# Air Quality Prediction Using Machine Learning

**Dr. K R Shylaja[1], Sushma C V[2]**

Professor, Dept of CS&E, Dr.AIT, Bengaluru, India[1]

M. Tech Student, Dept of CS&E, Dr.AIT, Bengaluru, India[2]

*Abstract: The environment is becoming increasingly contaminated as a result of human activities such as industrialization and urbanization. Air pollution is a combination of natural and man-made chemicals that has a variety of negative effects on humans and the environment. Toxic or hazardous contaminants such as SO2, NO2, CO, PM, and toxic organics are released in significant quantities by most manufacturing activities. This paper presents a Machine Learning based methodology for gauging air quality. It goes through several phases of pre-processing, learning, and evaluation. Based on the level of dataset obtained, the ML technique Random Forest (RF) is investigated to predict air quality. Precision, recall, and F1 score are used as performance assessment indicators.*

*Index Terms* - **Machine Learning, Random forest, Air pollution, Air quality index.**

## I. INTRODUCTION

Air contamination is really perhaps the most genuine ecological concerns. To guarantee better expectations for everyday comforts, air quality should be persistently controlled and estimated. Each living being requires clean air to endure. One of the principal issues affecting human wellbeing, farming harvests, woods species, and biological systems is air contamination. Air contamination has been related with improved grimness and mortality. Dirtied air contrarily affects the climate and living organic entities. Solid particles, fluid beads, or gases may all be utilized as the material. Toxins might be either normal or man-made. Essential foreign substances, for example, carbon monoxide gas from car fumes or Sulphur dioxide discharged by plants, are ordinarily the aftereffect of a technique. Essential foreign substances are not unequivocally radiated. Maybe, they structure noticeable all around because of the response or communication of essential pollutants. Ozone at ground level is a notable illustration of an optional toxin [1]. One such instrument for adequately dispersing air quality data to individuals is the Air Quality Index (AQI) [2]. As the Air Quality Index grows, a rising level of the populace is probably going to endure progressively genuine wellbeing fallouts. An air toxin focus from a sensor or model is needed to compute the AQI. The recipe for changing over air contamination focuses to AQI differs by toxin and by locale. The upsides of the air quality file are isolated into ranges, with each reach having its own descriptor and shading code. We utilized a managed learning approach in the proposed study. Direct Regression, Nearest Neighbor, SVM, bit SVM, Naive Bayes, and Random Forest are instances of regulated learning calculations. Since Random Forest calculation give more exactness thus, we picked it for our way to deal with precisely anticipate air contamination.

## II. RELATED WORK

### 2.1 Air contamination in Delhi

Delhi is one of the quickest developing urban communities on the planet, with a populace of more than 19.3 million individuals [3]. In ongoing many years, populace thickness and advancement, just as fast modern extension, have brought about hazardously significant degrees of air contamination, neglecting to furnish individuals with one of life's most fundamental necessities: clean air. According to the World Economic Forum, India has six of the world's ten most dirtied metropolitan regions, with Delhi being one of them [4]. As indicated by reports, contaminated air is one of the main sources of sudden passing [5], and the normal future is declining as air contamination levels increment [6]. Next to current transmissions, farm fires are eminent as one of the critical wellsprings of air defilement in Delhi [7].As per a United Nations concentrate from 2016, Delhi had a populace of 26 million individuals in the more prominent metro locale, which is required to develop to a day and a half by 2030, making it the world's second most crowded city after Tokyo [8]. The worldwide populace thickness is one of the greatest on the planet, representing extra difficulties to air quality and wellbeing. As shown by a report circulated in 2018 by the Indian Ministry of Earth Sciences, the

number of vehicles has extended fourfold since 2000, making it a huge ally of air defilement in India, which fuses PM2.5 and risky nitrogen oxide [9]. Among January and September of the previous five years, the air quality rundown in Delhi was generally moderate (101–200).

During October to December, the air quality record jumps to low (301–400), then to genuine (401–500) or risky (500+) levels due to an arrangement of factors [10]. To the extent the speed of unfamiliar substances and the control steps taken to restrict them, the air pollution condition in Delhi has changed fundamentally. Sulian et al. [11] give a proof-based gander at the current status of air tainting in Delhi, its impact on prosperity, and the controls that have been set up. The meteorological conditions, for example, provincial and brief meteorology are huge in deciding the air toxin fixations Lower wind speed (powerless scattering/ventilation) can bring about higher groupings of traffic poisons [12]. Regardless, strong breeze speed may outline dust storms and end up blowing the particles on the ground. Higher tenacity levels are connected with higher sums of air defilements like PM2.5, carbon monoxide (CO), nitrogen dioxide (NO2) and sulphur dioxide (SO2) [13]. Also, significant degrees of dampness frequently demonstrate precipitation occasions which bring about substantial wet statement prompting the bringing down of air contaminations.

## 2.2. AI approaches for air quality anticipating

We glance back at how AI has been utilized to foresee air quality throughout the long term. AI approaches for air quality determining have had impressive accomplishment in various regions. Notwithstanding the upsides of neural organizations over ordinary measurable methodologies in air quality determining, there is still space for development because of difficulties like computational expense, imperfect combination, over-fitting, and loud information. Besides, the setup of the organization geography and model boundaries, which influences forecast precision, is a test. Corani [14] utilized feedforward and pruned neural organizations to anticipate hourly PM-10 fixations dependent on information from the earlier day.

AI is a functioning space of study, and new strategies and procedures for more refined recreation of a specific issue emerge consistently. To foster standard neural organization models, Fu et al. [15] utilized a moving component and a dim model. Chang et al. [16] thought about the impacts of accumulated LSTM organizations to help vector relapse and inclination supported tree relapse. Diverse added substance backslides trees, significant feedforward neural associations, and a cross variety model ward on LSTM networks were used by Karimian et al. [17] to figure air quality given by PM-2.5 obsessions all through different time intervals, with the LSTM model being the best for gauging and controlling air tainting. To give fast reaction metropolitan air contamination conjectures and controls, Xiao et al. [18] parameterized a non-meddling decreased request model dependent on appropriate symmetrical disintegration for model decrease of contamination transport conditions.

## 2.3. Air contamination's impact

Information driven AI approaches empower specialists to assess the impact of various air poisons on wellbeing results simultaneously. There is mounting proof that early-life openness to indoor air contamination affects youngsters' neurodevelopment. An examination found that air contamination is a potential danger factor for corpulence in grown-ups with a higher BMI, which calls for additional investigation into other wellbeing impacts. In a deliberate investigation of younger students' affectability to air contamination during the day-by-day drive, scientists discovered examinations connecting schoolchildren's drive openness to antagonistic psychological results and outrageous wheeze in asthmatic kids. Moreover, kids' lung usefulness and other respiratory problems have been connected to encompassing air contamination. Aside from wellbeing, air contamination has detrimentally affected various different enterprises, including farming. As indicated by an examination in China, mechanical air contamination fundamentally affects farming efficiency, bringing about lower minor items and further adjusting work capital and different components [19]. Finally, air defilement influences monetary turn of events. An examination of the association between air defilement and stock returns uncovered that advanced air pollution decreases provincial creation's mechanical adequacy basically [20].

## III. METHODOLOGY

AI: Machine learning is a type of man-made reasoning (AI) that permits frameworks to adapt expressly and improve their presentation dependent on past experience. Learning begins with perceptions, practice, and guidance to discover patterns in information and improve perception and dynamic.
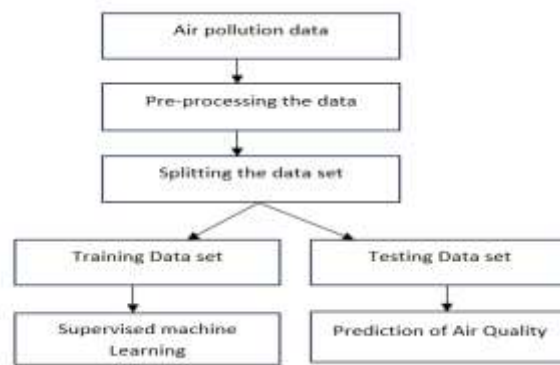
Figure 1: Architecture of proposed model

**Data Collection**: Data Collection is the way toward gathering and estimating data from an assortment of sources. It should be gathered and put away such that bodes well for the current issue. The dataset gathered from Kaggle it incorporates different convergence of toxins. The complete properties in the informational index are Temperature, $CH_4$ (Methane), CO (Carbon Monoxide), NMHC (Non-Methane Hydro-Carbons), PM2.5, RH (Relative Humidity), and $SO_2$ (Sulphur Dioxide).

**Pre-Processing**: The data we get from different sources may contain clashing data, missing characteristics and repeated data. To get authentic assumption result, the dataset ought to be cleaned, missing characteristics ought to be managed either by deleting or by stacking up with mean characteristics or some other system. Likewise, excess information should be taken out or wiped out to try not to inclination of the outcomes. Some dataset may have some anomaly or outrageous qualities which additionally must be eliminated to get great forecast precision.

**Building the characterization model:** In the first place, we need to partition the informational collection into preparing and testing set. The anticipating model is first prepared with the preparation dataset. Later it will be tried with the testing set. In the wake of testing the model, the exactness of the model is assessed by utilizing boundaries like location rate, exactness, review, F-Measure and generally precision.

**Model training:** Random forest, an administered learning group calculation, joins different choice trees to frame backwoods and the sacking standard, the last of which fuses haphazardness into the model structure measure. The individual is parted utilizing an arbitrary arrangement of highlights. For of choice tree, an irregular exhibit of cases is utilized to develop a preparation information subset. The variable from the irregular number of highlights is considered for the best split at every choice hub in any tree. Irregular woodlands would choose the most successive as its indicator if the objective trait is straight out. On the off chance that it's a mathematical question, the normal, everything being equal, will be chosen.
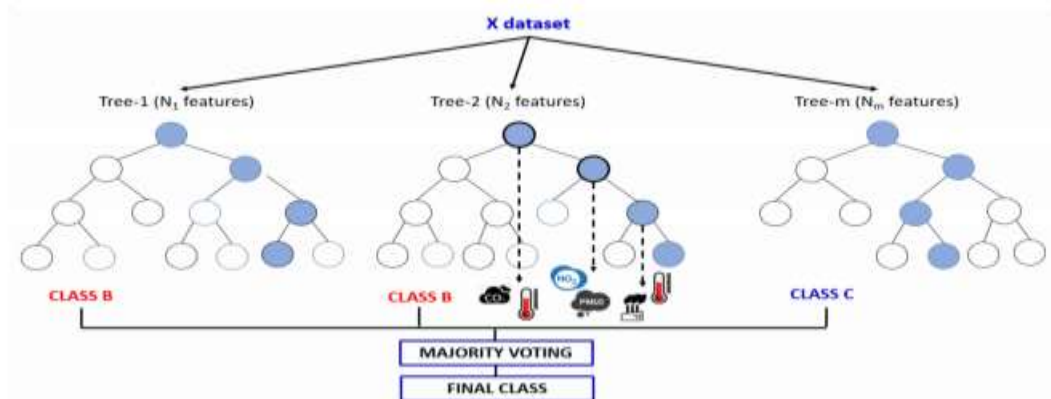


Figure 2: Random Forest Classifier

Random forest, can handle both classification and regression cases. Every test data point is passed through every decision tree in the forest for prediction. The trees then, at that point vote on a result, and the forecast is made dependent on a dominant part vote among the models, bringing about a more effective and stable single student. The forecast normal would inexact the ground truth (arrangement) or genuine worth, permitting arbitrary backwoods to determine the expectation change that every choice tree has (relapse).
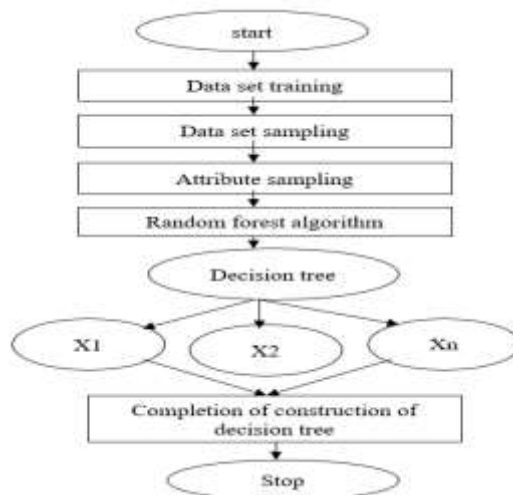
Figure 3: Flow chart of random forest

## IV. RESULTS

A result is the qualitative or quantitative expression of the end result of behavior or events. Performance analysis is a type of operational analysis that consists of a collection of fundamental quantitative relationships between performance variables. The delayed consequences of the examinations are discussed in this part to choose the capability of the proposed air quality assumption model.
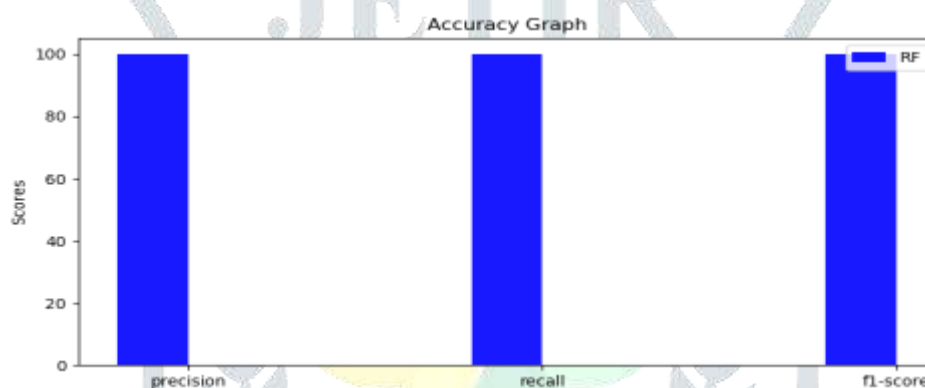


Figure 4: prediction of air quality using random forest algorithm

Negative tuples are any tuples other than sure tuples, with N being the quantity of negative tuples. The quantity of genuine positives is indicated by T P. Negative tuples that are erroneously named positives are known as bogus positives. The quantity of bogus positives is indicated by F P. Negative tuples that are effectively marked by the classifier are alluded to as obvious negatives. The quantity of genuine negatives is meant by T N. Positive tuples that are wrongly marked as negatives are known as bogus negatives. The quantity of bogus negatives is indicated by F N.

**Accuracy** - The most natural achievement metric is accuracy, which is basically the quantity of effectively anticipated that observations should all perceptions. One may accept that if our model is right, it is the awesome. Indeed, exactness is a helpful measurement, yet just when the datasets are symmetric and the upsides of bogus positives and bogus negatives are almost equivalent.

$$Accuracy = \frac{True_{positive} + True_{negative}}{True_{positive} + True_{negative} + False_{positive} + False_{negative}}$$

**Precision** - The proportion of accurately anticipated positive perceptions to add up to anticipated positive perceptions is known as exactness.

$$Precision = \frac{True_{positive}}{True_{positive} + False_{positive}}$$

**Recall**- The proportion of accurately anticipated that positive observations should all perceptions in the real class is known as Recall.

$$Recall = \frac{True_{positive}}{True_{positive} + False_{negative}}$$

### V. CONCLUSION

Air contamination is perhaps the most genuine ecological difficulties, and it is deteriorating as the planet turns out to be more urbanized and industrialized. Meteorological factors, for example, air wind speed, wind bearing, relative moistness, and temperature impact the grouping of air contaminations in encompassing air. The Air Quality Index (AQI) is a device for assessing air quality. The outcomes demonstrate that RF-based air quality expectation is promising, as demonstrated by the outcomes.

### REFERENCES

[1] D. Domanska, et al. "Explorative forecasting of air pollution". Atmospheric Environment 92 (2014) 19-30.

[2] http://pib.nic.in/newsite/PrintRelease.as px?relid=110654 Ministry of environment and forests

[3] "Unique identification authority of India," May 2020, [Online; accessed 16-July-2020]. [Online]. Available: https://uidai.gov.in/ images/state-wise-aadhaar-saturation.pdf

[4] "6 of the world's 10 most polluted cities is in India," August 2020, [Online; accessed 22-August-2020]. [Online]. Available: https://www.weforum.org/agenda/2020/03/ 6-of-the-world-s-10-most-polluted-cities-are-in-india/

[5] "7 million premature deaths annually linked to air pollution,"august 2020,[online;accessed 22-august-2020].[online].available:https://www.who.int/mediacentre/news/releases/2014/air-pollution/en/

[6] J. S. Apte, M. Brauer, A. J. Cohen, M. Ezzati, and C. A. Pope III, "Ambient pm2. 5 reduces global and regional life expectancy," Environmental Science & Technology Letters, vol. 5, no. 9, pp. 546–551, 2018.

[7] D. H. Cusworth, L. J. Mickley, M. P. Sulprizio, T. Liu, M. E. Marlier, R. S. DeFries, S. K. Guttikunda, and P. Gupta, "Quantifying the influence of agricultural fires in northwest india on urban air pollution in delhi, india," Environmental Research Letters, vol. 13, no. 4, p. 044018, 2018.

[8]" the worlds cities in 2016: data booklet," june 2016, [united nations, online; accessed 24-august2020].[online].available:https://www.un.org/en/development/desa/population/publications/index.asp

[9] "Usual suspects: Vehicles, industrial emissions behind foul play," August 2018, [Online; accessed24-August-2020]. [Online]. Available: https://timesofindia.indiatimes.com/city/delhi/ usual-suspects-vehicles-industrial-emissions-behind-foul-play-all-year/ articleshow/66228517.cms

[10] "Delhi breathed easier from january to april," June 2017, [Online; accessed 24-August-2020]. [Online]. Available: https://timesofindia.indiatimes.com/city/delhi/delhi-breathed-easier-from-january-to-april/articleshow/59011204.cms

[11] R. Suliankatchi, B. Nongkynrih, and S. Gupta, "Air pollution in delhi: Its magnitude and effects on health," Indian Journal of Community Medicine, vol. 38, pp. 4–8, 01 2013.

[12] A. Y. Watson, R. R. Bates, and D. Kennedy, "Atmospheric transport and dispersion of air pollutants associated with vehicular emissions," in Air Pollution, the Automobile, and Public Health. National Academies Press (US), 1988.

[13] R. Zalakeviciute, J. Lopez-Villada, and Y. Rybarczyk, "Contrasted effects of relative humidity and precipitation on urban pm2. 5 pollution in high elevation urban areas," Sustainability, vol. 10, no. 6, p. 2064, 2018.

[14] G. Corani, "Air quality prediction in milan: Feed-forward neural networks, pruned neural networks and lazy learning," Ecological Modelling, pp. 513–529, 07 2005

[15] M. Fu, W. Wang, Z. Le, and M. Safaei khorram, "Prediction of particular matter concentrations by developed feed-forward neural network with rolling mechanism and gray model," Neural Computing and Applications, 02 2015.

[16] Y.-S. Chang, H.-T. Chiao, S. Abimannan, Y.-P. Huang, Y.-T. Tsai, and K.-M. Lin, "A lstm-based aggregated model for air pollution forecasting," Atmospheric Pollution Research, vol.11,no.8,pp.1451 – 1463, 2020. [Online]. Available: http://www.sciencedirect.com/science/ article/pii/S1309104220301215

[17] H. Karimian, Q. Li, C. LI, L. Jin, J. Fan, and Y. Li, "An improved method for monitoring fine particulate matter mass concentrations via satellite remote sensing," Aerosol and Air Quality Research, vol. 16, pp. 1081– 1092, 01 2016

[18] D. Xiao, J. Zheng, C. Pain, and I. Navon, "Machine learning-based rapid response tools for regional air pollution modelling," Atmospheric Environment, vol. 199, 11 2018.

[19] Z. Wang, W. Wei, and F. Zheng, "Effects of industrial air pollution on the technical efficiency of agricultural production: Evidence from china," Environmental Impact Assessment Review, vol. 83, p. 106407, 2020.

[20] M. Xu, Y. Wang, and Y. Tu, "Uncovering the invisible effect of air pollution on stock returns: A moderation and mediation analysis," Finance Research Letters, p. 101646, 2020