

A comparison between geostatistical and machine learning models for spatio-temporal prediction of PM_{2.5} data

Zeinab Mohamed ^{*} Wenlong Gong [†]

September 16, 2025

Abstract

Ambient air pollution poses significant health and environmental challenges. Exposure to high concentrations of PM_{2.5} have been linked to increased respiratory and cardiovascular hospital admissions, more emergency department visits and deaths. Traditional air quality monitoring systems such as EPA-certified stations provide limited spatial and temporal data. The advent of low-cost sensors has dramatically improved the granularity of air quality data, enabling real-time, high-resolution monitoring. This study exploits the extensive data from PurpleAir sensors to assess and compare the effectiveness of various statistical and machine learning models in producing accurate hourly PM_{2.5} maps across California. We evaluate traditional geostatistical methods, including kriging and land use regression, against advanced machine learning approaches such as neural networks, random forests, and support vector machines, as well as ensemble model. Our findings enhanced the predictive accuracy of PM_{2.5} concentration by correcting the bias in PurpleAir data with an ensemble model, which incorporating both spatiotemporal dependencies and machine learning models.

Keywords: Air pollution; Geostatistical model; Spatio-temporal data; Machine learning model; Ensemble model

^{*}Department of Mathematics, Oberlin College and Conservatory

[†]Department of Mathematics and Statistics, University of Houston - Downtown

1 Introduction

The disparate distribution of air pollution underscores significant ecological justice concerns Hajat et al. (2015), particularly with particulate matter less than 2.5 micrometers ($PM_{2.5}$). As one of the most hazardous pollutants, it not only leads to significant health risks but also contributes to environmental issues Organization et al. (2018). In urban settings, the government typically installs a few EPA-certified air quality monitors that are updated hourly with an hourly average particle count. However, pollution levels exhibit considerable variability within small areas and over short periods. The advent of low-cost air pollution sensors, such as PurpleAir, has revolutionized monitoring by providing high spatiotemporal resolution data, thereby enhancing the detection of pollution hotspots and improving air quality indices during severe pollution events like wildfires Bi et al. (2020); Delp and Singer (2020); Morawska et al. (2018).

PurpleAir sensors, known for their affordability, ease of use and maintenance, have proliferated, significantly aiding public and scientific understanding of air quality dynamics Barkjohn et al. (2020). The PurpleAir sensor network offers a high spatial resolution and real-time particulate matter data at a 2-minute temporal resolution. The widespread use of PurpleAir sensors, gives us the opportunity to generate an hourly map of ambient $PM_{2.5}$ concentrations in California. Leveraging high spatiotemporal resolution data collected from 1110 PurpleAir sensors operating from January to December 2019, this study aims to assess various geostatistical and machine learning techniques in producing hourly $PM_{2.5}$ maps.

In recent years, considerable effort has been devoted to developing statistical methods that accurately predict $PM_{2.5}$ concentrations on the spatial, temporal and spatiotemporal domains. Traditional statistical methods like kriging and land use regression have been essential for epidemiological studies to estimate $PM_{2.5}$ levels. However, these methods often suffer from computational intensity and data availability limitations, respectively Alexeeff et al. (2015); Datta et al. (2016a); Hu et al. (2013). With the increase in data volume, computational challenges have prompted the adoption of machine learning models, which have demonstrated superior performance in capturing complex patterns in high-dimensional data. Datta et al. (2016a) introduced a hierarchical nearest-neighbor Gaussian process model for large geostatistical datasets, which used a hierarchical structure to reduce the computational complexity of the Gaussian process

while maintaining its accuracy. However, it requires spatiotemporal covariance structures and parametric assumptions.

More recently, machine learning (ML) algorithms have been explored for air pollution prediction, including random forests, support vector regression, and deep learning approaches. Hu et al. (2017) combined ground observations, satellite data, and meteorological variables to develop a model that could predict $\text{PM}_{2.5}$ concentrations at a high spatial resolution using a random forest (RF). The study was able to accurately estimate $\text{PM}_{2.5}$ concentrations across the United States, though they found that season and climate can affect the prediction accuracy. Mogollón-Sotelo et al. (2021) adopted support vector machines (SVM) to forecast $\text{PM}_{2.5}$ concentrations in Bogotá, Colombia. The study trained data for 12 monitoring stations using a radial basis kernel, and the model had better predictive accuracy compared with Bogotá Integrated Air Quality Modeling System, especially in areas with complex terrain, where traditional air quality models did not perform well. Gupta and Christopher (2009) used a neural network to produce more accurate predictions of $\text{PM}_{2.5}$ concentrations compared with other traditional regression models. The method involved training a neural network to predict $\text{PM}_{2.5}$ concentrations based on the input of three years of surface, satellite, and meteorological fields in the southeastern United States, resulting in significant improvements in accuracy compared to traditional regression-based methods.

Despite the variety of approaches available, there is no consensus on which produces the most accurate predictions. Several studies have evaluated traditional geostatistical methods with machine learning algorithms. Requía et al. (2019) compared ordinary kriging, hybrid interpolation, and random forests for estimating concentrations of 10 $\text{PM}_{2.5}$ components. And demonstrated that random forests surpassed both kriging and hybrid interpolation methods (empirical Bayesian kriging and land use regression). Wang et al. (2019) introduced a flexible spatial prediction approach using a Nearest-Neighbor Neural Network, which excelled over conventional geostatistical methods in handling non-normal data by innovatively utilizing neighboring information. Conversely, Berrocal et al. (2020) evaluated universal kriging and downscaler models against machine learning approaches for generating daily national maps of $\text{PM}_{2.5}$ concentrations in the United States, noting superior predictive performance in geostatistical methods under the assumption of constant spatial

covariance.

Although numerous studies have compared geostatistical and machine learning techniques, they often focus solely on spatial or temporal dependencies, neglecting their combined spatiotemporal dynamics. Leveraging the extensive network of PurpleAir sensors, this study aims to produce hourly maps for $\text{PM}_{2.5}$ concentrations in California. We explore the efficacy of both traditional geostatistical methods, such as universal kriging and fixed rank kriging, and advanced non-geostatistical approaches, including regression, random forests, support vector machines, and neural networks. Additionally, this study integrates geostatistical methods and nearest neighbor data with machine learning algorithms to enhance prediction accuracy and performance.

2 Data

PurpleAir sensors are one of the low-cost particulate matter (PM) sensors that have gained popularity for their affordability, accessibility, and ease of use in monitoring outdoor air quality. These sensors use laser-based optical sensors to count particles by sizes and uses the counts to calculate mass concentrations of $\text{PM}_{1.0}$, $\text{PM}_{2.5}$, and PM_{10} . Each sensor has two Plantower PMS5003 units, labeled as channels A and B, which operate alternately and provide 80-second averaged data values (and, from June 2019, 2-minute averaged vales) (Barkjohn et al., 2020). Widely adopted by communities, citizen scientists, and researchers, PurpleAir sensors are extensively used for local air quality monitoring, and their data has been utilized in various applications. The substantial rise in sensor counts from 251 in 2017 to 39266 in 2022 (API) as shown in Figure 1, reflects a growing interest in understanding air quality at a local level. This increase has important implications for air quality monitoring and decision-making. Figure 2 illustrates the rapid increase in PurpleAir installations across California from 2019 to 2023. This underscores the growing importance of incorporating these sensors into air quality monitoring efforts. With more sensors in operation, a higher density of data points can be obtained, enabling a more comprehensive and detailed assessment of air quality trends and patterns in specific areas.

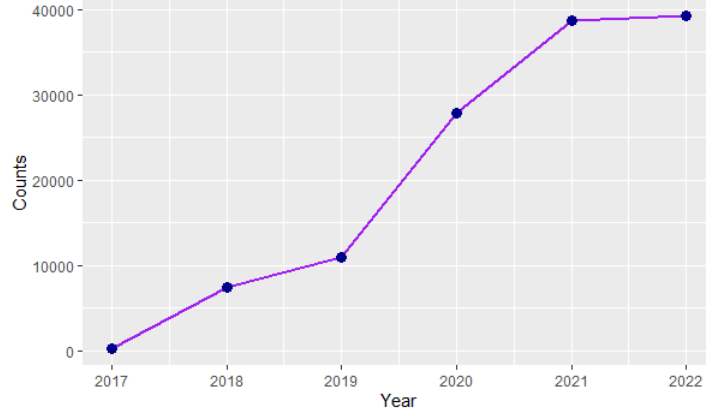


Figure 1: PurpleAir sensors counts from 2017-2019 in US

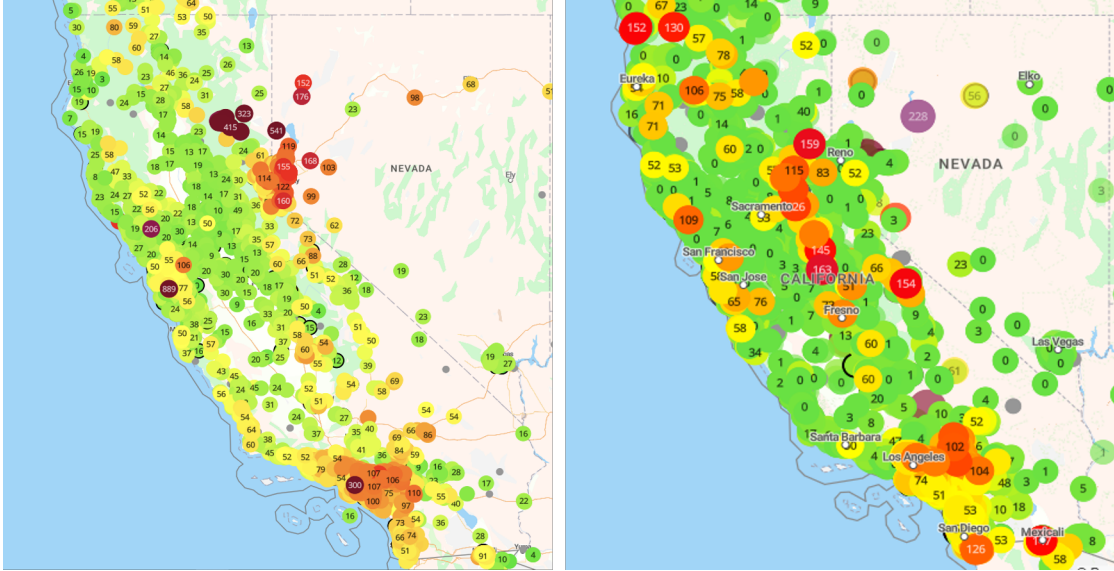


Figure 2: Real-time map for $PM_{2.5}$ concentration measured by PurpleAir sensors in 2019 (left) versus 2021 (right) PurpleAir (2021). (Green) indicates air quality is satisfactory, (yellow) air quality is acceptable, (orange) members from sensitive maybe affected, (red) the general public may experience health effects, and (purple) risk of health effects is increased for everyone.

However, it is essential to consider the limitations of PurpleAir sensors, such as potential issues with sensor accuracy, calibration, and data quality. As PurpleAir sensor counts continue to rise, it is crucial to ensure that proper quality control measures, calibration, and validation protocols are in place to ensure the reliability and accuracy of the data obtained from these sensors for robust air quality assessments and decision-making. To start with the most robust data we processed the PurpleAir data in three steps after

accessing and downloading the data in JSON format: 1) we removed observations with extremely high or low temperatures (i.e 2147483447 or -224) which indicates a communication error between PurpleAir micro-controller and temperature sensors. We also removed observations that have a relative humidity (RH) outside the range 0 – 100%. 2) We checked the consistency of the measurements for both channels A and B and then averaged the 2-minute or 80-second data hourly. 3) we applied the correction formula given by Barkjohn et al. (2020),

$$PM_{2.5corrected} = 0.524 * PA_{avg} - 0.0852 * RH + 5.72$$

where $PM_{2.5corrected}$ is the corrected $PM_{2.5}$, PA_{avg} is the hourly averaged $PM_{2.5}$ from channel A and B and RH is hourly relative humidity. For our study, We collected hourly measurements of $PM_{2.5}$ concentrations from January 2019 through to December 31, 2019, from 1015 PurpleAir sensors in California. Figure 3 illustrates an example of the data.

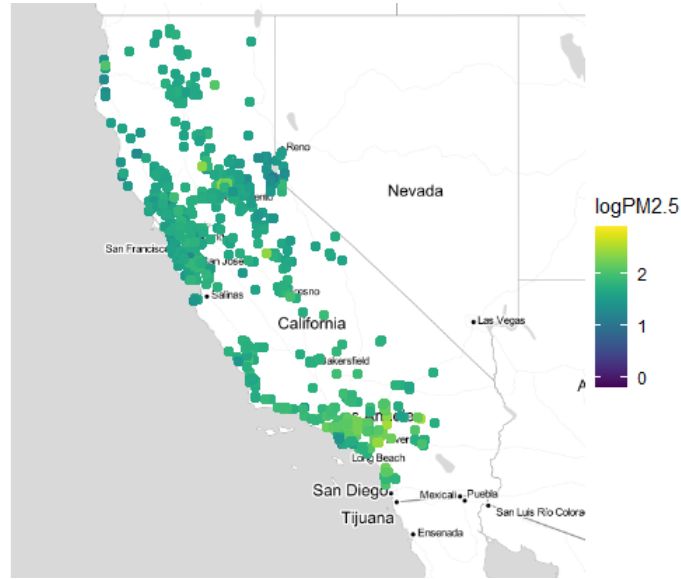


Figure 3: $PM_{2.5}$ concentration on the log scale on May 18th, 2019 at 2 pm.

3 Methods

3.1 Geostatistical Methods

Geostatistical methods are a set of statistical techniques that are used to analyze spatial/spatiotemporal data, explicitly take into account the spatial/spatiotemporal dependence of the data. These methods are widely used in various fields such as geology, hydrology, agriculture, environmental science, and many others (Maliva, 2016). A typical procedure in geostatistical techniques is to start with modeling the spatial/spatiotemporal correlation structure of the data, which describes how the values at different locations are related to each other in space and over time. Once we estimate the model parameter from the data, we can make predictions for the variable of interest at unsampled locations.

3.1.1 Universal Kriging (UK)

Let $Y(s, t)$ be the observed $\text{PM}_{2.5}$ (mg/m^3) from PurpleAir sensors at location $s \in \mathcal{D} \subset \mathbb{R}^2$ and time $t \in \mathbb{R}$. Then

$$Y(s, t) = \mu(s, t) + w(s, t) + \epsilon(s, t) \quad (1)$$

where, $\mu(s, t) = \mathbf{X}(s, t)\boldsymbol{\beta}$ is the mean with $\mathbf{X}(s, t) = [X_1(s, t), X_2(s, t), \dots, X_p(s, t)]^T$ are predictors and $\boldsymbol{\beta}$ is the p -dimensional vector of parameters. $w(s, t)$ represents a mean-zero Gaussian process with a spatiotemporal covariance function. $\epsilon(s, t)$ is the white noise measurement error with mean 0 and variance σ_ϵ^2 . We use a separable exponential covariance function,

$$\Gamma(\mathbf{h}; \tau, \boldsymbol{\theta}) = \sigma_s^2 \exp(-\|\mathbf{h}\|/\rho_s) \sigma_t^2 \exp(-|\tau|/\rho_t) \quad (2)$$

where $\mathbf{h} = s' - s$ is the spatial lag, $\tau = t' - t$ temporal lag and $\boldsymbol{\theta} = \{\sigma_s, \sigma_t, \rho_s, \rho_t\}$ is the set of parameters associated with this covariance function, where σ_s, σ_t are the standard deviations for space and time, and ρ_s, ρ_t are spatial and time range parameters respectively. To predict \mathbf{Y} at any unsampled location at (s_0, t_0) ,

we have

$$\hat{Y}(s_0, t_0) = \mathbf{X}(s_0, t_0)^T \hat{\boldsymbol{\beta}}_{gls} + c(\boldsymbol{\theta})^T \Sigma_y^{-1}(\boldsymbol{\theta})(\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{gls}), \quad (3)$$

where $c(\boldsymbol{\theta}) = \text{Cov}(Y(s_0, t_0), \mathbf{Y})$ and $\text{Var}(\mathbf{Y}) = \Sigma_y(\boldsymbol{\theta})$. The uncertainty can be obtained by

$$\text{Var}(\hat{Y}(s_0, t_0)) = c_{0,0} - c(\boldsymbol{\theta})^T \Sigma_y^{-1}(\boldsymbol{\theta}) c(\boldsymbol{\theta}) + \kappa, \quad (4)$$

where $c_{0,0} = \text{Var}(Y(s_0, t_0))$ and

$$\kappa = \left(\mathbf{X}(s_0, t_0) - \mathbf{X}^T \Sigma_y^{-1}(\boldsymbol{\theta}) c(\boldsymbol{\theta}) \right)^T \left(\mathbf{X}^T \Sigma_y^{-1}(\boldsymbol{\theta}) \mathbf{X} \right)^{-1} \left(\mathbf{X}(s_0, t_0) - \mathbf{X}^T \Sigma_y^{-1}(\boldsymbol{\theta}) c(\boldsymbol{\theta}) \right).$$

The above prediction method is called universal kriging. To get the prediction in (3) and prediction uncertainty in (4), we have to compute the estimates $\hat{\boldsymbol{\beta}}_{gls}$ and $\hat{\boldsymbol{\theta}}$ so that we get the estimated BLUP (EBLUP). A 95% prediction coverage is given by

$$\hat{Y}(s_0, t_0) \pm 1.96 \times \sqrt{\text{Var}(\hat{Y}(s_0, t_0))}. \quad (5)$$

In this study, we implemented universal kriging using the gstat package in R Gräler et al. (2016).

3.1.2 Nearest Neighbor Gaussian Process (NNGP)

The nearest neighbor Gaussian process extends Vecchia's approximation (Vecchia, 1988) to a process using conditional independence given information from neighboring locations (Datta et al., 2016a). For notation purposes, let

$$\mathbf{Y} = \left(Y(\mathbf{s}_1, t_1), Y(\mathbf{s}_2, t_1), \dots, Y(\mathbf{s}_m, t_1), \dots, Y(\mathbf{s}_1, t_T), \dots, Y(\mathbf{s}_m, t_T) \right)'. \quad (6)$$

The joint density of \mathbf{Y} can be written as the product of conditional densities, as

$$f(\mathbf{Y}) = \prod_{i=1}^N f(Y_i | Y_1, \dots, Y_{i-1}) \quad (7)$$

The idea for this method is to approximate $f(\mathbf{Y})$ given in (7) by

$$f(\mathbf{Y}) \approx \tilde{f}(\mathbf{Y}) = \prod_{i=1}^N f(\mathbf{Y}_i | \mathbf{Y}_{\mathcal{N}_i}) \quad (8)$$

where \mathcal{N}_i is the neighboring set for \mathbf{Y}_i and $\mathbf{Y}_{\mathcal{N}_i}$ are the observations in this set. Vecchia (1988) demonstrated that approximating the full conditional with a subset of m nearest neighbors provides an excellent approximation in the case of spatial data and Datta et al. (2016b) extended it to spatiotemporal data. In this study, we construct the m nearest neighbors using simple neighbor selection suggested by Datta et al. (2016b) and select the m nearest spatial neighbor with temporal lag 1. The NNGP is implemented with GpGp package in R (Guinness J, 2021). GpGp relies on ordering of the Gaussian process observations where the conditional distribution only conditions on a small subset of previous observations in the ordering. This results a sparse Cholesky factor of the precision matrix. After the estimation step, to predict at any unsampled location (s_0, t_0) ,

$$\hat{Y}(s_0, t_0) = \mathbf{X}(s_0, t_0)^T \hat{\boldsymbol{\beta}} + C_{s_0, \mathcal{N}_0} C_{\mathcal{N}_0}^{-1} (\mathbf{Y}_{\mathcal{N}_0} - \mathbf{X}(s_0, t_0)^T \hat{\boldsymbol{\beta}}), \quad (9)$$

where $C_{s_0, \mathcal{N}_0} = \text{Cov}(Y(s_0, t_0), \mathbf{Y}_{\mathcal{N}_0})$ and $C_{\mathcal{N}_0} = \text{Var}(\mathbf{Y}_{\mathcal{N}_0})$.

3.1.3 Fixed Rank Kriging (FRK)

Fixed Rank Kriging (FRK) is a spatial/spatiotemporal interpolation technique that combines the strengths of Kriging with a reduced rank approximation to address the computational challenges of large spatial datasets (Cressie and Johannesson, 2008). The Gaussian process here is modeled similarly to (1) except that $\mathbf{w}(s, t)$ is approximated by a linear combination of \mathbf{K} basis functions $\phi(\mathbf{s}, \mathbf{t}) = (\phi_1(\mathbf{s}, \mathbf{t}), \phi_2(\mathbf{s}, \mathbf{t}), \dots, \phi_{\mathbf{K}}(\mathbf{s}, \mathbf{t}))$ and the basis function coefficients $\mathbf{w}^* = (w_1^*, w_2^*, \dots, w_{\mathbf{K}}^*)$ (Wikle et al. (2019), Heaton et al. (2019)):

$$\mathbf{w}(s, t) \approx \tilde{\mathbf{w}}(s, t) = \sum_{i=1}^{\mathbf{K}} \phi_i(s, t) w_i^*. \quad (10)$$

This approximation reduces computational demands by ensuring that all estimations and predictions involve matrices of size $\mathbf{K} \times \mathbf{K}$, where $\mathbf{K} \ll \mathbf{N}$ and \mathbf{N} is the total number of observations. Moreover, $\phi(\mathbf{s}, \mathbf{t})$ can be composed at R different resolutions, thus

$$\tilde{\mathbf{w}}(s, t) = \sum_{r=1}^{\mathbf{R}} \sum_{i=1}^{\mathbf{K}_r} \phi_{ri}(s, t) w_{\mathbf{ri}}^*, \quad (11)$$

where $\phi_{ri}(s, t)$ is the i th spatiotemporal basis function at the r th resolution with basis function coefficients $w_{\mathbf{ri}}^*$ and the total number of basis functions is given by $\mathbf{K} = \sum_{r=1}^{\mathbf{R}} \mathbf{K}_r$. The coefficients $\mathbf{w}^* = (w_{\mathbf{ri}}^* : r = 1, \dots, R, i = 1, \dots, \mathbf{K}_r)$ can be modeled as spatial, temporal or spatiotemporal varying (Wikle et al., 2019). Here we consider the case where \mathbf{w}^* is spatial varying, thus $\mathbf{w}^* \sim \text{MVN}(0, \Sigma_{\mathbf{w}^*}(\boldsymbol{\theta}))$ and $\Sigma_{\mathbf{w}^*}(\boldsymbol{\theta})$ is chosen to be exponential covariance function, with $\boldsymbol{\theta} = (\sigma_s, \rho_s, \sigma_t = 1, \rho_t = 0)$ in equation (2). In this study, we construct $\{\phi_i(s, t), i = 1, \dots, \mathbf{K}\}$ by taking the tensor product of a spatial basis function with a temporal basis function (Wikle et al., 2019). Thus, the spatiotemporal basis functions $\phi(s, t) = \{\phi_{st,u} : u = 1, \dots, r_s r_t\} = \{\phi_p(s) \psi_q(t) : p = 1, \dots, r_s; q = 1, \dots, r_t\}$, where $r_s r_t$ are sets of spatial and temporal basis functions respectively. In this study, we choose bisquare basis functions for both spatial and temporal with two spatial resolutions and 1 temporal resolution (Heaton et al., 2019; Wikle et al., 2019)). We used 1600 basis functions with 80 basis functions for space and 20 basis functions for time and implemented this method using the FRK package in R (Zammit-Mangion and Cressie, 2021). With parameter estimates $\hat{\boldsymbol{\beta}}$, the FRK predictor at a new location (s_0, t_0) is given by

$$\hat{Y}(s_0, t_0) = \mathbf{X}(s_0, t_0) \hat{\boldsymbol{\beta}}' + c(s_0, t_0)^T \Sigma_y^{-1}(\hat{\boldsymbol{\theta}}) (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}), \quad (12)$$

where \mathbf{Y} as defined in (6), $c(s_0, t_0) = \text{Cov}(Y(s_0, t_0), \mathbf{Y})$ and $\text{Var}(\mathbf{Y}) = \Sigma_y(\boldsymbol{\theta})$. The uncertainty is given by

$$\text{Var}(\hat{Y}(s_0, t_0)) = \phi(s_0, t_0)^T \Sigma_{\mathbf{w}^*}(\hat{\boldsymbol{\theta}}) \phi(s_0, t_0) + \sigma_\epsilon^2 - c(s_0, t_0)^T \Sigma_y^{-1}(\hat{\boldsymbol{\theta}}) c(s_0, t_0). \quad (13)$$

3.2 Non-Geostatistical Methods Non-Geo

In this section will present methods that do not account for the spatiotemporal dependence in the data. We present the commonly used regression method and several machine learning (ML) algorithms. In recent years, ML algorithms have shown promise in modeling environmental data. Unlike traditional statistical models, machine learning algorithms don't rely on assumptions about the underlying distribution of the data and can identify complex patterns in the data and learn from them to make accurate predictions (Breiman, 2001; Di et al., 2016; Reid et al., 2015)). Furthermore, these algorithms can handle high-dimensional data and model complex non-linear relationships between variables (Athmaja et al., 2017). For the machine learning algorithm, we will discuss random forests, support vector regressions, and neural networks with a brief explanation for each algorithm.

3.2.1 Regression (Reg)

Consider a regression model where we assume that we will account for all spatiotemporal dependence in the predictors. The typical regression model is given by

$$Y(s, t) = \boldsymbol{\mu}(s, t) + \boldsymbol{\epsilon}(s, t), \quad (14)$$

where, $\boldsymbol{\mu}(s, t) = \mathbf{X}(s, t)\boldsymbol{\beta}$ is the mean with $\mathbf{X}(s, t) = [X_1(s, t), X_2(s, t), \dots, X_p(s, t)]^T$ are predictors and $\boldsymbol{\beta}$ is the p -dimensional vector of parameters and $\boldsymbol{\epsilon}(s, t)$ is the white noise measurement error with mean 0 and variance $\sigma_{\boldsymbol{\epsilon}}^2$. In this model, we will choose $\mathbf{X}(s, t)$ in terms of the nearest neighbor criteria. Although this model is simple to implement, it accounts for model error and allows us to obtain prediction error variance.

3.2.2 Random Forest (RF)

Random Forest Regression is a supervised learning algorithm that uses an ensemble of decision tree learning methods for regression.

First, the main idea of decision trees is to make tree that predicts a regression surface $\hat{f}(X|\hat{c}) =$

$\sum_{i=1}^M \hat{c}_i I(X \in R_i)$ for a partition of regions $\{R_i : i = 1, \dots, M\}$ and $\hat{c}_i \in \mathbb{R}$ for $i = 1, \dots, M$. In each decision tree, a random sample of m predictors is chosen as split candidates from the full set of p predictors (Breiman, 2001). The algorithm uses the method of least squares to minimize $RSS(c) = \sum_{i=1}^N (y_i - f(x_i|c))^2$. It results in $\hat{c}_m = \frac{1}{N_m} \sum_{i: x_i \in R_m} y_i$ where $N_m = |\{i : x_i \in R_m\}|$. The $\{R_i : i = 1, \dots, M\}$ is found using a greedy algorithm to grow the regression tree top-down. Decision trees method are not always a strong ML method and can be improved using ensembles.

Random forests are non-parametric ensemble machine learning algorithms that require the selection of two key parameters: the number of predictors in the random subset of each node (m) as the default value and the number of decision trees in the forest (n). This algorithm constructs multiple decision trees using bootstrapped training samples. At each node, a random sample of m features from the total p features is selected, and the algorithm identifies the optimal feature for creating a split. Predictions from the n trees are then aggregated using the mean to determine the final output. The error rate is assessed using predictions of out-of-bag samples (Ziegel, 2003). In this study, we used 500 trees to minimize the out-of-bag error, confirmed through five-fold cross-validation. Prediction coverage was evaluated by calculating the 2.5% and 97.5% quantiles across all 500 trees. This was implemented using the randomForest package in R (Liaw et al., 2002).

3.2.3 Support Vector Regression (SVR)

The groundwork for Support Vector Machines was provided by Vapnik and Chervonenkis (1974) and Vapnik (1999). A regression version of SVM, named Support Vector Regression (SVR), was introduced in Drucker et al. (1997). Based on selected kernel $k(x, y)$, a predictive model with input features \mathbf{x} is given by

$$F(x|w) = \sum_{i=1}^N (\alpha_i - \alpha_i^*) K(v_i, \mathbf{x}) \quad (15)$$

for vectors $(\alpha_1, \dots, \alpha_N)$ and $(\alpha_1^*, \dots, \alpha_N^*)$ and support vector v_i . Here we use the radial kernel $k(x, y) = \exp(-\frac{\sum_i (x_i - y_i)^2}{\gamma})$. The set of parameters is given by w . The primal objective function is

$$\lambda \sum_{i=1}^N \ell(y_i - F(x_i|w)) + \|w\|_2^2. \quad (16)$$

Note that the regularization parameter λ is introduced in front of the loss function and not the L_2 -regularization term. The best value of λ is to be found using cross-validation. Either α_i or α_i^* (but not both) will be non-zero based on the location of the observed point above or below the ϵ -tube, respectively. Both are zero if the point falls inside the tube. The ϵ -tube is induced by the ϵ -sensitive loss function defined so that it has a value of zero if the predicted value is within the tube.

The dual maximization problem is to

$$\max_{\alpha, \alpha^*} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N K(v_i, v_j) (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) - \sum_{i=1}^N y_i (\alpha_i - \alpha_i^*) - \epsilon \sum_{i=1}^N (\alpha_i + \alpha_i^*) \quad (17)$$

subject to

$$\sum_{i=1}^N \alpha_i = \sum_{i=1}^N \alpha_i^*, \quad 0 \leq \alpha_i, \alpha_i^* \leq \lambda \quad (18)$$

for $i = 1, 2, \dots, N$.

The dual maximization problem is equivalent to a quadratic programming problem

$$\min_{\beta} \frac{1}{2} \beta' Q \beta + \mathbf{c}' \beta \quad (19)$$

subject to

$$\sum_{i=1}^N \beta_i = \sum_{i=1}^N \beta_{i+N}, \quad 0 \leq \beta_i \leq \lambda \quad (20)$$

for $i = 1, 2, \dots, N$, where $\beta_i = \alpha_i^*$ and $\beta_{i+N} = \alpha_i$ for $i = 1, 2, \dots, N$ and $c_i = \epsilon - y_i$ and $c_{i+N} = \epsilon + y_i$ for $i = 1, 2, \dots, N$ while the block matrix $Q = \begin{bmatrix} D & -D \\ -D & D \end{bmatrix}$ with $D_{i,j} = K(v_i, v_j)$. We implemented SVR using

the e1071 package in R (Dimitriadou et al., 2008).

3.2.4 Neural Network (NN)

Neural networks algorithms are intended mimic human brain learning processes, see Figure 4. A neural network posses input data from an input layer through a linear map $\mathbf{L}_1 : X \mapsto X_{L_1}$ to the first hidden layer. The linear map $\mathbf{L}_1(X) = \mathbf{W}^{(1)}X + b^{(1)}$, where $\mathbf{W}^{(1)}$ is a $M_1 \times p$ matrix and $b^{(1)}$ is $M_1 \times 1$ bias vector. An activation function $\sigma_1 : X_{L_1} \mapsto Z^{(1)}$ activates each of the M_1 neurons in the first layer. The same process, using linear maps \mathbf{L}_i and activation functions σ_i , takes place between the $i - 1^{th}$ and i^{th} hidden layers with M_{i-1} and M_i neurons, respectively, for $i = 1, \dots, N$. The linear maps $\mathbf{L}_i(X) = \mathbf{W}^{(i)}X + b^{(i)}$, where $\mathbf{W}^{(i)}$ is a $M_i \times M_{i-1}$ matrix and $b^{(i)}$ is $M_i \times 1$ bias vector. A target linear map $\mathbf{T} : Z^{(N)} \mapsto T$ precedes the output function $g_k : T \mapsto Y_K$ at the output layer. In overall, it is a composition (\circ) of linear and activation functions

$$Y_k = g_k(\mathbf{T}(\circ_{i=1}^N \sigma_i(\mathbf{L}_i(Z^{(i-1)})))) =: f_k(X|\Theta)$$

with $Z^{(0)} := X$ and $f_k : X \mapsto Y_k$ is a nonlinear transformation of X into Y . Also, Θ is a vector of all parameters of the neural network that include the weights/entries in the matrices and the biases of the linear maps and the parameters of the activation and output functions. An optimization algorithm is used to minimize a loss function $\ell(\Theta) = \frac{1}{K} \sum_{i=1}^K \ell_i(y_k, f_k(X|\Theta))$. The loss function ℓ_i is usually the mean squared error $\ell_i(x, y) = \sum_{j=1}^n (x_j - y_j)^2$ in case of regression problems and cross-entropy $\ell_i(x, y) = -\sum_{j=1}^n x_j \log(y_j)$ in case of classification problems. An L_p -regularization is also used in which the objective function is $\ell(\Theta) + \lambda \|\Theta\|_p^p$ for $p = 1, 2$. The hyperparameter is tuned using cross-validation. The most common algorithm to find the optimal solution Θ is using back-propagation algorithm with the stochastic descent method.

There are different types of activation functions that can be used to activate the neurons

- Identity: $\sigma(x) = x$
- Sigmoid: $\sigma(x) = S(x) = \frac{1}{1+e^{-x}}$
- Hyperbolic tangent: $\sigma(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

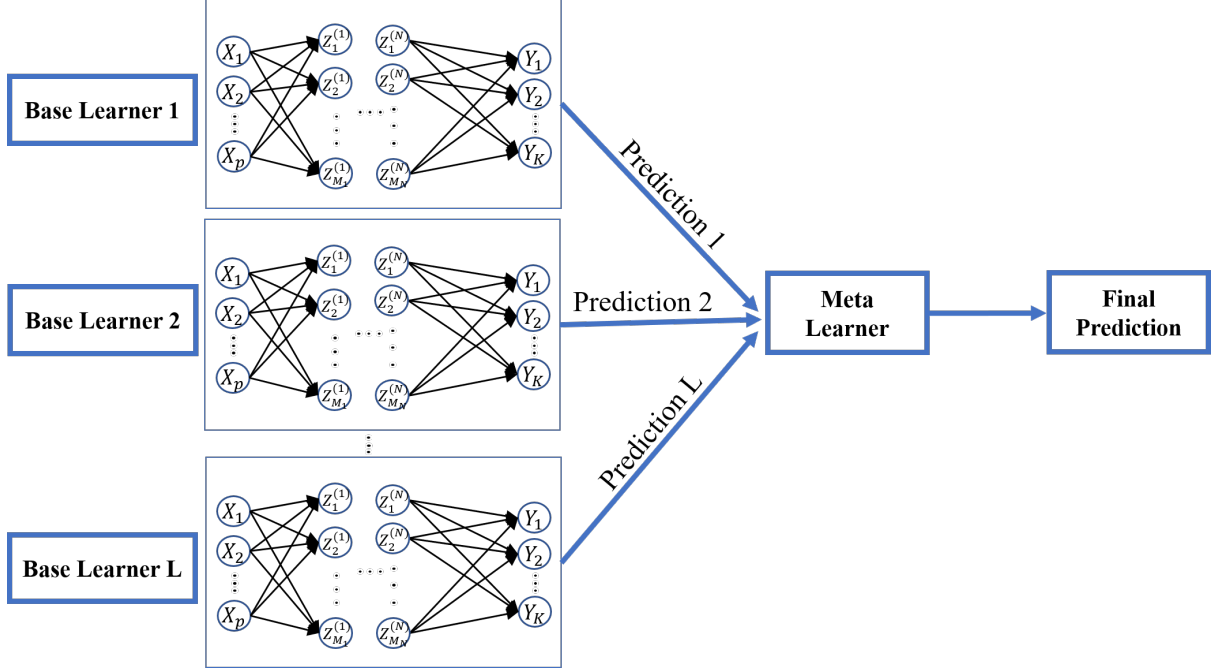


Figure 4: Illustration of a Multi-layer Neural Network with N hidden layers.

- Rectified Linear Unit (ReLU): $\sigma(x) = \text{ReLU}(x) = \max(x, 0) = x_+$
- Rectified softplus: $\sigma(x) = \text{ReSP}(x) = \log(1 + e^x)$

There are two main possible output functions g_k

- k^{th} element “identity” function: $g_k(T) = T_k$ which is used for regression problems
- Softmax function: $g_k(T) = \frac{e^{T_k}}{\sum_{\ell=1}^K e^{T_\ell}}$ which is used for classification problems.

Ensemble neural networks (ENN) are a type of ensemble method that use multiple neural network models, each with their own architectures, hyperparameters, and training data, to collectively make predictions. The idea behind ENNs is that by leveraging the complementary strengths and weaknesses of multiple models, we can achieve better performance than any individual model alone. There are several ways to construct an ENN. One common approach is to train several independent neural networks, called base models or base learners, on the training data and then combine their outputs through some form of averaging, in case of regression, or voting, in classification problems (Hansen and Salamon, 1990). Another approach is to use a single architecture and hyperparameters, but vary the initial weights and biases or the order in which

the data is presented during training (Ziegel, 2003). In this study, we constructed two neural networks, their predictions are combined to make a final omnibus prediction. Combining the predictions from the two neural networks can be done either by averaging or by assigning weights to each prediction according to the importance of the model, called stacking. The idea behind stacking is to use the predictions of multiple base models to train a higher-level “meta-model” that can make more accurate predictions than any of the individual models alone. We implemented ENN using Keras (Chollet et al., 2015) and Sklearn (Pedregosa et al., 2011) in Python (Van Rossum and Drake Jr, 1995).

4 Data Analysis

4.1 Feature Selection

Based on the set of features associated with each model and whether the model is considered to be a geostatistical model or non-geostatistical model, we divided all models into four groups. Group 1 contains geostatistical models namely: universal kriging (UK), nearest neighbor Gaussian process (NNGP), and fixed ranked kriging (FRK). Group 2 contains Non-geostatistical methods, namely regression (Reg), random forest (RF), support vector regression (SVR), and ensemble neural network (ENN). The set of features (p) for groups 1 and 2 will be the longitude and the latitude for each location. For group 3, we will add the set of the nearest neighbor observations (NNO) for each location and their longitude and latitude as features for Reg, RF, SVR, and ENN. Finally, for group 4 we will add the kriging prediction that we obtained from NNGP to be an additional feature in fitting Reg, RF, SVR, and ENN.

4.2 Performance Evaluation

To evaluate the performance of each model, we employed five-fold cross-validation. In each fold, PurpleAir (PA) sensors were randomly split spatially to mitigate spatial bias, with 80% for training and 20% for testing each model. There are two types of prediction: estimating PM_{2.5} concentrations at locations without PA sensors, which is our focus in this study; and forecasting future PM_{2.5} levels at locations with PA sensors.

The latter is not the subject of our current investigation. Model performance was assessed using five metrics: root mean squared error (RMSE) in (21), symmetric mean absolute percentage error (SMAPE) in (22), mean absolute deviation (MAD) in (23), the correlation (Cor) between predicted and observed PM_{2.5} concentrations, and 95% prediction coverage. For geostatistical methods, we constructed 95% prediction interval as illustrated in (5).

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}; \quad (21)$$

$$SMAPE = \frac{100\%}{n} \sum_{i=1}^n \frac{|Y_i - \hat{Y}_i|}{|Y_i| + |\hat{Y}_i|}; \quad (22)$$

$$MAD = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|, \quad (23)$$

where n denotes the number of test samples, Y_i represents the actual PM_{2.5} concentration, and \hat{Y}_i denotes the corresponding predicted value. These metrics are essential for assessing the accuracy of our predictive models.

4.3 Results

Table 1 presents the results of comparing geostatistical and non-geostatistical methods as detailed in Section 3. For the initial comparison of prediction performance, we considered geostatistical methods, designated as group 1. Out of the three geostatistical methods, universal kriging demonstrated superior performance across all six metrics, using longitude and latitude as predictors, followed by NNGP and FRK.

In group 2, focusing on non-geostatistical methods, SVR surpassed Reg, RF, and ENN for all metrics, using longitude and latitude as predictors. For group 3, we modified the predictors for the non-geostatistical methods (from group 2) to include the 10 nearest neighbor observations (NNO) of PM_{2.5} for each location with a one-time lag. Again, SVR outperformed all other methods in this group followed by RF, Reg, and ENN.

In group 4, we adjusted the set of predictors to include predictive values obtained from NNGP (the fastest

method from group 1), which significantly enhanced the performance for all machine learning algorithms. SVR maintained the best performance across all metrics.

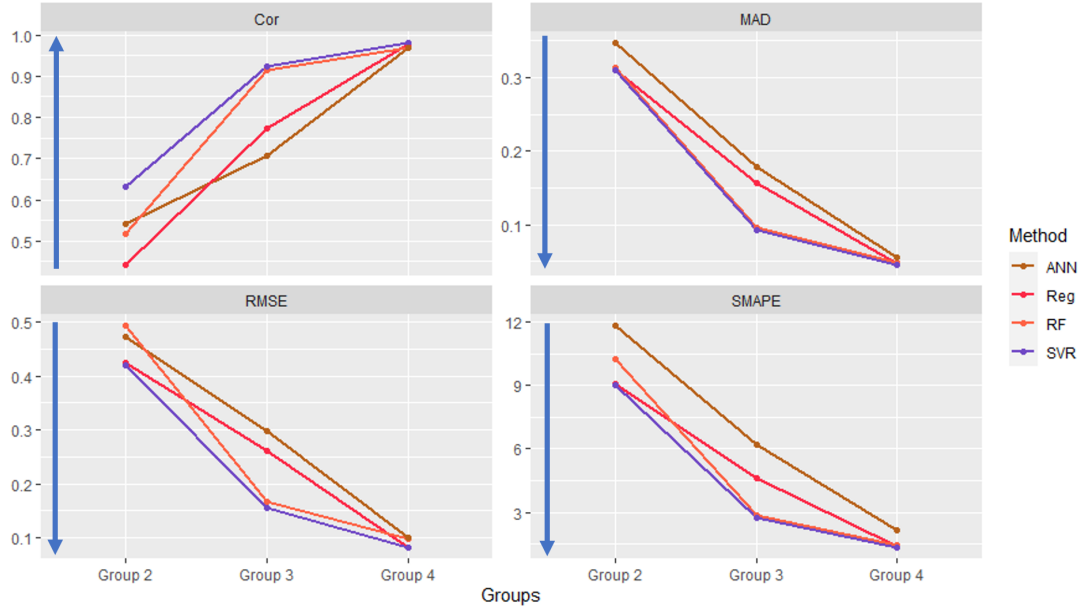


Figure 5: A comparison between all models in the Non-geostatistical groups with respect to root mean square error (RMSE), symmetric mean absolute percentage error (SMAPE), mean absolute deviation (MAD), and the correlation (Cor) between observed and predicted values.

Figure 5 illustrates how altering the set of predictors can improve RMSE, SMAPE, and Cor across groups 2,3, and 4. It is worth mentioning that, although the RF algorithm comes in second place after SVR, it offers a unique advantage: it allows for the construction of 95% empirical prediction coverage. Through five-fold cross-validation, RF has a predictive coverage of 95% with the addition of Nearest Neighbor Observations NNO and increased to 96% when incorporating the kriging output.

To generate an hourly map of $PM_{2.5}$ concentrations across California, we constructed a 50×50 grid over 24-hour period. Predictions were made using the top-performing model from each group: Universal Kriging (UK) from group 1, SVR from group 2, SVR with NNO from group 3, and SVR with NNO and Kriging from group 4. Figure 6 displays spatial predictions for $PM_{2.5}$ concentrations on June 15, 2019 at 2:00 pm. Particularly, panel (d) in 6, associated with the best performing method, reveals that southern California, especially along the coasts, experiences high $PM_{2.5}$ levels. Stowell et al. (2020) attribute these elevated $PM_{2.5}$

Method	RMSE	SMAPE	MAD	Cor	95% Cov	Time
Group 1: "Geostatical"						
UK	0.3730	7.802%	0.2630	0.8701	93%	3600
NNGP	0.4076	9.117%	0.3174	0.7925	89%	5
FRK	0.4227	9.477%	0.3537	0.7758	85%	420
Group 2: "Non-Geo"						
Reg	0.4239	9.098 %	0.3102	0.4416	40%	0.1025
RF	0.4927	10.278 %	0.3140	0.5186	67%	3.5726
SVR	0.4186	8.994%	0.3104	0.6333	-	1.5142
ENN	0.4725	11.825%	0.3464	0.5412	-	5.6514
Group 3: "Non-Geo+NNO"						
Reg+NNO	0.2608	4.625%	0.1567	0.7726	90%	0.5078
RF+NNO	0.1674	2.848%	0.0964	0.9133	95%	6.0145
SVR+NNO	0.1552	2.740%	0.0920	0.9249	-	5.1140
ENN+NNO	0.2983	6.208%	0.1785	0.7062	-	6.7755
Group 4: "Non-Geo+NNO+Krig"						
Reg+NNO+Krig	0.0826	1.390%	0.04697	0.9781	97%	6.0597
RF+NNO+Krig	0.0989	1.458%	0.0493	0.9694	96%	12.6580
SVM+NNO+Krig	0.0819	1.330%	0.0449	0.9792	-	10.9344
ENN+NNO+Krig	0.0994	2.149%	0.0552	0.9690	-	13

Table 1: A comparison among all models in the four groups with respect to root mean square error (RMSE), symmetric mean absolute percentage error (SMAPE), mean absolute deviation (MAD), the correlation between the observed and predicted values (Cor), 95% empirical coverage (95% Cov), and execution time in minutes (Time).

concentrations in southern California to Santa Ana winds, which can alter temperature, humidity, and wind speed, thereby escalating PM2.5 levels.

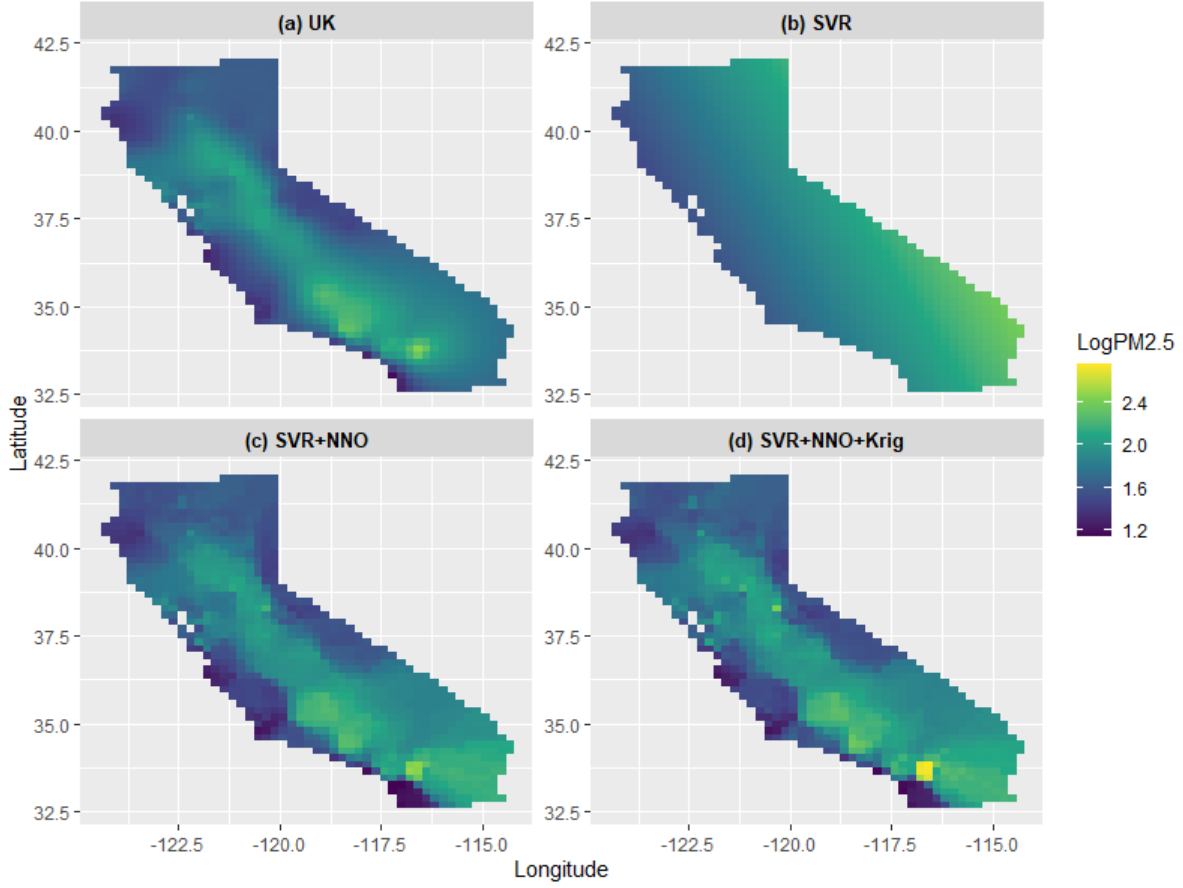


Figure 6: Predicted log PM_{2.5} concentration for June 15, 2019 at 2:00 pm using (a) Universal Kriging (UK), (b) Support Vector Regression (SVR) with longitude and latitude as predictors, (c) SVR with 10 NNO as predictors and (d) SVR with 10 NNO and Kriging as predictors.

5 Conclusions and Discussion

This study provides a comprehensive assessment of prevalent geostatistical and machine learning methods for estimating PM_{2.5} concentrations across spatial and temporal scales, utilizing an extensive network of PurpleAir sensors to enhance prediction accuracy for air quality in California. Among geostatistical methods, universal kriging shows best performance across all six metrics followed by Nearest Neighbor Gaussian Process (NNGP) and Fixed Rank Kriging (FRK). In machine learning methods, Support Vector Regression (SVR) surpassed Regression (Reg), Random Forest (RF), and Ensemble Neural Network (ENN).

While geostatistical methods effectively capture spatiotemporal dependencies, they are constrained by

high dimensionality and computational demands. For instance, Universal Kriging (UK) required up to 3600 minutes, compared to 420 minutes for Fixed Rank Kriging (FRK) and only 5 minutes for the Nearest Neighbor Gaussian Process (NNGP). Machine learning algorithms offer rapid implementation and robustness to non-linear relationships, but they generally overlook spatiotemporal dependencies and uncertainty. For the non-geostatistical algorithms, the maximum execution time was 8 minutes recorded for ensemble neural network (ENN).

Our findings highlight the synergistic potential of integrating geostatistical approaches with machine learning to enhance prediction accuracy without relying on assumptions of stationarity, linearity, or normality. Particularly, incorporating predictors from the fastest geostatistical method in our comparison, NNGP, into machine learning models significantly improved prediction quality, reduced RMSE, SMAPE, and MAD, and increased the correlation between observed and predicted PM_{2.5} values. Notably, Support Vector Regression (SVR) excelled across all metrics, demonstrating its capability in enhancing the spatiotemporal prediction of PM_{2.5} concentrations. Overall, our comparison addresses the often-neglected combined spatiotemporal dynamic, which are critical for understanding and managing air quality. More importantly, our findings reveal significant advantages when integrating geostatistical methods with machine learning algorithms.

Earlier research has confirmed that human activities, such as industrial processes, transportation, and agricultural practices significantly impact ground-level PM_{2.5} concentrations (Bao et al., 2016; Pérez et al., 2010)). Additionally, correlations between land cover/land use (LCLU) types and PM_{2.5} concentrations suggest that different land uses may reflect varied emission sources (Hoek et al., 2008). Therefore, incorporating meteorological variables, such as temperature, humidity, dew-point temperature, alongside these human factors could further refine predictive models, offering insights into pollution hotspots and assisting policymakers in optimizing air quality monitoring networks.

Looking ahead, future studies could explore the integration of more diverse data sources, including real-time traffic and industrial emissions data, to further contextualize the spatiotemporal dynamics of air pollution. Additionally, advancing algorithmic approaches to more effectively process large-scale environmental data could provide deeper insights into the complex interactions affecting air quality. Such advancements will not

only improve predictive accuracy but also enhance our understanding of the environmental and anthropogenic factors driving PM_{2.5} variability, supporting more targeted and effective air quality management strategies.

Acknowledgments

The authors would like to thank Yawen Guan for helpful discussions and we also appreciate the effort of the Editor, Associate Editor and Referees to improve quality of the manuscript.

References

- Airsensor. <https://api.purpleair.com/>. Accessed: 04 06, 2023.
- Stacey E Alexeeff, Joel Schwartz, Itai Kloog, Alexandra Chudnovsky, Petros Koutrakis, and Brent A Coull. Consequences of kriging and land use regression for PM_{2.5} predictions in epidemiologic analyses: insights into spatial variability using high-resolution satellite data. *Journal of exposure science & environmental epidemiology*, 25(2):138–144, 2015.
- S Athmaja, M Hanumanthappa, and Vasantha Kavitha. A survey of machine learning algorithms for big data analytics. In *2017 International conference on innovations in information, embedded and communication systems (ICIIECS)*, pages 1–4. IEEE, 2017.
- Chengzhen Bao, Pengfei Chai, Hongbo Lin, Zhenyu Zhang, Zhenhua Ye, Mengjia Gu, Huaichu Lu, Peng Shen, Mingjuan Jin, Jianbing Wang, et al. Association of PM_{2.5} pollution with the pattern of human activity: A case study of a developed city in eastern China. *Journal of the Air & Waste Management Association*, 66(12):1202–1213, 2016.
- KK Barkjohn, B Gantt, and AL Clements. Development and application of a United States wide correction for PM_{2.5} data collected with the PurpleAir sensor. *Atmos. Meas. Tech. Discuss.* <https://doi.org/10.5194/amt-2020-413>, 2020.
- Veronica J Berrocal, Yawen Guan, Amanda Muyskens, Haoyu Wang, Brian J Reich, James A Mulholland, and Howard H Chang. A comparison of statistical and machine learning methods for creating national daily maps of ambient PM_{2.5} concentration. *Atmospheric Environment*, 222:117130, 2020.
- Jianzhao Bi, Avani Wildani, Howard H Chang, and Yang Liu. Incorporating low-cost sensor measurements into high-resolution PM_{2.5} modeling at a large spatial scale. *Environmental Science & Technology*, 54(4):2152–2162, 2020.
- Leo Breiman. Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical science*, 16(3): 199–231, 2001.
- Francois Chollet et al. Keras, 2015. URL <https://github.com/fchollet/keras>.

- Noel Cressie and Gardar Johannesson. Fixed rank kriging for very large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(1):209–226, 2008.
- Abhirup Datta, Sudipto Banerjee, Andrew O Finley, and Alan E Gelfand. Hierarchical nearest-neighbor Gaussian process models for large geostatistical datasets. *Journal of the American Statistical Association*, 111(514):800–812, 2016a.
- Abhirup Datta, Sudipto Banerjee, Andrew O Finley, Nicholas AS Hamm, and Martijn Schaap. Nonseparable dynamic nearest neighbor Gaussian process models for large spatio-temporal data with an application to particulate matter analysis. *The annals of applied statistics*, 10(3):1286, 2016b.
- William W Delp and Brett C Singer. Wildfire smoke adjustment factors for low-cost and professional PM2.5 monitors with optical sensors. *Sensors*, 20(13):3683, 2020.
- Qian Di, Itai Kloog, Petros Koutrakis, Alexei Lyapustin, Yujie Wang, and Joel Schwartz. Assessing PM2.5 exposures with high spatiotemporal resolution across the continental United States. *Environmental science & technology*, 50(9):4712–4721, 2016.
- Evgenia Dimitriadou, Kurt Hornik, Friedrich Leisch, David Meyer, and Andreas Weingessel. Misc functions of the department of statistics (e1071), tu wien. *R package*, 1:5–24, 2008.
- Harris Drucker, Chris J.C. Surges, Linda Kaufman, Alex Smola, and Vladimir Vapnik. Support vector regression machines. In *Advances in Neural Information Processing Systems*, 1997. ISBN 0262100657.
- Benedikt Gräler, Edzer J Pebesma, and Gerard BM Heuvelink. Spatio-temporal interpolation using gstat. *R J.*, 8(1):204, 2016.
- Fahmy Y Guinness J, Katzfuss M. fast Gaussian process computation using vecchia’s approximation. *r package version 0.4.0*. 2021.
- Pawan Gupta and Sundar A Christopher. Particulate matter air quality assessment using integrated surface, satellite, and meteorological products: 2. a neural network approach. *Journal of Geophysical Research: Atmospheres*, 114(D20), 2009.
- Anjum Hajat, Charlene Hsia, and Marie S O’Neill. Socioeconomic disparities and air pollution exposure: a global review. *Current environmental health reports*, 2:440–450, 2015.
- Lars Kai Hansen and Peter Salamon. Neural network ensembles. *IEEE transactions on pattern analysis and machine intelligence*, 12(10):993–1001, 1990.
- Matthew J Heaton, Abhirup Datta, Andrew O Finley, Reinhard Furrer, Joseph Guinness, Rajarshi Guhaniyogi, Florian Gerber, Robert B Gramacy, Dorit Hammerling, Matthias Katzfuss, et al. A case study competition among methods for analyzing large spatial data. *Journal of Agricultural, Biological and Environmental Statistics*, 24:398–425, 2019.

- Gerard Hoek, Rob Beelen, Kees De Hoogh, Danielle Vienneau, John Gulliver, Paul Fischer, and David Briggs. A review of land-use regression models to assess spatial variation of outdoor air pollution. *Atmospheric environment*, 42(33):7561–7578, 2008.
- Xuefei Hu, Lance A Waller, Mohammad Z Al-Hamdan, William L Crosson, Maurice G Estes Jr, Sue M Estes, Dale A Quattrochi, Jeremy A Sarnat, and Yang Liu. Estimating ground-level PM2.5 concentrations in the southeastern us using geographically weighted regression. *Environmental research*, 121:1–10, 2013.
- Xuefei Hu, Jessica H Belle, Xia Meng, Avani Wildani, Lance A Waller, Matthew J Strickland, and Yang Liu. Estimating PM2.5 concentrations in the conterminous United States using the random forest approach. *Environmental science & technology*, 51(12):6936–6944, 2017.
- Andy Liaw, Matthew Wiener, et al. Classification and regression by randomforest. *R news*, 2(3):18–22, 2002.
- Robert G Maliva. Geostatistical methods and applications. *Aquifer Characterization Techniques: Schlumberger Methods in Water Resources Evaluation Series No. 4*, pages 595–617, 2016.
- Caroline Mogollón-Sotelo, Alejandro Casallas, Sergio Vidal, Nathalia Celis, Camilo Ferro, and Luis Belalcazar. A support vector machine model to forecast ground-level PM2.5 in a highly populated city with a complex terrain. *Air Quality, Atmosphere & Health*, 14:399–409, 2021.
- Lidia Morawska, Phong K Thai, Xiaoting Liu, Akwasi Asumadu-Sakyi, Godwin Ayoko, Alena Bartonova, Andrea Bedini, Fahe Chai, Bryce Christensen, Matthew Dunbabin, et al. Applications of low-cost sensing technologies for air quality monitoring and exposure assessment: How far have they gone? *Environment international*, 116:286–299, 2018.
- World Health Organization et al. Noncommunicable diseases country profiles 2018. 2018.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830, 2011.
- Noemí Pérez, Jorge Pey, Michael Cusack, Cristina Reche, Xavier Querol, Andrés Alastuey, and Mar Viana. Variability of particle number, black carbon, and PM10, PM2.5, and PM1 levels and speciation: influence of road traffic emissions on urban air quality. *Aerosol Science and Technology*, 44(7):487–499, 2010.
- PurpleAir. PurpleAir - real time air quality monitoring, 2021. URL <https://www.purpleair.com/>.
- Colleen E Reid, Michael Jerrett, Maya L Petersen, Gabriele G Pfister, Philip E Morefield, Ira B Tager, Sean M Raffuse, and John R Balmes. Spatiotemporal prediction of fine particulate matter during the 2008 northern california wildfires using machine learning. *Environmental science & technology*, 49(6):3887–3896, 2015.

- Weeberb J Requia, Brent A Coull, and Petros Koutrakis. Evaluation of predictive capabilities of ordinary geostatistical interpolation, hybrid interpolation, and machine learning methods for estimating PM2.5 constituents over space. *Environmental research*, 175:421–433, 2019.
- Jennifer D Stowell, Jianzhao Bi, Mohammad Z Al-Hamdan, Hyung Joo Lee, Sang-Mi Lee, Frank Freedman, Patrick L Kinney, and Yang Liu. Estimating PM2.5 in southern california using satellite data: factors that affect model performance. *Environmental Research Letters*, 15(9):094004, 2020.
- Guido Van Rossum and Fred L Drake Jr. *Python reference manual*. Centrum voor Wiskunde en Informatica Amsterdam, 1995.
- Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 1999.
- Vladimir Vapnik and Alexey Chervonenkis. Theory of pattern recognition, 1974.
- Aldo V Vecchia. Estimation and model identification for continuous spatial processes. *Journal of the Royal Statistical Society: Series B (Methodological)*, 50(2):297–312, 1988.
- Haoyu Wang, Yawen Guan, and Brain Reich. Nearest-neighbor neural networks for geostatistics. In *2019 international conference on data mining workshops (ICDMW)*, pages 196–205. IEEE, 2019.
- Christopher K Wikle, Andrew Zammit-Mangion, and Noel Cressie. *Spatio-temporal Statistics with R*. Chapman and Hall/CRC, 2019.
- Andrew Zammit-Mangion and Noel Cressie. Frk: An r package for spatial and spatio-temporal prediction with large datasets. *Journal of Statistical Software*, 98:1–48, 2021.
- Eric R Ziegel. The elements of statistical learning, 2003.