

Review

# Critical Review of Methods to Estimate PM<sub>2.5</sub> Concentrations within Specified Research Region

Guangyuan Zhang, Xiaoping Rui \* and Yonglei Fan

College of Resources and Environment, University of Chinese Academy of Sciences, Beijing 100049, China; zhangguangyuan16@mails.ucas.ac.cn (G.Z.); fanyonglei17@mails.ucas.ac.cn (Y.F.)

\* Correspondence: ruixpsz@163.com; Tel.: +86-136-7110-6220

Received: 15 June 2018; Accepted: 20 August 2018; Published: 7 September 2018



**Abstract:** Obtaining PM<sub>2.5</sub> data for the entirety of a research region underlies the study of the relationship between PM<sub>2.5</sub> and human spatiotemporal activity. A professional sampler with a filter membrane is used to measure accurate values of PM<sub>2.5</sub> at single points in space. However, there are numerous PM<sub>2.5</sub> sampling and monitoring facilities that rely on data from only representative points, and which cannot measure the data for the whole region of research interest. This provides the motivation for researching the methods of estimation of particulate matter in areas having fewer monitors at a special scale, an approach now attracting considerable academic interest. The aim of this study is to (1) reclassify and particularize the most frequently used approaches for estimating the PM<sub>2.5</sub> concentrations covering an entire research region; (2) list improvements to and integrations of traditional methods and their applications; and (3) compare existing approaches to PM<sub>2.5</sub> estimation on the basis of accuracy and applicability.

**Keywords:** PM<sub>2.5</sub> concentrations; spatial interpolation; remote sensing; air-quality model; CMAQ model; machine learning

## 1. Introduction

Air pollution has become a serious hazard for human health and public welfare in recent decades. Particularly problematic is PM<sub>2.5</sub>, the most well-known particulate pollutant, which is suspended particulate matter smaller than 2.5 µm in aerodynamic diameter. Concentrations of PM<sub>2.5</sub> have become a global problem, especially in England, China, India, and the United States [1]. The effects of PM<sub>2.5</sub> on public health have been well documented in the literature [2,3]. The World Health Organization (WHO) released the latest internationally-applicable air quality guidelines (AQG) on 6 October 2006. In these AQG, the WHO refers to the standard value of PM<sub>10</sub> (coarse particles with a diameter between 2.5 and 10 µm) and determines the standard value of PM<sub>2.5</sub> with reference to the mass concentration ratio of PM<sub>2.5</sub> and PM<sub>10</sub> to 0.5 [4]. The measurement unit of PM<sub>2.5</sub> is generally micrograms per cubic meter (µg/m<sup>3</sup>).

The acquisition of PM<sub>2.5</sub> data is the first step in research work of this nature. A professional sampler with a filter membrane is used to measure accurate values of PM<sub>2.5</sub> at single points in space [5]. However, many PM<sub>2.5</sub> sampling and monitoring facilities consider only vital representative data points, an approach which cannot give an accurate representation of data for the whole region of research interest. Therefore, researching the methods of estimation of particulate matter in areas having fewer monitors at a special scale is an important approach that is currently being extensively researched.

Different approaches have different levels of accuracy in estimating results. Accuracy levels are also dependent on different research regions, which may have complex geographic conditions, and on the spatiotemporal resolution of resource data. Approaches to predicting PM<sub>2.5</sub> are classified into two types: those involving ground-level monitor-based estimation, and those relying on satellite-based

(monitor-free) estimation [6]. The general methods involving ground-level monitor-based estimation include (a) land use regression (LUR) models; (b) generalized additive mixed models; (c) hierarchical models; and (d) geostatistical interpolation. The approaches that depend on satellite-based (monitor-free) estimation specifically refer to remote sensing techniques. However, in recent years, some new techniques have been pioneered for generating accurate PM<sub>2.5</sub> data. Researchers have improved the abovementioned traditional approaches using new methods specific to their research conditions. Additionally, there are manifold relative studies which combine two or more approaches to achieve precise results. For example, one of the most popular techniques—machine learning—is also applied to PM<sub>2.5</sub> estimations [7]. A nuanced understanding requires review, sorting, and classification of all these the approaches to estimating PM<sub>2.5</sub> concentrations covering entire research regions.

The aim of this study is to (1) reclassify and particularize the most frequently used approaches to estimate PM<sub>2.5</sub> concentrations covering the entire research region; (2) list improvements and integrations of traditional methods and their applications; and (3) compare the existing approaches of PM<sub>2.5</sub> estimations on the basis of accuracy and applicability.

The remainder of this paper is organized as follows. Section 2 introduces the most frequently used methods. Section 3 discusses the advantages and disadvantages of the methods outlined in Section 2, and presents the main upgraded methods of PM<sub>2.5</sub> estimation and current integrations, as well as introducing possible comparisons based on the existing research. Finally, Section 4 presents the conclusion of this study.

## 2. Most Frequently Used Approaches to the Estimation of PM<sub>2.5</sub> Concentrations Covering the Entire Research Region in Recent Years

This section may be divided by subheadings. It should provide a concise and precise description of the experimental results, their interpretation as well as the experimental conclusions that can be drawn.

### 2.1. Spatial Interpolation

Spatial interpolation is an important method to estimate unknown data by using known sample data. The commonly used spatial interpolation methods are inverse distance weighted (IDW) interpolation, ordinary kriging (OK) interpolation, and trend surface (TS) interpolation [8,9]. Additionally, there exist other methods such as collaborative kriging (CK) and the radial basis function, which are referenced by a number of researchers. Spatial interpolation distribution in different regions can employ various interpolation approaches to incorporate different interpolation effects and accuracy [10–12].

#### 2.1.1. OK Interpolation Method

The OK interpolation method assumes that sampling the point between the distance or direction can be used to illustrate the spatial correlation of the surface changes, where the mathematical function with the specified number of points or designated radius in all points is fitted to determine the location of each output value. The calculation formula is as follows:

$$Z_v^*(x) = \sum_{i=1}^n \lambda_i Z(x_i) \quad (1)$$

where  $Z(x)$  is the measurement of position  $i$ ,  $\lambda_i$  is the unknown weight of the measurement value at position  $i$ ,  $x$  is the predicted position, and  $n$  is the number of measurements.

### 2.1.2. IDW Interpolation Method

IDW interpolation is used to determine pixel values by a linear combination of a set of sampling points, which is assumed to be reduced by the distance between the mapped variables and the sampling locations. The calculation formula is as follows:

$$Z = \left[ \sum_{i=1}^n \frac{Z_i}{d_i^k} \right] / \left[ \sum_{i=1}^n \frac{1}{d_i^k} \right] \quad (2)$$

where  $Z$  is the estimated value of ,  $Z_i$  is the value of the control point  $i$ ,  $d_i$  is the distance between and  $i$ ,  $n$  is used in the estimation of the number of control points, and  $k$  is the power to be specified.

Keler and Krisp used the IDW method to create a visual overview for 36 established positions in Beijing over the period of one year (from 2 August 2013 until 2 August 2014) with varying PM<sub>2.5</sub> measurements in time [13].

### 2.1.3. TS Method

TS analysis is a statistical method based on known points in space, fitting a continuous mathematical surface, and studying the variation regularity of geological variables in the region and local area. The calculation formula is as follows:

$$Z = \beta_1 + \beta_2x + \beta_3y + \beta_4x^2 + \beta_5xy + \beta_6y^2 + \dots \quad (3)$$

where  $Z$  is the address variable,  $x$  and  $y$  are the coordinates of the observation points.

Ping Zhang and Taotao Shen [14] compared the accuracies of the IDW, TS, and OK interpolation methods based on GIS and Spearman correlation.

### 2.1.4. CK Method

One of the foundational components of geostatistics, kriging interpolation is a method of unbiased optimal estimation for regionalized variables, based on variation function theory and structural analysis [15]. The CK method uses one or more secondary variables to interpolate the variables of interest, which are related to the main variables, and assumes that the correlation between the variables can be used to improve the accuracy of the main predicted values.

Generally, there are several measured points that obey normal distribution. For any unknown point to be estimated, its estimator is expressed as a linear combination of the effective sample values, as follows:

$$\hat{Z}(s_i) = \sum_{j=1}^n \lambda_j Z(s_j) \quad (4)$$

where  $s_i$  is the predicted position,  $Z(s_j)$  is measured value of  $i$ th point, and  $n$  is the number of the measured points. The  $\lambda_j$  is unknown weight of  $Z(s_j)$ ,  $\lambda_j$  depends on the spatial relationship model for fitting the distances  $s_i-s_j$  and the measured values of  $s_i$ . In order to ensure that the model is an unbiased estimation,

$$\sum_{j=1}^n \lambda_j = 1 \quad (5)$$

where  $\hat{Z}(s_i)$  can be calculated under the condition of ensuring the kriging variance minimum.

Deng [16] gathered PM<sub>2.5</sub> concentration data from the Beijing City Environmental Protection Bureau from 32 observation sites over six days, with the intention of analyzing the distribution of PM<sub>2.5</sub> concentrations. Use the kriging interpolation method Deng predicted the concentrations of  $100 \times 100$  unknown points, and drew up PM<sub>2.5</sub> spatial distribution plots. The results show that PM<sub>2.5</sub> concentration has a gradient distribution pattern.

## 2.2. Remote Sensing Technique

Monitoring networks of monitors can only provide direct point-level observations at limited locations. At locations without monitors, a widely-employed approach to derive ground-level PM<sub>2.5</sub> concentrations uses observations of aerosol properties from satellite-based remote sensing, such as moderate resolution imaging spectroradiometer (MODIS) data, and aerosol optical depth (AOD) data [17].

AOD is defined as the integral of the aerosol extinction coefficient in the vertical direction for the entire atmosphere:

$$\text{AOD} = \int_0^{\infty} \sigma_{\text{ext},z} dz, \quad (6)$$

where  $\sigma_{\text{ext},z}$  is the extinction coefficient of aerosols at  $z$ :

$$\sigma_{\text{ext},z} = \pi \int Q_{\text{ext}}(m, r, \lambda) n(r) r^2 dr \quad (7)$$

Here,  $Q_{\text{ext}}$  is a function of the refractive index  $m$ , particle radius  $r$ , and wavelength  $\lambda$  and  $n(r)$  is the aerosol particle size distribution.

From the above two formulas, we can see that AOD is related to radiation wavelength, aerosol size, vertical profile, and particle size distribution. Early theoretical studies of multiangle imaging spectrometers mounted on the Terra satellite in the United States show that the particle sizes corresponding to the AOD inversion at the visible and near-infrared bands range from 0.1 to 2 nm, and are very close to the particle size range of PM<sub>2.5</sub>. This research provides an important theoretical basis for establishing the correlation between satellite remote sensing AOD and PM<sub>2.5</sub>. As such, the AOD obtained by satellite remote sensing can provide an effective means for monitoring PM<sub>2.5</sub> pollution.

The establishment of PM<sub>2.5</sub>-AOD relationships is affected by many factors such as the AOD vertical profile, humidity, temperature, and wind speed. Some land-use or geographical parameters such as area classification (urban/rural), road distribution, and forest cover are related to PM<sub>2.5</sub> emission sources. For the establishment of a PM<sub>2.5</sub> AOD model, these land-use parameters have very effective auxiliary variables [18]. Therefore, in most PM<sub>2.5</sub>-AOD advanced statistical models, various meteorological and land-use parameters are often implemented to improve accuracy.

Recently, more complex models have been proposed to describe the varied relationships between AOD and PM<sub>2.5</sub> levels. Ma Zongwei et al. presented geographically weighted regression models, examined by using adaptive bandwidths, and selected by the cross-validation (CV) method or used Akaike's information criterion (AIC) [19]. Lee et al. presented a linear mixed-effect model and established day-specific PM<sub>2.5</sub>-AOD relationships using the mixed-effects model to fully exploit satellite data [20]. Remote-sensing formulas are used worldwide [21–23]. Ma et al. developed a two-stage spatial statistical model using MODIS Collection 6 AOD and assimilated meteorology, land-use data, and PM<sub>2.5</sub> concentrations from China's recently established ground monitoring network. An inverse variance weighting approach was developed to combine the MODIS Dark Target and Deep Blue AOD methods to optimize data coverage; its evaluation model predicted PM<sub>2.5</sub> concentrations from the year 2004 to early 2014 using ground observations [24]. Chang et al. describe a statistical downscaling approach that combines (1) recent advances in PM<sub>2.5</sub> land-use regression models utilizing AOD; and (2) statistical data fusion techniques for combining air quality datasets that have different spatial resolutions [25]. Further, Lv [26] employed a Bayesian-based statistical down-scaler to model spatiotemporal linear AOD-PM<sub>2.5</sub> relationships.

In the next stage, the unmonitored PM<sub>2.5</sub> concentration covering the entire region is calculated by the relationships and all other factors, expecting unknown PM<sub>2.5</sub> values.

### 2.3. Air Quality Model Methods

Air quality models generally consider the following minimum set of atmospheric processes: emissions (anthropogenic and natural source emissions), transport (horizontal advection and vertical convection), diffusion (horizontal and vertical diffusion), chemical conversion (gas, liquid, and solid chemical reactions), and scavenging mechanism (dry and wet settlement). Three theoretical systems of turbulent diffusion inform research, specifically gradient transport theory (K theory), statistical theory, and similarity theory. At the core of statistical theory is the spatiotemporal probability distribution of the diffusion particle; in other words, the spatial distribution and time change of the concentration of the diffusion particle are described by the probability distribution function. As such, the air quality model can simulate the transmission process of PM<sub>2.5</sub> and predict its probability distribution on a spatial scale. The third generation of integrated air quality models, such as CMAQ, CAMx, and WRF-Chem, is the most widely used.

#### 2.3.1. Community Multiscale Air Quality (CMAQ) Model Aerosol Component

CMAQ was first released by the U.S. Environmental Protection Agency in June 1998. After more than ten years of research and development updates, it is currently at version 5.2.1 [27]. During simulation, CMAQ can integrate the influence of weather systems and small-scale meteorological processes on the transport, diffusion, transformation, and migration of pollutants. At the same time, the interaction of air pollutants with factors like the region and urban scale, as well as various chemical processes of pollutants in the atmosphere (including liquid chemistry, process effects, heterogeneous chemical processes, aerosol processes, and the effect of dry-wet deposition on concentration distribution) are taken into consideration [28].

The CMAQ model is composed of five main modules [29]. The core module is the chemical transmission module, CCTM (CMAQ Chemical-Transport Model Processor), which can simulate the transport process, chemical process, and settlement of pollutants. The Initial Conditions Processor module (ICON) and the Boundary Conditions Processor module (BCON), provide the initial and boundary fields for CCTM. The Photolysis Rate Processor module (JPROC) calculates the photochemical decomposition rate. The interface processor is the interface between the meteorological model and the CCTM and transforms meteorological data into a CCTM-identifiable data format.

The Models-3/CMAQ is a powerful tool for the study of air pollution and has been extensively applied in the estimation of PM<sub>2.5</sub> concentrations. It considers the complex physical and chemical processes involved, and describes the actual atmosphere in a holistic way, thereby obtaining not only the meteorological elements in the study area, but also spatiotemporal distribution, evolution of the pollutants, and other important factors [30].

#### 2.3.2. Comprehensive Air Quality Model with Extensions (CAMx)

The Comprehensive Air Quality Model with extensions (CAMx) is an integrated air quality model developed by the ENVIRON company based on the UAM-V mode. It synthesises the various technical features required by the “scientific” air-mass model into a single system that can be used to simulate air and particulate air pollutants on a variety of scales, such as cities and regions [31]. It is a state of the art photochemical grid model predicted on a “one-atmosphere” treatment approach to tropospheric air pollution (ozone, particulates, and air toxins) over spatial scales ranging from neighbourhoods to continents. In addition, it is an open-source system that is computationally efficient and flexible. The model’s FORTRAN source code is modular and well documented. The default input/output files are structured in a consistent FORTRAN binary format. Alternatively, output files may be optionally written in the Network Common Data Form (netCDF). Uncompressed netCDF output files are compatible with Models-3 I/O API software without the need to build CAMx with I/O API libraries. Meteorological fields are supplied to CAMx from separate weather prediction models. Emission inputs are supplied from external pre-processing systems [32].

The use of the CAMx model to simulate the source and spatial distribution of PM<sub>2.5</sub> has been widely studied. Wang used the CAMx model to analyze PM<sub>2.5</sub> sources and transboundary transport during a heavy pollution period in Shanghai [33]. Wu's model, and other CAMx models, were used to simulate the source and spatial distribution of PM<sub>2.5</sub> in Guangzhou's spring [34].

### 2.3.3. Weather Research and Forecasting Model with Chemistry (WRF-Chem)

The WRF-Chem model is the latest regional atmospheric dynamic chemical coupling model in the United States. Its biggest advantage is that the meteorological mode and chemical transmission mode are fully coupled in time and spatial resolution and achieve true online transmission. The model takes into account the processes of transport (including advection, diffusion, and convection), dry and wet sedimentation, gas phase chemistry, aerosol formation, radiation and photodecomposition rate, radiation produced by organisms, aerosol parameterization and photolysis [35]. The Weather Research and Forecasting (WRF) mode is a fully compressible non-hydrostatic model. It has different parameterization schemes for various physical processes, such as turbulence exchange, atmospheric radiation, cumulus cloud precipitation, cloud microphysics, and land surface type. It can provide the atmospheric flow field for the chemical model online. Details of the WRF mode physical process and integration scheme are shown in documents [36–38]. This model has been used to study the chemical reaction mechanism of urban compound pollution characteristics, PM<sub>2.5</sub>, O<sub>3</sub> and its precursor reactants (NO<sub>x</sub>, VOC, etc.) [39].

WRF-Chem has been widely used in the simulation of air pollution. Tie et al. used the WRF-Chem model to simulate the distribution and variation characteristics of ozone and its precursors in a Mexican urban area [40]. Molders et al. used mobile and fixed location measurements from the Fairbanks winter September 2008 field campaign for an operational evaluation of WRF/Chem's performance at high latitudes, which simulated PM<sub>2.5</sub> concentrations from 1 October 2008 to 1 April 2009 effectively [41].

### 2.4. Prediction of Spatial PM<sub>2.5</sub> Concentrations by Machine Learning Method

Machine learning provides a broad range of multivariate regression algorithms for empirically estimating PM<sub>2.5</sub> data when there is a set of useful observational data but no clear and complete theoretical description. The paper by Lary [42] describes the use of machine learning to estimate global daily PM<sub>2.5</sub> data from 1997 to 2014, using in situ hourly PM<sub>2.5</sub> observations from more than 8000 sites in 55 countries with approximately 100 parameters of comprehensive contextual data drawn from satellite data, meteorology, and demographics. Thus, it can be seen that machine learning can provide an extensive range of practical algorithms to facilitate the examination of the linear and Gaussian relationship between PM<sub>2.5</sub> abundance and meteorological variables [43]. There are plenty of machine learning methods; this paper lists some of the most frequently used approaches for predicting PM<sub>2.5</sub> concentrations on a spatial scale.

#### 2.4.1. BP Artificial Neural Network-Based Analysis

A BP artificial neural network is a neural network that interpolates high-dimensional space. Its structure is composed of three neuron layers, specifically the input layer, the hidden layer, and the output layer [7]. The three layers have a large number of simple neurons that are not connected to each other. The input layer neurons transfer the input data to the hidden layer, after which the data is activated by the hidden layer to the output layer and from the output layer to the output. The data, however, cannot be transmitted between neurons existing in the same layer. This is a forward transfer process; however, when the actual error exceeds the expected error, the error value propagates along the network to modify the connection weight and threshold between each neuron, and the training network repeats until the expected error is met, and thus the mapping relationship between the input and output is determined. The transfer function employed between the input and hidden layer is generally the S transform function, and the transfer function between the hidden and output layer is generally a pure linear transformation function.



The S function expression of the transfer function between the input and hidden layer of a BP artificial neural network is as follows:

$$f(x) = \frac{1}{1 + e^{-x}} \quad (8)$$

where the  $x$  represents the input of the neuron, and maps the input range of the neuron  $(-\infty, +\infty)$  to the interval  $(0, +1)$  to help the BP algorithm to train the neural network.

Chen uses surface meteorological observation data and the air pollution  $PM_{2.5}$  index of Wuhan City over the period from 1 November 2013 to 31 January 2014 to construct  $PM_{2.5}$  index forecasting models based on the use of a BP neural network [7].

#### 2.4.2. The Self-Organizing Map

It is often useful to apply an objective technique to classify large datasets into subclasses when using large datasets to characterize a problem. Self-organizing maps (SOMs) provide a method of performing such an unsupervised classification without any a priori assumption, a way to give the data “a voice” [44].

However, even with the assistance of unsupervised classification, high-dimensional data can still be challenging to visualize; this study dealt with a 36-dimensional space (the annual cycle was split into 36 ten-day windows). An SOM is a type of an artificial neural network for performing unsupervised classification. It is also a data visualization and unsupervised classification technique that can reduce the dimensions of high-dimensional data using self-organizing neural networks [45]. Similar to other forms of machine learning, an SOM operates in two phases, specifically (1) the training phase; and (2) the mapping phase. The map is built by training, using examples from the training dataset; the mapping determines the class for a new input vector [46]. An SOM consists of a 2-dimensional regular grid of components called nodes. Each node is associated with a weight vector of the same dimension and a position in the map. The procedure for converting a vector from the input data space to the map involves finding the node with a weight vector most similar to the input data space [47,48].

#### 2.4.3. Random Forest

Random forests were first introduced in 2001 by Leo Breiman [49]. They are a popular and efficient ensemble approach to statistical learning, useful for both classification and regression. A random forest is an ensemble of decision trees (hence the term forest). An ensemble approach allows more robust estimates, which are less prone to “over-learning.” The size of the “forest” ensemble is estimated by examining the estimated error as a function of the ensemble size [50]. In this study, the error rate plateaued at an ensemble size of approximately 30 trees. Therefore, an ensemble of 50 decision trees was used in our random forest. Random forests provide an objective way of highlighting the most important predictors, and ranking the relative importance of each predictor. To measure a variable’s importance, we first fit a random forest to the training dataset. This provides us with the so-called out-of-bag (OOB) error. If the importance of a specific variable  $X_i$  (where  $X_i$  is the  $i$ th predictor) is needed, the value  $X_i$  will be permuted, and the OOB error will be computed again for the permuted data. The importance of  $X_i$  will be the average of the difference between the OOB errors before and after all trees.

### 3. Discussion: Improvement and Integration of Traditional Methods with Their Comparisons

#### 3.1. Comparison between Traditional Methods

##### 3.1.1. Comparison between the Spatial Interpolation Methods for Estimating PM<sub>2.5</sub> Concentrations

The comparative accuracy of IDW, trend surface, and ordinary kriging interpolation is outlined in this section. Ping Zhang and Taotao Shen compared the accuracy of different spatial interpolation methods based on GIS and Spearman correlation [14]. The results show that the visualization of PM<sub>2.5</sub> concentrations' spatial distribution is best achieved using IDW interpolation, while that of TS is good, and that of ordinary kriging interpolation is the worst of the three methods, although still of standard quality. The IDW method is better than the other two methods, whether due to the error range of the forecast value, the accuracy of the predicting value relative to the observed value, or the sensitivity and reflection of the simulated value. The correlation coefficient of the simulated and observed values of ordinary kriging and trend surface are 0.62 and 0.67 respectively. That for IDW is 0.99 larger than the trend surface, indicating that the correlation coefficient of the trend surface and ordinary kriging methods have less precision than IDW.

Another study compares the accuracy of IDW, ordinary kriging interpolation, and cokriging interpolation methods, this time involving 100 data points in Zhaoxian of Shijiazhuang city which were used as sample points to conduct spatial interpolation [51]. Another 20 points in the study area were used to test interpolation accuracy. The results showed that the cokriging interpolation method had the highest precision; however, its visualization effect was not as good as the other two methods.

The lognormal ordinary kriging and regular ordinary kriging methods have also been compared. Liao's research used PM<sub>2.5</sub> data for the year 2000 with an aerodynamic diameter of  $\leq 10 \mu\text{m}$  for PM<sub>10</sub> and an aerodynamic diameter of  $\leq 2.5 \mu\text{m}$  for PM<sub>2.5</sub>, obtained from the U.S. Environmental Protection Agency [52]. Kriging estimations were performed at 94,135 geocoded addresses belonging to Women's Health Initiative study participants, using the ArcView geographic information system. The researchers developed a semiautomated program to enable large-scale daily kriging estimation and assessed the validity of semivariogram models using prediction error (PE), standardized prediction error (SPE), root mean square standardized (RMSS) error, and SE of the estimated PM<sub>2.5</sub>. National and regional scale kriging was performed satisfactorily, with the former emerging slightly better than the latter. The average PE, SPE, and RMSS of daily PM<sub>10</sub> semivariograms using regular ordinary kriging with a spherical model were 0.0629,  $-0.0011$ , and  $1.255 \mu\text{g}/\text{m}^3$  respectively; the average SE of the estimated residential-level PM<sub>10</sub> was  $27.36 \mu\text{g}/\text{m}^3$ . The values for PM<sub>2.5</sub> were 0.049, 0.0085, 1.389, and  $4.13 \mu\text{g}/\text{m}^3$  respectively. Lognormal ordinary kriging yielded a smaller average SE and effectively eliminated out-of-range predicted values compared to regular ordinary kriging. The study showed that semiautomated daily kriging estimations and semivariogram CVs are feasible on a national scale, and that lognormal ordinary kriging is more valid for estimating daily ambient PM<sub>2.5</sub> at geocoded residential addresses than regular ordinary kriging.

##### 3.1.2. Comparison between the CMAQ and CAMx Model for PM<sub>2.5</sub> Simulation

Shimadera et al. evaluated the year-long performance of the CMAQ model v5.0.1 and the comprehensive air quality model with extensions (CAMx) v6.00 [53]. They conducted year-long air quality simulations with common input meteorology, emission, and boundary concentration data in the Kinki region of Japan. They found that CAMx-simulated ground-level concentrations were generally higher by 10–20% than CMAQ-simulated values. Despite the systematic difference, comparisons with observed data proved that the overall year-long performances of the two models for simulating ground-level concentrations were similar. The two models approximately reproduced mass concentrations of PM<sub>2.5</sub>, but shared common difficulties in simulating PM<sub>2.5</sub> components.



### 3.1.3. The Hybrid Comparison

A comparison between geostatistical interpolation and remote sensing technique has been presented for the designated research area. Seung-Jae found that for most of the populated areas of the continental United States, geostatistical interpolation produced more accurate estimates than remote sensing [6]. The differences between the estimates resulting from the two methods, however, were relatively small. In areas with extensive monitoring networks, interpolation may provide more accurate estimates, but in many areas of the world without such monitoring, remote sensing can provide useful exposure estimates that perform equally well.

## 3.2. Improvement and Integration of the Current Methods

### 3.2.1. Improvement and Integration of Spatial Interpolation Methods

An IDW-based interpolation method of estimating  $PM_{2.5}$  concentrations has been developed.  $PM_{2.5}$  data interpolation is conducted in the continuous space-time domain by integrating space and time simultaneously, using the so-called extension approach [54]. Time values are calculated with the help of a factor under the assumption that spatial and temporal dimensions are equally important when interpolating a continuous changing phenomenon in the space-time domain. Various IDW-based spatiotemporal interpolation methods with different parameter configurations are evaluated by CV. Additionally, this study explores computational issues (computer processing speed) faced during the implementation of spatiotemporal interpolation for huge datasets. Parallel programming techniques and a k-d tree advanced data structure are adapted in this paper to address the computational challenges, with significant computational improvement achieved. Finally, a web-based spatiotemporal IDW-based interpolation application is designed and implemented, wherein users can visualize and animate spatiotemporal interpolation results.

The consideration of wind-field path distance can be also added into traditional methods to improve the accuracy. Li developed an interpolation method based on the shortest path distance to characterize the impact of complex urban wind-fields on the distribution of the particulate matter concentration [54]. In this method, the wind-field is incorporated by first interpolating the observed wind-field from a meteorological-station network, and then using this continuous wind-field to construct a cost surface based on a Gaussian dispersion model, calculating the shortest wind-field path distances between locations, and finally replacing the Euclidean distances typically used in IDW with the shortest wind-field path distances. This proposed methodology is used to generate daily and hourly estimation surfaces for the particulate matter concentration in the urban area of Beijing in May 2013, and results demonstrate that wind-fields can be incorporated into an interpolation framework using the shortest wind-field path distance. This leads to a remarkable improvement in both the prediction accuracy and the visual reproduction of the wind-flow effect, both of which are of great importance for the assessment of the effects of pollutants on human health.

### 3.2.2. Improvement and Integration of Remote Sensing Methods

There are several studies on achieving improvements to remote sensing methods of estimating  $PM_{2.5}$  concentrations covering the entire research region. Lary et al. used a suite of remote sensing and meteorological data products together with ground-based observations of  $PM_{2.5}$  from 8329 measurement sites in 55 countries taken between the years 1997 and 2014 to train a machine learning algorithm to estimate the daily distributions of  $PM_{2.5}$  from 1997 to the present [55]. Using ground-based observations of particulate matter together with a suite of remote sensing and meteorological data products to train a machine learning algorithm to estimate the daily distributions of  $PM_{2.5}$  demonstrates a new approach. The new  $PM_{2.5}$  daily global data product reproduces global observations and spans an unprecedented 16 years from 1997 to the present. The correlation coefficient for each of the five training datasets is 0.96 or greater, and the correlation coefficient for each of the

independent validation datasets is 0.52 or greater. This implies that the  $PM_{2.5}$  abundances inferred using machine learning agree well with actual field conditions determined from in situ observations.

The comparison between the AOD measured from the ground-based Aerosol Robotic Network (AERONET) system and the satellite MODIS instruments at 550 nm shows that a bias exists between the two data products. A comprehensive search was performed to explore possible factors which may be contributing to this [56]. The analysis used several measured variables, including the MODIS AOD as an input, in order to train a neural network in regression mode to predict the AERONET AOD values. This not only allowed us to obtain an estimate, but also allowed us to infer the optimal sets of variables that significantly influenced the prediction results. In addition, we applied machine learning to infer the global abundance of ground-level  $PM_{2.5}$  from the AOD data and other ancillary satellite and meteorology products. This research is part of our goal to provide air quality information, which can also be useful for global epidemiology studies.

### 3.2.3. The Improvement and Integration of air Quality Model Methods

The general approach to air quality modelling is the CMAQ model combined with the dynamic linear regression method. The Models-3/CMAQ air quality modelling system was applied to forecast  $PM_{2.5}$  concentrations in Shanghai [57]. Observation data from ten monitoring sites were chosen to evaluate the model performance. The results indicate that the CMAQ model can simulate the variation of  $PM_{2.5}$  concentrations satisfactorily. However, the simulated  $PM_{2.5}$  concentrations are underestimated by 25% under normal conditions. During a high pollution episode, the underestimation, which could be caused by the uncertainty of the emission inventory, can rise as high as 32%. The dynamic linear regression method is used in order to improve the accuracy of the  $PM_{2.5}$  forecast. The statistical results show that after a revised forecast the accuracy improves from 76.4% to 79.3%, and the crisis success index improves from 56.4% to 72.1%, proving the value of this method.

Hourly  $PM_{2.5}$  concentrations at 252 environmental monitoring stations in China during the period January–December 2014 forecast by the real-time running fifth-generation Penn State/NCAR mesoscale model (MM5)-CMAQ model system are corrected using the dynamical-statistical method based on the CMAQ model, and by adapting the partial least square regression technique [58]. Temporal and spatial variations of  $PM_{2.5}$  concentrations before and after correction are analyzed with a focus on the applicability of the dynamical-statistical method in different areas and seasons in China. Furthermore, the method presented in this study is applicable to the correction of  $PM_{2.5}$  forecasts for both polluted and clean days in China. The correction is more effective during polluted processes in the Beijing-Tianjin-Hebei region; additionally, correction effects are better during clean processes than on polluted days in the other three regions. Results of this study will provide a scientific basis and corresponding new technique for improving air quality forecasting, and for early warning and prevention of heavy haze weather.

In addition, some techniques such as data assimilation (DA) and model output statistics can effectively decrease uncertainties resulting from the uncertainties of atmospheric chemical models. To improve the initialization of  $PM_{2.5}$  in CMAQ, Kumar et al. developed a new capability in the community Gridpoint Statistical Interpolation (GSI) system to assimilate MODIS AOD retrievals in CMAQ [59]. Specifically, they developed new capabilities within GSI to read/write CMAQ data, a forward operator that calculates AOD at 550 nm from CMAQ's aerosol chemical composition, and an adjoint of the forward operator that translates the changes in AOD to aerosol chemical composition. A generalized background error covariance program (GEN\_BE) has been extended to calculate background error covariance using CMAQ output. The background error variances are generated using a combination of both emissions and meteorological perturbations to better capture sources of uncertainties in  $PM_{2.5}$  simulations. They used the CMAQ-GSI system to perform daily 24 hourly  $PM_{2.5}$  forecasts with and without DA from 15 July to 14 August 2014 and compared the resulting forecasts against AirNOW  $PM_{2.5}$  measurements at 550 stations across the U.S. Results

indicated that the assimilation of MODIS AOD retrievals improved the initialization of the CMAQ model in terms of improved correlation coefficient and reduced bias.

### 3.2.4. Improvement and Integration of Machine Learning Methods

Combining LUR models with machine learning can be a powerful approach. Beckerman created a model to predict ambient particulate matter less than 2.5 microns in aerodynamic diameter ( $PM_{2.5}$ ) across the contiguous United States, to be applied in health-effects modelling [60]. The authors developed a hybrid approach combining a selected LUR model with a machine learning method, and Bayesian Maximum Entropy (BME) interpolation of the LUR space-time residuals. The  $PM_{2.5}$  dataset included 104,172 monthly observations at 1464 monitoring locations with approximately 10% of locations reserved for CV. LUR models were based on remote sensing estimates of  $PM_{2.5}$ , land use, and traffic indicators. Normalized cross-validated  $R^2$  values for LUR were 0.63 and 0.11 with and without remote sensing respectively, suggesting that remote sensing is a strong predictor of ground-level concentrations. In models including the BME interpolation of the residuals, cross-validated  $R^2$  was 0.79 for both configurations; the model without remotely sensed data described a more fine-scale variation than the model including remote sensing. The results suggest that the modelling framework can predict ground-level concentrations of  $PM_{2.5}$  at multiple scales over contiguous areas of the U.S.

## 4. Conclusions

The  $PM_{2.5}$  concentration of a certain point can be measured by appropriate facilities with accurate data; however, it is hard to obtain values representing conditions over a full research region in which setting up the necessary equipment is not feasible. Therefore, an increasing number of methods for estimating  $PM_{2.5}$  concentrations in areas without special sensors were presented by the researchers. The main approaches are divided into four categories namely, (a) spatial interpolation methods; (b) remote sensing techniques; (c) air quality model methods; and (d) machine learning methods. However, they usually overlap in terms of processing and source data. For example, machine learning can also be a part of the remote sensing process, the purpose of which is to find the relationship between  $PM_{2.5}$  concentrations and the AOD value of each grid. However, as machine learning has emerged as the most popular technique during recent years, it is worth special mention. This is because there is a trend in considering most relationship issues and a relative exact approach to predict the unknown value in time as well as spatial scale. In this study, we reviewed relevant recently published papers, as well as classical papers in the field of methods and application of  $PM_{2.5}$  concentration estimation. The main findings include:

- (1) There are various levels of accuracy that are determined by a large number of factors including circumstances specific to the research area, resolution of the source data, parameters chosen by specific models, and the details in the process, used by different methods to estimate the  $PM_{2.5}$  concentrations.
- (2) The most convenient and time-efficient method is spatial interpolation. However, the accuracy of this method is relatively low in comparison with other approaches. The more complex of the above methods always relate to machine learning, which can predict the unknown data of  $PM_{2.5}$  on spatiotemporal scales. There are many improvements and integrations of different methods which can provide more accurate results under specific circumstances.
- (3) The traditional approaches for estimation of  $PM_{2.5}$  concentrations are outdated because their results are not as accurate as the results obtained using newer techniques. Furthermore, there are an increasing number of integration models and methods that can be applied to various conditional applications, rendering the use of only traditional approaches obsolete.
- (4) Presently, because of the rise in computing power and urban computing, rapid development in machine learning has become a major research focus area for estimating and predicting environmental

problems through continuous acquisition, integration, and analysis of a variety of heterogeneous and large data in cities. Combining physical models with machine learning holds obvious potential for estimating spatial-temporal dynamic distribution of urban PM<sub>2.5</sub> concentrations.

The estimation of PM<sub>2.5</sub> concentrations can be simulated by a computer, with advances in information technology and research contributions from both individuals and organizations. However, there are not many general models or methods for estimating data with high accuracy, mainly because there is less source data and little approach to take all factors in different situations into account. Therefore, it is necessary to apply different approaches with improvements and integrations, selecting these as is appropriate for various different scenarios. Contemporary growth in machine learning and artificial intelligence technology is likely to result in more accurate methods of PM<sub>2.5</sub> estimation in spatial scale that will be adaptable to a variety of research types.

**Author Contributions:** G.Z. is considered to be someone who has made substantive intellectual contributions to this review. He reviewed and sorted through most of the references and drew the main conclusions. X.R. summarized the framework of the article and revised critically for important intellectual content. And Y.F. analyzed and improved formulas, checked the English editing as well as the format of references of the final version to be published.

**Funding:** This research was funded by the National Key Research and Development Program (Grant No. 2017YFB0503600), the National Natural Science Foundation of China (Grant No. 41771478) and the Beijing Natural Science Foundation (Grant No. 8172046).

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. Tang, M.; Wu, X.; Agrawal, P.; Pongpaichet, S.; Jain, R. Integration of diverse data sources for spatial PM<sub>2.5</sub> data interpolation. *IEEE Trans. Multimed.* **2017**, *19*, 408–417. [[CrossRef](#)]
2. Hanzalova, K.; Rossner, P., Jr.; Sram, R.J. Oxidative damage induced by carcinogenic polycyclic aromatic hydrocarbons and organic extracts from urban air particulate matter. *Mutat. Res. Genet. Toxicol. Environ. Mutagen.* **2010**, *696*, 114–121. [[CrossRef](#)] [[PubMed](#)]
3. Valavanidis, A.; Fiotakis, K.; Vlachogianni, T. Airborne particulate matter and human health: Toxicological assessment and importance of size and composition of particles for oxidative damage and carcinogenic mechanisms. *J. Environ. Sci. Health Part C Environ. Carcinog. Ecotoxicol. Rev.* **2008**, *26*, 339–362. [[CrossRef](#)] [[PubMed](#)]
4. World Health Organization (OMS). *Who Air Quality Guidelines Global Update 2005*; OMS: Copenhagen, Denmark, 2005.
5. Zheng, M.; Zhan, Y.; Yan, C.; Zhu, X.; Schauer, J.; Zhang, Y. Review of PM<sub>2.5</sub> Source Apportionment Methods in China. *Acta Scientiarum Naturalium Universitatis Pekinensis* **2014**, *50*, 1141–1154.
6. Lee, S.-J.; Serre, M.L.; van Donkelaar, A.; Martin, R.V.; Burnett, R.T.; Jerrett, M. Comparison of geostatistical interpolation and remote sensing techniques for estimating long-term exposure to ambient PM<sub>2.5</sub> concentrations across the continental United States. *Environ. Health Perspect.* **2012**, *120*, 1727–1732. [[CrossRef](#)] [[PubMed](#)]
7. Chen, Y.; Qin, H.; Zhou, Z.G. A comparative study on multi-regression analysis and bp neural network of PM<sub>2.5</sub> index. In Proceedings of the 10th International Conference on Natural Computation, Xiamen, China, 19–21 August 2014; pp. 155–159.
8. Hueglin, C.; Gehrig, R.; Baltensperger, U.; Gysel, M.; Monn, C.; Vonmont, H. Chemical characterisation of PM<sub>2.5</sub>, PM<sub>10</sub> and coarse particles at urban, near-city and rural sites in Switzerland. *Atmos. Environ.* **2005**, *39*, 637–651. [[CrossRef](#)]
9. Van Donkelaar, A.; Martin, R.V.; Spurr, R.J.; Burnett, R.T. High-resolution satellite-derived PM<sub>2.5</sub> from optimal estimation and geographically weighted regression over North America. *Environ. Sci. Technol.* **2015**, *49*, 10482–10491. [[CrossRef](#)] [[PubMed](#)]
10. Cordero, L.; Malakar, N.; Wu, Y.H.; Gross, B.; Moshary, F. Assessing surface PM<sub>2.5</sub> estimates using data fusion of active and passive remote sensing methods. *Br. J. Environ. Clim. Chang.* **2013**, *3*, 547–565. [[CrossRef](#)] [[PubMed](#)]

11. Garcia, C.A.; Yap, P.S.; Park, H.Y.; Weller, B.L. Association of long-term PM<sub>2.5</sub> exposure with mortality using different air pollution exposure models: Impacts in rural and urban California. *Int. J. Environ. Health Res.* **2016**, *26*, 145–157. [CrossRef] [PubMed]
12. Zou, B.; Luo, Y.; Wan, N.; Zheng, Z.; Sternberg, T.; Liao, Y. Performance comparison of lur and ok in PM<sub>2.5</sub> concentration mapping: A multidimensional perspective. *Sci. Rep.* **2015**, *5*, 8698. [CrossRef] [PubMed]
13. Keler, A.; Krisp, J.M. Spatiotemporal visualization of interpolated particulate matter (PM<sub>2.5</sub>) in Beijing. *GI\_Forum J. Geogr. Inf. Sci.* **2015**, *1*, 464–474.
14. Zhang, P.; Shen, T. Comparison of different spatial interpolation methods for atmospheric pollutant PM<sub>2.5</sub> by using GIS and Spearman correlation. *J. Chem. Pharm. Res.* **2015**, *7*, 452–469.
15. Ji, Q.; Yu, M. Parameters selection of the annual mean temperature spatial interpolation method based on collaborative kriging interpolation. *J. Cap. Norm. Univ.* **2010**, *31*, 81–87.
16. Deng, L. Estimation of PM<sub>2.5</sub> spatial distribution based on kriging interpolation. In Proceedings of the First International Conference on Information Sciences, Machinery, Materials and Energy, Chongqing, China, 11–13 April 2015.
17. Donkelaar, A.V.; Martin, R.V.; Levy, R.C.; Silva, A.M.D.; Krzyzanowski, M.; Chubarova, N.E.; Semutnikova, E.; Cohen, A.J. Satellite-based estimates of ground-level fine particulate matter during extreme events: A case study of the Moscow fires in 2010. *Atmos. Environ.* **2011**, *45*, 6225–6232. [CrossRef]
18. Liu, Y. New Directions: Satellite driven PM<sub>2.5</sub> exposure models to support targeted particle pollution health effects research. *Atmos. Environ.* **2013**, *68*, 52–53. [CrossRef]
19. Ma, Z.; Hu, X.; Huang, L.; Bi, J.; Liu, Y. Estimating ground-level PM<sub>2.5</sub> in China using satellite remote sensing. *Environ. Sci. Technol.* **2014**, *48*, 7436–7444. [CrossRef] [PubMed]
20. Lee, H.J.; Liu, Y.; Coull, B.A.; Schwartz, J.; Koutrakis, P. A novel calibration approach of modis aod data to predict PM<sub>2.5</sub> concentrations. *Atmos. Chem. Phys.* **2011**, *11*, 9769–9795. [CrossRef]
21. Li, Z.; Zhang, Y.; Shao, J.; Li, B.; Hong, J.; Liu, D.; Li, D.; Wei, P.; Li, W.; Li, L. Remote sensing of atmospheric particulate mass of dry PM<sub>2.5</sub> near the ground: Method validation using ground-based measurements. *Remote Sens. Environ.* **2016**, *173*, 59–68. [CrossRef]
22. Lin, C.; Li, Y.; Yuan, Z.; Lau, A.K.H.; Li, C.; Fung, J.C.H. Using satellite remote sensing data to estimate the high-resolution distribution of ground-level PM<sub>2.5</sub>. *Remote Sens. Environ.* **2015**, *156*, 117–128. [CrossRef]
23. Zhang, Y.; Li, Z. Remote sensing of atmospheric fine particulate matter (PM<sub>2.5</sub>) mass concentration near the ground from satellite observation. *Remote Sens. Environ.* **2015**, *160*, 252–262. [CrossRef]
24. Ma, Z.; Hu, X.; Sayer, A.M.; Levy, R.; Zhang, Q.; Xue, Y.; Tong, S.; Bi, J.; Huang, L.; Liu, Y. Satellite-based spatiotemporal trends in PM<sub>2.5</sub> concentrations: China, 2004–2013. *Environ. Health Perspect.* **2016**, *124*, 184–192. [CrossRef] [PubMed]
25. Chang, H.H.; Hu, X.; Liu, Y. Calibrating modis aerosol optical depth for predicting daily PM<sub>2.5</sub> concentrations via statistical downscaling. *J. Expo. Sci. Environ. Epidemiol.* **2014**, *24*, 398–404. [CrossRef] [PubMed]
26. Lv, B.; Hu, Y.; Chang, H.H.; Russell, A.G.; Cai, J.; Xu, B.; Bai, Y. Daily estimation of ground-level PM<sub>2.5</sub> concentrations at 4km resolution over Beijing-Tianjin-Hebei by fusing modis aod and ground observations. *Sci. Total Environ.* **2017**, *580*, 235–244. [CrossRef] [PubMed]
27. National Exposure Research Laboratory (NERL), Ecological Exposure Research Division, United States Environmental Protection Agency (US EPA). CMAS: Community Modeling and Analysis System. Available online: <http://www.cmascenter.org> (accessed on 1 June 2018).
28. Wang, Z.; Li, X.; Wang, Z.; Wu, X.; Che, F.; Nie, P. Application Status of Models-3/CMAQ in Environmental Management. *Environ. Sci. Technol.* **2013**, *36*, 386–391.
29. Byun, D.; Schere, K.L. Review of the governing equations, computational algorithms, and other components of the models-3 community multiscale air quality (CMAQ) modeling system. *Appl. Mech. Rev.* **2006**, *59*, 51–77. [CrossRef]
30. Nie, T.; Li, X.; Wang, Z.; Qi, J.; Zhou, Z. Spatial and temporal distribution and process analysis of PM<sub>2.5</sub> pollution over Beijing during APEC. *China Environ. Sci.* **2016**, *36*, 349–355.
31. Wagstrom, K.M.; Pandis, S.N. Contribution of long range transport to local fine particulate matter concerns. *Atmos. Environ.* **2011**, *45*, 2730–2735. [CrossRef]
32. ENVIRON. User's Guide Comprehensive Air Quality Model with Extensions Version 6.5. Available online: [http://www.camx.com/files/camxusersguide\\_v6-50.pdf](http://www.camx.com/files/camxusersguide_v6-50.pdf) (accessed on 1 June 2018).



33. Wang, Y.; Li, L.; Chen, C.; Huang, H.; Feng, J.; Wang, S.; Wang, H.; Zhang, G.; Zhou, M.; Cheng, P.; et al. Source apportionment of fine particulate matter during autumn haze episodes in Shanghai, China. *J. Geophys. Res. Atmos.* **2014**, *119*, 1903–1914. [[CrossRef](#)]
34. Wu, D.; Fung, J.C.H.; Yao, T.; Lau, A.K.H. A study of control policy in the pearl river delta region by using the particulate matter source apportionment method. *Atmos. Environ.* **2013**, *76*, 147–161. [[CrossRef](#)]
35. Grell, G.A.; Schmitz, P.R.; Mckeen, S.A.; Frost, G.; Skamarock, W.C.; Eder, B. Fully coupled ‘online’ chemistry within the wrf model. *Atmos. Environ.* **2005**, *39*, 6957–6975. [[CrossRef](#)]
36. Janjic, Z.I. The NCEP WRF Core. In Proceedings of the 84th AMS Annual Meeting, Seattle, WA, USA, 10–15 January 2004.
37. McCaslin, P.T.; Smart, J.R.; Shaw, B.; Jamison, B.D. A graphical user interface to prepare the standard initialization for WRF. In Proceedings of the 84th AMS Annual Meeting, Seattle, WA, USA, 10–15 January 2004.
38. Klemp, J.B. Weather Research and Forecasting Model: A technical Overview. In Proceedings of the 84th AMS Annual Meeting, Seattle, WA, USA, 10–15 January 2004.
39. Geng, F.; Zhao, C.; Tang, X.; Lu, G.; Tie, X. Analysis of ozone and vocs measured in Shanghai: A case study. *Atmos. Environ.* **2007**, *41*, 989–1001. [[CrossRef](#)]
40. Tie, X.; Madronich, S.; Li, G.; Ying, Z.; Weinheimer, A.; Apel, E.; Campos, T. Simulation of mexico city plumes during the MIRAGE-Mex field campaign using the WRF-Chem model. *Atmos. Chem. Phys.* **2009**, *9*, 4621–4638. [[CrossRef](#)]
41. Mölders, N.; Tran, H.N.Q.; Cahill, C.F.; Leelasakultum, K.; Tran, T.T. Assessment of wrf/chem pm 2.5 forecasts using mobile and fixed location data from the fairbanks, alaska winter 2008/09 field campaign. *Atmos. Pollut. Res.* **2012**, *3*, 180–191. [[CrossRef](#)]
42. Lary, D.J.; Faruque, F.S.; Malakar, N.; Moore, A.; Roscoe, B.; Adams, Z.L.; Eggelston, Y. Estimating the global abundance of ground level presence of particulate matter (PM<sub>2.5</sub>). *Geospat. Health* **2014**, *8*, S611–S630. [[CrossRef](#)] [[PubMed](#)]
43. Wu, D.; Zewdie, G.K.; Liu, X.; Kneen, M.A.; Lary, D.J. Insights into the morphology of the East Asia PM<sub>2.5</sub> annual cycle provided by machine learning. *Environ. Health Insights* **2017**, *11*. [[CrossRef](#)] [[PubMed](#)]
44. Hastie, T.; Friedman, J.; Tibshirani, R. The elements of statistical learning. *J. R. Stat. Soc.* **2010**, *45*, 267–268.
45. Kohonen, T. Self-organized formation of topologically correct feature maps. *Biol. Cybern.* **1982**, *43*, 59–69. [[CrossRef](#)]
46. Von der Malsburg, C. Self-organization of orientation sensitive cells in the striate cortex. *Kybernetik* **1974**, *14*, 85–100. [[CrossRef](#)]
47. Kakuda, N.; Miwa, T.; Nagaoka, M.; Kohonen, T. The self-organizing map. *Neurocomputing* **1998**, *21*, 1–6.
48. Vesanto, J.; Alhoniemi, E. Clustering of the self-organizing map. *IEEE Trans. Neural Netw.* **2002**, *11*, 586–600. [[CrossRef](#)] [[PubMed](#)]
49. Breiman, L. Random forests. *Mach. Learn.* **2001**, *45*, 5–32. [[CrossRef](#)]
50. Pal, M. Random forest classifier for remote sensing classification. *Int. J. Remote Sens.* **2005**, *26*, 217–222. [[CrossRef](#)]
51. Li, J.; Zhai, L.; Sang, H.; Zhang, Y.; Yuan, J. Comparison of different spatial interpolation methods for PM<sub>2.5</sub>. *Sci. Surv. Mapp.* **2016**, *41*, 50–54.
52. Liao, D.; Peuquet, D.J.; Duan, Y.; Whitsel, E.A.; Dou, J.; Smith, R.L.; Lin, H.M.; Chen, J.C.; Heiss, G. Gis approaches for the estimation of residential-level ambient pm concentrations. *Environ. Health Perspect.* **2006**, *114*, 1374–1380. [[CrossRef](#)] [[PubMed](#)]
53. Shimadera, H.; Kojima, T.; Kondo, A.; Inoue, Y. Performance comparison of cmaq and camx for one-year PM<sub>2.5</sub> simulation in Japan. *Int. J. Environ. Pollut.* **2015**, *57*, 146–161. [[CrossRef](#)]
54. Li, L.; Gong, J.; Zhou, J. Spatial interpolation of fine particulate matter concentrations using the shortest wind-field path distance. *PLoS ONE* **2014**, *9*, e96111. [[CrossRef](#)] [[PubMed](#)]
55. Lary, D.J.; Lary, T.; Sattler, B. Using machine learning to estimate global PM<sub>2.5</sub> for environmental health studies. *Environ. Health Insights* **2015**, *9*, 41–52. [[CrossRef](#)] [[PubMed](#)]
56. Malakar, N.K.; Lary, D.J.; Moore, A.; Gencaga, D. Estimation and bias correction of aerosol abundance using data-driven machine learning and remote sensing. In Proceedings of the 2012 Conference on Intelligent Data Understanding, Boulder, CO, USA, 24–26 October 2012; pp. 24–30.
57. Wang, Q.; Wu, J.; Lin, Y. Implementation of a dynamic linear regression method on the CMAQ forecast of PM<sub>2.5</sub> in Shanghai. *Acta Sci. Circumst.* **2015**, *35*, 1651–1656.



58. Cheng, X.; Diao, Z.; Hu, J.; Xu, X.; Zhang, J.; Li, D. Dynamical-statistical forecasting of PM<sub>2.5</sub> concentration based on CMAQ model and adapting partial least square regression method in China. *Acta Sci. Circumst.* **2016**, *36*, 2771–2782.
59. Kumar, R.; Delle Monache, L.; Alessandrini, S.; Saide, P.; Lin, H.C.; Liu, Z.; Pfister, G.; Edwards, D.P.; Baker, B.; Tang, Y.; et al. Improving short-term air quality predictions over the U.S. Using chemical data assimilation. In Proceedings of the AGU Fall Meeting, New Orleans, LA, USA, 11–13 December 2017.
60. Beckerman, B.S.; Jerrett, M.; Serre, M.; Martin, R.V.; Lee, S.J.; Van, D.A.; Ross, Z.; Su, J.; Burnett, R.T. A hybrid approach to estimating national scale spatiotemporal variability of PM<sub>2.5</sub> in the contiguous United States. *Environ. Sci. Technol.* **2013**, *47*, 7233–7241. [[CrossRef](#)] [[PubMed](#)]



© 2018 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).