground-level PM<sub>2.5</sub>. *Remote Sens. Environ.*  
424 156, 117-128. <https://doi.org/10.1016/j.rse.2014.09.015>.

425 Khan, M.B., Masiol, M., Forementon, G., Gilio, A.D., de Gennaaro, G., Agostinelli, C., and  
426 Pavoni, B, 2016. Carboneous PM2.5 and secondary organic aerosol across the Veneto region (NE  
427 Italy). *Sci. Total Environ.* 542, 172-181, doi:10.1016/j.scitotenv.2015.10.103.

428 Khosravi, K ; Mao, L; Kisi, O; Yaseen, Z. M; Shahid, S. Quantifying hourly suspended sediment load  
429 using data mining models: case study of a glacierized Andean catchment in Chile. *J. Hydrol.*  
430 2018, 567, 165-179.

431 Kong, W., Z. Y. Dong, Y. Jia, D. J. Hill, Y. Xu, and Y. Zhang. “Short-term residential load forecasting  
432 based on LSTM recurrent neural network”, *IEEE Trans. Smart Grid*, vol. 10, no.1, pp. 841-851,  
433 Jan. 2017.

434 Kuremoto, T., Kimura, S., Kobayashi, K., and Obayashi, M., 2014. Time series forecasting using a  
435 deep belief network with restricted Boltzmann machines, *Neurocomputing*, 137, 47-56.

436 Liu, B.,; Philip, S. Y.; Top 10 algorithms in data mining. *Knowl. Inf. Syst.* 2008, 14, 1-37.

437 Liu, Y., Sarnat, J.A., Kilaru, V., Jacob, D.J., Koutrakis, P., 2005. Estimating ground-level PM2.5 in  
438 the eastern United States using satellite remote sensing, *Environ. Sci. Technol.*, 39, 3269-3278.

439 Malik, A.; Kumar. A.; Kisi, O. Daily pan evaporation estimation using heuristic methods with gamma  
440 test. *J. Irrig. Drain. Eng.* 2018, 144, 4018023.

441 Masiol, M., Benetello, F., Harrisom, R.M., Fornenton, G., Gaspari, F.D., and Pavoni, B., 2015.  
442 Spatial, seasonal trends and trans-boundary transport of PM2.5 inorganic ions in the Veneto  
443 region (northeastern Italy), *Atmos. Environ.*, 117, 19-31, doi:10.1016/j.atmosenv.2015.06.044.

444 Meng, X., Garay, M.J., Diner, D.J., Kalashnikova, O.V., Xu, J., Liu, Y., 2018. Estimating PM2.5  
445 speciation concentrations using prototype 4.4 km resolution misr aerosol properties over Southern  
446 California, *Atmos. Environ.*, 181, 70-81.

447 Nash J. E., Sutcliffe, J. V. River flow forecasting through conceptual models part I – A discussion of  
448 principles. *J. Hydrol.* 1970, 10, 282-290.

449 Nury, A.H.; Hasan, K.; Alam, M. J. Bin comparative study of wavelet-ARIMA and wavelet-ANN  
450 models for temperature time series data in northeastern Bangladesh. *J. King. Saud. Univ. Sci.*  
451 2017, 29, 47-61.

452 Ong, B.T., Sugiura, K., and Zettsu, K., 2016. Dynamically pre-trained deep recurrent neural networks  
453 using environmental monitoring for predicting PM2.5, *Neural Comput. Appl.*, 27, 6, 1553-1566.

454 Park, Y., Kwon, B., Heo, J., Hu, X., Liu, Y., Moon, T. 2019. Estimating PM2.5 concentration of the  
455 conterminous Unites states via interpretable convolutional neural networks. *Environ. Pollut.*  
456 113395.

457 Pietrogrande, M.C., Bacco, D., Ferrari, S., Ricciardelli, I., Scotto, F., Trentini, A., and Visentin, M.:  
458 2016. Characteristics and major sources of carbonaceous aerosols in PM2.5 in Emilia Romagna  
459 Region (Northern Italy) from four-year observations. *Sci. Total Environ.*, 553, 172-183,  
460 doi:10.1016/j.scitotenv.2016.02.074.

461 Santhi, C; Arnold, J. G.; Williams, J. R.; Dugas, W. A; Srinivasan, R.; and Hauck, L. M. Validation  
462 of the swat model on a large river basin with point and non-point sources, *JAWRA. J. Am. Water*  
463 *Resour. Assoc.*, 2001, 37, 1169-1188.

464 Schapire, R. (1990). The strength of weak learnability. *Machine Learning*, 5 (2), 197-227.

465 Freund, Y., Schpire, R (1997). A decision-theoretic generalisation of on-line learning and an  
466 application of boosting. *J. Computer and System Sciences*, 55 (1), 119-139.

467 Soni, M., Payra, S., Verma, S., 2018. Particulate matter estimation over a semi-arid region Jaipur,  
468 India using satellite AOD and meteorological parameters. *Atmospheric Pollution Research* 9 (5),  
469 949-958.

470 Stajkowski, S; Kumar, D; Samui, P; Bonakdari, H; and Gharabaghi, B, Genetic algorithm-optimized  
471 sequential model for water temperature prediction, *Sustainability*, 12, 13, 5374, 2020.

472 Rumelhart, D. E., G. E. Hinton, and R. J. Williams, “learning representations by back-propagating  
473 errors,” *Nature*, Vol. 323, no.6088, pp. 533-536, 1986.

474 Taylor, K.E., Summarizing multiple aspects of model performance in a single diagram, *J. Geophys.*  
475 *Res.*, 106, 7183-7192, 2001.

476 Van Donkelaar, A., Martin, R.V., Brauer, M., Kahn, R., Levy, R., Verduzco, C., Villeneuve, P.J.,  
477 2010. Global estimates of ambient fine particulate matter concentrations from satellite-based  
478 optical depth: development and application. *Environ. Health Perspect.* 118 (6), 847-855.

479 Van Liew, M. W; Arnold, J. G.; Garbrecht, J. D. Hydrologic simulation on agricultural watersheds:  
480 Choosing between two models. *Trans. ASAE* 2003, 56, 1539.

481 Vutukuru, S., Dabdub, D., 2008. Modeling the effects of ship emissions on coastal air quality: a case  
482 study of Southern California. *Atmos. Environ.* 42, 3751-3764.

483 Wang, J., Christopher, S.A., 2003. Intercomparison between satellite-derived aerosol optical  
484 thickness and PM<sub>2.5</sub> mass: Implications for air quality studies. *Geophys. Res. Lett.* 30 (21).

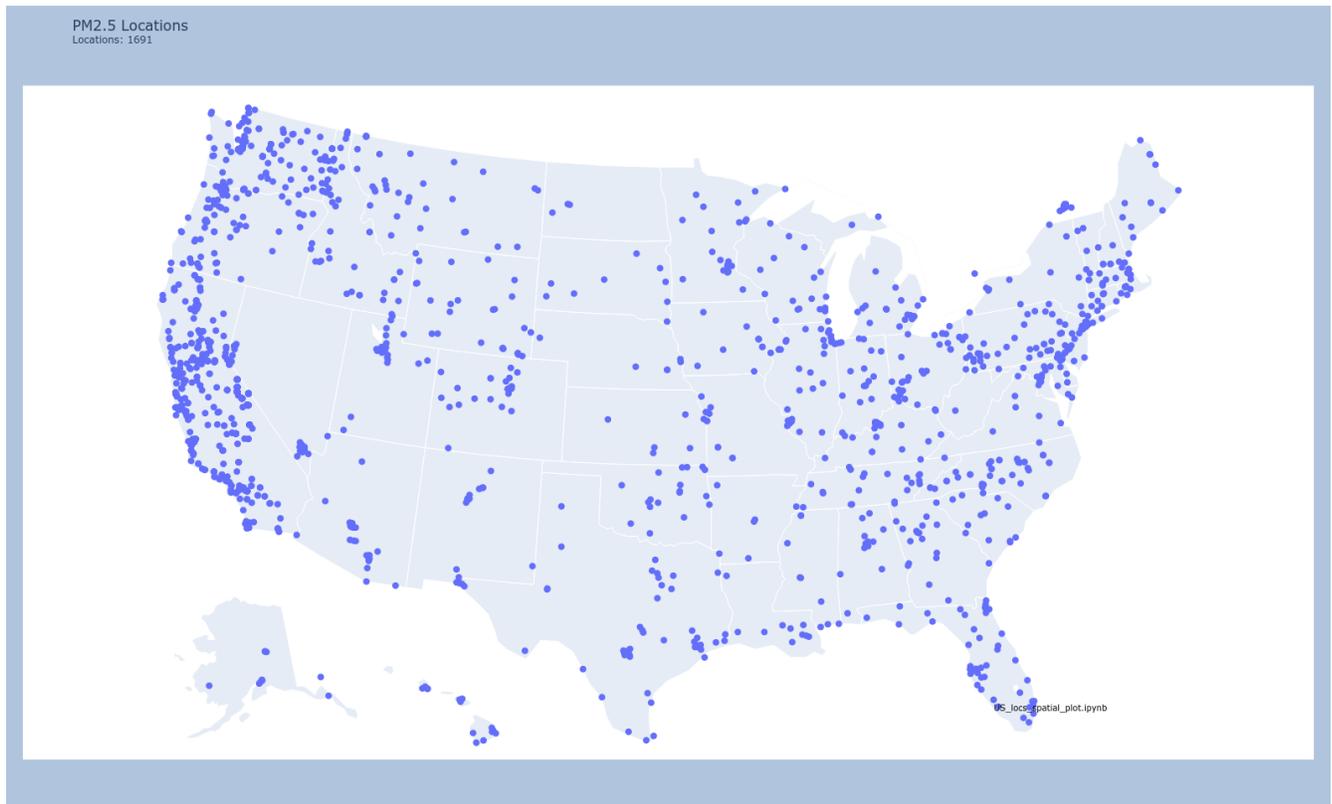
485 Wei, J., Huang, W., Li,, Z., Xue, W., Peng, Y., Sun, L., Gribb, M., 2019. Estimating 1-km resolution  
486 PM<sub>2.5</sub> concentrations across China using space-time random forest approach. *Rem. Sens.*  
487 *Environ.* 231, 111221.

488 World Health Organization, media centre (2016). Air pollution levels are rising in many of the  
489 world’s poorest cities: <http://www.int/mediacentre/news/releases/2016/air-pollution-raising/>.

490 Wu, X.; Kumar, V.; Quinlan, J. R.; Ghosh, J.; Yang, Q.; Motoda, H.; McLachlan, G. J.; Ng, A.;

491 Yi. L., Mengfan, T., Kun, Y., Yu, Z., Xiaolu, Z., Miao, Z., Yan, S., 2019. Research on PM<sub>2.5</sub>  
492 estimation and prediction method and changing characteristics analysis under long temporal and  
493 large spatial scale – a case study in China typical regions. *Sci. Total Environ.* 696, 133983,  
494 <https://doi.org/10.1016/j.scitotenv.2019.133983>.

495 Zhang, Y., Cao,F., 2015. Fine particle matter (PM<sub>2.5</sub>)in China at a city level. *Sci. Rep.*, 5, 14884.



497

498

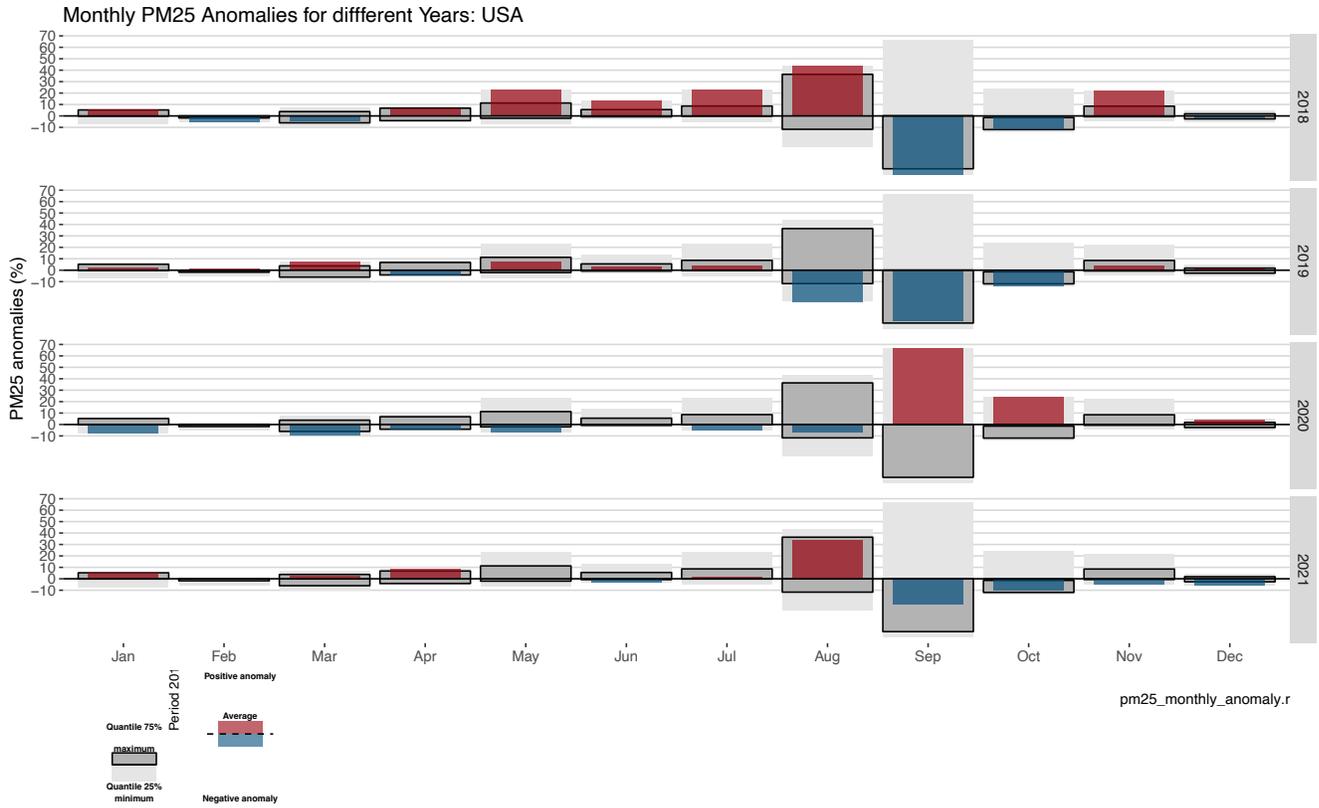
**Figure 1.** Locations of PM<sub>2.5</sub> monitoring sites over USA

499

500

501

502



503

504 **Figure 2.** Monthly anomalies and quantiles for the observed period (2018-2021) using daily PM<sub>2.5</sub>

505 values over United States.

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

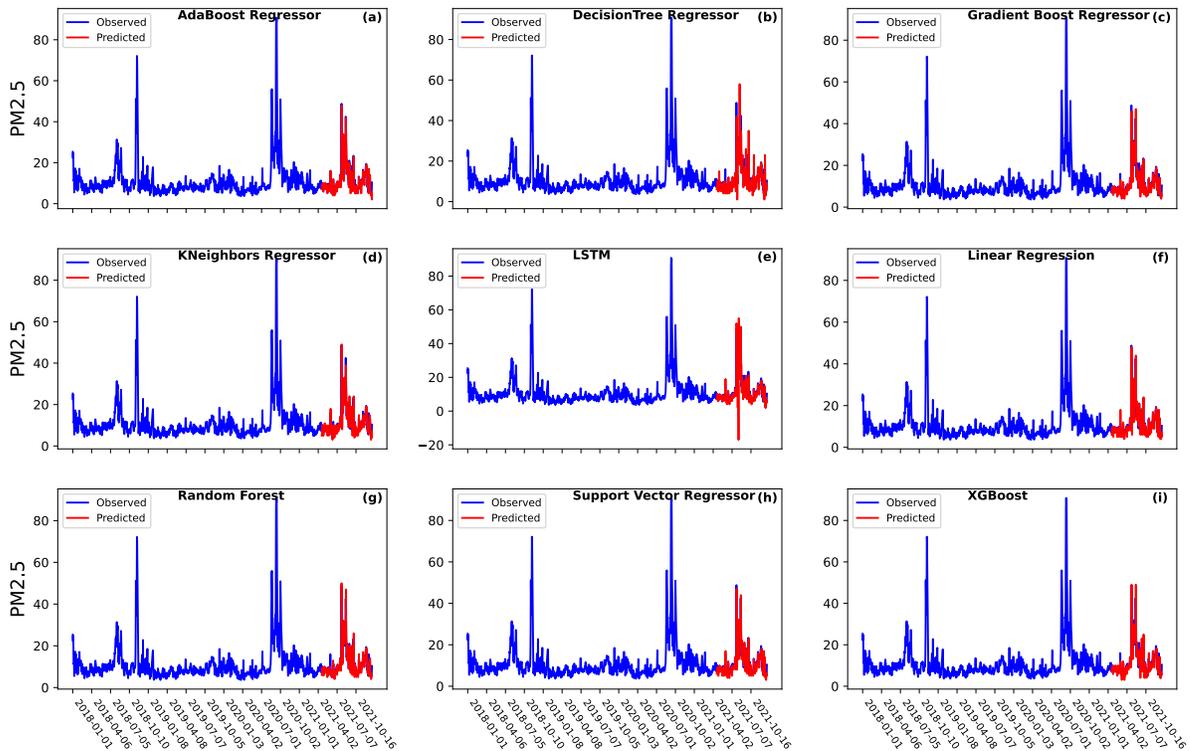
522

523

524

525

526



PM25\_dly\_dfrntmdl\_obs\_pred\_lineplt.ipynb

527

528

529

530

531

532

533

534

535

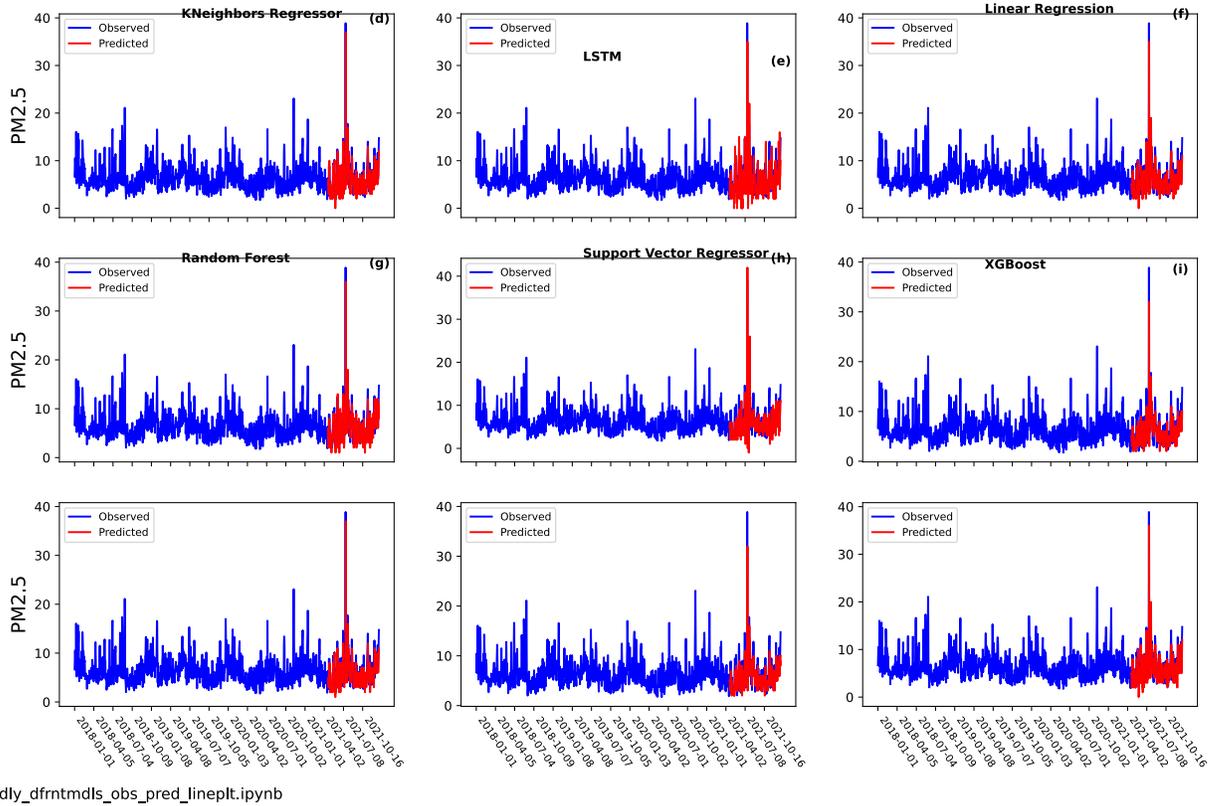
536

537

538

539

**Figure 3.** The comparison of the time series of estimated and observed  $PM_{2.5}$  concentrations over California using different machine learning models: (a) AdaBoost regressor, (b) Decision Tree regression, (c) Gradient Boost regression, (d) K-neighbors regression (e) LSTM, (f) Linear regression, (g) Random Forest, (h) Support Vector regression, and (I) XGBoost.



540

541 **Figure 4.** The comparison of the time series of estimated and observed  $PM_{2.5}$  concentrations over  
 542 New York using different machine learning models: (a) AdaBoost regressor, (b) DecisionTree  
 543 regression, (c) Gradient Boost regression, (d) Kneighbors regression (e) LSTM, (f) Linear regression,  
 544 (g) Random Forest, (h) Support Vector regression, and (I) XGBoost.

546

547

548

549

550

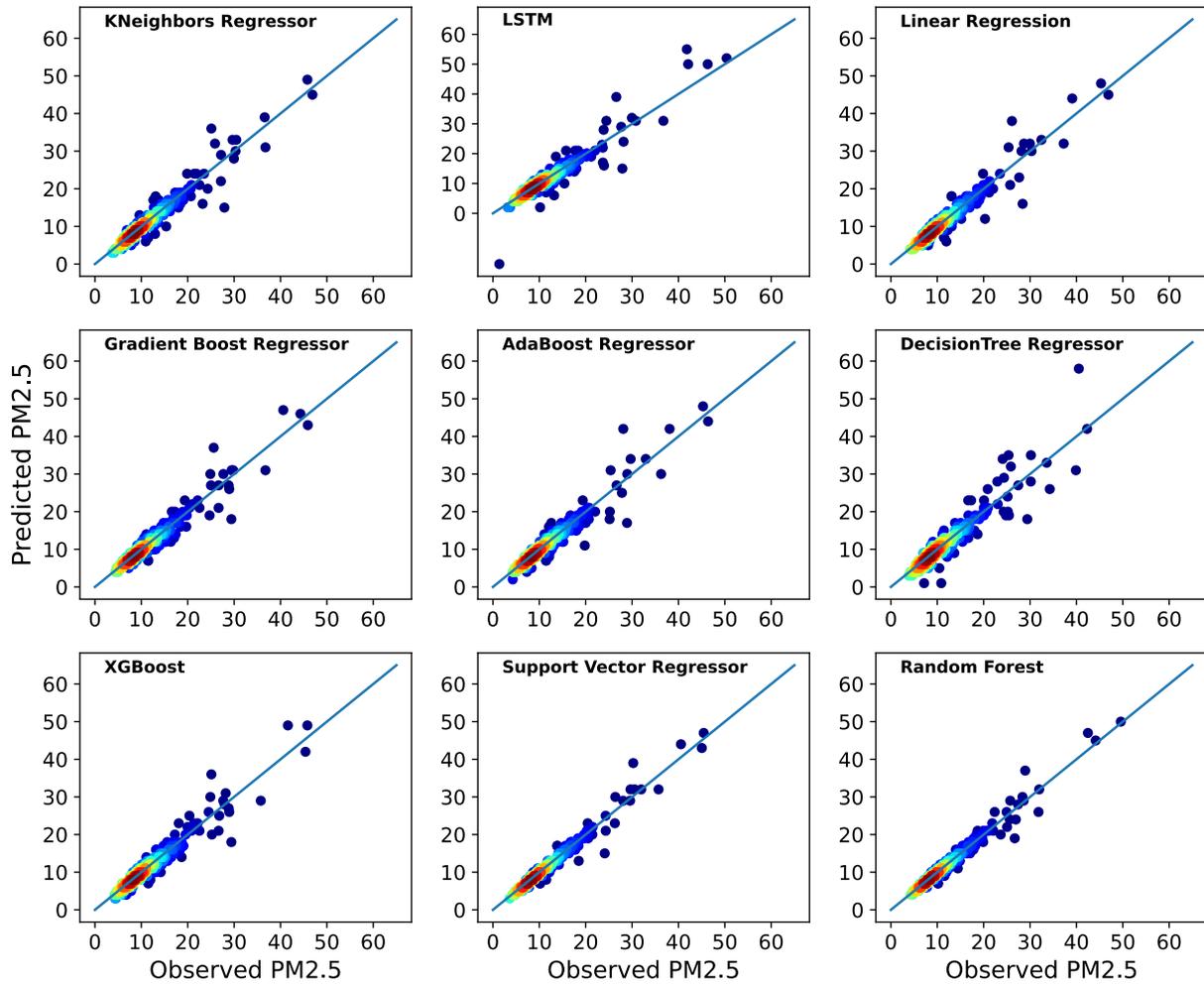
551

552

553

554

555



556

557 **Figure 5.** Scatter plots of observed and estimated daily  $PM_{2.5}$  concentrations over California using  
558 different machine learning models: (a) AdaBoost regressor, (b) DecisionTree regression, (c) Gradient  
559 Boost regression, (d) Kneighbors regression (e) LSTM, (f) Linear regression, (g) Random Forest, (h)  
560 Support Vector regression, and (I) XGBoost.

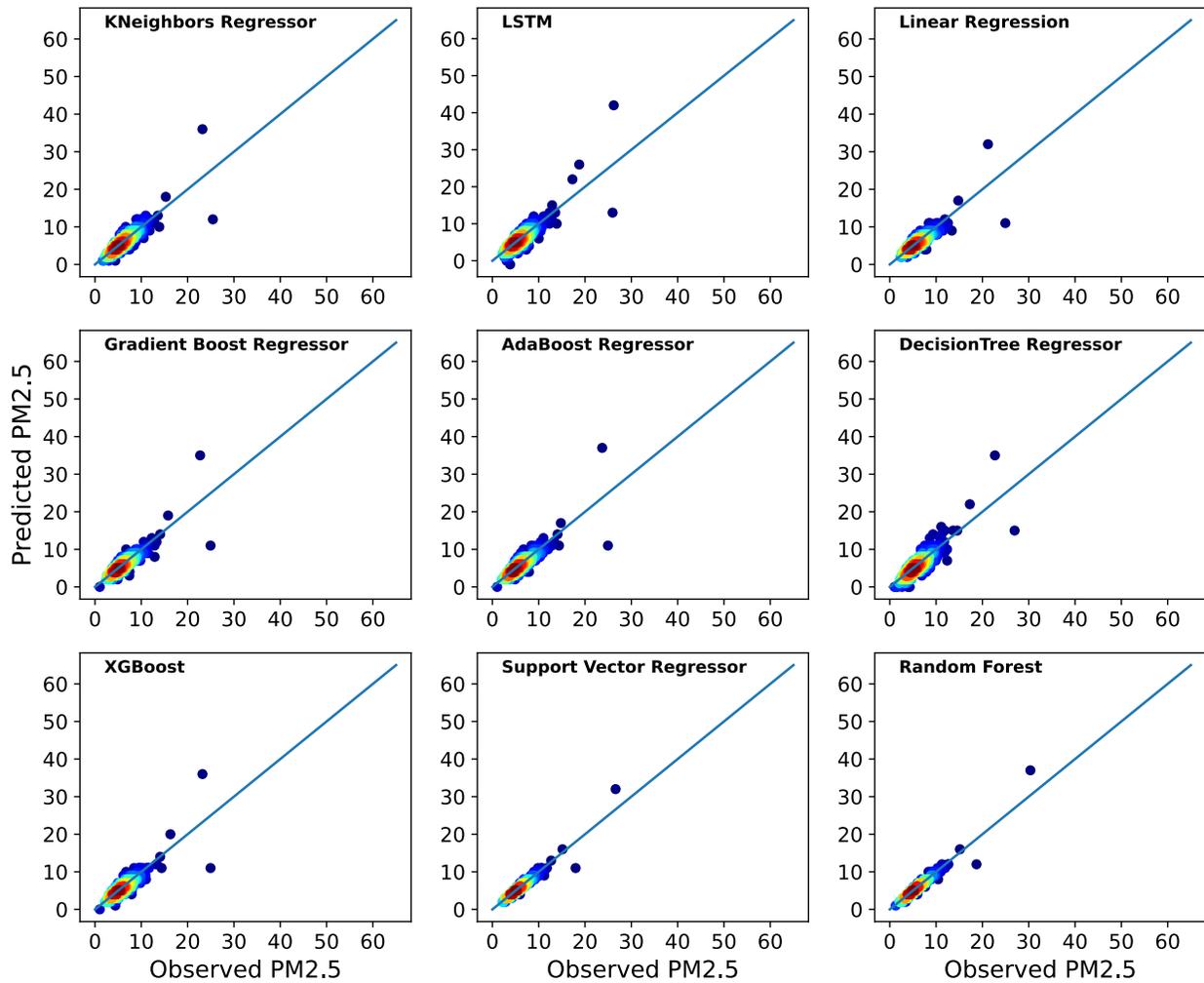
561

562

563

564

565



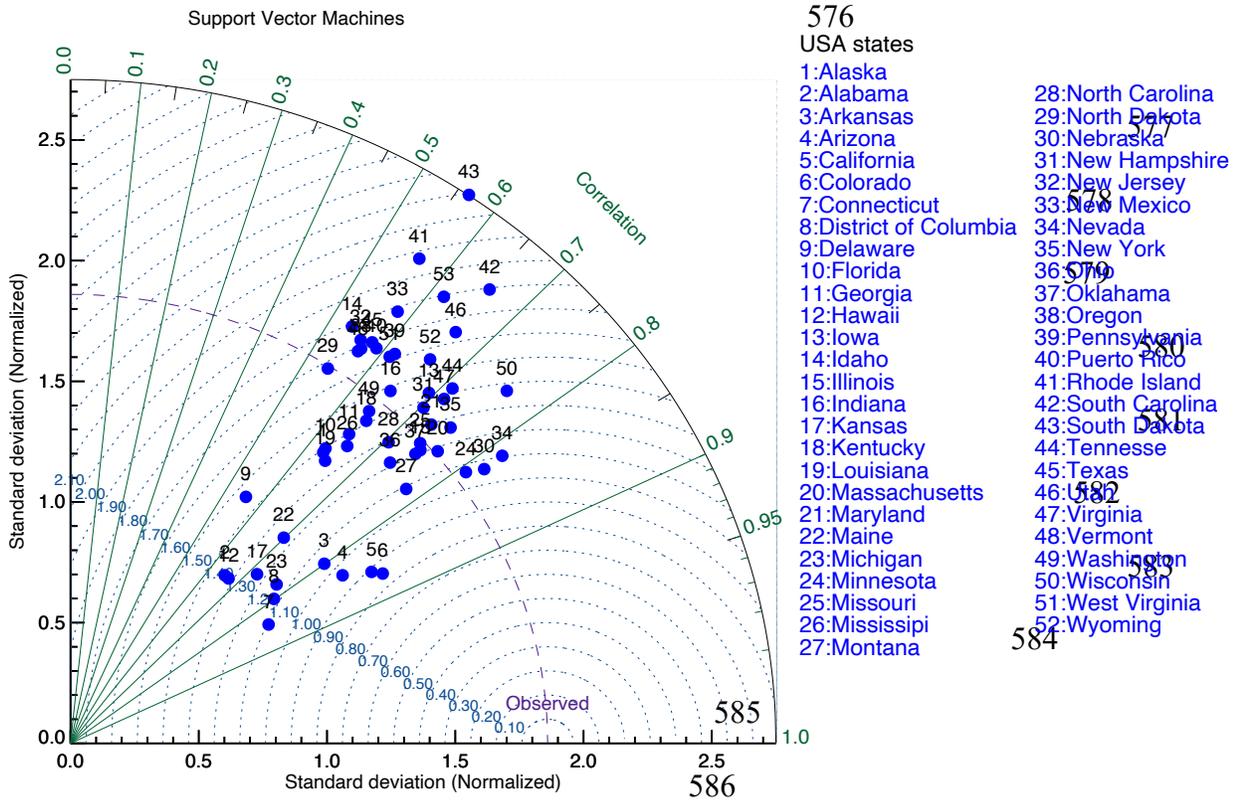
567

568 **Figure 6.** Scatter plots of observed and estimated daily  $PM_{2.5}$  concentrations over New York using  
 569 different machine learning models: (a) AdaBoost regressor, (b) DecisionTree regression, (c) Gradient  
 570 Boost regression, (d) Kneighbors regression (e) LSTM, (f) Linear regression, (g) Random Forest,  
 571 (h) Support Vector regression, and (I) XGBoost.  
 572

573

574

575



576

USA states

- 1:Alaska
- 2:Alabama
- 3:Arkansas
- 4:Arizona
- 5:California
- 6:Colorado
- 7:Connecticut
- 8:District of Columbia
- 9:Delaware
- 10:Florida
- 11:Georgia
- 12:Hawaii
- 13:Iowa
- 14:Idaho
- 15:Illinois
- 16:Indiana
- 17:Kansas
- 18:Kentucky
- 19:Louisiana
- 20:Massachusetts
- 21:Maryland
- 22:Maine
- 23:Michigan
- 24:Minnesota
- 25:Missouri
- 26:Mississippi
- 27:Montana
- 28:North Carolina
- 29:North Dakota
- 30:Nebraska
- 31:New Hampshire
- 32:New Jersey
- 33:New Mexico
- 34:Nevada
- 35:New York
- 36:Ohio
- 37:Oklahoma
- 38:Oregon
- 39:Pennsylvania
- 40:Puerto Rico
- 41:Rhode Island
- 42:South Carolina
- 43:South Dakota
- 44:Tennessee
- 45:Texas
- 46:Utah
- 47:Virginia
- 48:Vermont
- 49:Washington
- 50:Wisconsin
- 51:West Virginia
- 52:Wyoming

taylor\_pm25mod.pro

587

588

589 **Figure 7.** Taylor diagram of the Support Vector Machines (SVM) over each state of the United States.

590

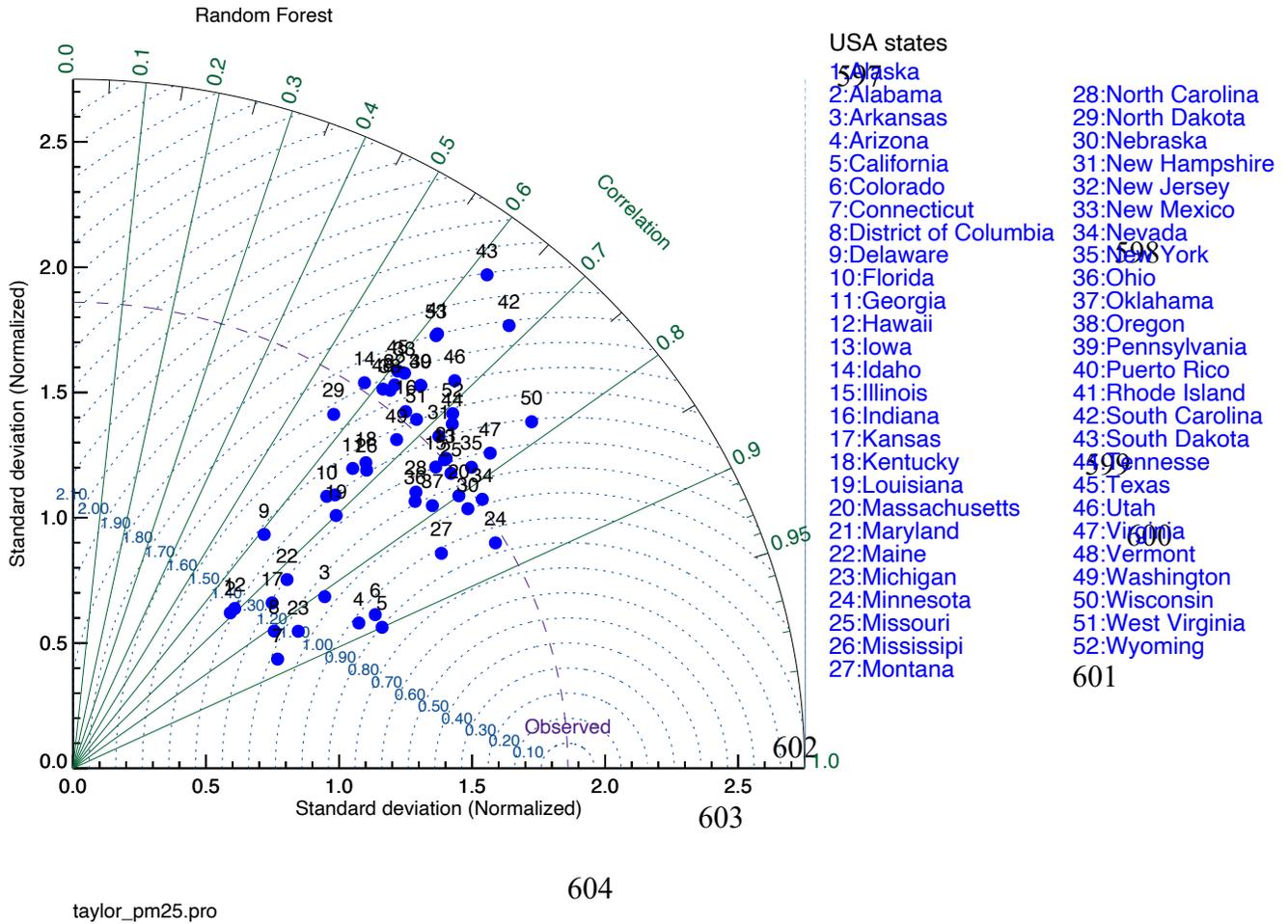
591

592

593

594

595



605

606 **Figure 8.** Taylor diagram of the Random Forest (RF) over each state of the United States.

607

608

609

610

611

612 **Table 1: Different Model Metrics for New York State**

New York								
Model	RMSE	MAE	MAPE	R2	NSE	NORM	PBIAS	RSR
<b>Linear Regression</b>	3.883	2.309	0.285	0.688	0.613	60.156	11.24	0.561
<b>Decision Tree</b>	5.136	3.109	0.254	0.454	0.533	79.58	13.44	0.691
<b>Gradient Boost Regressor</b>	3.822	2.394	0.545	0.698	0.683	59.207	8.210	0.546
<b>AdaBoost Regressor</b>	3.961	2.316	0.188	0.676	0.683	61.369	9.653	0.576
<b>XG Boost</b>	3.898	2.501	0.202	0.686	0.681	60.393	8.342	0.559
<b>KNeighbors Regressor</b>	3.919	2.379	0.195	0.683	0.677	60.711	7.515	0.562
<b>LSTM</b>	7.487	3.359	0.218	0.158	0.455	115.991	6.020	0.812
<b>Random Forest</b>	3.121	2.122	0.182	0.899	0.811	38.671	2.989	0.331
<b>SVM</b>	3.125	2.145	0.183	0.857	0.820	39.161	3.011	0.338

- 613
- 614 RMSE = Root mean squared error
- 615 MAE = Mean absolute error
- 616 MAPE = Mean absolute percentage error
- 617  $R^2$  = The coefficient of determination
- 618 NSE = Nash-Sutcliffe efficiency
- 619 PBIAS = Percent Bias
- 620 RSR = root mean square error ratio

621  
622  
623 **Table 2: Different Model Metrics for California State**

California								
Model	RMSE	MAE	MAPE	R <sup>2</sup>	NSE	NORM	PBIAS	RSR
<b>Linear Regression</b>	3.695	2.599	0.326	0.43	0.694	57.243	12.086	0.932
<b>Decision Tree</b>	5.481	3.743	0.467	0.23	0.576	84.917	19.901	0.732
<b>Gradient Boost Regressor</b>	4.051	2.736	0.340	0.28	0.461	62.758	16.891	1.017
<b>AdaBoost Regressor</b>	3.804	2.636	0.342	0.33	0.435	58.938	17.532	0.969
<b>XG Boost</b>	4.271	2.972	0.372	0.17	0.438	66.178	18.726	1.075
<b>KNeighbors Regressor</b>	4.394	3.062	0.392	0.22	0.286	68.071	17.076	1.106
<b>LSTM</b>	5.025	3.252	0.339	0.46	0.309	77.853	18.027	0.618
<b>Random Forest</b>	3.051	2.233	0.315	0.77	0.817	46.894	7.022	0.355
<b>SVM</b>	3.714	2.618	0.320	0.71	0.897	47.853	7.027	0.424

- 624
- 625 RMSE = Root mean squared error
- 626 MAE = Mean absolute error
- 627 MAPE = Mean absolute percentage error
- 628  $R^2$  = The coefficient of determination
- 629 NSE = Nash-Sutcliffe efficiency
- 630 PBIAS = Percent Bias
- 631 RSR = root mean square error ratio
- 632

# Predicting PM<sub>2.5</sub> Concentrations Across USA Using Machine Learning

P. Preetham Vignesh<sup>1</sup>, Jonathan H. Jiang<sup>2</sup>, P. Kishore

<sup>1</sup> University of California, Los Angeles, USA

<sup>2</sup> Jet Propulsion Laboratory, California Institute of Technology, Pasadena, California, USA.

<sup>3</sup> Retired, University of California, Irvine, USA

Copyright ©2022, All Rights Reserved.

Correspondence: Jonathan.H.Jiang@jpl.nasa.gov

Keywords: Surface Temperature, Climate Model, Global Warming Projection

## Abstract:

Fine particulate matter with a size less than 2.5  $\mu\text{m}$  (PM<sub>2.5</sub>) is increasing due to economic growth, air pollution, and forest fires in some states in the United States. Although previous studies have attempted to retrieve the spatial and temporal behavior of PM<sub>2.5</sub> using aerosol remote sensing and geostatistical estimation methods the coarse resolution and accuracy limit these methods. In this paper the performance of machine learning models on predicting PM<sub>2.5</sub> is assessed with Linear Regression (LR), Decision Tree (DT), Gradient Boosting Regression (GBR), AdaBoost Regression (ABR), XG Boost (XGB), k-nearest neighbors (KNN), Long Short-Term Memory (LSTM), Random Forest (RF), and support vector machine (SVM) using PM<sub>2.5</sub> station data from 2017-2021. To compare the accuracy of all the nine machine learning models the coefficient of determination ( $R^2$ ), root mean square error (RMSE), Nash-Sutcliffe efficiency (NSE), root mean square error ratio (RSR), and percent bias (PBIAS) were evaluated. Among all nine models the RF and SVM models were the best for predicting PM<sub>2.5</sub> concentrations. Comparison of the PM<sub>2.5</sub> performance metrics displayed that the models had better predictive behavior in the western United States than that in the eastern United States.

42

43 **1. Introduction:**

44

45 Air pollution has had negative effects on human health and has interfered with social functions;  
46 particles with diameters less than  $2.5 \mu\text{m}$  ( $\text{PM}_{2.5}$ ) have especially been the primary pollutants in  
47 many cities in the USA. Among air pollutants,  $\text{PM}_{2.5}$  is among the most harmful and can easily cross  
48 the human defense barrier, enter the lungs, and cause human disease and even death because of its  
49 small particle size and potential for long-term exposure (Wu et al., 2018; Chen et al., 2019c; Wei et  
50 al., 2019). The  $\text{PM}_{2.5}$  observations were from environmental monitoring stations, however, the  
51 quantity of available  $\text{PM}_{2.5}$  data presented regional differences due to the uneven station distribution.  
52 He et al. (2016) conducted research that indicates the  $\text{PM}_{2.5}$  pollution index was positively correlated  
53 with the emergency admission rate of female acute myocardial infarction and with the increased  
54 incidence of diabetes and hypertension. According to the latest urban air quality database, 98% of  
55 low and middle income countries with more than 100,000 inhabitants do not meet the World Health  
56 Organization (WHO) air quality guidelines [2].

57 Several researchers have used satellite remote sensing data for spatial monitoring coverage in  
58 their studies to estimate  $\text{PM}_{2.5}$  concentrations (Fang et al., 2016; Hu et al., 2017; Park et al., 2019).  
59 One way of using remote sensing satellites for estimating  $\text{PM}_{2.5}$  levels is through the aerosol optical  
60 depth (AOD) parameter, which refers to the solar radiation attenuation due to the scattering and  
61 absorption characteristics of aerosols within the atmosphere (Hutschison et., 2005; Van Donkelaar et  
62 al., 2010; Soni et al., 2018). Wang and Christopher (2003) was the first estimated  $\text{PM}_{2.5}$  using AOD  
63 measurements from Moderate Resolution Imaging Spectrometer (MODIS). Several researchers noted  
64 that satellite AOD as well as monitoring sources and transport of aerosols are key variables in  
65 estimating  $\text{PM}_{2.5}$  and air quality (Gupta and Christopher, 2009). Most have used linear regression  
66 models to correlate AOD and  $\text{PM}_{2.5}$  (Gupta and Christopher, 2009). Grahremanloo et al., 2021

67 examined seasonal behavior of PM<sub>2.5</sub> over Texas using the Random Forest model. Liu et al. (2005)  
68 studied PM<sub>2.5</sub> levels in three different areas such as urban, suburban, and county in the Eastern United  
69 States using multiple linear regression (MLR). They concluded that the model performance may  
70 decrease since the satellite images have a relatively coarse spatial resolution since each pixel  
71 represents a large area on the ground.

72 The design of a model for time series prediction focuses on the application of algorithms to predict  
73 future events based on past trends. The model captures the variables with certain assumptions and  
74 represents the existing dynamic relations, summarizing them to better understand the process that  
75 produced the past data to better predict the future. Most of the above studies have used linear and  
76 non-linear regressions to correlate various parameters with PM<sub>2.5</sub> concentrations over a particular  
77 region. In our study we focused on the entire United States and predicted PM<sub>2.5</sub> concentrations over  
78 various regions using different machine learning models.

79 Recently, due to an increase in the application of machine learning models to various fields  
80 in order to increase the accuracy of predictions, machine learning has also been used to predict particle  
81 concentrations (Kuremoto et al., 2014; Ong et al., 2016; Gui et al., 2020). However, the data mining  
82 does not only differ from one study to another but also in terms of classification algorithms and used  
83 features. The regression, boosting models, and deep learning-based methods display remarkable  
84 performance in time-series data processing to make predictions (Hochreiter and Schmidhuber, 1997).  
85 The estimation using traditional statistical methods requires a large amount of historical data to  
86 construct the relationship between explanatory variables and target variables (Breiman, 2001b). Since  
87 machine learning is a very promising tool to forecast pollution, we proposed applying this approach  
88 to predict PM<sub>2.5</sub> concentrations in the USA. The model predictions based on ML algorithms were  
89 checked by cross-validation and evaluated using appropriate metrics such as root mean square  
90 (RMSE) and mean absolute error (MAE).

91 Earlier studies used a limited number of statistical models, but in our study, we used nearly six  
92 machine learning models to find the best accuracy of predictions. In addition to this, our research  
93 paper took a novel approach in PM<sub>2.5</sub> concentration research by exploring concentrations over USA  
94 as opposed to China where many existing PM<sub>2.5</sub> studies have already been conducted. The purpose of  
95 this paper is to present the predictions of PM<sub>2.5</sub> over different states over the USA. The data collection  
96 and different machine learning techniques applied in the context of time series predictions are adopted  
97 for the present study as described in Section 2. Results and discussion are given in Section 3 and  
98 finally the overall conclusions are drawn from the present study presented in Section 4.

## 99 **2. Datasets:**

### 100 **2.1 Ground PM<sub>2.5</sub> Measurements:**

101 Daily PM<sub>2.5</sub> observational data was collected from January 2015 to December 2021 from the  
102 openaq air quality database (<https://openaq.org/>). These datasets are available from nearly 1081  
103 stations around the USA. The PM<sub>2.5</sub> concentrations of ground sites were taken as the dependent  
104 variable of the model. In this paper, the daily PM<sub>2.5</sub> concentration data of 1081 ground monitoring  
105 stations were sorted in to monthly and seasonal data from January 2015 to December 2021, and the  
106 data integrity exceeded 97%. The datasets were calibrated and quality-controlled according to  
107 national standards. Figure 1 shows the ground-level monitoring site coverage over the United States;  
108 these sites collected 7 years of daily continuous observations. From this figure, we can see that PM<sub>2.5</sub>  
109 monitoring sites are greater in number in the eastern part than in the western part of USA. We  
110 observed small data gaps and therefore applied linear interpolation for filling the gaps of PM<sub>2.5</sub>  
111 datasets. However, stations are sparsely located, therefore ground level PM<sub>2.5</sub> monitoring sites face  
112 difficulties in meeting the data requirements (Lin et al., 2015). As expected, the PM<sub>2.5</sub> concentrations  
113 were much lower at remote sites compared to urban areas, mainly due to the absence of anthropogenic  
114 sources.

115 This study aims to achieve the best statistical comparison of nine machine learning models: Linear  
116 Regression, K-Nearest Neighbors Regressor, Logistic Regression, Gradient Boosting Regressor, Ada  
117 Boost Regressor, Decision Tree Regressor, XG Boost, Support Vector Regressor, Random Forest,  
118 Support Vector Machine, and LSTM for estimating the PM<sub>2.5</sub> concentrations over the specified  
119 period. The datasets are split into 80% and 20% as training and testing datasets, respectively. The  
120 training datasets are used to build the model, and the testing dataset is used to verify the model  
121 performance of the trained model.

## 122 **2.2 K Nearest Neighbors (K-NN):**

123 The K-NN model is one of the earliest ML models (reference). The K-NN model categorizes each  
124 unknown instance in the training set by choosing the majority class label among its k nearest  
125 neighbors. Its performance is also crucially dependent on the Euclidean distance metric used to define  
126 the most immediate neighbors. After determining the Euclidean distance between the data, the  
127 database samples are sorted in ascending order from the least distance (maximal similarity) to  
128 maximum distance (minimal similarity) [Wu et al. 2008]. The k nearest distances are looked at, and  
129 the highest occurring class label of these k nearest points to the instance is decided to be the class  
130 label of the previously unknown instance in the training set. Selecting an optimal value of k becomes  
131 challenging since too low of a value for k can result in overfitting while a larger value of k can cause  
132 the opposite to occur.

## 133 **2.3 Random Forest (RF):**

134 RF is a machine learning algorithm and was proposed by Breiman (2001); it integrates multiple  
135 trees through the idea of ensemble learning, utilizes classification and regression tree (CART) as  
136 learning algorithms of decision trees. The RF is a set of decision trees, where the structure of each  
137 one, and the space of the variables is divided into smaller subspaces so that the data in each region is  
138 as uniform as possible [Hastie et al., 2005 and Breiman, 2001]. It uses the bootstrap resampling

139 technique to randomly extract  $k$  samples (with replacement) from the original training set to generate  
140 new training samples. RF uses multiple base classifiers to obtain higher accuracy classification results  
141 by voting or averaging. RF excels because of its ability to leverage several different independent  
142 decision trees in order to classify better, thereby reducing the error from using a single decision tree  
143 because oftentimes viewing classification in independent directions can lead to lower error than a  
144 single decision tree's direction.

#### 145 **2.4 XGBoost:**

146 This is a highly efficient and optimized distributed gradient boosting algorithm. XGBoost  
147 supports a range of different predictive modeling problems such as classification and regression. It is  
148 trained by minimizing the loss of an objective function against a dataset, and the loss function is a  
149 critical hyperparameter which is tied directly to the type of problem being solved. Regular gradient  
150 boosting, stochastic gradient boosting, and regularized gradient boosting are the three main forms of  
151 gradient boosting. For efficiency, the system features include parallelization, distributed computing,  
152 out-of-core computing, cache optimization, and optimization of data structures to achieve the best  
153 global minimum and run time.

#### 154 **2.5 Long Short-Term Memory (LSTM):**

155 LSTM is well suited for prediction based on time-series data, with better performance, to learn  
156 long-term dependency, and it deals with exploding and vanishing gradient problems [Alahi et al.,  
157 2016, Kong et al., 2017]. LSTM is superior to traditional ML methods in processing large input data  
158 and is a type of Recurrent Neural Network (RNN) [Rumelhart et al., 1986], that has been proposed  
159 to predict future outputs using past inputs. LSTM is great at processing time-series data because the  
160  $PM_{2.5}$  concentrations are time-dependent, and it can better predict future air pollution concentrations  
161 by learning features contained in past air pollution concentration time-series data.

#### 162 **2.6 Decision Tree (DT):**

163 Decision Trees are one of the most commonly used machine learning models in classification and  
164 regression problems. To split a node into two or more sub-nodes DT uses mean squared error (MSE).  
165 It is a tree structure with three types of nodes. The root node is the initial node, which may get split  
166 into further nodes of the branched tree that finally leads to a terminal node (leaf node) that represents  
167 the prediction or final outcome of the model. The interior nodes and branches represent features of  
168 a data set and decision rules respectively. The final prediction is the average of the value of the  
169 dependent variable in that particular leaf node.

### 170 **2.7 Gradient Boosting Regression (GBR):**

171 The type of boosting that combines simple models called weak learners into a single composite  
172 model. Gradient boosting involves optimizing the loss function and a weak learner which makes  
173 predictions. Generally, the gradient descent procedure is used to minimize a set of parameters, such  
174 as coefficients in a regression equation or weights in a neural network. After estimating loss or error,  
175 the weights are updated to minimize that error. Gradient Boosting algorithms minimize the bias error  
176 of the model. The Gradient Boosting algorithm predicts the target variable using a regressor and Mean  
177 Square Error (MSE) as the cost function (for regression problems) or predicts the target variable with  
178 a classifier using a Log Loss cost function (for classification problems).

### 179 **2.8 Support Vector Regression (SVR):**

180 The SVR model is widely applied to time series prediction problems. It is a novel forecasting  
181 approach, which is trained independently based on the same training data with different targets. The  
182 SVR can be used with functions that are linear or non-linear (called kernel functions). The linear  
183 function is used for the linear regression model and evaluates results with metrics such as Root Mean  
184 Square Error (RMSE) and Mean Absolute Error (MAE) to estimate the performance of the model.

### 185 **2.9 AdaBoost Regressor (ABR):**

186 AdaBoost (Adaptive Boosting) is a popular technique, as it combines multiple weak classifiers to  
187 build one strong classifier. The boosting approach is a class of ensembles of ML algorithms and is  
188 described by Schapire (1990). Generally, the boosting approach requires a large amount of training  
189 data which is not possible for many cases, and one way of mitigating this issue is by using AdaBoost  
190 (Freund and Schapire, 1997). The main difference of AdaBoosting from most of the other boosting  
191 approaches is in computing loss functions using relative error rather than absolute error. AdaBoost  
192 regressor fits the data set and adjusts the weights according to the error rate of the current prediction,  
193 and reduces the bias as well as the variance for supervised learning.

#### 194 **2.10 Linear Regression:**

195 Linear Regression is a great statistical tool that achieves to model and predict variables by fitting  
196 the predicted values to the observed values with a straight line or surface. This fitting process is  
197 implemented by reducing the average perpendicular distance from the straight line/surface (which  
198 are the predictions) to the observed values which oftentimes are scattered. The lower this  
199 perpendicular distance, the better the line of best fit; based on this line of best fit's equation future  
200 values can be predicted. In this case, the line of best fit's equation uses the  $PM_{2.5}$  values as the  
201 dependent and output variable whereas time is the independent variable.

#### 202 **3.0 Results and Discussion:**

203 Before proceeding to apply machine learning models on the  $PM_{2.5}$  data we will first discuss the  
204  $PM_{2.5}$  concentrations monthly mean structures, a common method of data exploration to better  
205 understand the data and potentially adjust hyperparameters of the models. Figure 2 shows the USA  
206 monthly anomalies and quantiles for four years using daily  $PM_{2.5}$  values. The monthly anomalies are  
207 in percent form, so we subtracted 100 to set the average value to zero. In addition, we estimated the  
208 anomaly to be positive or negative. Using anomalies we estimated the minimum, maximum values,  
209 the 25%, 75% quantiles, and the interquartile ranges for each month of the entire time period, and the

210 resultant plot is shown in Figure 2. During 2018, in USA, the highest levels of PM<sub>2.5</sub> were observed  
211 in the inland locations and they declined nearly 20% in the year 2019. In the inland areas, PM<sub>2.5</sub>  
212 concentrations are primarily influenced by the secondary particles' formation resulting from the  
213 oxidation of gaseous precursors (NO<sub>x</sub>, SO<sub>x</sub>, and NH<sub>3</sub>) (South Coast Air Quality Management  
214 District, 2017). PM<sub>2.5</sub> concentrations show a drastic change before and during pandemic years. Before  
215 pandemic years the PM<sub>2.5</sub> concentrations are higher in the spring and summer months especially  
216 towards the end of summer (August) and early fall (September) during summer years.

217 The monthly PM<sub>2.5</sub> concentrations are greatest in 2018 when compared to other years. The  
218 positive anomalies are observed on a higher frequency in August 2018 whereas negative anomalies  
219 are observed more in September 2018. This indicates that before COVID-19 the PM<sub>2.5</sub> concentrations  
220 were a little higher than in other years throughout the USA. PM<sub>2.5</sub> values were also higher in the  
221 Eastern USA than in Western USA (Figure not shown). The decrease was moderate (in absolute and  
222 relative terms) in urban areas and progressively became lower from the urban to the rural sites. From  
223 our review of recent sources, primary traffic emissions are highest at traffic sites in absolute and  
224 relative terms (Masiol et al., 2015; Khan et al., 2016, Pietrogrande et al., 2016). Before proceeding  
225 with applying machine learning models to the data, a preliminary statistical analysis was performed  
226 for each state's PM<sub>2.5</sub> values and all time series values were freed of trend and outliers. This was done  
227 because otherwise the time-series data values would give rise to several issues during training like  
228 overfitting or significantly decreasing the performance of the model. The seasonal and annual  
229 variations were removed from all states' time series data points from the entire time period. This  
230 ensured stationarity in the time series data, which is a preprocessing prerequisite before applying  
231 different machine learning algorithms. This is because it is better to observe statistical properties of  
232 a time series which do not change over time, since statistical properties would have to be averaged  
233 for the entire time period, which is not as accurate.

### 234 **3.1 Evaluation Parameters:**

235 For model evaluation, the errors between the estimated and true values were evaluated using  
236 several evaluation indices (Chadalawada & Babovic 2017; Shahid et al., 2018; Yi et al., 2019). The  
237 statistical metrics selected for comparing the performance of the models and error-values between  
238 computed and observed data are evaluated by Root Mean Square Error (RMSE): square root of the  
239 mean squared differences between observed and predicted, and suggests the dispersion of the sample.  
240 Smaller RMSE indicates better performance, and as performance decreases, the RMSE increases.  
241 The coefficient of determination ( $R^2$ ) indicates the collinearity (relationship) between the observed  
242 and predicted data. The  $R^2$  value ranges from 0 to 1 (Santhi et al., 2001 and Van Liew et al., 2003).  
243 Mean absolute error (MAE): average of the absolute differences between the observed and predicted  
244 values where a small value of MAE indicates better performance. Mean absolute percentage error  
245 (MAPE): this index indicates the ratio between errors and observations, the lower the MAPE the  
246 higher the accuracy (Chen et al., 2018). Root mean square error ratio (RSR): the ratio of the RMSE  
247 to the standard deviation of measured data (Stajkowski et al., 2020). RSR is classified into four  
248 intervals: very good ( $0.0 \leq RSR \leq 0.50$ ), good ( $0.50 < RSR \leq 0.60$ ), acceptable ( $0.60 < RSR \leq 0.70$ ),  
249 and unacceptable ( $RSR > 0.70$ ), respectively (Khosravi et al., 2018). Nash-Sutcliffe efficiency (NSE):  
250 is a normalized statistical metric to determine the relative magnitude of the residual variance relative  
251 to the variance or noise (Nash and Sutcliffe 1970). NSE performance ratings are very good ( $0.75 <$   
252  $NSE \leq 1.0$ ), good ( $0.65 < NSE \leq 0.75$ ), satisfactory ( $0.50 < NSE \leq 0.65$ ), and unsatisfactory ( $NSE \leq$   
253  $0.50$ ). Percent bias (PBIAS): it measures the average percent of the predicted value that is smaller or  
254 larger than the observed value (Malik et al., 2018; Nury et al., 2017). The PBIAS is classified into  
255 four ranges, very good ( $PBIAS < \pm 10$ ), good ( $\pm 10 \leq PBIAS < \pm 15$ ), satisfactory ( $\pm 15 \leq PBIAS <$   
256  $\pm 25$ ), and unsatisfactory ( $PBIAS \geq \pm 25$ ).

257 
$$MSE = \frac{\sum_{i=1}^n (x_{oi} - x_{pi})^2}{N}$$

258 
$$MAE = \frac{1}{N} \sum_{i=1}^n |x_{oi} - x_{pi}|$$

259

260 
$$R^2 = 1 - \frac{\sum_{i=1}^n (x_{oi} - x_{pi})^2}{\sum_{i=1}^n (x_{oi} - x_{mean})^2}$$

261

262

263 
$$RSR = \frac{RMSE}{STDEV_{obj}} = \frac{\sqrt{\sum_{i=1}^n (x_{oi} - x_{pi})^2}}{\sqrt{\sum_{i=1}^n (x_{oi} - x_{mean})^2}}$$

264

265

266 
$$PBIAS = \left| \frac{\sum_{i=1}^n (x_{oi} - x_{pi})}{\sum_{i=1}^n x_{oi}} \right| * 100$$

267

268

269 
$$NORM = \sqrt{\sum_{i=1}^n (x_{oi} - x_{pi})^2}$$

270

271

272 
$$MAPE = \frac{\sum_{i=1}^n \frac{|x_{oi} - x_{pi}|}{x_{oi}}}{N} * 100\%$$

273

274 
$$NSE = 1 - \left[ \frac{\sum_{i=1}^n (x_{oi} - x_{pi})^2}{\sum_{i=1}^n (x_{oi} - x_{mean})^2} \right]$$

275

276 where N refers to the number of data points,  $x_{oi}$ ,  $x_{pi}$  are the observed and predicted daily  $PM_{2.5}$   
 277 concentrations, respectively.

278 The nine machine learning models can describe daily variations of observed and estimated values  
 279 of  $PM_{2.5}$  concentrations as shown in Figure 3 and Figure 4, in which the blue curve represents the  
 280 observed  $PM_{2.5}$  concentrations, while the red curve represents the estimated  $PM_{2.5}$  concentrations.  
 281 We generated time series plots for all states but we showed one state from the western side of the  
 282 USA: California (Figure 3) and another state from east USA: New York (Figure 4). All nine machine  
 283 learning models show that the seasonal variability of  $PM_{2.5}$  concentration is lower in the spring and

284 summer and higher in autumn and winter, maybe due to atmospheric circulation of autumn and  
285 winter. The  $PM_{2.5}$  concentrations in the autumn and winter are less accurate because air pollution is  
286 more severe than that in spring and summer. The SVM and RF models give better agreement with  
287 observed  $PM_{2.5}$  concentrations. However, the California  $PM_{2.5}$  estimations are less accurate than  
288 those of the New York because pollution is more severe due to forest fires in the summer. Sulfate  
289 concentrations may reflect regional influences of  $PM_{2.5}$ ; these concentrations decreased from east to  
290 west but with higher amounts in California (Meng et al. 2018).

291         Figures 5 and 6 display California and New York's scatter plots of the observed vs estimated  
292 daily  $PM_{2.5}$  concentrations during the period of observations using different machine learning models  
293 respectively. The scatter plot of the two variables suggests a positive linear relationship between  
294 them. All points on the scatter plot lie on a straight line; this indicates the differences are zero and  
295 suggest a strong correlation between the observed and estimated  $PM_{2.5}$  concentrations. Tables 1 and  
296 2 indicate the performance and statistical metrics as estimated for New York and California. The  
297 metrics of all models in Table 1 are for New York: Random Forest with  $R^2 = 0.899$ ,  $MAE = 2.122$ ,  
298 and  $RMSE = 3.121$  has less error than the other models. The next model with the lowest error is  
299 Support Vector Machine with  $R^2 = 0.857$ ,  $MAE = 2.145$ , and  $RMSE = 3.125$ .

300         The performance of the models at different states are good at most sites, as 73% of them show an  
301  $R^2 > 0.62$  and 10% show an  $R^2$  less than 0.3. Moreover, an average  $RMSE$  less than  $4.5 \text{ Mg/m}^3$  in  
302 70% of the states and more than  $5 \text{ Mg/m}^3$  in rest of the states demonstrates good performance.  $PM_{2.5}$   
303 estimations are lower and higher than observations with high and low  $PM_{2.5}$  concentration scenarios  
304 respectively, indicating that estimation accuracy will decline in extreme cases in both states. Zhan et  
305 al. (2017) also found similar behavior using  $PM_{2.5}$  concentration in some parts of China. This may be  
306 due to the model's lack of performance caused by a smaller amount of training data, especially  
307 during extreme  $PM_{2.5}$  concentrations. Ghahremanloo et al. 2021 observed  $PM_{2.5}$  levels in Texas are

308 maximal in the summer and are attributed to higher temperatures and humidity that accelerate the  
309 formation of nitrate and sulfate from NO<sub>2</sub> and SO<sub>2</sub> (Lin et al., 2019). Overall, the performance of RF  
310 is reasonable, with California's R<sup>2</sup>, RMSE, and MAE values of 0.77, 3.051 mg/m<sup>3</sup>, and 2.233 mg/m<sup>3</sup>,  
311 respectively. New York's R<sup>2</sup>, RMSE, and MAE values were 0.899, 3.121 mg/m<sup>3</sup>, and 2.12 mg/m<sup>3</sup>,  
312 respectively. Comparing California's to New York's results, we observe that the California PM<sub>2.5</sub>  
313 concentration values and biases were slightly higher. Overall, the average error values are slightly  
314 lower in the Eastern states than in the Western states. Each state's R<sup>2</sup>, RMSE, MAE, and bias values  
315 are estimated for each model and we observed RF and SVM models produce better estimates than  
316 the other models. On average, the R<sup>2</sup> of the SVM model is 5% higher than that of the RF model. The  
317 biases are 15% lower in the Eastern states than in the Western states of the USA. The high sulfate  
318 concentrations around Los Angeles and Long Beach may be due to the ship emissions, since these  
319 two areas combined have one-fourth of all container cargo traffic in the United States  
320 (<http://www.dot.ca.gov>) (Vutukuru and Dabdu, 2008). However, the PM<sub>2.5</sub> estimations in the  
321 autumn and winter are less accurate because air pollution is more severe than that present in the spring  
322 and summer. Among the nine machine learning models, only the SVM and RF models give desirable  
323 results in the mildest air pollution cases. The LSTM model performs the outperformed among all  
324 models, which can neither reflect the variations of PM<sub>2.5</sub> concentrations significantly nor estimate the  
325 PM<sub>2.5</sub> concentrations accurately.

326         A Taylor diagram can display multiple metrics in a single plot and can be used to summarize  
327 the relative skill with several states' PM<sub>2.5</sub> model outputs. The Taylor diagram characterizes the  
328 statistical relationship between two fields (Taylor, 2001). In this paper, observed is representing the  
329 values based on observations, and predicted indicates that the values were simulated by a machine  
330 learning model. Figures 7 and 8 illustrate the Random Forest and Support Vector Machine of standard  
331 deviation and correlation of all states of USA. Metrics of RF and SVM were computed at each state,

332 and a number was assigned to each state considered. The position of each number appearing on the  
333 plot quantifies how closely model  $PM_{2.5}$  values matches with different states. Consider state 50, for  
334 example and its correlation is about 0.78. The centered standard deviation difference between the  
335 observed and predicted patterns is proportional to the point on the x-axis identified as observed. The  
336 dotted line contours indicate the normalized standard deviation values, and it can be seen that in the  
337 case of state 50 it is centered at about 1.65. Predicted patterns that agree well with observed test data  
338 will lie nearest to the observed marked point. The state values lie near or on the observed dotted line,  
339 and it indicates a small predicted pattern difference. Some of the state values are slightly further from  
340 the observed value, it also shows that the predicted values are larger than the observed.

#### 341 **4. Conclusion:**

342 In this paper, we present the prediction of  $PM_{2.5}$  concentrations over USA using various machine  
343 learning algorithms with the goal of improving our understanding of the differences among them.  
344 Machine learning algorithms are new approaches for analyzing large datasets due to the  
345 computational speed and easy implementation for massive data. In this paper we studied and  
346 examined nine machine learning models (Linear Regression, Decision Tree, Gradient Boost, Ada  
347 Boost, XG Boost, K-Nearest Neighbors, LSTM, Random Forest, and SVM) and their performance  
348 in predicting  $PM_{2.5}$  concentrations.

349 The obtained machine learning-based methods' accuracies vary in all of USA's states, but the  
350 performance of RF (California:  $R^2=0.77$ ,  $NSE = 0.817$ ,  $PBIAS=7.022$ , and  $RSR=0.355$ ; New York:  
351  $R^2=0.899$ ,  $NSE=0.811$ ,  $PBIAS=2.989$ , and  $RSR=0.331$ ) and SVM (California:  $R^2=0.71$ ,  $NSE=0.897$ ,  
352  $PBIAS=7.027$ , and  $RSR=0.424$ ; New York:  $R^2= 0.857$ ,  $NSE=0.280$ ,  $PBIAS=3.011$ , and  $RSR=0.338$ )  
353 were better than the other examined methods. Moreover, it should be noted that the accuracy and  
354 performance of these machine learning methods are not constant in different climates and regions.

355 Both RF and SVM models'  $R^2$  scores were between 0.71 and 0.899, RMSE scores ranged between  
356 3.05 to 3.714, NSE values ranged between 0.811 to 0.899, PBIAS ranged between 2.989-7.027, and  
357 RSR scores ranged between 0.331-0.424 for California and New York states. These metrics revealed  
358 high model reliability and performed well for both RF and SVM and larger datasets produced better  
359 prediction results.

360 Our study can also contribute to limiting human health exposure risks and helping future  
361 epidemiological studies of air pollution. With the improved computational efficiency, machine  
362 learning models improved prediction performance and served as a better scientific tool for decision-  
363 makers to make sound  $PM_{2.5}$  control policies. Real-time measurements of the chemical composition  
364 of  $PM_{2.5}$  taken as regulatory air quality measurements are needed in the future.

365 Several parameters affect  $PM_{2.5}$  concentrations; in the future, it is possible to improve the  
366 performance of our machine learning models with GDP per capita, urbanization data, and other  
367 atmospheric parameters which would be investigated for model development. In the United States  
368 more extensive ground monitoring is needed, as the total number of stations is 1000, suggesting the  
369 network of stations is too sparse for a large nation (See Figure 1). This becomes much more apparent  
370 in some states as also displayed in Figure 1. However, understanding the spatial and temporal  
371 distribution of each region over the United States is helpful, especially over rural areas. Considering  
372 these areas, a larger amount of data for these locations and other ground-based locations would  
373 enhance predicting  $PM_{2.5}$  concentrations. Furthermore, the machine learning models can always be  
374 updated to yield better results as new data becomes available, therefore, the expansion of sources of  
375 data becomes even more important as models can be updated.

376 **Acknowledgements:** The first author (PPV) acknowledges the Jet Propulsion Laboratory (JPL) for  
377 providing him the opportunity with their summer internship program. Author JHJ conducted research

378 at the Jet Propulsion Laboratory and California Institute of Technology under contract by NASA. We  
379 sincerely acknowledge the open air quality group for providing PM<sub>2.5</sub> station data used in this study.  
380 **Data availability:** All PM<sub>2.5</sub> data used for this study can be downloaded from the public website  
381 <https://openaq.org>. For additional questions regarding the data sharing, please contact the  
382 corresponding author at [Jonathan.H.Jiang@jpl.nasa.gov](mailto:Jonathan.H.Jiang@jpl.nasa.gov).

### 383 **References:**

- 384 Alahi, A., K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, “Social LSTM:  
385 Human trajectory prediction in crowded spaces”, in Proc. IEEE Conf. Comput.Vis. Pattern  
386 Recognit., Jun, 2016, pp.961-971.
- 387 Breiman, L. Random Forests, Mach. Learn. 2001, 45, 5-32. [https://doi.org/10.1023/A:](https://doi.org/10.1023/A:1010933404324)  
388 1010933404324.
- 389 Breiman, L., 2001b, Statistical modeling: the two cultures. Stat. Sci., 16 (3), 199-215,  
390 <https://doi.org/10.1214/ss/1009213726>.
- 391 Chadalawada, J., and Babovic, V., 2017. Review and comparison of performance indices for  
392 automatic model induction. J. of Hydroinformatics, 21, 13-31,  
393 <https://doi.org/10.2166/hydro.2017.078>.
- 394 Chen, S., Li, D.C., Zhang, H.Y., Yu, D.K., Chen, R., Zhang, B., Tan, Y.F. et al., 2019c. The  
395 development of a cell-based model for the assessment of carcinogenic potential upon long-term  
396 PM<sub>2.5</sub> exposure. Environ. Int. 131, <https://doi.org/10.1016/j.envint.2019.104943>.
- 397 Fang, X., Zou, B., Liu, X., Sternberg, T., Zhai, I., 2016. Satellite-based ground PM<sub>2.5</sub>  
398 estimation using timely structure adaptive modeling. Rem. Sens. Environ, 186, 152-163.
- 399 Freund, Y., Schapire, R.E., 1997. A decision-theoretic generalization of on-line learning and an  
400 application to boosting, J. computer and System Sciences, 55 (1), 119-139.

401 Ghahremanloo, M., Lops, Y., Choi, Y., Mousavinezhad, S., 2021. Impact of the COVID-19  
402 outbreak on air pollution levels in East Asia. *Sci. Total Environ.* 142226.

403 Gui, K., Che, H., Zeng, Z., Wang, Y., Zhai, S., Wang, Z., Luo, M., Zhang, L., Liao, T., Zhao,  
404 H., Li, L., Zheng, Y., Zhang, X., 2020. Construction of a virtual PM<sub>2.5</sub> observation network in  
405 China based on high-density surface meteorological observations using the extreme gradient  
406 boosting model. *Environ. Int.* 141, 105801. <https://doi.org/10.1016/j.envint.2020.105801>.

407 Gupta, P., Christopher, S.A., 2009. Particulate matter air quality assessment using integrated  
408 surface, satellite, and meteorological products: 2. A neural network approach. *J. Geophys. Res.*  
409 *Atmosphere* 114 (D20).

410 Hastie, T.; Tibshirani, R.; Friedman, J.; Franklin, J. The elements of statistical learning: Data mining,  
411 inference and prediction. *Math. Intell.* 2005, 27, 83-85.

412 He, X. N., Chen, P., Zhang, C., Chen, J.Y. Study on the correlation between PM<sub>2.5</sub> and onset of acute  
413 myocardial infarction among female patients. *Child Care China* 31, 22, 4626-4629, 2016.

414 Hu, X., Belle, J.H., Meng, X., Wildani, A., Waller, L.A., et al. Estimating PM<sub>2.5</sub> concentrations in  
415 the conterminous United States using the Random Forest approach. *Environ, Sci., Technol.* 2017,  
416 51, 6936-6944. <https://doi.org.10.1021/acs.est.7b01210>.

417 Hochreiter, S., and Schmidhuber, J., 1997. Long short-term memory, *Neural Comput.*, 9, 8, 1735-  
418 1780.

419 Hutschison, K.D., Smith, S., Faruqui, S.J., 2005. Correlating MODIS aerosol optical thickness data  
420 with ground-based PM<sub>2.5</sub> observations across Texas for use in a real time air-quality prediction  
421 system. *Atmos. Environ.* 39 (37), 7190-7203.

422 Lin, C., Li, Y., Yuan, Z., Lau, A.K.H., Li, C., Fung, J.C.H., 2015. Using satellite remote sensing  
423 data to estimate the high-resolution distribution of ground-level PM<sub>2.5</sub>. *Remote Sens. Environ.*  
424 156, 117-128. <https://doi.org/10.1016/j.rse.2014.09.015>.

425 Khan, M.B., Masiol, M., Forementon, G., Gilio, A.D., de Gennaaro, G., Agostinelli, C., and  
426 Pavoni, B, 2016. Carboneous PM2.5 and secondary organic aerosol across the Veneto region (NE  
427 Italy). *Sci. Total Environ.* 542, 172-181, doi:10.1016/j.scitotenv.2015.10.103.

428 Khosravi, K ; Mao, L; Kisi, O; Yaseen, Z. M; Shahid, S. Quantifying hourly suspended sediment load  
429 using data mining models: case study of a glacierized Andean catchment in Chile. *J. Hydrol.*  
430 2018, 567, 165-179.

431 Kong, W., Z. Y. Dong, Y. Jia, D. J. Hill, Y. Xu, and Y. Zhang. “Short-term residential load forecasting  
432 based on LSTM recurrent neural network”, *IEEE Trans. Smart Grid*, vol. 10, no.1, pp. 841-851,  
433 Jan. 2017.

434 Kuremoto, T., Kimura, S., Kobayashi, K., and Obayashi, M., 2014. Time series forecasting using a  
435 deep belief network with restricted Boltzmann machines, *Neurocomputing*, 137, 47-56.

436 Liu, B.,; Philip, S. Y.; Top 10 algorithms in data mining. *Knowl. Inf. Syst.* 2008, 14, 1-37.

437 Liu, Y., Sarnat, J.A., Kilaru, V., Jacob, D.J., Koutrakis, P., 2005. Estimating ground-level PM2.5 in  
438 the eastern United States using satellite remote sensing, *Environ. Sci. Technol.*, 39, 3269-3278.

439 Malik, A.; Kumar. A.; Kisi, O. Daily pan evaporation estimation using heuristic methods with gamma  
440 test. *J. Irrig. Drain. Eng.* 2018, 144, 4018023.

441 Masiol, M., Benetello, F., Harrisom, R.M., Fornenton, G., Gaspari, F.D., and Pavoni, B., 2015.  
442 Spatial, seasonal trends and trans-boundary transport of PM2.5 inorganic ions in the Veneto  
443 region (northeastern Italy), *Atmos. Environ.*, 117, 19-31, doi:10.1016/j.atmosenv.2015.06.044.

444 Meng, X., Garay, M.J., Diner, D.J., Kalashnikova, O.V., Xu, J., Liu, Y., 2018. Estimating PM2.5  
445 speciation concentrations using prototype 4.4 km resolution misr aerosol properties over Southern  
446 California, *Atmos. Environ.*, 181, 70-81.

447 Nash J. E., Sutcliffe, J. V. River flow forecasting through conceptual models part I – A discussion of  
448 principles. *J. Hydrol.* 1970, 10, 282-290.

449 Nury, A.H.; Hasan, K.; Alam, M. J. Bin comparative study of wavelet-ARIMA and wavelet-ANN  
450 models for temperature time series data in northeastern Bangladesh. *J. King. Saud. Univ. Sci.*  
451 2017, 29, 47-61.

452 Ong, B.T., Sugiura, K., and Zettsu, K., 2016. Dynamically pre-trained deep recurrent neural networks  
453 using environmental monitoring for predicting PM2.5, *Neural Comput. Appl.*, 27, 6, 1553-1566.

454 Park, Y., Kwon, B., Heo, J., Hu, X., Liu, Y., Moon, T. 2019. Estimating PM2.5 concentration of the  
455 conterminous Unites states via interpretable convolutional neural networks. *Environ. Pollut.*  
456 113395.

457 Pietrogrande, M.C., Bacco, D., Ferrari, S., Ricciardelli, I., Scotto, F., Trentini, A., and Visentin, M.:  
458 2016. Characteristics and major sources of carbonaceous aerosols in PM2.5 in Emilia Romagna  
459 Region (Northern Italy) from four-year observations. *Sci. Total Environ.*, 553, 172-183,  
460 doi:10.1016/j.scitotenv.2016.02.074.

461 Santhi, C; Arnold, J. G.; Williams, J. R.; Dugas, W. A; Srinivasan, R.; and Hauck, L. M. Validation  
462 of the swat model on a large river basin with point and non-point sources, *JAWRA. J. Am. Water*  
463 *Resour. Assoc.*, 2001, 37, 1169-1188.

464 Schapire, R. (1990). The strength of weak learnability. *Machine Learning*, 5 (2), 197-227.

465 Freund, Y., Schpire, R (1997). A decision-theoretic generalisation of on-line learning and an  
466 application of boosting. *J. Computer and System Sciences*, 55 (1), 119-139.

467 Soni, M., Payra, S., Verma, S., 2018. Particulate matter estimation over a semi-arid region Jaipur,  
468 India using satellite AOD and meteorological parameters. *Atmospheric Pollution Research* 9 (5),  
469 949-958.

470 Stajkowski, S; Kumar, D; Samui, P; Bonakdari, H; and Gharabaghi, B, Genetic algorithm-optimized  
471 sequential model for water temperature prediction, *Sustainability*, 12, 13, 5374, 2020.

472 Rumelhart, D. E., G. E. Hinton, and R. J. Williams, “learning representations by back-propagating  
473 errors,” *Nature*, Vol. 323, no.6088, pp. 533-536, 1986.

474 Taylor, K.E., Summarizing multiple aspects of model performance in a single diagram, *J. Geophys.*  
475 *Res.*, 106, 7183-7192, 2001.

476 Van Donkelaar, A., Martin, R.V., Brauer, M., Kahn, R., Levy, R., Verduzco, C., Villeneuve, P.J.,  
477 2010. Global estimates of ambient fine particulate matter concentrations from satellite-based  
478 optical depth: development and application. *Environ. Health Perspect.* 118 (6), 847-855.

479 Van Liew, M. W; Arnold, J. G.; Garbrecht, J. D. Hydrologic simulation on agricultural watersheds:  
480 Choosing between two models. *Trans. ASAE* 2003, 56, 1539.

481 Vutukuru, S., Dabdub, D., 2008. Modeling the effects of ship emissions on coastal air quality: a case  
482 study of Southern California. *Atmos. Environ.* 42, 3751-3764.

483 Wang, J., Christopher, S.A., 2003. Intercomparison between satellite-derived aerosol optical  
484 thickness and PM<sub>2.5</sub> mass: Implications for air quality studies. *Geophys. Res. Lett.* 30 (21).

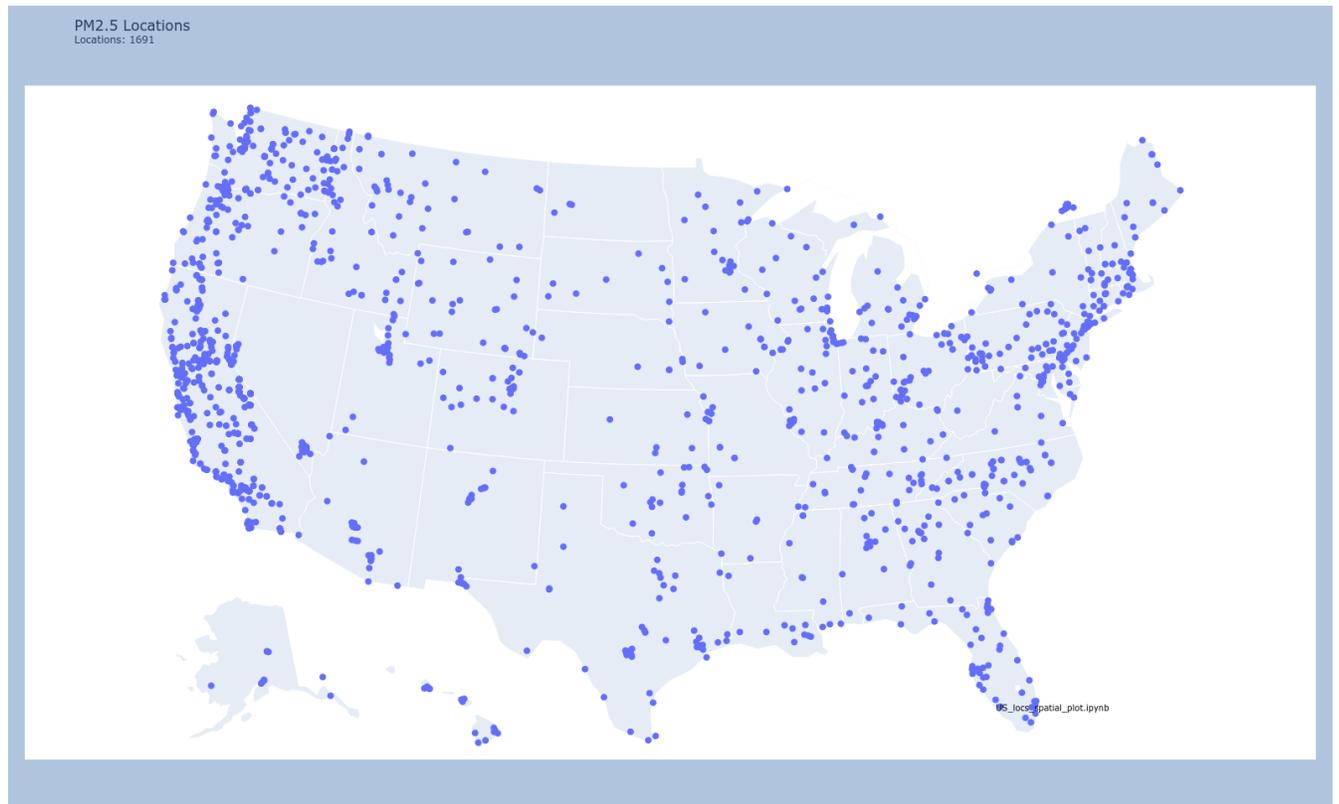
485 Wei, J., Huang, W., Li,, Z., Xue, W., Peng, Y., Sun, L., Gribb, M., 2019. Estimating 1-km resolution  
486 PM<sub>2.5</sub> concentrations across China using space-time random forest approach. *Rem. Sens.*  
487 *Environ.* 231, 111221.

488 World Health Organization, media centre (2016). Air pollution levels are rising in many of the  
489 world’s poorest cities: <http://www.int/mediacentre/news/releases/2016/air-pollution-raising/>.

490 Wu, X.; Kumar, V.; Quinlan, J. R.; Ghosh, J.; Yang, Q.; Motoda, H.; McLachlan, G. J.; Ng, A.;

491 Yi. L., Mengfan, T., Kun, Y., Yu, Z., Xiaolu, Z., Miao, Z., Yan, S., 2019. Research on PM<sub>2.5</sub>  
492 estimation and prediction method and changing characteristics analysis under long temporal and  
493 large spatial scale – a case study in China typical regions. *Sci. Total Environ.* 696, 133983,  
494 <https://doi.org/10.1016/j.scitotenv.2019.133983>.

495 Zhang, Y., Cao,F., 2015. Fine particle matter (PM<sub>2.5</sub>)in China at a city level. *Sci. Rep.*, 5, 14884.



497

498

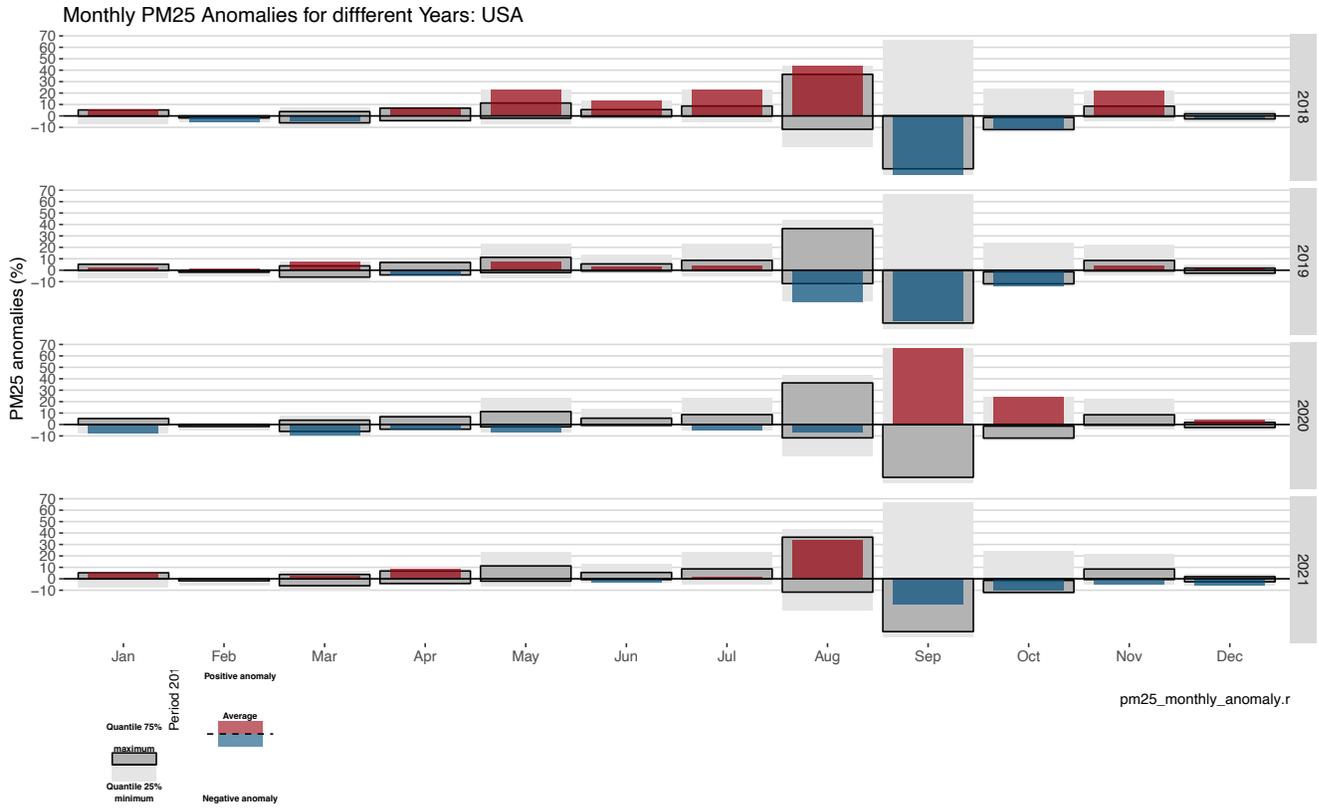
**Figure 1.** Locations of PM<sub>2.5</sub> monitoring sites over USA

499

500

501

502



503

504 **Figure 2.** Monthly anomalies and quantiles for the observed period (2018-2021) using daily PM<sub>2.5</sub>

505 values over United States.

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

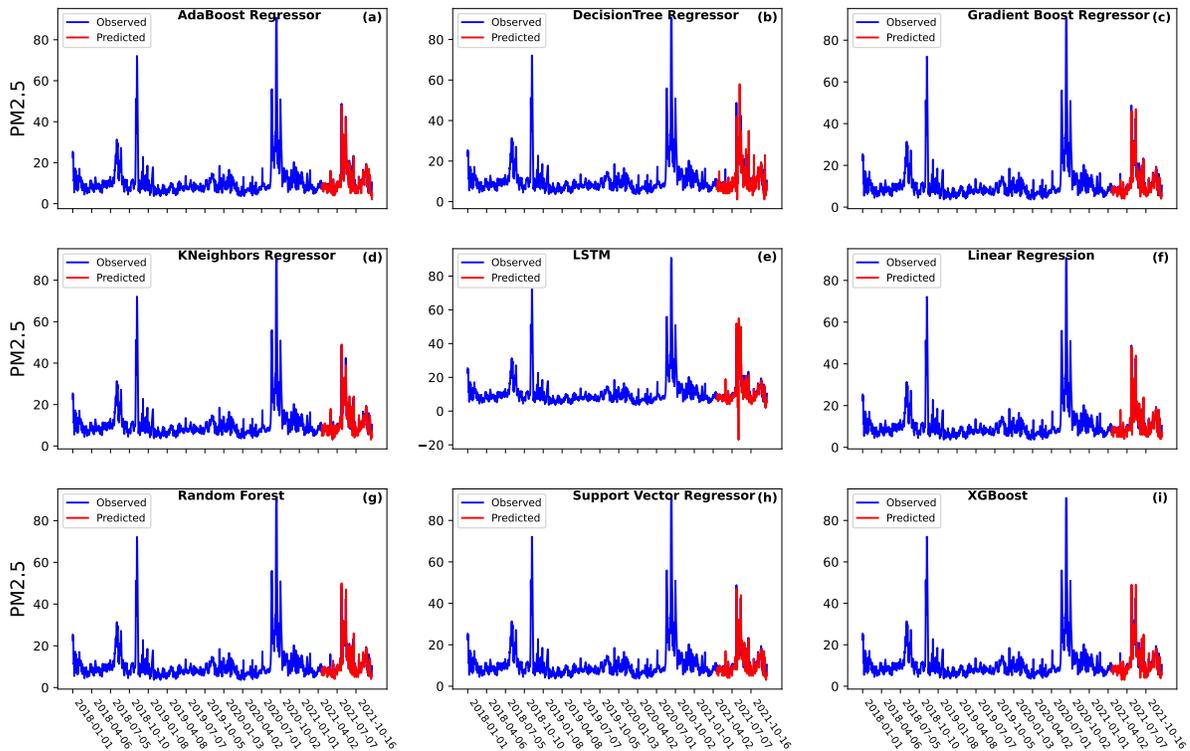
522

523

524

525

526



PM25\_dly\_dfrntmdl\_obs\_pred\_lineplt.ipynb

527

528

529

530

531

532

533

534

535

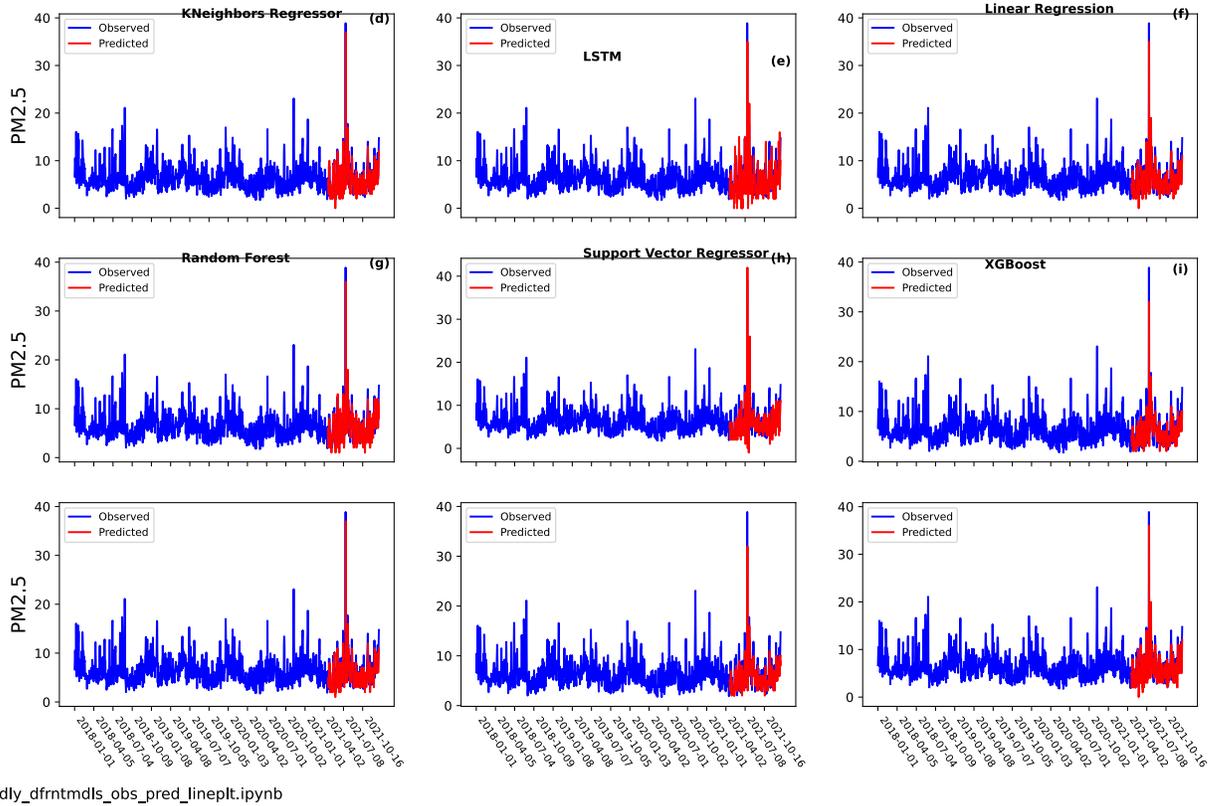
536

537

538

539

**Figure 3.** The comparison of the time series of estimated and observed  $PM_{2.5}$  concentrations over California using different machine learning models: (a) AdaBoost regressor, (b) Decision Tree regression, (c) Gradient Boost regression, (d) K-neighbors regression (e) LSTM, (f) Linear regression, (g) Random Forest, (h) Support Vector regression, and (I) XGBoost.



540

541 **Figure 4.** The comparison of the time series of estimated and observed  $PM_{2.5}$  concentrations over  
 542 New York using different machine learning models: (a) AdaBoost regressor, (b) DecisionTree  
 543 regression, (c) Gradient Boost regression, (d) Kneighbors regression (e) LSTM, (f) Linear regression,  
 544 (g) Random Forest, (h) Support Vector regression, and (I) XGBoost.

546

547

548

549

550

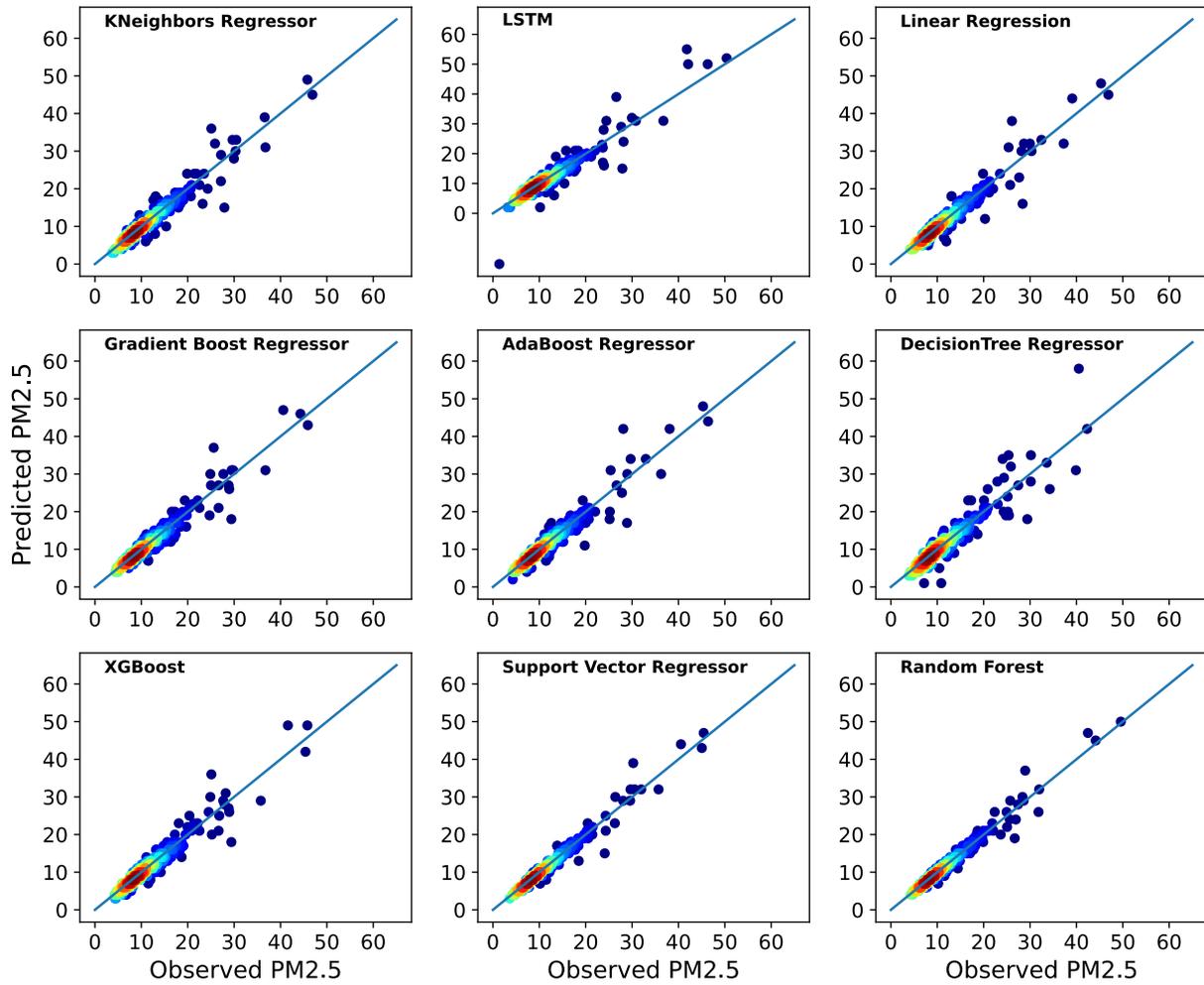
551

552

553

554

555



556

557 **Figure 5.** Scatter plots of observed and estimated daily PM<sub>2.5</sub> concentrations over California using  
558 different machine learning models: (a) AdaBoost regressor, (b) DecisionTree regression, (c) Gradient  
559 Boost regression, (d) Kneighbors regression (e) LSTM, (f) Linear regression, (g) Random Forest, (h)  
560 Support Vector regression, and (I) XGBoost.

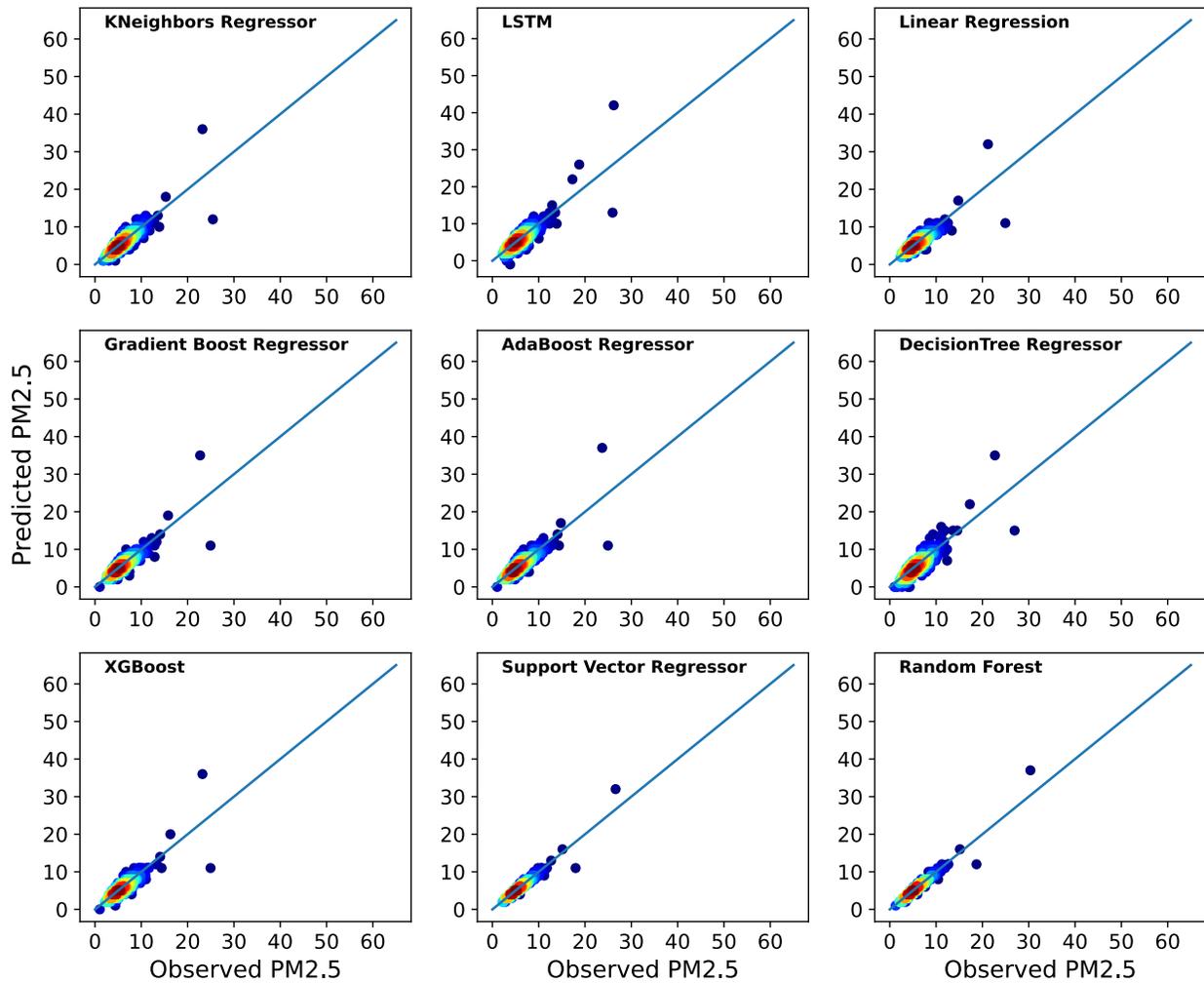
561

562

563

564

565



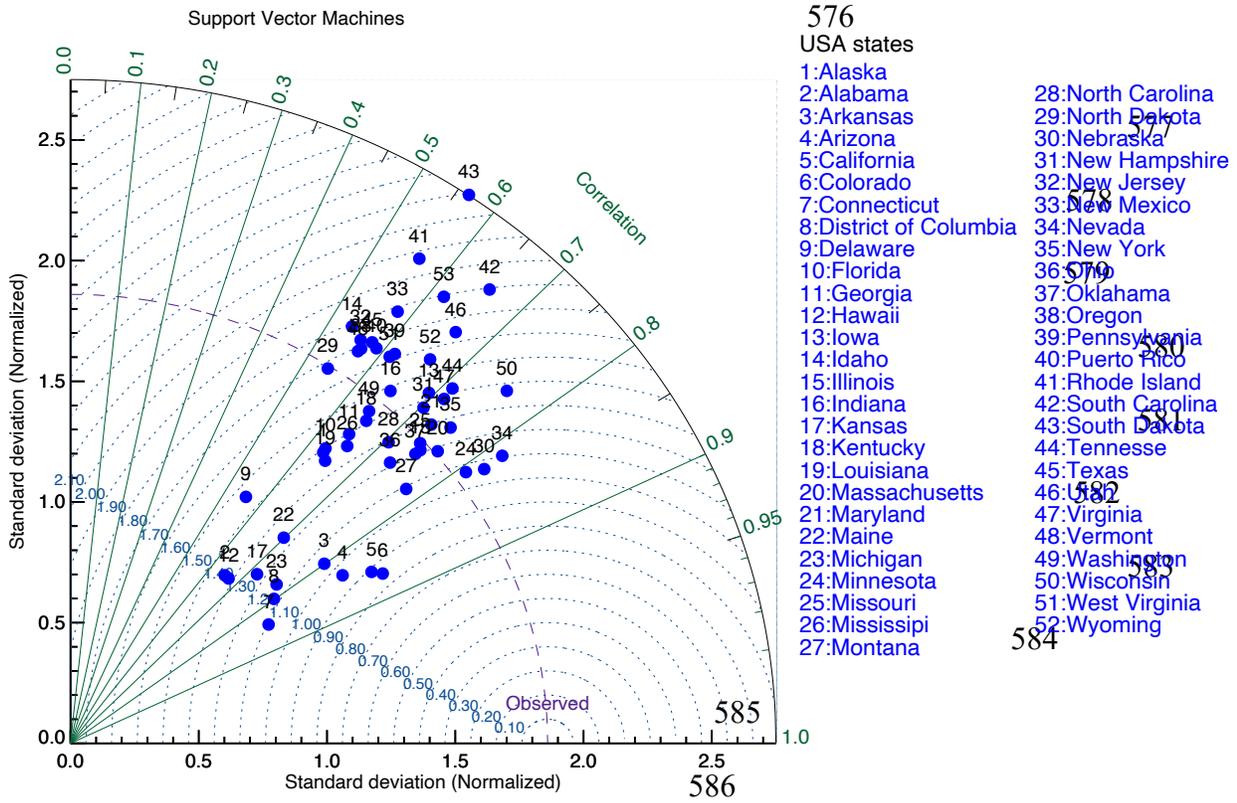
567

568 **Figure 6.** Scatter plots of observed and estimated daily  $PM_{2.5}$  concentrations over New York using  
 569 different machine learning models: (a) AdaBoost regressor, (b) DecisionTree regression, (c) Gradient  
 570 Boost regression, (d) Kneighbors regression (e) LSTM, (f) Linear regression, (g) Random Forest,  
 571 (h) Support Vector regression, and (I) XGBoost.  
 572

573

574

575



576

USA states

- 1:Alaska
- 2:Alabama
- 3:Arkansas
- 4:Arizona
- 5:California
- 6:Colorado
- 7:Connecticut
- 8:District of Columbia
- 9:Delaware
- 10:Florida
- 11:Georgia
- 12:Hawaii
- 13:Iowa
- 14:Idaho
- 15:Illinois
- 16:Indiana
- 17:Kansas
- 18:Kentucky
- 19:Louisiana
- 20:Massachusetts
- 21:Maryland
- 22:Maine
- 23:Michigan
- 24:Minnesota
- 25:Missouri
- 26:Mississippi
- 27:Montana
- 28:North Carolina
- 29:North Dakota
- 30:Nebraska
- 31:New Hampshire
- 32:New Jersey
- 33:New Mexico
- 34:Nevada
- 35:New York
- 36:Ohio
- 37:Oklahoma
- 38:Oregon
- 39:Pennsylvania
- 40:Puerto Rico
- 41:Rhode Island
- 42:South Carolina
- 43:South Dakota
- 44:Tennessee
- 45:Texas
- 46:Utah
- 47:Virginia
- 48:Vermont
- 49:Washington
- 50:Wisconsin
- 51:West Virginia
- 52:Wyoming

taylor\_pm25mod.pro

587

588

589 **Figure 7.** Taylor diagram of the Support Vector Machines (SVM) over each state of the United States.

590

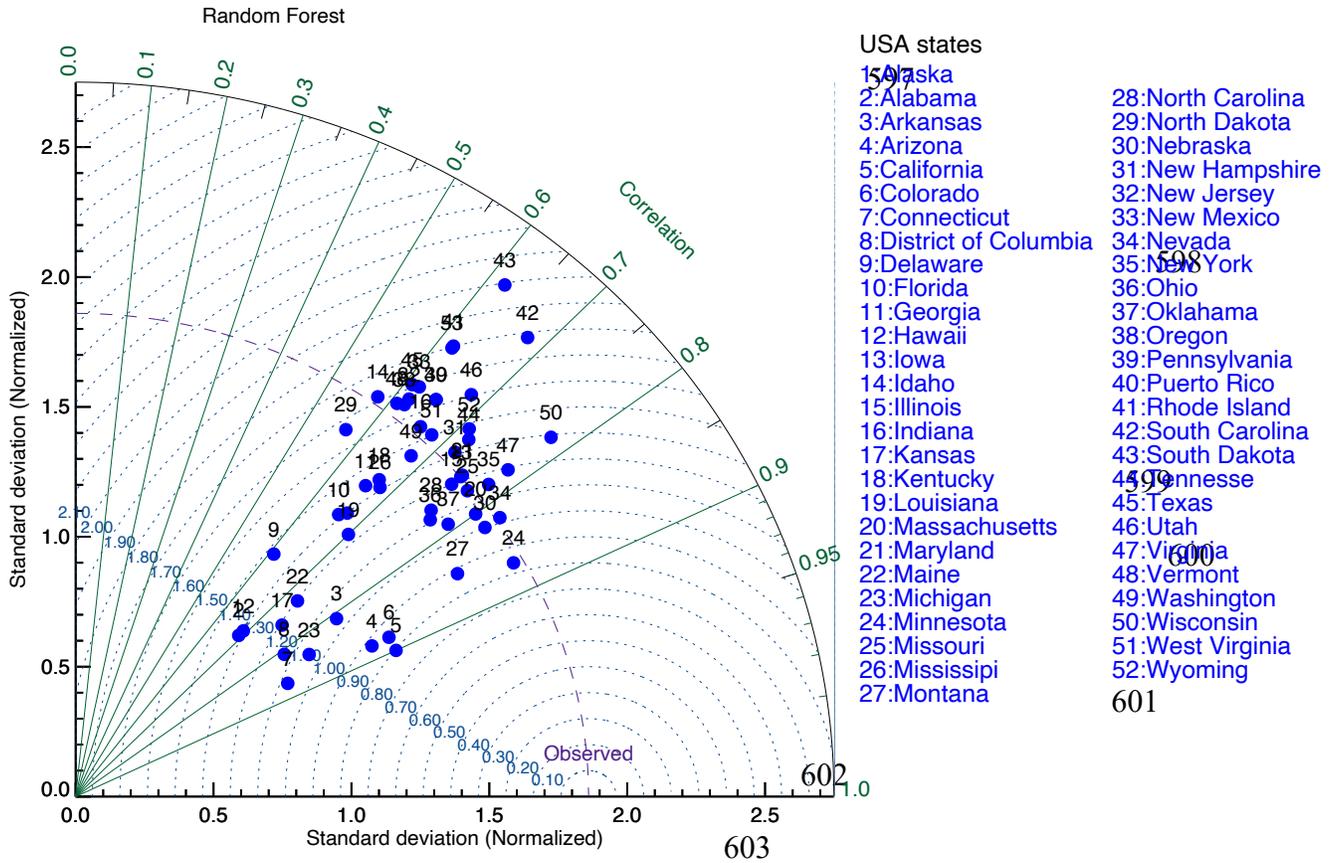
591

592

593

594

595



604

taylor\_pm25.pro

605

606 **Figure 8.** Taylor diagram of the Random Forest (RF) over each state of the United States.

607

608

609

610

611

612 **Table 1: Different Model Metrics for New York State**

New York								
Model	RMSE	MAE	MAPE	R2	NSE	NORM	PBIAS	RSR
<b>Linear Regression</b>	3.883	2.309	0.285	0.688	0.613	60.156	11.24	0.561
<b>Decision Tree</b>	5.136	3.109	0.254	0.454	0.533	79.58	13.44	0.691
<b>Gradient Boost Regressor</b>	3.822	2.394	0.545	0.698	0.683	59.207	8.210	0.546
<b>AdaBoost Regressor</b>	3.961	2.316	0.188	0.676	0.683	61.369	9.653	0.576
<b>XG Boost</b>	3.898	2.501	0.202	0.686	0.681	60.393	8.342	0.559
<b>KNeighbors Regressor</b>	3.919	2.379	0.195	0.683	0.677	60.711	7.515	0.562
<b>LSTM</b>	7.487	3.359	0.218	0.158	0.455	115.991	6.020	0.812
<b>Random Forest</b>	3.121	2.122	0.182	0.899	0.811	38.671	2.989	0.331
<b>SVM</b>	3.125	2.145	0.183	0.857	0.820	39.161	3.011	0.338

- 613
- 614 RMSE = Root mean squared error
- 615 MAE = Mean absolute error
- 616 MAPE = Mean absolute percentage error
- 617  $R^2$  = The coefficient of determination
- 618 NSE = Nash-Sutcliffe efficiency
- 619 PBIAS = Percent Bias
- 620 RSR = root mean square error ratio

621  
622  
623 **Table 2: Different Model Metrics for California State**

California								
Model	RMSE	MAE	MAPE	R <sup>2</sup>	NSE	NORM	PBIAS	RSR
<b>Linear Regression</b>	3.695	2.599	0.326	0.43	0.694	57.243	12.086	0.932
<b>Decision Tree</b>	5.481	3.743	0.467	0.23	0.576	84.917	19.901	0.732
<b>Gradient Boost Regressor</b>	4.051	2.736	0.340	0.28	0.461	62.758	16.891	1.017
<b>AdaBoost Regressor</b>	3.804	2.636	0.342	0.33	0.435	58.938	17.532	0.969
<b>XG Boost</b>	4.271	2.972	0.372	0.17	0.438	66.178	18.726	1.075
<b>KNeighbors Regressor</b>	4.394	3.062	0.392	0.22	0.286	68.071	17.076	1.106
<b>LSTM</b>	5.025	3.252	0.339	0.46	0.309	77.853	18.027	0.618
<b>Random Forest</b>	3.051	2.233	0.315	0.77	0.817	46.894	7.022	0.355
<b>SVM</b>	3.714	2.618	0.320	0.71	0.897	47.853	7.027	0.424

- 624
- 625 RMSE = Root mean squared error
- 626 MAE = Mean absolute error
- 627 MAPE = Mean absolute percentage error
- 628  $R^2$  = The coefficient of determination
- 629 NSE = Nash-Sutcliffe efficiency
- 630 PBIAS = Percent Bias
- 631 RSR = root mean square error ratio
- 632