



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

IJIRSET

International Journal of Innovative Research in
SCIENCE | ENGINEERING | TECHNOLOGY



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 13, Issue 11, November 2024

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 8.524

 9940 572 462

 6381 907 438

 ijirset@gmail.com

 www.ijirset.com

Air Quality Prediction in Tumakuru City: A Comparative Analysis of Machine Learning Models with Visualizations

Himaja S, Gagana A N, Mahesh N, Suhas K C

U.G. Student, Department of Computer Science and Engineering, Channabasaveshwara Institute of Technology, Gubbi, Tumkur, Karnataka, India

U.G. Student, Department of Computer Science and Engineering, Channabasaveshwara Institute of Technology, Gubbi, Tumkur, Karnataka, India

Assistant Professor, Department of Computer Science and Engineering, Channabasaveshwara Institute of Technology, Gubbi, Tumkur, Karnataka, India

Assistant Professor, Department of Computer Science and Engineering, Channabasaveshwara Institute of Technology, Gubbi, Tumkur, Karnataka, India

ABSTRACT: Predicting air quality has become more important in cities as health risks from pollution exposure have grown. This study attempts to predict Tumakuru City's air quality with the goal of determining which model is the most accurate. To create machine learning models that predict pollution levels, key air quality indices such as particulate matter (PM_{2.5}, PM₁₀), CO, NO₂, O₃ and SO₂ were studied, visualized and used to calculate the Air Quality Index (AQI). The dataset was used to train and verify a number of machine learning methods, such as Random Forest, XGBoost, RNN, LSTM and Bi-LSTM models. Comprehensive data visualizations were conducted to understand pollutant trends and the spatial distribution of pollution across Tumakuru. These visuals provided essential context, highlighting patterns such as high-risk locations and peak pollution periods. Based on prediction accuracy and consistency, the most effective model was Random Forest. Our study demonstrates the effectiveness of machine learning in managing urban air quality and offers a framework for the adoption of data-driven strategies for pollution control and public health protection in cities facing comparable environmental difficulties.

KEYWORDS: Air quality prediction, Machine learning, Data visualization, Air quality index.

I. INTRODUCTION

As urbanization continues to accelerate, cities worldwide are facing escalating challenges related to air quality and pollution. Tumakuru City, like many rapidly developing urban centres, has experienced a significant increase in pollutants that pose serious health risks to its residents. Poor air quality has been linked to various respiratory and cardiovascular diseases, underscoring the urgent need for effective monitoring and management. Recently, machine learning has emerged as a powerful tool for the analysis and prediction of environmental data, providing the capability to uncover complex relationships within large datasets. The ability to predict air quality accurately can inform public health initiatives, guide regulatory policies, and enable proactive measures to mitigate pollution exposure.

The incorporation of the Air Quality Index (AQI) into our analysis allows for a standardized assessment of air quality levels, making it easier for public administrators and the public to understand the implications of pollution data. The AQI serves as an essential tool for communicating the health risks associated with various pollutant concentrations, enabling timely responses to pollution events.

In light of the escalating health concerns associated with urban air pollution, this study investigated the air quality dynamics in Tumakuru City, Karnataka, using advanced machine learning techniques. By focusing on key air quality indices, including PM_{2.5}, PM₁₀, CO, NO₂, O₃, and SO₂, we developed a comprehensive dataset to facilitate the

accurate prediction of air quality. The study highlights the crucial role of data pre-processing, which involved cleaning the dataset to remove inconsistencies, handling missing values through techniques like imputation.

This study employed an array of machine learning algorithms, including Random Forest, XGBoost, RNN, LSTM, and Bi-LSTM models, to systematically assess their performance against various evaluation metrics. Effective visualization techniques were also emphasized to understand pollutant trends and spatial distributions. Ultimately, this work not only advances the application of machine learning in environmental science but also contributes to the development of strategic frameworks for managing urban air quality and protecting public health.

The paper is organized as follows. Section II discusses the methodology for predicting air quality in Tumakuru City. In Section III, we describe the implementation of various machine learning algorithms used in the prediction model. Section IV presents the experimental results, showcasing the predictive accuracy of each model. Section V concludes the paper, summarizing the findings, discussing their significance in urban air quality monitoring, and outlining potential avenues for future research in air quality prediction using advanced machine learning techniques.

II. LITERATURE SURVEY

In this work [1], the authors highlight the critical role of monitoring air quality in real time and show that among various models, including Support Vector Machine (SVM), Logistic Regression, Linear Regression, and Random Forest Regression, the Random Forest Classifier yields the most accurate results. Their future plans involve expanding this research into developing an application for tracking air pollution and implementing the Extended Kalman Filter algorithm to estimate the air quality index (AQI).

In this paper [2], the author's investigation highlighted the increasing significance of Machine Learning in forecasting the Air Quality Index (AQI) and its crucial role in safeguarding public health from pollution. They identified a knowledge gap in current research regarding the elements affecting air quality and the model parameters that improved prediction precision. Analysing data from Manali and Velachery, two regions in Chennai, they concluded that XGBoost demonstrated superior performance compared to other models in AQI prediction. The researchers recommended that future studies should explore the integration of sophisticated neural networks like LSTM and time-series models such as SARIMA. These proposed improvements could potentially enhance the accuracy and efficacy of AQI forecasting substantially.

The authors developed an innovative approach to predict the Air Quality Index (AQI) in urban areas, particularly in Delhi, India, by integrating ensemble learning techniques in this paper [3]. They utilized a Meta Classifier that combined the strengths of multiple models, including Multi-Layer Perceptron (MLP), XGBoost, and Random Forest, to achieve superior predictive performance. Their findings indicated that the Stacking Classifier reached an impressive accuracy of 98.25%, showcasing its effectiveness in categorizing air quality conditions. The high F1 Score of 0.98 reflected a well-balanced performance between precision and recall. The research emphasized the potential of combining predictive models with health advisory systems to provide timely alerts to vulnerable populations, ultimately contributing to improved air quality monitoring and urban sustainability.

The researchers emphasized the crucial importance of monitoring air quality, given the severe health consequences of pollution, including brain damage and lung cancer in this study [4]. To evaluate air quality, they employed various machine learning approaches, such as Linear Regression, Random Forest, Decision Tree, Artificial Neural Network, and Support Vector Machine. The research examined datasets from Kaggle and air quality monitoring locations to determine the most effective predictive model. The findings revealed that Decision Tree regression outperformed Random Forest and Gradient Boosting methods, achieving an MAE of 20% and an RMSE of 0.0725.

The study [5] underscored the critical need for air quality monitoring, given the detrimental health effects of pollution, including neurological damage and pulmonary carcinoma. To assess air quality, the investigators utilized diverse machine learning techniques, including linear regression, random forest, decision tree, artificial neural network, and support vector machines. This study analysed datasets from Kaggle and air quality monitoring sites to identify the most accurate predictive model. The results indicated that Decision Tree regression surpassed Random Forest and Gradient Boosting methods in terms of performance, yielding an MAE of 20% and an RMSE of 0.0725.

The paper [6] tackles air pollution in smart cities using machine learning techniques like Decision Tree and Random Forest regression. Cloud computing reduced execution times while maintaining model accuracy. Random Forest algorithm outperformed others.

In this paper [7], three machine learning algorithms—SARIMAX, Prophet, and LSTM—are compared in this work for univariate time series forecasting in Delhi, India, the most polluted city in the world. Better accuracy may be achieved by LSTM with more training, but SARIMAX outperformed its rivals. SARIMAX could be improved by modifying parameters, and the models should be tested on a variety of datasets to increase performance. In order to record the levels of air pollution in other places, the study also recommends broadening the reach outside Delhi.

With an emphasis on PM2.5, the article [8] addresses air pollution control in developing and heavily populated countries. By examining characteristics including transmission, sky clarity, and colour variance, it suggests a technique for detecting PM levels using image processing technologies. The study emphasises how crucial precise picture processing is for evaluating the surroundings.

III. METHODOLOGY

The methodology employed in this study for predicting air quality using machine learning is depicted in Figure 1 and comprises several critical phases. It begins with data preparation and collection, followed by addressing class imbalances to ensure model reliability. The final phase involves training the selected machine learning models on the processed dataset. This structured approach is designed to enhance the accuracy and effectiveness of air quality predictions.

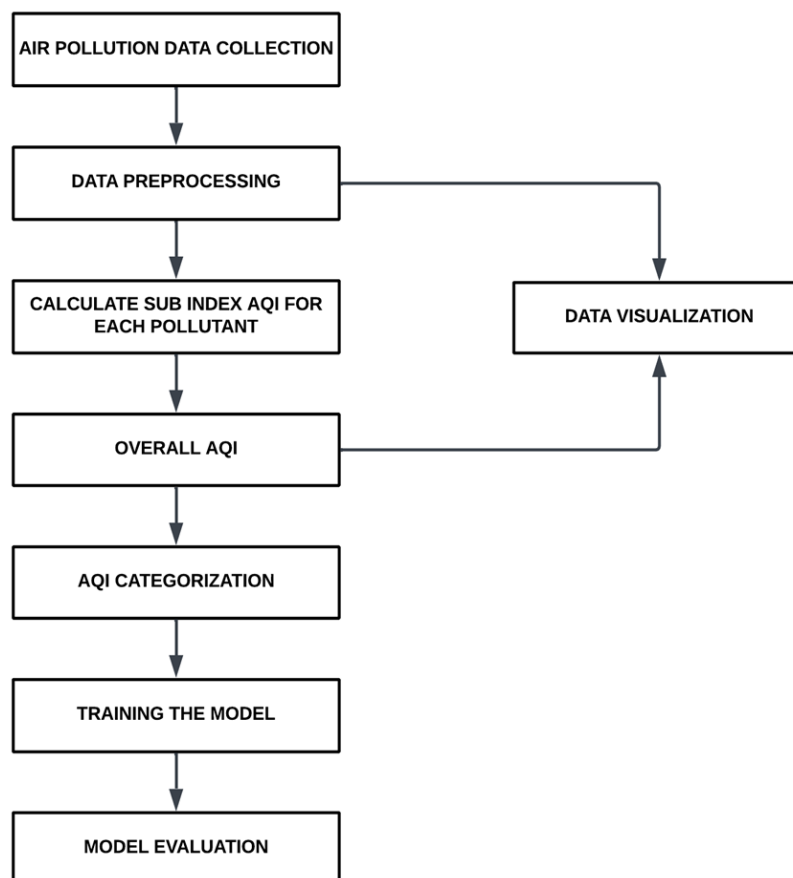


Figure 1. Air Quality Prediction Methodology

- Data Collection:** This study utilized air quality data from the Integrated Command and Control Centre (ICCC) in Tumakuru, where advanced sensor technology continuously monitors pollutant levels in real time, as summarized in Table 1. Data were acquired covering the period from May 2024 to September 2024.¹

This dataset includes key air quality indices necessary for understanding the pollution levels in Tumakuru City. Upon acquiring the dataset, several alterations were performed to tailor it to the research requirements. This involved the removal of irrelevant columns that did not contribute to our analysis, ensuring that the dataset was streamlined and focused on essential air quality metrics.

- Data Pre-processing:** In the data pre-processing phase, the dataset was examined for any missing or null values across all columns because handling these values was essential for maintaining the integrity of the data. To address missing values, an imputation strategy was applied by filling them with the mean value specific to each location and hour, ensuring that location-based variations and temporal patterns were preserved without introducing bias from other locations or timeframes. Categorical variables such as "Month" and "Location" were transformed into numerical formats to facilitate integration with machine learning algorithms. This conversion preserved essential information while optimizing the dataset for numerical analysis. These pre-processing steps enhanced the dataset, ensuring it was ready for accurate air quality predictions and providing a solid foundation for further model development.

By using the Adaptive Synthetic Sampling (ADASYN) technique, class imbalance in air quality prediction was addressed. In order to balance the distribution of classes, this oversampling technique creates synthetic samples for the minority class. The prediction performance of the machine learning model is improved by increasing the representation of underrepresented groups. When ADASYN was applied to the training dataset, artificial samples that closely matched the actual data points from a minority class were added. This improved the model's learning capacity and strengthened its comprehension of the minority class by balancing the dataset.

Table 1: Data Collection Parameters

Column Name	Description
Date Time	Date and time in the format YYYY-MM-DD HH:MM
CO (mg/m ³)	Concentration of carbon monoxide in milligrams per cubic meter
NO2 (µg/m ³)	Concentration of nitrogen dioxide in micrograms per cubic meter
O3 (µg/m ³)	Concentration of ozone in micrograms per cubic meter
SO2 (µg/m ³)	Concentration of sulphur dioxide in micrograms per cubic meter
PM2.5 (µg/m ³)	Concentration of particulate matter less than 2.5 µm in micrograms per cubic meter.
PM10 (µg/m ³)	Concentration of particulate matter less than 10 µm in micrograms per cubic meter
Temperature (°C)	Temperature in degrees Celsius
Humidity	Relative humidity percentage
CO2 (ppm)	Concentration of carbon dioxide in parts per million
Month	Month of the recorded measurements
Year	Year of the recorded measurements
Location	Specific location where the measurements were taken.

- AQI Calculation:**

Each pollutant's AQI is initially established using the concentration ranges and breakpoints specified by standard air quality recommendations. Depending on their observed quantities, pollutants such as PM2.5, PM10, CO, NO2, O3, and SO2 each give a unique AQI score. Each pollutant's contamination level is represented by these distinct AQIs, which are determined separately. The highest individual AQI among all pollutants is then chosen to indicate the dominant pollutant that has the biggest influence on air quality, and this is how the total AQI is calculated. In order to reflect the

¹ Data obtained through email correspondence with ICCC, Tumakuru, as part of an internship.

severity of air pollution and direct public health responses based on the possible health risks associated with each category, this Overall AQI value is then divided into standardized levels, such as Good, Moderate, Unhealthy for Sensitive Groups, Unhealthy, Very Unhealthy, and Hazardous as shown in table 2.

The general formula to calculate AQI for a pollutant is as follows:

$$AQI_p = \frac{(I_{high} - I_{low})}{(C_{high} - C_{low})} \times (C_p - C_{low}) + I_{low}$$

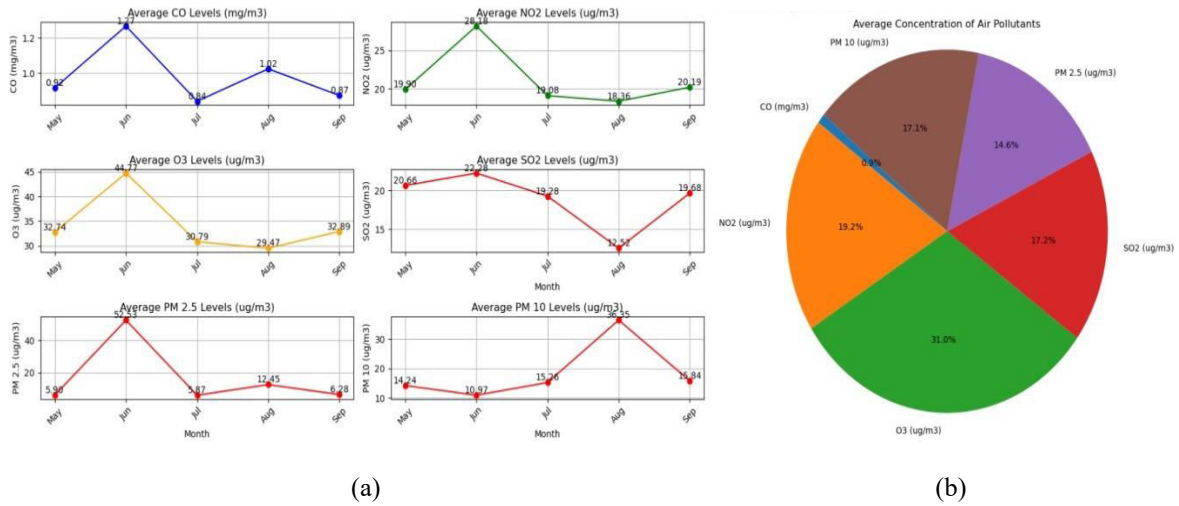
- AQI_p : AQI for a specific pollutant.
- C_p : Observed concentration of the pollutant.
- C_{low} : The breakpoint concentration that is the closest value less than or equal to C_p .
- C_{high} : the breakpoint concentration that is the closest value greater than C_p .
- I_{low} : The AQI value corresponding to C_{low} .
- I_{high} : The AQI value corresponding to C_{high} .

Table 2: Air Quality Index Categorization

AQI Range	Category
0-50	Good
51-100	Moderate
101-150	Unhealthy for sensitive groups
151-200	Unhealthy
201-300	Very Unhealthy

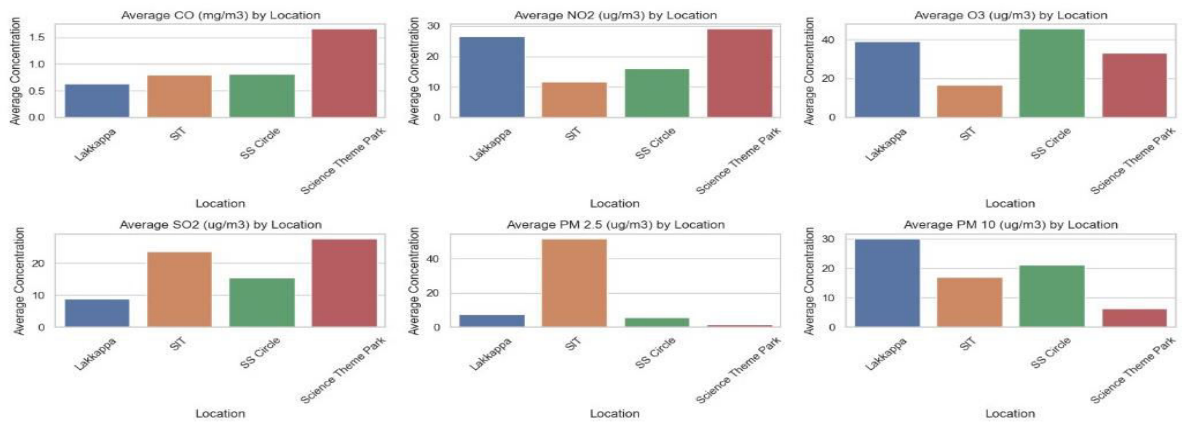
- **Data Visualization:** Data visualizations are graphical depictions of data that facilitate the comprehension and analysis of complex information. Visualizations of data, such as charts, graphs, or maps, highlight trends, patterns, and insights that may be overlooked in unprocessed data. Good visualizations facilitate interactive data exploration, make decision-making easier, and effectively convey findings. They are crucial resources for enabling a broad audience to receive data-driven insights. The figure 2(a) displays average pollutant levels for each month from May to September, capturing shifts in air quality over time. It helps to identify which pollutants are most prevalent in specific months. The figure 2(b) presents the average concentration levels of major air pollutants, offering a clear view of their overall impact on air quality. It highlights the relative intensity of each pollutant, aiding in the identification of primary pollution sources. The figure 2(c) shows the average concentration of pollutants across different locations, allowing for comparison of air quality levels regionally. It helps to identify areas with higher pollutant concentrations and provides insights into localized pollution sources. The figure 3(a) represents heatmap illustrating overall AQI across different locations and times offers a clear visual representation of air quality variations, with colour gradients indicating levels of pollution. The figure 3(b) displays the average AQI for each junction, providing a comparative overview of air quality at different traffic points. It highlights areas with higher pollution levels, aiding in identifying critical zones for air quality management. The figure 3(c) illustrates the average overall AQI by hour and day of the week, revealing patterns in air quality fluctuations throughout the week. It enables the identification of peak pollution times and days, aiding in understanding how traffic and other factors influence air quality on a temporal basis. The figure 3(d) presents the distribution of AQI categories across various locations, showcasing the proportion of each category (e.g., Good, Moderate, Unhealthy) at each site. It highlights areas with frequent high pollution levels and facilitates a better understanding of regional air quality variations, helping inform public health initiatives and environmental policies.

Monthly Average Pollutant Levels from May to September



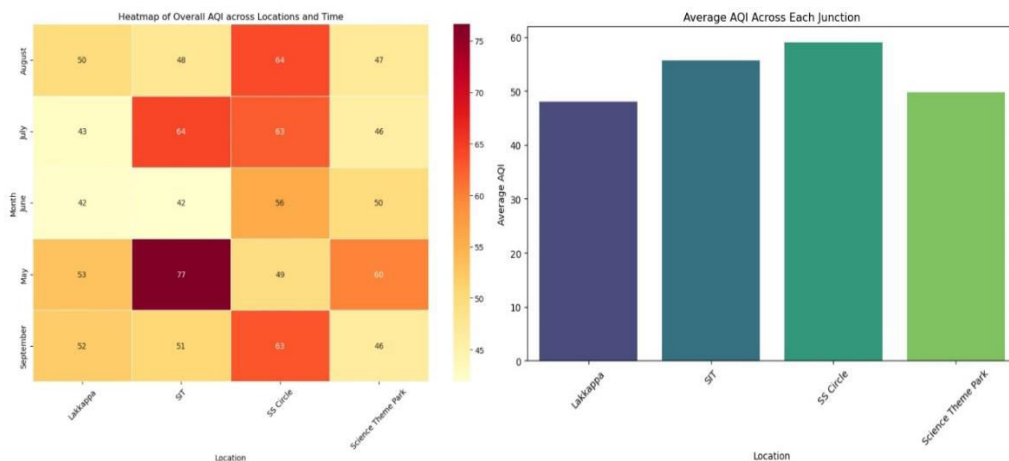
(a)

(b)



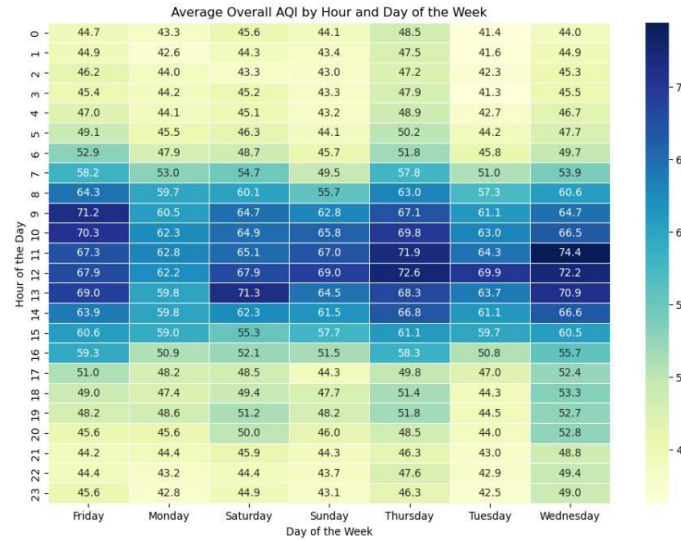
(c)

Figure 2. Pollutant Concentrations.

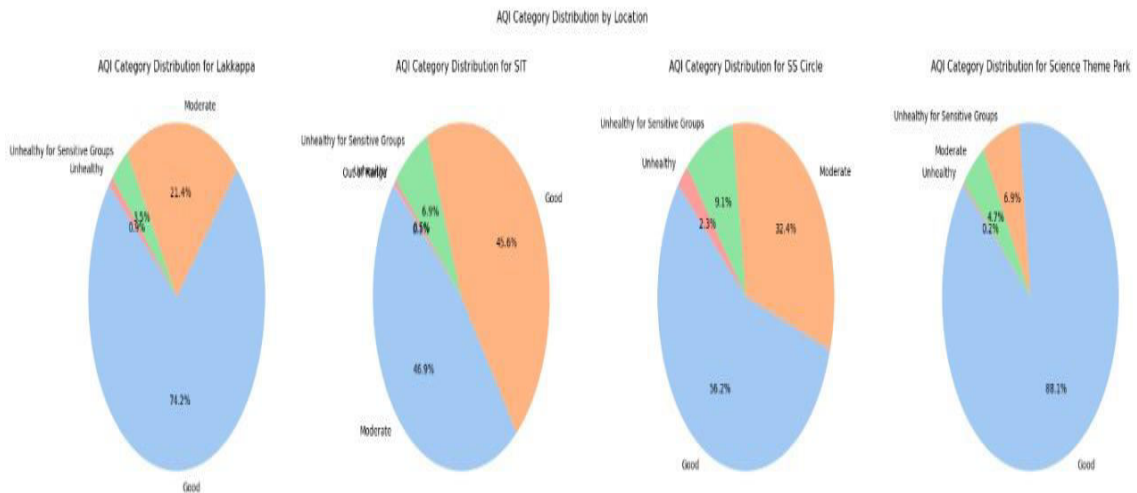


(a)

(b)



(c)



(d)

Figure 3. AQI Values

• **Model Training:**

In the model training phase, the dataset was divided into training and testing sets. This split allowed to train various machine learning models on the training subset while setting aside the test subset for model evaluation. By reserving a portion of the data for testing, the models could be validated on previously unseen data, reducing the risk of overfitting. This approach helped improve the accuracy and reliability of our air quality predictions.

- a) **Random forest:** Random Forest is an ensemble learning method that constructs multiple decision trees on different subsets of the data, where each tree learns from a random sample of the training data and uses a random subset of features. During prediction, each tree in the forest votes, and the majority vote (for classification) or average prediction (for regression) is chosen as the final output. This process helps to reduce overfitting compared to individual decision trees and makes Random Forest a robust model for high-dimensional data. Random Forest is known for its accuracy, resilience to noisy data, and ability to handle both categorical and numerical features,

although it can be computationally expensive on very large datasets and lacks the interpretability of simpler models.

- b) **XGBoost:** Extreme Gradient Boosting, or XGBoost, is a scalable, highly optimized gradient boosting method that prioritizes speed and performance. It creates an ensemble of trees, each of which uses gradient descent to fix the residual mistakes of the ones before it. It is perfect for huge datasets because of its key characteristics, which include parallel processing, handling missing data, regularization to avoid overfitting, and efficient memory consumption. XGBoost also supports both regression and classification and is commonly utilized in data science competitions for its accuracy and efficiency. Model performance can be improved by adjusting hyper parameters like learning rate, maximum depth, and tree count.
- c) **LSTM:** Long Short-Term Memory (LSTM) is a type of recurrent neural network (RNN) specifically designed to model long-term dependencies in sequential data by incorporating memory cells that can retain information for long periods. Unlike traditional RNNs, LSTM networks are equipped with gates (input, forget, and output gates) that control the flow of information, allowing them to overcome the vanishing gradient problem and learn patterns across time steps. This makes LSTMs particularly useful for tasks involving time series, text processing, and other sequential data, where long-term context is important. Despite their complexity and higher training times, LSTMs are widely used for applications like language modelling, speech recognition, and forecasting.
- d) **Bi-LSTM:** Bidirectional Long Short-Term Memory (Bi-LSTM) is an extension of the LSTM model that processes sequential data in both forward and backward directions to capture past and future context simultaneously. By having two separate LSTMs—one reading the sequence from beginning to end and the other from end to beginning—Bi-LSTMs gain a deeper understanding of the data, especially in cases where context from both directions is essential, such as natural language processing tasks. This bidirectional approach helps improve model accuracy in tasks like sentiment analysis, machine translation, and named entity recognition, though it comes with increased computational costs and longer training times.
- e) **RNN:** Recurrent Neural Networks (RNNs) are a class of neural networks designed for processing sequential data by allowing connections between nodes to form cycles, which enables them to maintain a memory of previous inputs. This memory mechanism makes RNNs useful for tasks that involve time series or ordered data, such as language processing and speech recognition. However, standard RNNs face limitations with long sequences due to the vanishing gradient problem, which affects their ability to learn long-term dependencies. Despite this, RNNs serve as the foundation for more advanced architectures like LSTM and GRU, which are better equipped to handle complex sequence-based tasks.

IV. EXPERIMENTAL RESULTS

The evaluation of the air quality dataset from Tumakuru City revealed that the Random Forest model achieved an impressive accuracy of 100%, showcasing its effectiveness in classifying complex air quality indicators. This level of accuracy indicates the model's ability to manage the intricate relationships among various pollutants, making it a robust tool for monitoring air quality. The XGBoost Classifier closely followed, with an accuracy of 99.93%, demonstrating its strength in enhancing predictions through boosting techniques tailored to handle diverse data patterns.

Deep learning models, including Bi-LSTM and LSTM, also showed strong performance, with accuracies of 99.65% and 99.53%, respectively. The marginally higher accuracy of the Bi-LSTM suggests its capability to better capture temporal relationships in the data due to its bidirectional processing approach. The RNN model, while achieving a respectable accuracy of 96.49%, indicated some limitations in effectively managing long-term dependencies, which could impact its utility in air quality analysis.

Overall, the results underscore the effectiveness of both ensemble and deep learning models for predicting air quality metrics. The high accuracies achieved by Random Forest and XGBoost highlight their potential for providing valuable insights into air quality dynamics. These models can play a significant role in informing strategies for air quality management and public health initiatives, contributing to better environmental outcomes.

Table 3: Results

Sl. No	Model	Accuracy
1.	Random Forest	100 %
2.	XGBoost Classifier	99.93 %
3.	RNN	96.49 %
4.	LSTM	99.53 %
5.	Bi-LSTM	99.65%

V. CONCLUSION

In conclusion, this study evaluated the performance of multiple machine learning models for predicting air quality in Tumakuru City, with Random Forest and XGBoost achieving exceptional accuracy rates of 100% and 99.93%, respectively. These results highlight the effectiveness of these models in analysing the intricate relationships among various pollutants and their impact on air quality in the city. The findings provide valuable insights into Tumakuru City's pollution dynamics and indicate that these models can serve as reliable tools for forecasting air quality, supporting urban planning initiatives and public health strategies.

Furthermore, future research could benefit from integrating additional data sources to enhance model accuracy and applicability. By incorporating real-time traffic information, industrial emissions data, meteorological variables, and satellite observations, a more comprehensive understanding of the factors influencing air quality could be achieved. Additionally, exploring the implementation of these models within an automated real-time monitoring system could facilitate timely air quality updates for residents and local authorities through mobile applications or web platforms. Evaluating hybrid models that combine ensemble techniques with advanced deep learning methods may further improve AQI forecasting capabilities. With these advancements, this research could play a pivotal role in improving air quality management in Tumakuru City, equipping policymakers with actionable insights for effective environmental and public health interventions.

REFERENCES

- [1] H. Srivastava, G. K. Sahoo, S. K. Das and P. Singh, "Performance Analysis of Machine Learning Models for Air Pollution Prediction," 2022 International Conference on Smart Generation Computing, Communication and Networking (SMART GENCON), Bangalore, India, 2022, pp. 1-6, doi: 10.1109/SMARTGENCON56628.2022.10084037.
- [2] R, S. C., & S, S. (2024). Unveiling Hidden Air Quality Patterns: A Data-Driven Predictive Regression Analysis for Sustainable Air Management.
- [3] B. Ida Seraphim, E. Nagarajan and S. Pathak, "Predicting Air Quality: Ensemble Learning Approach," 2024 Second International Conference on Data Science and Information System (ICDSIS), Hassan, India, 2024, pp. 1-8, doi: 10.1109/ICDSIS61070.2024.10594204.
- [4] S. Goyal, A. Kaur, V. Sharma, N. Batra and P. Das, "AQMP: Air Quality Monitoring and Prediction using Machine Learning," 2024 2nd International Conference on Disruptive Technologies (ICDT), Greater Noida, India, 2024, pp. 775-777, doi: 10.1109/ICDT61202.2024.10488943.
- [5] R. Gupta, K. Khandal and M. Kandan, "Air Quality Prediction in Smart Cities Using Regression Techniques," 2024 2nd International Conference on Networking and Communications (ICNWC), Chennai, India, 2024, pp. 1-6, doi: 10.1109/ICNWC60771.2024.10537556.

- [6] N. Niveshitha, F. Amsaad and N. Z. Jhanjhi, "Air Quality Prediction in Smart Cities Using Cloud Machine Learning," 2023 Second International Conference on Smart Technologies for Smart Nation (SmartTechCon), Singapore, Singapore, 2023, pp. 1115-1119, doi: 10.1109/SmartTechCon57526.2023.10391685.
- [7] S. Singh, V. Kumar, Z. Ahmed and K. Mittal, "Delhi Air Pollution Prediction: A Comparative Analysis using Time Series Forecasting," 2023 International Conference on Disruptive Technologies (ICDT), Greater Noida, India, 2023, pp. 604-608, doi: 10.1109/ICDT57929.2023.10151445.
- [8] Kumar, R., Sharma, B., Singhal, S., Dhaou, I. B., & Shekhar, S. (2023, January 23). Real Time Prediction Model for Air Pollution and Air Quality Index based on Machine Learning. <https://doi.org/10.1109/icaisc56366.2023.10085379>



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  ijirset@gmail.com



www.ijirset.com

Scan to save the contact details