

# Machine Learning for Air Quality Analysis and Prediction

Bommanna K<sup>1</sup>, Dr Radha H R<sup>2</sup>, Praveen Math<sup>3</sup>, Dr Yuvaraja Naik<sup>4</sup>, Dr A Hareesh<sup>5</sup>

<sup>1</sup>Ph.D Research Scholar, Department of Mechanical Engineering, Visvesvaraya Technological University, Bangalore, India  
[bommanna\\_kta@rediffmail.com](mailto:bommanna_kta@rediffmail.com)

<sup>2</sup>Research Scientist, Department of Chemistry, Sarojaayudh Innovation cell, Bangalore, India  
[radha.sudhakar99@gmail.com](mailto:radha.sudhakar99@gmail.com)

<sup>3</sup>Assistant Professor, Department of Mechanical Engineering, Reva University, Bangalore, India  
[praveenmath@reva.edu.in](mailto:praveenmath@reva.edu.in)

<sup>4</sup>Assistant Professor, Department of Mechanical Engineering, Presidency University, Bangalore, India  
[yuvarajanaik@presidencyuniversity.ac.in](mailto:yuvarajanaik@presidencyuniversity.ac.in)

<sup>5</sup>Associate Professor, Department of Mechanical Engineering, ACS College of Engineering, Bangalore, India  
[hareeshcm@gmail.com](mailto:hareeshcm@gmail.com)

**Abstract:** One of the most serious ecological problems is air pollution. It can have negative health consequences such as cancer, cardiovascular disease, and a high death rate. High population density is a major contributor to air pollution in urban areas and urbanized regions, as well as having a negative impact on climate. The climate changes dramatically as more area is used for agriculture or dwelling. Every year, over 2,000,000 people die prematurely as a result of the effects of polluted air, according to the World Health Organization. Air pollution is a huge problem that affects everyone, not just those who live in brown haze-engulfed cities: it may affect us all through things like unusual weather changes and ozone layer damage. According to the Global Burden of Disease Report, outside air pollution is the fifth leading cause of death in India. A study conducted by scientists from the University of Chicago, Harvard, and Yale found that high Particulate Matter (PM) fixation is responsible for reducing the life expectancy of 660 million Indians living in urban areas by 3.2 years.

**Keywords**— Sensor Networks, synchronization, Noise factor, Linear regression

## 1. INTRODUCTION

Air pollution has recently elicited basic measurements, and the air quality in most Indian metropolitan communities that monitor open air pollution fails to meet WHO guidelines for safe levels. PM2.5 and PM10 (airborne particles with a width of less than 2.5 micrometres and a measurement of less than 10 micrometres) levels, as well as groupings of cancer-causing substances like Sulfur Dioxide (SO<sub>2</sub>) and Nitrogen Dioxide (NO<sub>2</sub>), have reached alarming levels in most Indian cities, putting people at increased risk of respiratory illnesses and other medical conditions. Furthermore, the problem of indoor air pollution has put women and children in grave danger.

According to data available from the Central Electricity Authority, coal-controlled nuclear power plants account for 60.72 percent of India's total installed capacity (CEA). Coal plants are one of the primary sources of SO<sub>2</sub> and NO<sub>2</sub>. Biomass is used by 87 percent of provincial families and 26% of metropolitan families for cooking.

Biomass consumption is a major source of indoor air pollution, causing respiratory and pneumonic illnesses in nearly 400 million Indians. In Indian cities, the number of vehicles is increasing. Private and commercial vehicles account for 66.28 percent of total diesel usage. Low standards for car emissions and fuel have resulted in increased nitrogen levels. Biomass is used by 87 percent of country families and 26 percent of city families for cooking. Biomass use is a major source of indoor air pollution, causing respiratory and

aspiratory difficulties in around 400 million Indians. The percentage of rural families who rely on lamp oil as a primary source of energy for lighting is about 30%. Lamp oil lights, which are commonly used in rural areas, are a major source of dark carbon residue emissions and have a significant impact on women's and children's health.

Indoor (family) and outdoor air pollution both have an impact on people's health as well as the economy. The negative effects of air pollution are not limited to urban areas; they also affect rural areas, where a large portion of the population relies on lamp oil and biomass consumption for lighting and cooking, respectively.

Contrary to popular belief, population growth has a negative impact on the climate. The climate changes dramatically as more land is used for horticulture or dwelling reasons. As the population of individual urban towns or provincial regions grows, more resources should be used to keep up with the population's wealth. Numerous dwelling areas are being eliminated as the strain on available resources grows. People are consuming more assets, and the economy and nature are unable to replenish those assets quickly enough to meet our needs. Population growth has a negative impact on the climate as well.

## 2. EASE OF USE

### 2.1 Machine Learning

The assessment of encouraging PCs to act without being unambiguously changed is known as AI. AI has brought us

self-driving cars, quick voice confirmation, practical web search, and a massively enhanced understanding of the human genome in the previous decade. Today, AI is so pervasive that you probably utilize it on a daily basis without even realizing it. Several experts believe it is the most effective strategy to move toward human-level AI.

Bunching is a Machine Learning approach that involves grouping data centers together. We can utilize a Clustering algorithm to organize each datum point into a specific get-together given a game plan of data centers. On a basic level, data centers in similar social gatherings should have similar properties and also incorporates, whilst data centers in other gatherings should have considerably diverse properties and also incorporates. Bunching is a tool for self-study and a common framework for analyzing quantitative data in a variety of professions.

## 2.2 K-Means

K-Means appears to be the easiest bunching calculation to grasp. It's covered in a large percentage of first-year data science and AI curricula. It's simple to understand and use in code; see the example below for an example of data analysis.

K-Means has the preferred standpoint that it's truly quick, as all we're truly doing is figuring the separations among focuses and bunch focuses; not very many calculations, it in this way has a direct complexity  $O(n)$ .

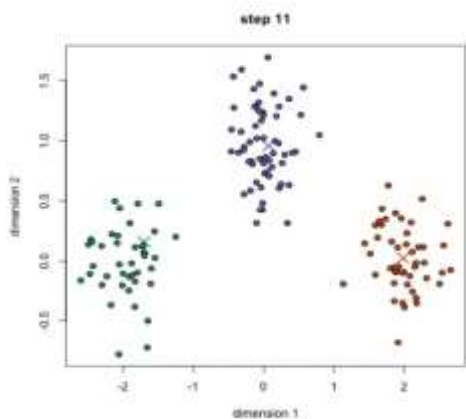


Fig 1: K-Means illustrational values

## 2.3 Linear Regression

A regression is a method for displaying a target's respect while accounting for independent pointers. The majority of the time, this method is utilised to determine and find conditions and final products associated with factors. Backslide techniques are often compared based on the number of independent components and the type of relationship between the free and ward elements.

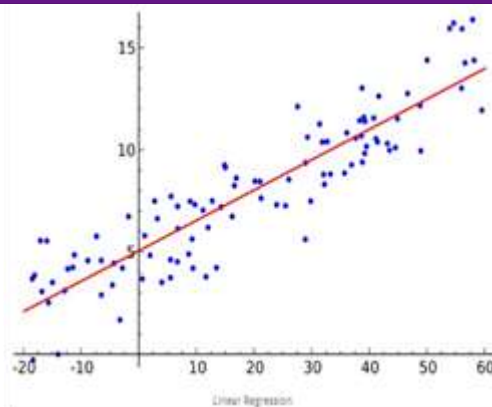


Fig 2: Linear regression value graph

Straight relapse is a type of backslide test in which the number of independent components is one and the independent(x) and dependent(y) variables have an immediate relationship. In the diagram above, the red line is inferred to be the best-fit straight line. We attempt to plot a line that best mimics the given data centres, taking into account the given data centres. The line can be depicted as follows when considering the straight condition.

$$y = a_0 + a_1 * x \quad (1)$$

## 2.4 R is a Programming Language

The R Foundation for Statistical Computing supports R, a programming language and free programming for real figuring and plans. The R programming language is often used by investigators and data miners for quantitative programming and data analysis. R's popularity has grown dramatically in recent years, according to overviews, polls of data diggers, and investigations of quickly constructing data sets. R is ranked eighteenth in the TIOBE document, which measures the reputation of programming dialects, as of August 2018.

## 2.5 R studio

R Studio is a free and open-source integrated development environment (IDE) for R, a computer language for calculating enrollment and planning. JJ Allaire, the creator of the ColdFusion programming language, founded RStudio. RStudio's Chief Scientist is Hadley Wickham. RStudio is well-versed in the C++ programming language, and its graphical user interface is built on the Qt framework. The higher level of the code is written in Java, and JavaScript is also one of the languages used.

## 3. PROPOSED MODEL

In the field of air registration, the insertion, expectation, and component investigation of air quality are three important topics. The responses to these questions can provide extremely useful information for air pollution mitigation, resulting in remarkable cultural and specialized

consequences. A lot of the current research tackles these difficulties in different ways using different models. We present a generic and persuasive technique to deal with the three difficulties in a single model in this study.

Our suggested framework consists of putting together an AI model that can predict air quality. To put our plan into action, we'll use unaided learning to obtain comparative knowledge and administered learning to predict outcomes. Whenever we used k-implies bunching to build groups based on their proximity and named these data for the next stage. We are using Collaborative separating AI computations to predict air quality, and we have the surrounding air quality statistics from the [www.data.gov.in](http://www.data.gov.in) site.

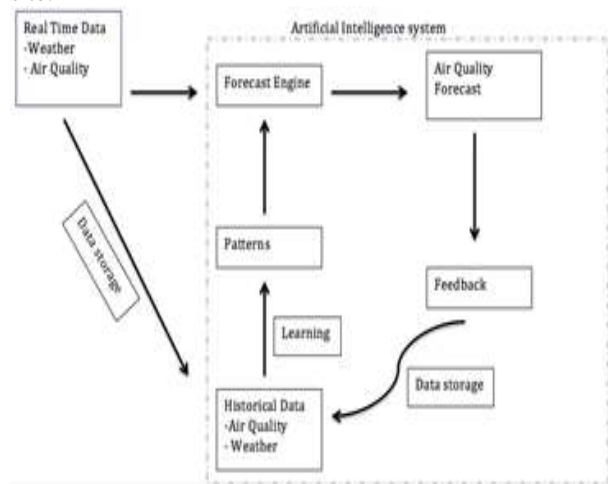


Fig 3: Proposed Model

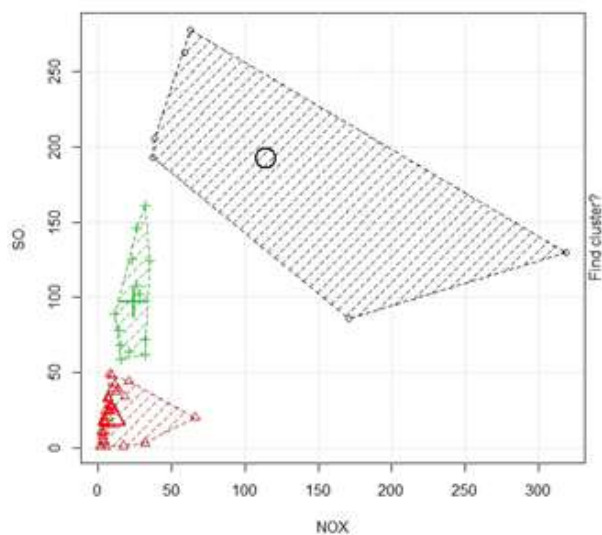


Fig 4: Air Quality values

#### 4. CONCLUSIONS

We used a probabilistic model to discover how settings are appropriated in metropolitan regions in terms of a few components in this study. We adapt to the incited sparsity by changing the scanty demonstrating approach of nearby information because most residents of a city do not visit the great majority of the accessible settings. We demonstrate the pieces of knowledge provided by a particularly unique methodology by applying our model to a vast informational collection of over 11 million checking in 40 urban communities all around the world. To begin, we calculated the likelihood of appropriation of a single component moulded on a certain location in the city using the removed model events. This allowed us to create a hotness guide for that element, highlighting what element esteems are unquestionably and specifically at various locations around a metropolis. We also demonstrated a principled approach to assessing the importance of various components inside the produced models. While all provisions help, we discovered that the guests of scenes are the most distinguishing feature of the sections disclosed by the model. This study suggests that a closer examination of client behaviour is a possible avenue for eliciting more experiences. Third, after focusing on the various districts of a single city, we used the deleted model instances to locate the most two comparable localities across two urban communities, a task that had previously been attempted using a more heuristic manner. This time, we use the strong hypothetical foundations of probabilistic models to specify a principled proportion of closeness, and we depict a strategy for covetously discovering two locations enhancing measure. In a few metropolitan societies, we outline this collaborating with method using narrated proof. Finally, we compare our methodology to other ways that provide equivalent results, demonstrating that our locations are both more dependable with information (in terms of predictive execution) and have more well specified attributes, making them easier to distinguish from one another. An examination of recent related works in the Urban Computing field suggests that, while the region is dynamic and understanding metropolitan exercises is a commendable undertaking that benefits from geo-labeled data, it could very well be aided by the use of probabilistic models, as such models have tremendous interpretative power. The collected centroids are used to group and label informational indexes. The K-implies grouping algorithm is implemented to the given informational collection esteems, and the bunches are classified into three classes based on the acquired cetroid values. When compared to the outcomes from the relapse results, the Linear Regression Model yield demonstrates that the predicted qualities are even closer to the actual qualities.

#### 5. REFERENCES

[1] Y. Zheng, L. Capra, O. Wolfson, and H. Yang (2014). "Urban Computing: Concepts, Methodologies, and

- Applications,” *ACM Transaction on Intelligent Systems and Technology*, vol. 5, no. 3, pp. 38:1–38:55.
- [2] J. Eisenstein, A. Ahmed, and E. P. Xing (2016). “Sparse Additive Generative Models of Text,” in *ICML*, Seattle, WA, pp. 1041–1048.
- [3] J. Cranshaw, J. I. Hong, and N. Sadeh (2017). “The livelihoods project: Utilizing social media to understand the dynamics of a city,” in *ICWSM*, pp. 58–65.
- [4] A. X. Zhang, A. Noulas, S. Scellato, and C. Mascolo, (2016). “Hoodsquare: Modeling and recommending neighborhoods in location-based social networks,” in *ASE/IEEE SocialCom*, pp. 69–74.
- [5] V. Frias-Martinez and E. Frias-Martinez (2017). “Spectral clustering for sensing urban land use using Twitter activity,” *Engineering Applications of Artificial Intelligence*, vol. 35, pp. 237–245.
- [6] G. Le Falher, Gionis Aristides, and M. Mathioudakis, (2015). “Where Is the Soho of Rome? Measures and Algorithms for Finding Similar Neighborhoods in Cities,” in *ICWSM*, Oxford.
- [7] J. L. Toole, M. Ulm, M. C. González, and D. Bauer (2018). “Inferring Land Use from Mobile Phone Activity,” in *UrbComp*, New York, NY, USA, pp. 1–8.
- [8] J. R. Hipp, R. W. Faris, and A. Boessen (2017). “Measuring ‘neighborhood’: Constructing network neighborhoods,” *Social Networks*, vol. 34, no. 1, pp. 128–140.
- [9] J. O. Berger and D. Sun (2017). “Objective priors for the bivariate normal model,” *The Annals of Statistics*, pp. 963–982.
- [10] S. C. Deerwester, S. T. Dumais, T. K. Landauer, G. W. Furnas, and R. A. Harshman (2018). “Indexing by latent semantic analysis,” *JAsIs*, vol. 41, no. 6, pp. 391–407.
- [11] N. J. Yuan, Y. Zheng, X. Xie, Y. Wang, K. Zheng, and H. Xiong (2015). “Discovering Urban Functional Zones Using Latent Activity Trajectories,” *Knowledge and Data Engineering, IEEE Transactions on*, vol. 27, no. 3, pp. 712–725.
- [12] Z. Cao, S. Wang, G. Forestier, A. Puissant, and C. F. Eick (2016). “Analyzing the Composition of Cities Using Spatial Clustering,” in *UrbComp*, New York, NY, USA, pp. 14:1–14:8.