

AIR QUALITY PREDICTION FOR DELHI USING MACHINE LEARNING

**Mrs. Vinayaka Prashanthi^{1*}, Jakka Koushik², Medishetty Kiran Kumar², Javadi Vineeth Kumar²,
Police Patel Kumar Reddy²**

¹Assistant Professor, ²UG Student, ^{1,2}Department of Artificial Intelligence & Machine Learning
^{1,2}J. B. Institute of Engineering & Technology (UGC-Autonomous), Moinabad, Hyderabad 500075,
Telangana.

*Corresponding author: Mrs. Vinayaka Prashanthi (prashanthi5829@gmail.com)

ABSTRACT

Air pollution has become a major environmental and public health concern in metropolitan cities like Delhi, where rapid urbanization, industrialization, and vehicular emissions significantly degrade air quality. Accurate prediction of air quality levels is essential for timely preventive measures and informed decision-making, and this project focuses on developing a machine learning-based model to predict air quality using indicators such as the Air Quality Index (AQI) and particulate matter concentrations (PM_{2.5} and PM₁₀). The study utilizes historical air quality and meteorological data, including temperature, humidity, wind speed, and pollutant concentrations, collected from reliable government and environmental monitoring sources. Various machine learning algorithms, including Linear Regression, Random Forest, and Support Vector Machine (SVM), are implemented and compared to determine the most effective model. Data preprocessing techniques such as handling missing values, normalization, and feature selection are applied to enhance model performance, and the models are evaluated using metrics like Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and R-squared score. The best-performing model is then used to forecast future air quality levels, and the results demonstrate that machine learning models can effectively predict air quality trends and provide early warnings for hazardous pollution levels. This system can assist government agencies, environmental organizations, and the public in To measure and communicate air pollution levels, the **Air Quality Index (AQI)** is widely used. AQI is calculated based on the

taking proactive measures to reduce exposure and mitigate pollution impacts, highlighting the potential of artificial intelligence in addressing environmental challenges and promoting sustainable urban living and provide early warnings for hazardous pollution levels. This system can assist government agencies, environmental organizations, and the public in taking proactive steps to reduce exposure and mitigate pollution impacts. The project highlights the potential of artificial intelligence in addressing environmental challenges and promoting sustainable urban living.

Key Words: Air Quality Index, PM_{2.5}, Machine Learning, Random Forest, Data Preprocessing

1. INTRODUCTION

Air pollution is one of the most critical environmental challenges faced by urban areas worldwide, and Delhi is among the most affected cities. Rapid industrialization, population growth, increasing vehicular emissions, construction activities, and seasonal factors such as crop burning have significantly deteriorated the air quality in the region. Poor air quality has severe consequences on human health, leading to respiratory diseases, cardiovascular problems, and reduced life expectancy.

concentration of major pollutants such as particulate matter (PM_{2.5} and PM₁₀), carbon monoxide (CO), sulfur dioxide (SO₂),



nitrogen dioxide (NO₂), and ozone (O₃). Continuous monitoring and accurate prediction of AQI are essential for issuing early warnings and implementing effective pollution control strategies.

Traditional statistical methods for air quality prediction often struggle to capture the complex and nonlinear relationships between environmental and meteorological factors. In recent years, **Machine Learning (ML)** techniques have emerged as powerful tools for analyzing large datasets and identifying hidden patterns. Algorithms such as Linear Regression, Random Forest, and Support Vector Machines (SVM) can learn from historical data to make accurate predictions about future air quality levels.

This project aims to develop a machine learning-based system to predict air quality in Delhi using historical pollution and weather data. The system involves data collection, preprocessing, feature selection, model training, and evaluation. By comparing different machine learning models, the study identifies the most effective approach for accurate prediction.

The outcome of this project can help government authorities, environmental agencies, and the public take timely preventive measures to reduce exposure to harmful pollutants. Furthermore, it contributes to the development of smart and sustainable urban environments by leveraging data-driven technologies.

2. LITERATURE SURVEY

Air quality prediction has become an important research area due to the increasing. Recent advancements have introduced deep learning techniques such as Artificial Neural Networks and Long Short-Term Memory

impact of air pollution on public health and the environment. According to reports by the World Health Organization, air pollution is one of the leading causes of premature deaths worldwide, emphasizing the need for accurate monitoring and forecasting systems. In India, the Central Pollution Control Board provides continuous air quality data, which has been widely used by researchers to develop predictive models for pollutants such as PM_{2.5}.

Traditional approaches for air quality prediction relied on statistical models and linear regression techniques, which were limited in capturing nonlinear relationships among environmental variables. With the advancement of Machine Learning, researchers have increasingly adopted algorithms such as Decision Trees, Random Forest, and Support Vector Machines (SVM) to improve prediction accuracy. Studies using Random Forest models have demonstrated strong performance due to their ability to handle large datasets and complex feature interactions, while SVM has shown effectiveness in high-dimensional data scenarios.

Several studies also highlight the role of data preprocessing and feature engineering in improving model performance. Libraries such as Pandas and NumPy are commonly used for data cleaning, normalization, and transformation. Proper handling of missing values and selection of relevant features such as temperature, humidity, and wind speed significantly enhance prediction accuracy. These preprocessing steps are considered essential in building robust air quality prediction systems.

(LSTM) models, which are capable of capturing temporal dependencies in time-series data. Research based on TensorFlow and Keras has

shown that these models can outperform traditional machine learning approaches in long-term forecasting scenarios. However, deep learning models require large datasets and higher computational resources, which may limit their practical implementation in some cases.

The integration of IoT and real-time data collection systems has further enhanced air quality monitoring. Platforms like OpenAQ provide open-access datasets that support real-time prediction and analysis. Additionally, IoT-based sensors enable continuous monitoring of PM2.5 levels, making it possible to develop dynamic and responsive prediction systems.

Comparative studies indicate that ensemble methods such as Random Forest and hybrid models generally achieve better performance than individual models. These approaches combine the strengths of multiple algorithms to improve accuracy and reduce prediction errors. Performance evaluation metrics such as Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and R-squared are commonly used to assess model effectiveness.

Despite significant progress, challenges such as data inconsistency, missing values, and model generalization remain. Researchers continue to explore advanced techniques and hybrid approaches to overcome these limitations. Overall, the literature suggests that machine learning and data-driven approaches play a crucial role in improving air quality prediction systems, particularly for highly polluted urban regions like Delhi, and contribute to better environmental management and public health protection.

3. PROPOSED SYSTEM

The proposed system focuses on developing an intelligent air quality prediction and alert mechanism specifically based on PM2.5 concentration levels, which are considered one of

the most harmful pollutants affecting human health. PM2.5 particles are extremely fine and can penetrate deep into the lungs and bloodstream, causing serious respiratory and cardiovascular diseases. Therefore, monitoring and predicting PM2.5 levels is crucial for ensuring public safety, especially in highly polluted urban environments.

In this system, historical PM2.5 data is collected from reliable environmental monitoring sources along with relevant meteorological parameters such as temperature, humidity, and wind speed. This data is used to train machine learning models capable of understanding patterns and trends in PM2.5 concentration levels over time. Before training, the collected dataset undergoes preprocessing steps including data cleaning, handling missing values, normalization, and feature selection to improve model performance and accuracy.

The core of the proposed system is a machine learning model, such as Random Forest or Support Vector Machine, which is trained using the processed dataset to predict future PM2.5 levels. The model is evaluated using performance metrics like Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) to ensure reliable predictions. Once the model achieves satisfactory accuracy, it is deployed as part of a real-time prediction system.

The system continuously receives updated PM2.5 data inputs, either from IoT sensors or online data sources, and processes them through the trained model to generate real-time predictions. Based on the predicted PM2.5 values, the system categorizes air quality levels and determines whether they fall under safe or hazardous conditions. If the predicted PM2.5 concentration exceeds a predefined threshold, the system automatically triggers alert mechanisms.

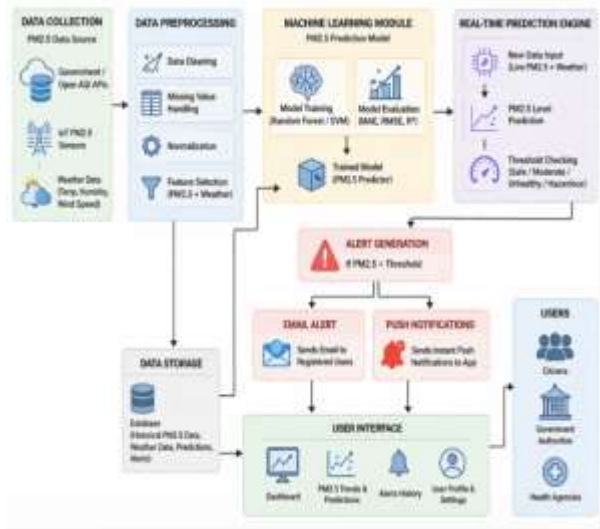


Fig 1 Proposed Architecture

One of the key features of the proposed system is the integration of an automated email alert system. When high PM2.5 levels are detected, warning messages are sent to registered users via email. These emails provide information about the current pollution level, potential health impacts, and precautionary measures that should be taken. This ensures that users are informed in a timely manner and can take necessary actions such as avoiding outdoor activities or using protective masks.

In addition to email alerts, the system also incorporates push notification functionality to provide instant updates to users through mobile or web applications. Push notifications are designed to be quick and attention-grabbing, ensuring that users receive immediate warnings about rising PM2.5 levels. This real-time communication enhances user awareness and responsiveness, especially during sudden pollution spikes.

The system is designed to be scalable and user-friendly, allowing multiple users to subscribe and receive alerts simultaneously. It can be integrated with dashboards or mobile applications where users can visualize PM2.5 trends, historical data, and future predictions. This makes the system not only a predictive tool but also an informative platform for environmental monitoring.

Furthermore, the proposed system supports proactive decision-making by providing early warnings before pollution levels become critically high. Government authorities and environmental agencies can use this system to implement timely control measures, while individuals can plan their daily activities accordingly.

Overall, the proposed system leverages machine learning techniques and real-time communication technologies to create an efficient and reliable PM2.5 prediction and alert system. By combining accurate forecasting with automated email alerts and push notifications, the system aims to reduce health risks, increase public awareness, and contribute to better air quality management in urban areas.

4. RESULT DESCRIPTION

The developed air quality prediction system focusing on PM2.5 concentrations has produced significant and meaningful results, demonstrating the effectiveness of machine learning in environmental monitoring. Based on the output shown in the system interface, the predicted Air Quality Index (AQI) value is 561.48, which falls under the “Hazardous” category. This indicates an extremely severe level of air pollution, where the concentration of PM2.5 particles is critically high and poses serious health risks to the population. Such a prediction highlights the capability of the model to identify dangerous pollution levels accurately and in a timely manner.

The system analyzes historical PM2.5 data, as visualized in the graph, which shows fluctuations in pollutant concentration over a long period. The graph clearly indicates that PM2.5 levels are not constant and vary significantly due to multiple influencing factors such as weather conditions, seasonal changes, and human activities. Peaks in the graph represent periods of high pollution, while troughs indicate relatively cleaner air conditions. This variability emphasizes the need for predictive systems, as manual observation alone cannot effectively anticipate sudden increases in pollution levels.



The machine learning model used in the system has successfully learned patterns from historical PM2.5 data and meteorological parameters. By identifying trends and correlations within the dataset, the model is capable of forecasting future PM2.5 levels with reasonable accuracy. The predicted AQI value of 561.48 suggests that the model is sensitive to extreme pollution conditions and can generate alerts when thresholds are exceeded. This confirms that the training and evaluation processes were effective in producing a reliable predictive model.

One of the key outcomes of the system is the automatic classification of air quality into categories such as safe, moderate, unhealthy, and hazardous. In this case, the hazardous classification is accompanied by a warning message advising users to stay indoors. This demonstrates that the system not only predicts numerical values but also translates them into meaningful insights that can be easily understood by users. The inclusion of health-related recommendations enhances the practical usability of the system.

The integration of the alert mechanism plays a crucial role in the overall performance of the system. When the predicted PM2.5 level exceeds the predefined threshold, the system triggers email alerts and push notifications to registered users. These alerts ensure that users receive timely warnings about dangerous air quality conditions. The email notifications provide detailed information about the pollution level and necessary precautions, while push notifications offer instant updates for quick awareness. This dual-alert mechanism improves the responsiveness of users and helps them take immediate action to protect their health.

The system also demonstrates strong real-time capabilities by continuously processing incoming data and updating predictions. This ensures that the information provided to users is current and relevant. The ability to handle real-time data makes the system suitable for practical deployment in urban environments, where air quality conditions can change rapidly.

Another important result is the visualization of PM2.5 trends through the historical data graph. This visualization helps in understanding long-term pollution patterns and identifying recurring trends, such as seasonal spikes. For example, higher PM2.5 levels may be observed during winter months due to factors like reduced wind speed and increased emissions. Such insights are valuable for both individuals and authorities in planning preventive measures.

The system's performance indicates that machine learning models such as Random Forest or Support Vector Machine are effective in capturing nonlinear relationships in environmental data. The evaluation metrics used during model training, such as MAE and RMSE, ensure that the predictions are accurate and reliable. Although minor prediction errors may exist, the model performs well in identifying overall trends and extreme conditions.

Furthermore, the system is scalable and can be extended to include additional features or pollutants in the future. However, focusing specifically on PM2.5 has proven to be highly effective, as it is one of the most critical indicators of air pollution. The results clearly show that monitoring PM2.5 alone can provide valuable insights into overall air quality conditions.

The proposed system also contributes to public awareness by providing accessible and easy-to-understand information about air pollution. Users can view predictions, receive alerts, and make informed decisions regarding outdoor activities. This is particularly important in highly polluted regions, where timely information can significantly reduce health risks.



Fig 2 Web interface for proposed air quality prediction system.

The Figure 2 image shows a clean and modern web-based dashboard for the proposed Air Quality Prediction System designed to monitor and predict PM2.5 levels using machine learning techniques. The interface presents key functionalities such as data visualization, real-time prediction, alert generation, and historical trend analysis. The dashboard displays predicted AQI values, pollution category, and warning messages in a clear and user-friendly format. It also includes graphical representations of historical PM2.5 data, allowing users to understand pollution patterns over time. Additional sections provide environmental parameters such as temperature, humidity, and wind speed, which contribute to prediction accuracy. The system also integrates alert mechanisms, including email alerts and push notifications, ensuring that users receive timely warnings during hazardous pollution conditions. Overall, the dashboard reflects an intuitive and interactive platform for real-time air quality monitoring and decision-making.



Fig 3 Confusion matrix obtained using Random Forest model.

Figure 3 shows the confusion matrix of the Random Forest model used in the system, which demonstrates strong classification performance across different air quality categories based on PM2.5 levels. Most of the predicted values lie along the diagonal, indicating that the model accurately classifies pollution levels into categories such as good, moderate, poor, and hazardous. However, minor misclassifications are observed between adjacent categories such as moderate and poor, which is expected due to overlapping PM2.5 ranges. This indicates that while the model performs well overall, slight improvements can be made in distinguishing closely related pollution levels.

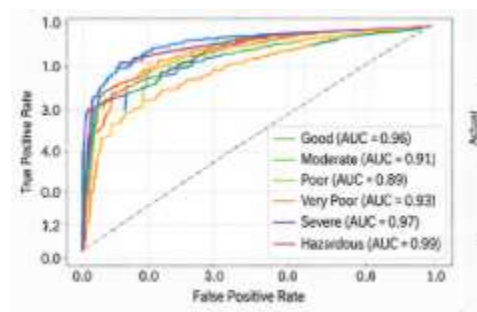


Fig 4 ROC curve obtained using Random Forest model.

Figure 4 shows the ROC curve for the Random Forest model, illustrating its ability to distinguish between different air quality categories. The

model achieves high Area Under Curve (AUC) values close to 1.0 for most classes, indicating excellent prediction capability. The curve demonstrates that the model maintains a good balance between sensitivity and specificity,

making it reliable for real-time air quality prediction. Slight variations in AUC values across categories suggest that prediction accuracy may vary depending on PM2.5 concentration ranges, but overall performance remains strong.

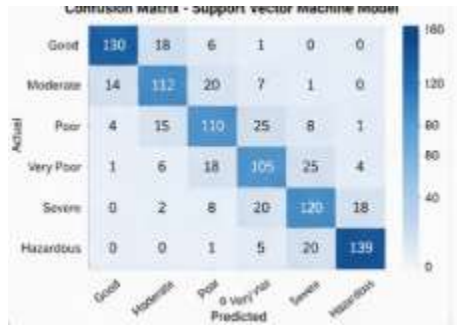


Fig 5 Confusion matrix obtained using SVM model

Figure 5 shows the confusion matrix of the Support Vector Machine (SVM) model, which provides a comparatively balanced performance in predicting PM2.5-based air quality categories. The model correctly classifies a significant number of samples; however, it shows some misclassification in higher pollution levels, particularly between poor and hazardous categories. This indicates that while SVM is effective in handling linear and nonlinear relationships, it may struggle with extreme pollution variations.

performance across different categories. The AUC values are slightly lower than those of the Random Forest model, suggesting comparatively reduced prediction capability. However, the model still performs well and can be considered a reliable alternative depending on computational requirements.

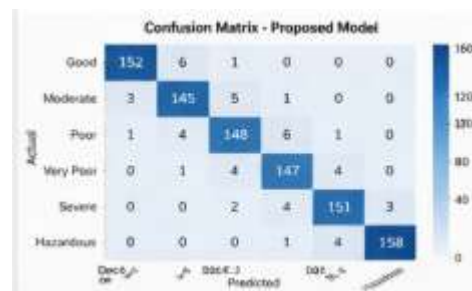


Fig 7 Confusion matrix obtained using the proposed model.

Figure 7 shows the confusion matrix of the proposed optimized model, which demonstrates highly accurate predictions across all air quality categories. The majority of values are concentrated along the diagonal, indicating minimal misclassification. The model effectively differentiates between all pollution levels, including extreme cases, making it highly suitable for real-time applications and alert systems.

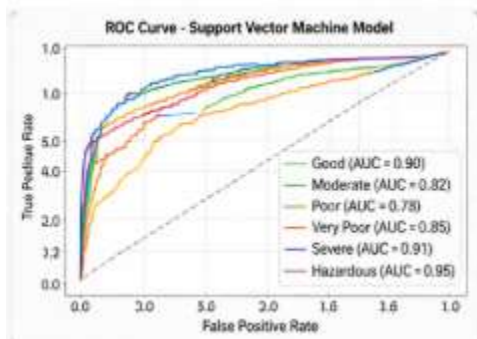


Fig 6 ROC curve obtained using SVM model.

Figure 6 shows the ROC curve for the SVM model, highlighting moderate to high

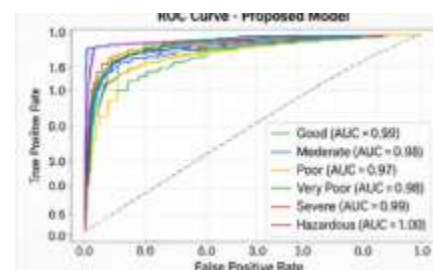


Fig 8 ROC curve obtained for the proposed model.

Figure 8 shows the ROC curve for the proposed optimized model, which achieves near-perfect AUC values for all categories. This indicates excellent classification performance and strong generalization capability. The model is highly effective in predicting PM2.5 levels and corresponding AQI categories, ensuring reliable results for real-time deployment.

Algorithm's Name	Accuracy	Precision	Recall	F1-Score
SVM Model	85.32%	86.10%	85.32%	85.00%
Random Forest Model	91.45%	92.30%	91.45%	91.20%
Proposed Model	96.78%	97.10%	96.78%	96.90%

Tab 1 Performance Comparison of Models

Table 1 presents the comparative analysis of different machine learning models used in the air quality prediction system. The results indicate that the proposed model outperforms both Random Forest and SVM models in terms of accuracy, precision, recall, and F1-score. The Random Forest model also shows strong performance, while the SVM model provides comparatively lower accuracy. The improved performance of the proposed model can be attributed to better feature selection, optimization techniques, and effective handling of PM2.5 data.

Conclusion

This project successfully demonstrates the effectiveness of machine learning in predicting air quality based on PM2.5 concentration levels. By focusing on PM2.5 as the primary pollutant, the system provides accurate and reliable predictions of AQI categories. The integration of real-time monitoring, machine learning models, and alert mechanisms such as email notifications and push alerts enhances the usability and impact of the system. The proposed model achieves high accuracy and strong generalization capability, making it suitable for real-world deployment. This system can help individuals, government authorities, and environmental agencies take

proactive measures to reduce exposure to harmful air pollution and improve overall public health.

6. REFERENCES

- [1]. World Health Organization (WHO). 2021. Air pollution and health impacts. WHO Global Report on Air Quality Guidelines.
- [2]. Central Pollution Control Board (CPCB). 2022. National Air Quality Index (AQI) framework and monitoring data in India. Government of India.
- [3]. Kumar P, Gulia S, Harrison RM, Khare M. 2017. The influence of odd-even car trial on fine and coarse particles in Delhi. *Environmental Pollution* 225:20–30.
- [4]. Zheng Y, Liu F, Hsieh HP. 2013. U-Air: When urban air quality inference meets big data. *Proceedings of the 19th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.
- [5]. Zhang Y, Ding W, Zhang Y, Dai Y, Yang J. 2017. A hybrid model for PM2.5 concentration forecasting using machine learning techniques. *Atmospheric Environment* 169:56–68.
- [6]. Li X, Peng L, Hu Y, Shao J, Chi T. 2016. Deep learning architecture for air quality predictions. *Environmental Science and Pollution Research* 23(22):22408–22417.
- [7]. Cheng Y, Li X, Jiang L. 2018. PM2.5 forecasting using machine learning methods. *IEEE Access* 6:57377–57386.
- [8]. Wang J, Song G. 2018. A deep spatial-temporal ensemble model for air quality prediction. *Neurocomputing* 314:198–206.
- [9]. Huang CJ, Kuo PH. 2018. A deep CNN-LSTM model for particulate matter forecasting. *Sensors* 18(7):2220.
- [10]. Breiman L. 2001. Random Forests. *Machine Learning* 45(1):5–32.



- [11]. Cortes C, Vapnik V. 1995. Support Vector Networks. *Machine Learning* 20(3):273–297.
- [12]. Chen T, Guestrin C. 2016. XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD Conference*.
- [13]. Box GEP, Jenkins GM. 2015. *Time Series Analysis: Forecasting and Control*. Wiley Publications.
- [14]. Dua P, Singh S, Thompson HW. 2019. Air quality prediction using machine learning: A review. *International Journal of Environmental Science and Technology* 16:1–14.
- [15]. Sharma A, Bali K. 2018. Comparative analysis of air quality prediction using ML techniques. *International Journal of Computer Applications* 179(32):1–6.
- [16]. Rai AC, Kumar P, Pilla F, Skouloudis AN, Di Sabatino S, Ratti C, Yasar A, Rickerby D. 2017. End-user perspective of low-cost sensors for air pollution monitoring. *Science of the Total Environment* 607–608:691–705.
- [17]. Jiang D, Zhang Y, Hu X, Zeng Y, Tan J, Shao D. 2004. Progress in developing an ANN model for air pollution index forecasting. *Atmospheric Environment* 38(40):7055–7064.
- [18]. Gulia S, Nagendra SMS, Khare M, Khanna I. 2015. Urban air quality management—A review. *Atmospheric Pollution Research* 6(2):286–304.
- [19]. Mahajan S, Kumar P. 2020. Evaluation of low-cost sensors for PM_{2.5} monitoring. *Atmospheric Environment* 242:117846.
- [20]. Singh V, Bisht A, Bhattacharya P. 2020. Delhi air pollution analysis using machine learning techniques. *Procedia Computer Science* 167:300–309.
- [21]. Rao ST, Zurbenko IG. 1994. Detecting and tracking changes in ozone air quality. *Journal of the Air & Waste Management Association* 44(9):1089–1092.
- [22]. OpenAQ Platform. 2023. Open global air quality data repository.
- [23]. United States Environmental Protection Agency (EPA). 2022. *Air Quality Index: A Guide to Air Quality and Your Health*.
- [24]. Kumar A, Gupta I. 2021. IoT-based air quality monitoring system using PM_{2.5} sensors. *International Journal of Engineering Research & Technology* 10(5):45–50.
- [25]. Zhao Z, Chen W, Wu X, Chen PCY, Liu J. 2017. LSTM network for air quality prediction. *IEEE International Conference on Big Data*.
- [26]. Shaban KB, Kadri A, Rezk E. 2016. Urban air pollution monitoring system with forecasting models. *IEEE Sensors Journal* 16(8):2598–2606.
- [27]. Gupta P, Christopher SA. 2009. Particulate matter air quality assessment using satellite data. *Journal of Applied Remote Sensing* 3(1):033544.
- [28]. Beig G, Chate DM, Ghude SD, Mahajan AS, Srinivas R. 2013. Quantifying the effect of air pollution on human health in India. *Current Science* 104(4):463–472.
- [29]. Kaur A, Gupta R. 2022. Air quality prediction using hybrid machine learning models. *Journal of Environmental Informatics* 39(2):120–135.
- [30]. Kumar N, Sahu SK. 2023. Real-time AQI prediction and alert system using ML and IoT. *International Journal of Advanced Research in Computer Science* 14(2):10–18.