

---

# Application of Machine Learning for PM<sub>2.5</sub> estimation using low-cost sensor network

---

Prateek Sethi  
Virginia Tech  
Arlington, VA  
prateek20@vt.edu

Kuldeep Dixit  
Virginia Tech  
Arlington, VA  
kdixit@vt.edu

## Abstract

1 Fine particulate matter or PM<sub>2.5</sub> is known to harm human health. The monitoring  
2 of PM<sub>2.5</sub> is a challenging task due to the high cost of PM<sub>2.5</sub> monitors. The low-cost  
3 sensors are increasingly being employed to have more spatial monitoring coverage  
4 but, these sensors are less reliable due to lower accuracy. In this study, we have  
5 shown that machine learning models can be employed to predict accurate PM<sub>2.5</sub>  
6 and perform better than the traditional linear regression method used to correct the  
7 low-cost sensor data. Most of our machine learning models performed better than  
8 the baseline linear regression model used by Environmental Protection Agency's  
9 AirNow portal.

## 10 1 Introduction

11 High concentration of PM<sub>2.5</sub> is known to cause respiratory and cardiovascular health problems. This  
12 study[2] finds PM<sub>2.5</sub> attributed to having caused 8.42 million deaths in 2016. Although there is an  
13 urgent need for universal monitoring of PM<sub>2.5</sub>, it is not possible because of the high cost of the  
14 monitors. To address this problem, the usage of low-cost sensors is growing. The Purple air sensor  
15 has become a popular low-cost sensor. It is increasingly employed to estimate the PM<sub>2.5</sub> concentration  
16 using meteorological data and PM signals from these sensors.

## 17 2 Methodology

18 A popular standard definition of Machine Learning by Tom Mitchell is: "A computer program is said  
19 to learn from experience E with respect to some class of tasks T and performance measure P if its per-  
20 formance at tasks in T, as measured by P, improves with experience E." Machine learning algorithms  
21 can make decisions automatically, and these decisions can be improved through experience/data.  
22 The logic is not hardcoded but is programmed by the machine itself. This study uses the popular  
23 algorithms - CART - Classification and Regression Trees, KRR - Kernel Ridge Regression, KNN - K  
24 - Nearest Neighbors, and ENS - Ensemble methods - Bagging, Voting Regressor, Random Forrest.

### 25 2.1 Data

26 In our study we have used two data sets namely the Environmental Protection Agency's PM2.5  
27 dataset and the low cost sensor dataset. We extracted the data from the Purple air sensors which are  
28 within 50 meters of the EPA monitor and merged both the EPA data and Purple air data to feed into  
29 the models. This process is similar to the baseline study[1] that we have used to compare our machine

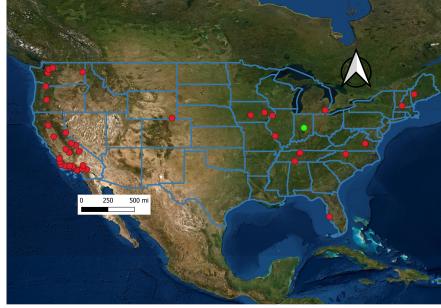


Figure 1: Location of 46 EPA monitors used in our study

30 learning models. We have used the temperature, humidity and the PM signal from low cost sensor as  
 31 our features and we have used the EPA monitor's PM<sub>2.5</sub> concentration as the label. We have removed  
 32 the invalid values of temperature and humidity and also compared the two channel readings of the  
 33 purple air sensor and removed the values showing more than 3 percentage deviation. We followed the  
 34 process of Lu et al[2] here, lastly we used data from 46 EPA monitors and 96 purple air sensors  
 35 in our study. The location of EPA monitors used is given in Figure 1.

## 36 2.2 Classification and Regression Trees

37 It is the process of predicting the class of given data points. Classes are sometimes called targets  
 38 or labels or categories. Classification predictive modeling is the task of approximating a mapping  
 39 function ( $f$ ) from input variables ( $X$ ) to discrete output variables ( $y$ ). In contrast regression analysis  
 40 consists of an output variable which is continuous or a real value. Regression techniques try to fit  
 41 the best hyperplane which goes through the data points. The CART algorithm can be used to tackle  
 42 decision making problems relying on a number of features in a systematic order. CART techniques  
 43 have been a major contributor to a lot of health studies [7][8]. Different steps in the CART process  
 44 are explained and carried out in [3][4] gives an example of another use case in the health and wellness  
 45 industry. CART can handle both classification and regression tasks. This algorithm uses a metric  
 46 named gini index to create decision points for classification tasks

$$Gini = 1 - \sum_{i=1}^C (p_i)^2 \quad (1)$$

## 47 2.3 Kernel Ridge Regression

48 Kernel Ridge Regression, is a very simple special case of Support Vector Regression. The main  
 49 formula of the method is identical to a formula in Bayesian statistics, but Kernel Ridge Regression  
 50 has performance guarantees that have nothing to do with Bayesian assumptions.[3] “Kernel Ridge  
 51 Regression” was coined in 2000 by Cristianini and Shawe-Taylor [7] to refer to a simplified version of  
 52 Support Vector Regression; this was an adaptation of the earlier “ridge regression in dual variables”

$$\begin{aligned} \min \quad & \|w\|^2 + C \sum_{t=1}^T ((\xi_t)^k + (\xi'_t)^k) \\ \text{s.t.} \quad & (w \cdot x_t + b) - y_t \leq \epsilon + \xi_t, t = 1, \dots, T \\ & y_t - (w \cdot x_t + b) \leq \epsilon + \xi'_t, t = 1, \dots, T \end{aligned} \quad (2)$$

53 where  $(x^t; y^t) \in R_n \times R$  are the training examples,  $w$  is the weight vector,  $b$  is the bias term,  $\xi'_t$  and  $\xi_t$   
 54 are the slack variables, and  $T$  is the size of the training set  $\epsilon$ ;  $C > 0$  and  $k$  eta 1,2 are the parameters.  
 55 Simplify the problem by ignoring the bias term  $b$  (it can be partially recovered by adding a dummy

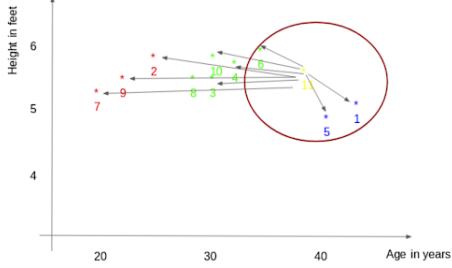


Figure 2: KNN Algorithm.

Table 1: Different distance metric for KNN

Metric	Formula
Euclidean Distance:	$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$ (4)
Manhattan Distance:	$\sum_{i=1}^k  x_i - y_i $ (5)
Hamming Distance	$\min D_H = \sum_{i=1}^k  x_i - y_i $ $x = y \Rightarrow D = 0$ $x \neq y \Rightarrow D = 1$ (6)

56 attribute 1 to all  $x_t$  ), setting  $\epsilon := 0$ , and setting  $k := 2$ . The optimization problem becomes

$$\min a\|w\|^2 + \sum_{t=1}^T (y_t - w \cdot x_t)^2 \quad (3)$$

57 (where  $a := \frac{1}{c}$ ), the usual Ridge Regression problem. And Vapnik's usual method ([6], Sect. 11.3.2)  
58 then gives the prediction.

59 In KRR, as in Support Vector Regression in general, the kernel is not supposed to reflect any  
60 knowledge or beliefs about reality, and the usual approach is pragmatic: one consults standard  
61 libraries of kernels and uses whatever works.[3]

## 62 2.4 K Nearest Neighbor Regression

63 The KNN algorithm uses ‘feature similarity’ to predict the values of any new data points. This means  
64 that the new point is assigned a value based on how closely it resembles the points in the training set.

65 Consider the figure 3 from [8], to make a prediction the distance to the existing data points is  
66 calculated and based on the K closest data points a prediction is made. For the calculation of distance  
67 different methods can be used depending on the use case. Some of these are indicated in Table 1.

## 68 2.5 Ensemble Methods

69 Ensemble methods are techniques that create multiple models and then combine them to produce  
70 improved results. The most commonly used methods are the voting and the averaging ensemble  
71 methods. Where multiple classification or regression models are created using a dataset. This can be  
72 done by using the same algorithm or by using different algorithms. The prediction of these models  
73 can be combined in various ways to come to a final conclusion. Weights can be added to improve the  
74 quality of the model and hence use Weighted Voting or Weighted Averaging. It is very important  
75 that the ensemble classifier results in a prediction which is better than its individual classifiers. A  
76 necessary and sufficient condition for an ensemble of classifiers to be more accurate than any of

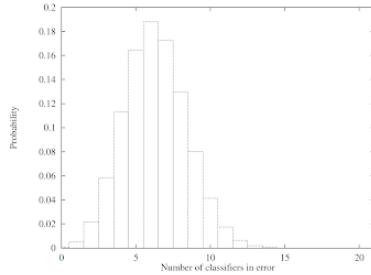


Figure 3: The Probability that exactly 1 (of 21) hypothesis will make an error, assuming each hypothesis has an error rate of 0.3 and makes its error independently of the other

77 its individual members is if the classifiers are accurate and diverse [9]. To see why accuracy and  
 78 diversity are good, imagine that we have an ensemble of three classifiers:  $h_1, h_2, h_3$  and consider a  
 79 new case  $x$ . If the three classifiers are identical (i.e., not diverse), then when  $h_1(x)$  is wrong,  $h_2(x)$   
 80 and  $h_3(x)$  will also be wrong. However, if the errors made by the classifiers are uncorrelated, then  
 81 when  $h_1(x)$  is wrong,  $h_2(x)$  and  $h_3(x)$  may be correct, so that a majority vote will correctly classify  
 82  $x$ . More precisely, if the error rates of  $L$  hypotheses  $h$  are all equal to  $p < 1/2$  and if the errors  
 83 are independent, then the probability that the majority vote will be wrong will be the area under  
 84 the binomial distribution where more than  $L/2$  hypotheses are wrong. Figure 1 shows this for a  
 85 simulated ensemble of 21 hypotheses, each having an error rate of 0.3. The area under the curve  
 86 for 11 or more hypotheses being simultaneously wrong is 0.026, which is much less than the error  
 87 rate of the individual hypotheses. Additionally [7] discusses how ensemble methods result in highly  
 88 accurate classifiers by combining less accurate ones. Also it discusses the three fundamental reasons  
 89 why ensemble methods are able to out-perform any single classifier within the ensemble. For our  
 90 experimentation we considered the following methods: Voting Regression, Bootstrap Aggregation  
 91 (Bagging) and Random Forest Regression

### 92 3 Results and Observations

93 To compare the performance of the linear regression model[1] with our machine learning models  
 94 we considered five metrics, the mean squared error, mean absolute error, r2 distance, fit time, and  
 95 the score time. For each of these metrics, we performed the 10 Fold cross-validation to get a robust  
 96 performance estimation. The results and comparison is shown in Table 1 and Figure 4. Table 1 shows  
 97 the median values obtained from 10-fold cross-validation.

98 **Mean Squared Error** From Table 1 we can observe that almost all the algorithms perform  
 99 comparable to the linear regression model used in the reference study.

100 **Mean Absolute Error** According to this metric, the performance of most of the models is better  
 101 than the linear regression model. As can be observed that the Bagging using KNN and the Voting  
 102 regression outperform the linear regression model.

103 **Fit time and Score time** This metric was not available for the base model but, we have used  
 104 this to compare the machine learning models that we implemented in this study. Most machine  
 105 learning models took almost equal time to fit and score but, the KNN and the voting regression took  
 106 comparatively more time.

#### 107 3.1 Model Explanation

108 We also used SHapley Additive exPlanations or, SHAP to explain our models. We applied SHAP  
 109 method on the best performing voting regression method. The most important feature for prediction

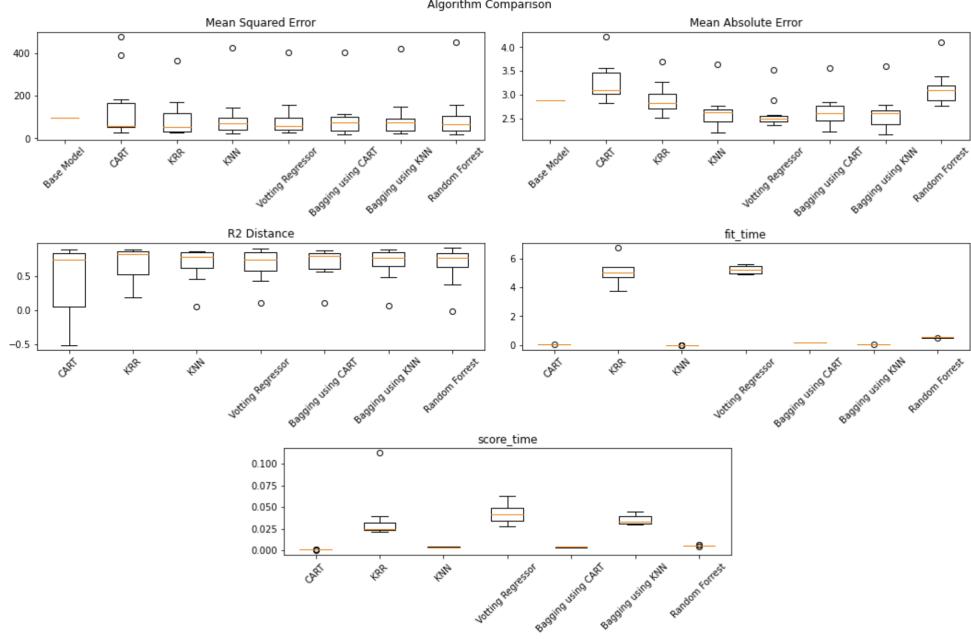


Figure 4: ML models comparison

Table 2: The comparison of models

Metric	Base Model	CART	KRR	KNN	Votting Regressor	Bagging using CART	Bagging using KNN	Random Forrest
Mean Squared Error	96.44	107.96	54.3	69.06	59.88	73.41	74.1	66.99
Mean Absolute Error	2.88	3.13	2.82	2.63	2.49	2.61	2.6	3.09
R <sup>2</sup>	0.69	0.72	0.82	0.78	0.79	0.8	0.77	0.78
Fit Time	NA	0.06	6.36	0.01	8.35	0.54	0.1	1.6
Score Time	NA	0	0.04	0.02	0.05	0.01	0.09	0.02

110 PM<sub>2.5</sub> turns out to be low cost sensor's PM signal followed by humidity and temperature. Figure 5  
111 shows the mean SHAP values for the voting regression model.

## 112 4 Conclusion

113 In our study most of the machine learning algorithms showed either comparable or better performance  
114 than the baseline linear regression model. The metrics do not indicate a clear winner but, according  
115 to the metric mean absolute error, voting regression outperforms, and according to mean squared error  
116 and R<sup>2</sup> Kernal Ridge Regression (KRR) outperforms all other models.

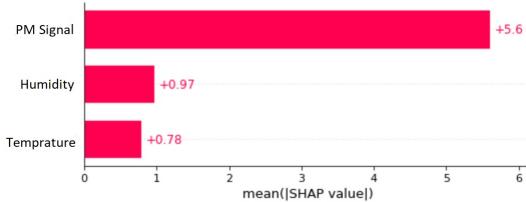


Figure 5: The average SHAP values for the voting regression model

117 **References**

- 118 [1] Barkjohn, K. K., Gantt, B., and Clements, A. L.: Development and application of a United States-  
119 wide correction for PM2.5 data collected with the PurpleAir sensor, *Atmos. Meas. Tech.*, 14, 4617–4637,  
120 <https://doi.org/10.5194/amt-14-4617-2021>, 2021.
- 121 [2] Lu T, Marshall JD, Zhang W, Hystad P, Kim SY, Bechle MJ, Demuzere M, Hankey S. National Empirical  
122 Models of Air Pollution Using Microscale Measures of the Urban Environment. *Environ Sci Technol*. 2021 Nov  
123 16;55(22):15519-15530. doi: 10.1021/acs.est.1c04047. Epub 2021 Nov 5. PMID: 34739226.
- 124 [3] Roger J. Lewis, M.D., Ph.D. Department of Emergency Medicine Harbor-UCLA Medical Center Torrance,  
125 California - “*An Introduction to Classification and Regression Tree (CART) Analysis*” Presented at the 2000  
126 Annual Meeting of the Society for Academic Emergency Medicine in San Francisco, California.
- 127 [4] N. Speybroeck - “*Classification and regression trees*” - Received: 9 February 2011 / Revised: 2 September  
128 2011 / Accepted: 8 October 2011 / Published online: 21 October 2011.
- 129 [5] Festschrift in Honor of Vladimir N. Vapnik - “*Empirical Inference*”
- 130 [6] Vapnik, V.N.: *Statistical Learning Theory*. Wiley, New York (1998)
- 131 [7] Cristianini, N., Shawe-Taylor, J.: *An Introduction to Support Vector Machines and Other Kernel-Based*  
132 *Methods*. Cambridge University Press, Cambridge (2000)
- 133 [8] A Practical Introduction to K-Nearest Neighbors Algorithm for Regression (with Python code) - ‘Aishwarya  
134 Singh — August 22, 2018’ - <https://www.analyticsvidhya.com/blog/2018/08/k-nearest-neighbor-introduction-regression-python/>
- 136 [9] Thomas G. Dietterich, Oregon State University, Corvallis, Oregon, USA - ‘*Ensemble Methods in Machine*  
137 *Learning*’