

# Using Big Data to Battle Air Pollution

Karthik Vegi

Indiana University Bloomington  
2619 East 2nd Street, Apt 11  
Bloomington, IN 47401, USA  
kvegi@iu.com

## ABSTRACT

We have come a long way from the stone age to build large scale industries, big cities, bullet trains, and a booming automobile industry. Technological and industrial advances are making our cities smarter by the day and yet a nagging side-effect is air pollution. Air pollution is not only creating local health hazards like respiratory and heart problems, but also directly leading to an increase in temperatures and contributing to global warming. We show how the advances in *Big Data*, *Cloud Computing*, and *Internet of Devices* can be used to combat air pollution.

## KEYWORDS

i523, hid231, big data, environment, air pollution, global warming

## 1 INTRODUCTION

Air pollution is no longer a local problem. It is a global environmental issue which involves individual countries to come together and device measures to combat it [5]. It is causing about 3.7 million premature deaths worldwide from cardiovascular and respiratory diseases and also ruins the crops that feed the world [5]. Air pollution also has a direct effect on a number of environmental issues like global warming, depletion of ozone layer, acid rains, and impacts wild-life [5].

Back in the year 1990, the job of a typical air quality scientist was to develop atmospheric dispersion models to evaluate the air pollution caused by industries and make sure that it is within the permissible level suggested by the *Environmental Protection Agency* [2]. These models gather historic data of many years from airports and weather balloons to predict the pollution with the help of meteorology theory [2]. Although the methods used to derive the values were good enough, the limitations with respect to the technology posed a real challenge which took weeks to run the simulations, only to be cut-off in the middle due to power and storage issues [2]. The data processing engine was built on Sun-Solaris workstations with tapes handling the data storage [2]. The work-stations set up in major points in the country would communicate using a very slow network connection [2]. The data processing would be done locally and later written to all the servers which would then be split and distributed among many machines and consolidated in the end [2]. "If only we had that much more data and that much more ability to handle it, we could iterate through the model at a much finer scale. Real-time data processing remained a pipe dream" [2].

## 2 AIR POLLUTION AS A BIG DATA PROBLEM

The advent of *Big Data* and the technological advances changed the way the data is ingested and analyzed [2]. The network speeds have increased, wide range of sensors are available to collect data with a lot of precision which would feed the high speed data processing systems. Batch processing has become easier with *Hadoop* and *Map-Reduce*. The storage mechanisms have become cheaper and more disaster proof.

IBM is helping Beijing combat air pollution by analyzing huge amounts of data using a data analysis platform *Green Horizons* [3]. IBM has signed up partnerships with different cities in China and India to deploy *Big Data Analytics*, *Machine Learning* and *Internet of Things* to improve traffic, keep a check on the pollution from industrial machines, and other pollution causing agents [3]. IBM will deploy sensors in various places to collect data in real-time and analyze previous weather forecasts, and build improved iterative models over time [3]. The system continuously streams data from the sensors and improves the forecast by learning over time using *Machine Learning* algorithms [3]. Figure 1 shows the Green Horizons air quality management for Beijing.

[Figure 1 about here.]

IBM is collaborating with the United Nations to push the use of technological advances by every country for the common good of the world [3]. More and more cities and countries are opening air quality data to public where you can get reports in real time [1]. The *BreezoMeter* is the first mobile application that provides real-time information of the street's air quality information using geo-location maps [1]. *Copernicus* is another monitoring service that ingests data from satellites and on site sensors on land, air and sea to provide continuous information to the users [1]. *Open Data Week* is an intergovernmental organization where 34 states come together to bring reforms and discuss how to use technology and services like *Copernicus* that use *Big Data* to test prototypes of new products to ensure they operate within the permissible levels of pollution [1].

While these initiatives help bring awareness about the seriousness of the issue, each state and country should take strict measures to bring out reforms that will help eradicate pollution. *Big Data* might never replace the environmental responsibility but it will help to plan the vision for environmental awareness and its tools make it easier to achieve the vision [1]. These tools can also be used to gauge the alternative sources of energy and the feasibility of tapping into other natural resources ensuring responsible consumption of energy [1]. For example, IBM *Bluemix* analyzed data from a steel industry and the analysis uncovered an interesting insight that the furnace wastes a lot of energy to offset the temperature of the smoke which resulted in optimizing its operation [2].

### 3 BIG DATA TECHNIQUES TO COMBAT POLLUTION

#### 3.1 Random Forest Approach for predicting air quality in Urban Sensing Systems

Air pollution in an urban setting is very important to monitor because of the population density. Air quality in these areas varies a lot in various parts of the city owing to traffic and presence of industries [6]. A random forest approach ingests data from meteorology, urban sensors, road information, and real-time traffic and predicts the air quality with utmost precision [6]. Real-time air quality information consists of measuring the concentration of  $PM_{2.5}$ ,  $PM_{10}$ , and  $NO_2$  [6].

The *Air Quality Index* AQI is the measure that is used to understand how polluted the air is [6]. AQI is measured by reading the concentration of 6 pollutant gases namely, sulfur dioxide  $SO_2$ , nitrogen dioxide  $NO_2$ , air particles smaller than  $10\ \mu m$   $PM_{10}$ , air particles smaller than  $2.5\ \mu m$   $PM_{2.5}$ , carbon monoxide  $CO$ , and ozone  $O_3$  [6]. Based on the level of AQI, the air quality is classified as shown in Figure 2.

[Figure 2 about here.]

*Traffic Congestion Status* TCS, explains the traffic status at the current hour [6]. Figure 3 shows how colors are used to represent the traffic congestion [6].

[Figure 3 about here.]

#### 3.2 RAQ Algorithm

The RAQ algorithm collects data from air monitoring station AQI, meteorology data MD, traffic congestion TCS, road information RI, and point of interest POI which is the specific location that someone is interested to visit [6]. The data refresh rate is one hour and the data is collected from different parts of the city which are divided in grids from  $G_1$  to  $G_n$  [6]. The data is divided into training and testing data sets to train the model and evaluate the model [6]. Figure 4 shows the structure of the data.

[Figure 4 about here.]

A decision tree is used to split and classify the data and the results are aggregated by collecting the data from all the sub-trees [6]. Figure 5 illustrates the procedure of RAQ.

[Figure 5 about here.]

The *Random Forest* algorithm is employed using the tree type classifier to recursively partition the dataset and generate sub-trees and finally aggregate the results of each sub-tree [6]. Each sub-tree is constructed using *Bootstrap Aggregating* where each data set is divided into different buckets by using statistical samples [6]. Once the trees are constructed, each subset of data is fed into a decision tree and the estimated AQI index is calculated [6]. The final AQI index is determined as the maximum value out of all the individual values [6]. Figure 6 shows the step-by-step RAQ algorithm.

[Figure 6 about here.]

### 4 MACHINE LEARNING MODELS

*Machine Learning* deals with augmenting computers with the ability to learn from data and program themselves [4]. These algorithms can be used to evaluate the air quality [4].

#### 4.1 Artificial Neural Network Model

Artificial Neural Network Model tries to solve the problem by simulating the functioning of brain and neurons [4]. The model architecture is a function of a sigmoid [4]. For this experiment, the air quality data was divided into training, test, and validation data with split of 60, 20, and 20 with a back propagation network of two hidden layers [4]. To ensure consistency, the air quality data for the training and test sets are derived from the same season [4]. The air quality is forecast by looking at the historic data where the input and output are represented by the air quality data measured at different times [4]. The model turns out to be reliable with a good prediction accuracy with the lowest mean square error of  $3.7 \times 10^{-4}$  [4]. The Artificial Neural Network Model is combined with *Markov Chains* to develop a new improved model with improved prediction accuracy where the ANN computes the primary values and the results are re-computed and improved by the markov transitional probability matrices [4]. Figure 7 shows the Artificial Neural Network Model with two hidden layers.

[Figure 7 about here.]

#### 4.2 Least squares Support Vector Machine Model

Least squares support vector machine is a supervised learning model used for classification and regression analysis which arrives at the solution by solving the data represented in the form of linear equations [4]. For this model, the sample data was collected from 100 sensor points in different intervals of time and at different geographical locations that ranged from urban areas with population, areas near the airport, water surface areas like lakes, and sewage processing areas [4]. The sample data was a good split with 80 percent collected from urban sewage area and the other data collected from air surface areas [4]. The fluorescence content in the air was analyzed by a portable air quality measuring device developed in-house by Zhejiang University [4]. The fluorescence data captured using the device is highly dimensional and non-linear and therefore data pre-processing is essential to bring the dimensions down to a manageable level [4]. This eliminates the ambient noise and the temperature drift from the data [4]. The algorithm predicts the regression model by looking at the training data for each cluster [4]. Finally, the vector cosine distance is used to classify the sample into clusters and the performance criterion such as *Root Mean Square Error* and *Mean Absolute Error* are computed which demonstrate the efficiency of the algorithm [4]. Figure 8 shows the pictorial representation of the algorithm.

[Figure 8 about here.]

### 5 CONCLUSION

While the new age technologies have a big role to play in measuring, tracking, and keeping air pollution in check, each person should have individual environmental responsibility to make the world a better place to live in. *Internet of Things* and *Machine Learning* are augmenting the *Big Data* capabilities like never before. This ensures that we have more data points to work in a given time and continuous data streaming means more accurate real-time analytics with efficient *Machine Learning* algorithms. All these

three technologies will continue to work in tandem to keep a check on air pollution and the imminent threats.

## ACKNOWLEDGMENTS

The author would like to thank Dr. Gregor von Laszewski and the teaching assistants for their support and suggestions in writing this paper.

## REFERENCES

- [1] Ferrovia Blog. 2017. Big data will control pollution in your city. Webpage. (April 2017). <http://blog.ferrovial.com/en/2017/04/big-data-pollution-control-in-cities/>
- [2] Jay Hardikar. 2017. Environmental analysis in the era of cloud and big data platforms. Webpage. (Jan. 2017). <https://www.ibm.com/blogs/bluemix/2017/01/environmental-analysis-era-cloud-big-data-platforms/>
- [3] Alexander Howard. 2015. How IBM Is Using Big Data To Battle Air Pollution In Cities. Webpage. (Sept. 2015). <https://www.ibm.com/blogs/bluemix/2017/01/environmental-analysis-era-cloud-big-data-platforms/>
- [4] Gaganjot Kaur Kang, Jerry Gao, Sen Chiao, Shengqiang Lu, and Gang Xie. 2017. Air Quality Prediction: Big data and Machine Learning Approaches. *International Conference on Sustainable Environment and Agriculture* 1 (10 2017).
- [5] Research Applications Laboratory. 2016. Air Pollution: A Global Problem. Webpage. (April 2016). <https://ral.ucar.edu/pressroom/features/air-pollution-a-global-problem>
- [6] Ruiyun Yu, Yu Yang, Leyou Yang, Guangjie Han, and Oguti Ann Move. 2016. RAQ: A Random Forest Approach for Predicting Air Quality in Urban Sensing Systems. *Sensors* 16 (2016), 1–86. <http://www.mdpi.com/1424-8220/16/1/86>

## LIST OF FIGURES

1	Green Horizons air quality management for Beijing [3]	5
2	AQI classification [6]	5
3	Traffic Congestion[6]	5
4	Structure of RAQ data[6]	6
5	Procedure of RAQ [6]	6
6	RAQ Algorithm [6]	6
7	Artificial Neural Network(ANN) Model [4]	6
8	Least squares Support Vector Machine Model [4]	7

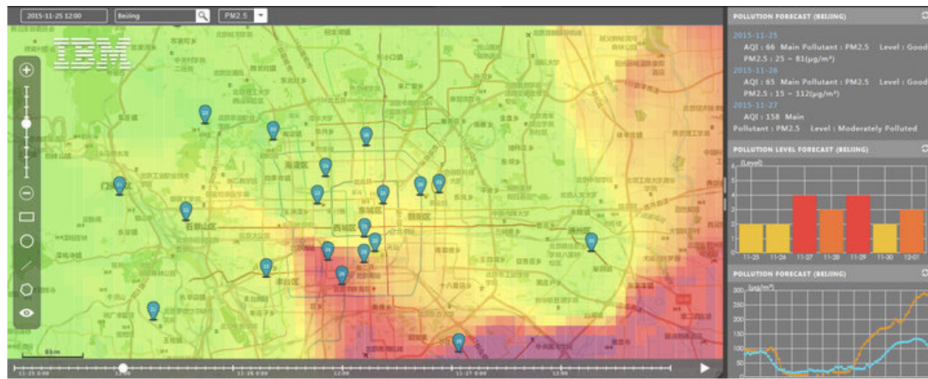


Figure 1: Green Horizons air quality management for Beijing [3]

AQI	Air Pollution Level
0-50	Excellent
51-100	Good
101-150	Lightly Polluted
151-200	Moderately Polluted
201-300	Heavily Polluted
300+	Severely Polluted

Figure 2: AQI classification [6]

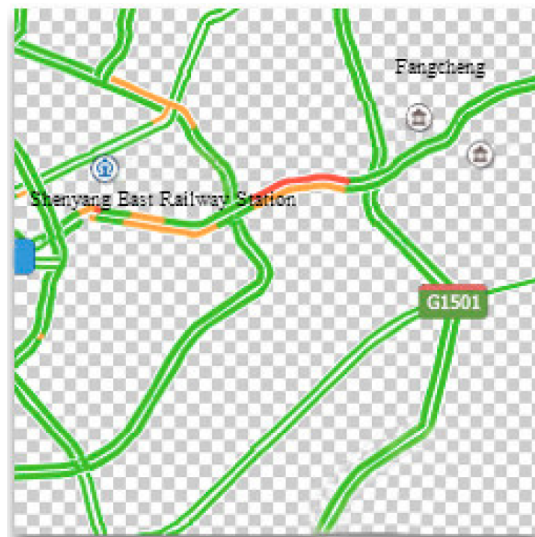


Figure 3: Traffic Congestion[6]

temperature	humidity	pressure	wind	visibility	road_length	tfs	poi_number	aqi
Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Nominal
5.5	89.0	758.1	2.0	14.0	2185.0	2371.0	63.0	excellent

Figure 4: Structure of RAQ data[6]

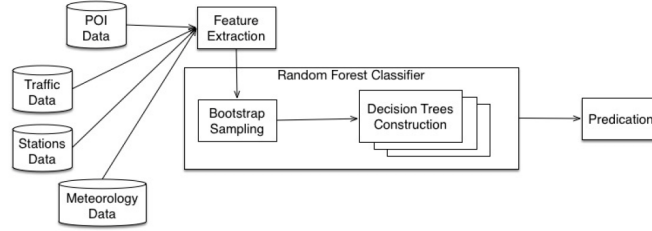


Figure 5: Procedure of RAQ [6]

Algorithm 1. RAQ	
<b>Input:</b>	A dataset $S$ with features: $F_{mt}, F_{nh}, F_{mp}, F_{mw}, F_{mv}, F_{ti}, F_{ics}, F_{pi}$ and labeled AQI level; unlabeled dataset $U$ ; trees quantity $T$ ; features quantity $m$ ;
<b>Output:</b>	AQI level
1	for $T$ trees
2	randomly select $m$ features from $S$ ;
3	for $m$ features in each node
4	calculate information gain by Equation (3);
5	choose maximum gain to split the dataset in the node;
6	remove used feature from feature candidates;
7	input unlabeled data into trees;
8	get predicted AQI level according to Equations (5) and (6);

Figure 6: RAQ Algorithm [6]

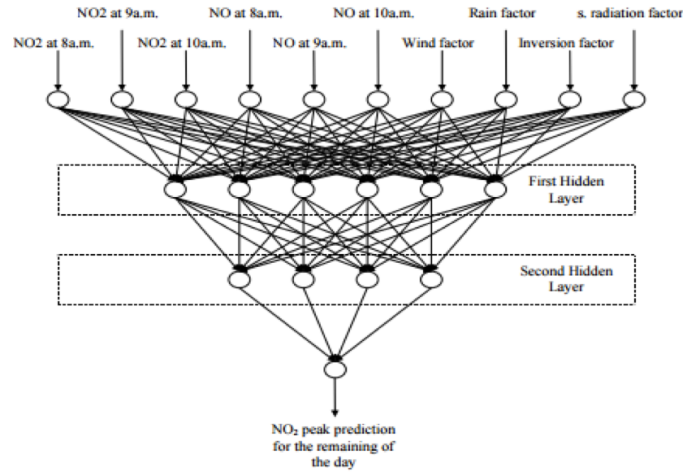


Figure 7: Artificial Neural Network(ANN) Model [4]

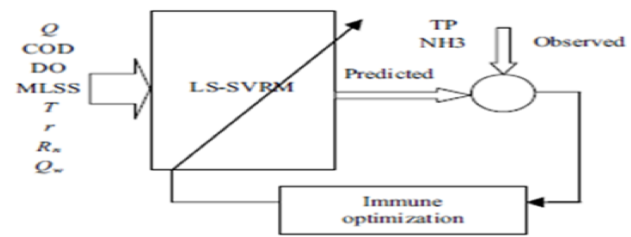


Figure 8: Least squares Support Vector Machine Model [4]