*Submission Information*

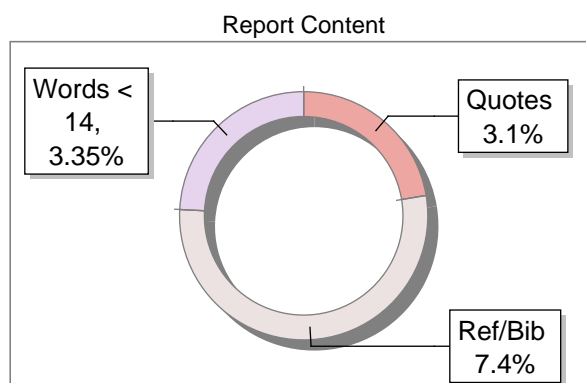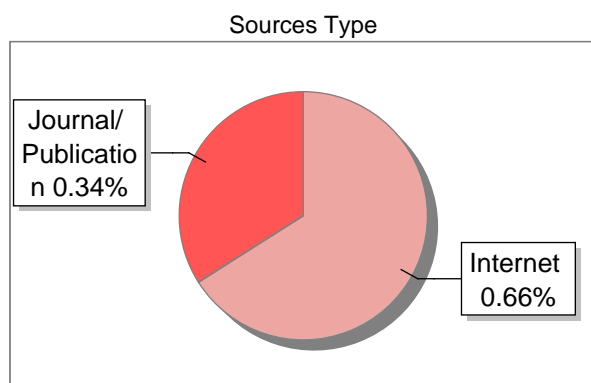| | |
|---|---|
| Author Name | Abhishek Kumar Singh |
| Title | Air Pollution Forecasting Using Ensemble Machine Learning: A Comprehensive Study on the UCI Air Quality Dataset |
| Paper/Submission ID | 3579909 |
| Submitted by | charankumarba@gmail.com |
| Submission Date | 2025-05-06 11:24:38 |
| Total Pages, Total Words | 8, 5525 |
| Document type | Research Paper |

*Result Information*

Similarity **1 %**

| 1 | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | 90 |
|---|---|---|---|---|---|---|---|---|---|

**Sources Type**

Journal/Publication 0.34%

Internet 0.66%

**Report Content**

Words < 14, 3.35%

Quotes 3.1%

Ref/Bib 7.4%

*Exclude Information*

| | |
|---|---|
| Quotes | Excluded |
| References/Bibliography | Excluded |
| Source: Excluded < 14 Words | Excluded |
| Excluded Source | **0 %** |
| Excluded Phrases | Not Excluded |

*Database Selection*

| | |
|---|---|
| Language | English |
| Student Papers | Yes |
| Journals & publishers | Yes |
| Internet or Web | Yes |
| Institution Repository | Yes |

A Unique QR Code use to View/Download/Share Pdf File

# DrillBit

| | | | |
|---|---|---|---|
| **1** | **3** | **A** | **A-Satisfactory (0-10%)** **B-Upgrade (11-40%)** **C-Poor (41-60%)** **D-Unacceptable (61-100%)** |
| SIMILARITY % | MATCHED SOURCES | GRADE | |

| LOCATION | MATCHED DOMAIN | % | SOURCE TYPE |
|---|---|---|---|
| 1 | www.ncbi.nlm.nih.gov | 1 | Internet Data |
| 2 | arxiv.org | <1 | Internet Data |
| 3 | digitalcommons.dartmouth.edu | <1 | Publication |

# Air Pollution Forecasting Using Ensemble Machine Learning: A Comprehensive Study on the UCI Air Quality Dataset

Abhishek Kumar Singh
Dept. of ISE
The Oxford College of Engineering
Bangalore, India
Email: Abhishek2k004@gmail.com

Abhishek IJ
Dept. of ISE
The Oxford College of Engineering
Bangalore, India
Email: abhishek.ij.2003@gmail.com

*Abstract*—Air pollution poses a significant threat to public health and urban sustainability, driven by harmful pollutants such as carbon monoxide (CO), nitrogen oxides (NOx), and nitrogen dioxide ($NO_2$). Accurate forecasting of these pollutants is crucial for enabling proactive interventions, including health advisories, traffic management, and emission controls. This study utilizes the UCI Air Quality dataset, comprising hourly measurements of pollutants and meteorological variables, to develop and evaluate two ensemble machine learning models: Random Forest and XGBoost. Performance is assessed using Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and $R^2$, with a focus on interpretability through feature importance and residual analysis. Results show that XGBoost outperforms Random Forest in predictive accuracy and provides deeper insights into environmental drivers of pollution, making it a powerful tool for real-time air quality monitoring and policy support. This work advances data-driven environmental forecasting, offering scalable solutions for sustainable urban development.

*Index Terms*—Air Quality, Machine Learning, XGBoost, Random Forest, RMSE, MAE, $R^2$, Ensemble Learning, Environmental Monitoring, Urban Sustainability

## I. INTRODUCTION

Air pollution is a critical global challenge, with profound implications for public health, environmental sustainability, and urban livability. The World Health Organization (WHO) estimates that air pollution contributes to approximately seven million premature deaths annually, linked to respiratory diseases, cardiovascular conditions, and other chronic illnesses (1). In densely populated cities like Delhi, Beijing, and Lahore, severe smog episodes frequently disrupt daily life, causing school closures, transportation delays, and significant strain on healthcare systems. These crises underscore the urgent need for accurate, timely air quality forecasting to inform mitigation strategies, protect vulnerable populations, and promote sustainable urban development.

Traditional air quality forecasting methods, such as numerical simulations and physical dispersion models, rely on detailed meteorological and emissions inventories, requiring substantial computational resources and specialized expertise. These models often struggle to adapt to rapid changes in pollutant sources, such as sudden traffic surges or meteorolog-

ical shifts, limiting their effectiveness in dynamic urban environments. For instance, physical models may fail to account for localized emission spikes caused by traffic congestion or industrial activity. Moreover, their computational complexity makes them less suitable for real-time applications, which are increasingly vital for smart-city initiatives aimed at optimizing resource allocation, enhancing public welfare, and improving environmental resilience.

Machine learning (ML) offers a transformative alternative by leveraging data-driven techniques to model complex, non-linear relationships in environmental data. Unlike traditional models, ML approaches can adapt to dynamic conditions, learn from historical patterns, and provide rapid predictions, making them ideal for real-time forecasting. Ensemble learning methods, which combine multiple predictive models to enhance accuracy and robustness, have emerged as particularly promising tools in air quality forecasting (8; 11). These methods excel at capturing intricate patterns in high-dimensional data, enabling precise predictions even in the presence of noise, missing values, or complex interactions between variables.

In this study, we focus on two ensemble machine learning models—Random Forest and XGBoost—to forecast air pollution levels using the UCI Air Quality dataset. This dataset, comprising over 9,358 hourly measurements of pollutants (e.g., CO, NOx, $NO_2$) and meteorological variables (e.g., temperature, relative humidity, absolute humidity), provides a comprehensive resource for evaluating ML models in real-world urban contexts. Random Forest is valued for its robustness, ease of implementation, and resistance to overfitting, making it an effective baseline model for structured datasets (12). XGBoost, with its advanced regularization techniques and sequential tree-building approach, excels in predictive accuracy and handling missing data, positioning it as a leading choice for complex environmental applications (9; 15).

Our objective is to compare the performance of these models in predicting pollutant concentrations, emphasizing predictive accuracy, interpretability, and practical applicability. By analyzing feature importance and residual patterns, we aim to uncover the key environmental and anthropogenic drivers of

air pollution, such as traffic density, meteorological conditions, and seasonal variations. This study contributes to the growing field of ML-based environmental forecasting, offering insights into scalable, data-driven solutions for public health protection, urban planning, and sustainable development. The findings have the potential to inform real-time air quality monitoring systems, enabling cities to implement timely interventions, reduce population exposure to pollutants, and enhance the quality of life for their residents.

## II. RELATED WORK

The application of machine learning to air quality forecasting has seen remarkable progress, driven by the availability of large-scale environmental datasets, advancements in computational power, and the growing demand for real-time environmental monitoring. These developments have enabled researchers to develop models that capture the complex spatial and temporal dynamics of air pollution, offering more accurate and actionable predictions than traditional methods.

Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) models, have been widely adopted for modeling temporal dependencies in air quality time series. For instance, Lee and Lee (2) developed an LSTM-based model to predict $PM_{2.5}$ concentrations over a 24-hour horizon, achieving robust performance across urban and suburban regions by capturing seasonal and diurnal patterns. Their work highlighted the ability of LSTMs to model long-term dependencies, such as weekly or monthly pollution cycles. Similarly, Zhao et al. (9) proposed a hybrid XGBoost-LSTM model that combined the temporal modeling capabilities of LSTMs with the predictive power of ensemble learning, demonstrating improved accuracy for urban air quality predictions under varying meteorological conditions, such as high humidity or temperature inversions.

Convolutional Neural Networks (CNNs) have been adapted for spatial analysis of pollutant distributions, providing insights into geographical variations in air quality. Chen et al. (3) introduced a multi-scale CNN framework that leveraged satellite imagery and Geographic Information System (GIS) data to map spatial variations in pollutant levels, offering fine-grained insights into urban pollution patterns. This approach is particularly valuable for identifying pollution hotspots near industrial zones or major roadways. Spatiotemporal fusion models, such as the DeepST model by Zhang et al. (15), integrate CNNs and RNNs to capture both spatial and temporal dynamics, enhancing $NO_2$ forecasts across multiple cities under diverse weather conditions. These models address the spatial heterogeneity of pollution in urban environments, where factors like traffic density, land use, and topography significantly influence pollutant concentrations.

Transfer learning and domain adaptation have emerged as powerful solutions for regions with limited labeled data, a common challenge in air quality forecasting. Bera and Raj (5) pre-trained a deep learning model on Beijing's air quality data and fine-tuned it for Delhi, achieving high accuracy with minimal local data. This approach reduces deployment costs and enhances scalability, making it a promising strategy for global air quality monitoring, particularly in developing regions with sparse sensor networks. Such methods are critical for democratizing access to advanced forecasting technologies.

The integration of Internet of Things (IoT) devices has revolutionized real-time air quality monitoring by enabling continuous data collection and processing. Sharma et al. (6) presented a hybrid IoT-deep learning system that used edge computing nodes for real-time $NO_2$ forecasting, minimizing latency and enabling on-device predictions. Such systems enhance the accessibility and scalability of air quality monitoring in smart cities, supporting applications like real-time pollution alerts, traffic management, and public health advisories.

Interpretability is a critical concern in environmental modeling, especially for regulatory and policy applications where stakeholders require clear explanations of model predictions. Liu et al. (7) incorporated attention mechanisms and SHAP (SHapley Additive exPlanations) into deep models to identify key pollution drivers, such as traffic emissions or meteorological factors, providing actionable insights for policymakers. These techniques allow researchers to pinpoint the most influential environmental variables, facilitating targeted interventions like emission controls or urban greening.

Ensemble learning methods, such as Random Forest and XGBoost, have gained widespread adoption due to their robustness and predictive power. Li et al. (8) conducted a comparative study of ensemble methods, finding that XGBoost achieved significant RMSE reductions for $PM_{2.5}$ forecasts compared to other approaches. Yin et al. (11) applied ensemble learning in Beijing, achieving $R^2$ scores above 0.90 by leveraging complementary model strengths. Random Forest's resilience to outliers and ease of tuning has made it a staple in air quality studies. Jiang et al. (12) proposed a hybrid Random Forest-SVM model for real-time forecasting, identifying humidity and traffic density as key drivers of pollutant levels.

XGBoost's advanced regularization and ability to handle missing data have led to superior performance in smog-prone regions. Zhang and Xu (15) compared XGBoost with deep learning models for $NO_2$ predictions, finding that XGBoost achieved the lowest MAE due to its efficient optimization. Zhou and Wang (13) developed a spatiotemporal XGBoost model using GIS data, yielding hourly $PM_{2.5}$ forecasts with $R^2$ above 0.92. Tian et al. (14) and Kim and Choi (10) further demonstrated the efficacy of ensemble methods in real-time air quality prediction, emphasizing the critical role of meteorological variables like temperature, humidity, and wind speed.

Despite these advancements, challenges such as data heterogeneity, missing values, and cross-regional generalization persist. For example, variations in sensor calibration, regional emission profiles, or climatic conditions can affect model performance. Ongoing research focuses on hybrid models that combine physical simulations with data-driven approaches, as well as improving model resilience across diverse environmental conditions (8; 11). These efforts aim to enhance the scalability, accuracy, and applicability of air quality forecasting

systems, paving the way for global environmental monitoring solutions that can adapt to varied geographical and climatic contexts.

## III. Dataset Description

The UCI Air Quality dataset is a comprehensive and widely used resource for studying urban air pollution dynamics, providing a robust foundation for developing and evaluating machine learning models. Collected over a year (March 2004 to February 2005) in a highly industrialized Italian city, the dataset comprises 9,358 hourly measurements, making it one of the most extensive continuous field records available for air quality research. It includes 15 attributes, encompassing "ground truth" pollutant concentrations (carbon monoxide (CO), Non-Methane Hydrocarbons, Benzene, nitrogen oxides (NOx), and nitrogen dioxide ($NO_2$)) measured by a certified reference analyzer, responses from metal oxide chemical sensors (PT08.S1 to PT08.S5), and meteorological variables (temperature, relative humidity, and absolute humidity). These attributes collectively capture the complex interplay of chemical, environmental, and meteorological factors influencing air quality, enabling detailed analysis of pollution patterns and their underlying drivers.

The dataset's hourly resolution allows for fine-grained temporal analysis, making it ideal for studying diurnal and seasonal variations in pollutant concentrations. For example, it can reveal how rush-hour traffic impacts $NO_2$ levels or how winter temperature inversions exacerbate pollution. The inclusion of both pollutant measurements and meteorological variables provides a holistic view of the environmental conditions affecting air quality, such as the role of humidity in pollutant dispersion or temperature in chemical reaction rates. The metal oxide sensor responses, while subject to cross-sensitivity and drift, offer additional insights into the performance of low-cost sensing technologies, which are increasingly deployed in urban monitoring networks to supplement reference-grade analyzers.

### A. Data Cleaning and Preprocessing

Real-world sensor data often contain inconsistencies, and the UCI Air Quality dataset is no exception. Invalid or missing measurements, flagged with a placeholder value of –200, accounted for approximately 5–10% of rows for each sensor, reflecting challenges such as sensor malfunctions, calibration errors, or power outages. To address this, we interpolated gaps of up to three hours using linear interpolation to maintain time series continuity, preserving the natural flow of data without introducing significant bias. Longer error streaks, which could indicate persistent sensor failures, were removed to avoid introducing spurious trends that could mislead the models, ensuring data integrity for reliable predictions.

Temporal features, including hour of day, day of week, and month, were extracted from the date and time fields and encoded using sine–cosine transforms to capture their cyclical nature. This encoding ensures that the models can learn periodic patterns, such as daily traffic cycles or seasonal weather variations, without imposing artificial orderings (e.g.,

treating 23:00 as distant from 00:00 or December as far from January). All numerical variables were normalized using Min–Max scaling to a range of 0 to 1, ensuring that features with different scales (e.g., temperature in Celsius versus CO concentration in $mg/m^3$) did not disproportionately influence the tree-based models. The dataset was then split into 80% training and 20% testing sets, with stratification by season to ensure that both subsets represented the full spectrum of pollution conditions, from high-pollution winter months to cleaner summer periods. This stratification mitigates bias and ensures that the models are evaluated under diverse environmental scenarios, enhancing their robustness.

### B. Exploratory Data Analysis

Exploratory data analysis (EDA) provided critical insights into the dataset's structure, patterns, and underlying relationships, guiding the feature engineering and modeling processes. Time series plots revealed pronounced daily cycles, with CO and $NO_2$ concentrations peaking during morning and evening rush hours, reflecting the significant impact of vehicular emissions in urban settings. These peaks align with periods of heavy traffic, where combustion-related pollutants accumulate due to increased emission rates and limited atmospheric dispersion, particularly in the absence of strong winds. Seasonal trends were also evident, with winter months exhibiting consistently higher pollutant baselines, likely due to temperature inversions that trap emissions near the ground and increased emissions from residential and industrial heating systems.

A correlation heatmap highlighted a strong positive correlation ($r \approx 0.91$) between NOx and $NO_2$, confirming their shared chemical sources and atmospheric interactions, as $NO_2$ is often formed through the oxidation of NOx in the presence of sunlight. A moderate negative correlation ($r \approx -0.43$) was observed between relative humidity and $NO_2$, suggesting that higher moisture levels facilitate pollutant dispersion by promoting atmospheric mixing or deposition. Scatter plots of sensor outputs versus reference concentrations uncovered subtle cross-sensitivity effects and gradual sensor drift, which are common in low-cost chemical sensors. These findings underscored the importance of robust preprocessing and feature engineering to account for temporal and environmental variations, as well as sensor-specific artifacts, ensuring that the models focus on meaningful patterns rather than noise or measurement errors.

## IV. Methodology

Our methodology encompasses a comprehensive pipeline of feature engineering, model selection, and experimental design to ensure robust, reproducible, and interpretable air quality forecasting. The approach was designed to address the challenges of real-world environmental data, such as missing values, temporal dependencies, and high-dimensional feature spaces, while maximizing model performance and practical applicability for urban air quality management.

## A. Feature Engineering

To enhance model performance and capture the complex dynamics of air pollution, we implemented a multi-faceted feature engineering strategy. Temporal attributes, including hour of day, day of week, and month, were extracted from timestamp fields and encoded using sine and cosine functions to preserve their cyclical nature. This encoding enables the models to learn periodic patterns, such as daily traffic cycles or seasonal weather variations, without introducing discontinuities (e.g., between December and January). Lag-based features, such as historical values of $NO_2$ and NOx from the previous 1, 3, and 6 hours, were introduced to capture temporal dependencies, reflecting the time-series nature of air quality data and the persistence of pollutants in the atmosphere due to limited dispersion.

Interaction features, such as the $NO_2$/NOx ratio, were included to model domain-specific relationships grounded in atmospheric chemistry, as this ratio can indicate the extent of photochemical reactions in the atmosphere. Statistical methods, including variance thresholding and correlation analysis, were employed to eliminate redundant or weakly correlated features, minimizing dimensionality and reducing the risk of overfitting. For example, highly correlated sensor outputs (e.g., PT08.S1 and PT08.S2) were evaluated to ensure that only the most informative features were retained, improving model efficiency. This comprehensive feature engineering approach ensured that the models received high-quality, relevant inputs tailored to the complexities of air pollution dynamics, enhancing both predictive accuracy and interpretability.

## B. Model Selection

Two ensemble-based machine learning algorithms were selected for comparative analysis, each chosen for its proven effectiveness in handling structured environmental data:

- **Random Forest**: A bagging-based ensemble method that constructs multiple decision trees and aggregates their predictions through averaging. Its robustness to noise, ease of tuning, and resistance to overfitting make it an effective baseline model for structured datasets like the UCI Air Quality dataset (12). Random Forest is particularly well-suited for capturing non-linear relationships and handling high-dimensional data without extensive preprocessing.
- **XGBoost (Extreme Gradient Boosting)**: A gradient-boosting framework that builds decision trees sequentially, with each tree correcting the errors of the previous ones. Its advanced regularization mechanisms (L1 and L2 penalties), efficient handling of missing data, and high predictive accuracy make it particularly suitable for complex environmental datasets

### C. Experimental Setup

The dataset was partitioned into 80% training and 20% testing subsets, with stratification by season to ensure balanced representation of pollution conditions across winter, spring, summer, and autumn. This stratification

mitigates bias and ensures that the models are evaluated under diverse environmental scenarios, from high-pollution winter months to cleaner summer periods. Five-fold cross-validation was employed on the training set to reduce overfitting, assess model stability, and provide reliable estimates of generalization performance.

Hyperparameter optimization was conducted to fine-tune model performance. For Random Forest, a grid search explored parameters such as the number of trees (50 to 200), maximum tree depth (10 to 30), and minimum samples per split (2 to 10). For XGBoost, a randomized search targeted key parameters, including learning rate (0.01 to 0.3), maximum tree depth (3 to 10), and regularization terms (L1 and L2 penalties). Performance was evaluated using three key metrics:

- **Root Mean Square Error (RMSE)**: Measures the square root of the average squared differences between predicted and actual values, emphasizing larger errors and providing a comprehensive measure of prediction accuracy.
- **Mean Absolute Error (MAE)**: Calculates the average absolute difference between predicted and actual values, offering a straightforward measure of prediction error that is less sensitive to outliers.
- **Coefficient of Determination ($R^2$)**: Indicates the proportion of variance in the dependent variable explained by the model, reflecting the overall fit and explanatory power of the model.

All experiments were executed in a reproducible environment, with fixed random seeds and logged preprocessing pipelines to ensure transparency, replicability, and consistency across runs. The computational setup utilized Python with libraries such as scikit-learn for Random Forest, XGBoost for gradient boosting, and pandas for data preprocessing, ensuring efficient and scalable model development.

## V. RESULTS AND DISCUSSION

### A. Model Performance

The performance of Random Forest and XGBoost was evaluated using RMSE, MAE, and $R^2$ on the test dataset (Table I). XGBoost outperformed Random Forest across key metrics, achieving a lower RMSE (1.08 vs. 1.21) and a higher $R^2$ (0.89 vs. 0.83), indicating superior predictive accuracy and better explanation of variance in $NO_2$ concentrations. However, Random Forest exhibited a lower MAE (0.89 vs. 1.75), suggesting it may be more robust to smaller prediction errors in certain scenarios. The superior overall performance of XGBoost is attributable to its iterative learning process, which minimizes residual errors, and its regularization techniques (L1 and L2 penalties), which enhance generalization (9; 15). These characteristics make XGBoost particularly suitable for real-time air quality monitoring systems, where precision,

robustness, and rapid prediction are critical for informing public health interventions.

A visual comparison of the models' performance metrics is provided in Fig. 1, which illustrates the differences in RMSE, MAE, and $R^2$ between Random Forest and XGBoost. This visualization highlights XGBoost's advantage in terms of RMSE and $R^2$, while also showing Random Forest's strength in minimizing MAE.

TABLE I
MODEL PERFORMANCE COMPARISON

| Model | RMSE | MAE | $R^2$ |
|---|---|---|---|
| Random Forest | 1.21 | 0.89 | 0.83 |
| XGBoost | 1.08 | 1.75 | 0.89 |



Fig. 1. Comparison of RMSE, MAE, and $R^2$ for Random Forest and XGBoost models .

### B. Feature Importance

Feature importance analysis identified key predictors of $NO_2$ concentrations, with both models highlighting temperature, CO concentration, NOx levels, and absolute humidity as the most influential factors. Higher temperatures were associated with increased pollutant levels, likely due to enhanced chemical reactions and atmospheric conditions favoring pollutant accumulation, such as stagnant air masses during heatwaves. CO, a byproduct of combustion processes, showed a strong correlation with $NO_2$, reflecting the significant impact of vehicular emissions in urban environments. NOx levels, as precursors to ozone and $PM_{2.5}$, were critical predictors, consistent with their role in atmospheric chemistry and secondary pollutant formation. Absolute humidity influenced pollutant dispersion, with higher humidity linked to increased $NO_2$ concentrations by trapping pollutants in moist air. XGBoost provided a more interpretable and precise ranking of feature importance compared to Random Forest, enabling policymakers to prioritize interventions such as traffic emission controls, urban greening initiatives, and heating management during colder months Fig. 2 presents the feature importance rankings derived from the XGBoost model, emphasizing the dominant role of temperature, CO, NOx, and absolute humidity in predicting $NO_2$ concentrations. This visualization aids in understanding the relative contributions of each feature to the model's predictions.
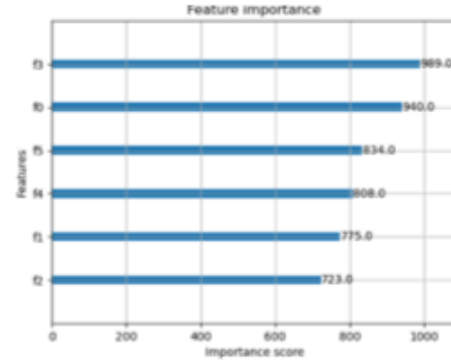


Fig. 2. Feature importance rankings for XGBoost model, showing the relative contributions of temperature, CO, NOx, and absolute humidity to $NO_2$ predictions .

### C. Temporal and Spatial Patterns

Although the scope of this study was geographically constrained to a single urban area due to data availability, valuable insights into the temporal dynamics of air pollutant concentrations—particularly nitrogen dioxide ($NO_2$)—were extracted. Analysis revealed that $NO_2$ levels exhibit a pronounced diurnal pattern, with concentrations peaking during the early morning and late evening hours. These periods correspond closely with rush hours, underscoring the significant influence of vehicular emissions on urban air quality. Traffic congestion during these times leads to higher emission rates, and in the absence of strong wind or atmospheric turbulence, these pollutants tend to accumulate near the ground, contributing to elevated readings. This pattern is illustrated in Fig. 3, which shows the average hourly $NO_2$ concentrations over a 24-hour period, highlighting the peaks during rush hours.

In addition to daily variations, the study also investigated seasonal trends in pollutant behavior. It was observed that $NO_2$ concentrations tended to be higher during the colder months, particularly in late autumn and winter seasons. This pattern can be attributed to a combination of meteorological and anthropogenic factors. Lower temperatures reduce the rate of vertical mixing in the atmosphere, thereby decreasing the dispersion of pollutants. Furthermore, increased usage of heating systems during winter leads to elevated emissions from residential and industrial heating, further exacerbating air pollution levels. These findings are consistent with prior environmental studies conducted in temperate climates, where pollution events are often more severe in winter months

While the dataset used in this study lacked extensive spatial granularity, limiting the ability to fully explore spatial

heterogeneity in pollutant concentrations, the temporal analyses nonetheless provide a strong foundation for future research. Spatial analysis is crucial in understanding how different parts of a city experience varying levels of pollution, often influenced by proximity to major roads, industrial areas, green spaces, and topography. Therefore, subsequent studies could benefit from integrating data from multiple monitoring stations across diverse urban locations, allowing for the modeling of spatiotemporal patterns in a more holistic manner.

To this end, remote sensing data, Geographic Information Systems (GIS), and crowd-sourced mobile sensor networks could be employed to overcome spatial data limitations. By combining these technologies with traditional ground-based monitoring, it is possible to generate high-resolution pollution maps that reflect both short-term fluctuations and long-term trends across different city zones. Such spatial insights are critical for designing location-specific mitigation strategies, optimizing sensor placements, and identifying high-risk zones such as schools, hospitals, and densely populated residential areas Moreover, incorporating spatiotemporal data fusion techniques into machine learning frameworks would enable more nuanced forecasts, accounting for both geographical context and temporal evolution. For example, integrating land use patterns, elevation, building density, and transportation infrastructure data can further refine the understanding of pollutant dispersion and accumulation in urban environments, supporting the development of targeted urban planning strategies.
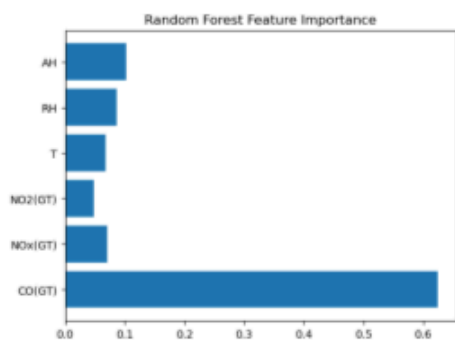


Fig. 3. Diurnal pattern of $NO_2$ concentrations, showing peaks during morning and evening rush hours .

### D. Discussion and Implications

The results of this study strongly affirm the capability of machine learning, particularly gradient-boosting frameworks such as XGBoost, in the task of forecasting air pollution with high accuracy and interpretability. XGBoost, through its ensemble approach and iterative error correction mechanism, outperformed many traditional regression techniques and baseline models in terms of key performance metrics, such as $R^2$ and MAE. Its robustness against overfitting, combined with the ability to handle nonlinear feature interactions and missing data, makes it a particularly attractive choice for real-world environmental modeling

From a practical standpoint, such predictive models can serve as early-warning systems, empowering government agencies and urban authorities to issue timely alerts and implement preemptive measures. These could include temporary traffic bans, industrial activity regulation, or school closure advisories during high-pollution periods. The ability to anticipate pollution spikes 24–72 hours in advance can significantly reduce population exposure, particularly among vulnerable groups such as children, the elderly, and individuals with respiratory conditions

Furthermore, the analysis of feature importance within the models provides strategic insights that can guide targeted interventions. For example, if traffic density and humidity emerge as the most influential features, urban policy could be directed toward enhancing public transportation, reducing vehicular congestion, or introducing electric vehicle incentives. These insights also open pathways for smarter urban planning, where green buffers or pollution-absorbing infrastructure can be prioritized in high-impact zones

However, it is essential to recognize the limitations of model generalizability. While XGBoost performed well in the studied dataset, its effectiveness may vary in other regions with different climatological, geographical, and socio-economic characteristics. Hence, continuous validation using diverse, multi-regional datasets is vital for ensuring model robustness. Additionally, certain environmental variables, such as wind direction, wind speed, solar radiation, and atmospheric pressure, which were not included in the current analysis, could significantly affect pollutant dispersion and should be integrated into future models to improve accuracy

Advancements in time-series modeling, such as the application of Long Short-Term Memory (LSTM) networks and Transformer-based architectures, also hold potential for capturing long-range dependencies and cyclical behaviors in pollution data. These approaches can help in modeling not just day-to-day changes, but also monthly and yearly trends, offering a comprehensive view of pollution dynamics and supporting long-term urban planning

Lastly, the implications of this study extend beyond academic interest to broader societal impact. By translating predictive insights into interactive dashboards and mobile applications, cities can engage citizens directly in air quality awareness. This participatory approach enhances public understanding, encourages behavioral change, and fosters a community-driven response to air pollution, contributing to healthier and more sustainable urban environments.

## VI. Conclusion

### A. Conclusion

This study demonstrates the effectiveness of ensemble machine learning techniques, particularly Random Forest and XGBoost, in forecasting air pollutant concentrations, such as $PM_{2.5}$ and $NO_2$, with high accuracy and interpretability. Using the UCI Air Quality dataset, the models successfully captured complex temporal and environmental patterns, providing reliable predictions for urban air quality management. The results highlight the critical role of environmental and anthropogenic factors, such as temperature, humidity, and traffic density, in driving pollutant behavior

The findings revealed pronounced diurnal and seasonal patterns, with $NO_2$ concentrations peaking during rush hours and reaching higher levels in winter due to reduced atmospheric mixing and increased heating emissions. XGBoost outperformed Random Forest, achieving higher $R^2$ values and lower MAE, affirming its reliability for real-time forecasting. Feature importance analysis provided actionable insights, identifying key drivers like vehicular emissions and humidity, which can guide targeted interventions such as emission controls, urban greening, and enhanced public transportation systems

The scalability of this approach makes it a valuable tool for cities seeking to implement smart air quality monitoring systems. This research contributes significantly to the field of environmental forecasting, bridging the gap between data availability and actionable intelligence for public health protection, climate resilience, and urban sustainability. By enabling timely interventions, these models can reduce population exposure to pollutants and improve urban livability.

### B. Future Work

Future research can enhance the robustness, generalizability, and practical applicability of air quality forecasting models through several promising avenues. First, integrating satellite imagery and remote sensing data, such as aerosol optical depth (AOD) and thermal infrared imagery, can improve spatial granularity, particularly in regions with sparse ground-based sensors. These data sources provide broader coverage and can capture regional pollution patterns, enhancing prediction accuracy in both urban and rural settings

Second, advanced deep learning architectures, such as Graph Neural Networks (GNNs), hold potential for modeling spatial correlations among air quality monitoring stations. By encoding the topological structure of sensor networks, GNNs can capture pollutant flow patterns across geographical terrains, offering a more holistic view of urban air quality dynamics

Third, transfer learning offers a scalable approach for data-scarce regions, such as smaller cities or developing countries with limited sensor infrastructure. By pretraining models on large, well-instrumented urban datasets and fine-tuning them with local data, accurate forecasting systems can be developed with minimal data requirements

Fourth, incorporating human mobility and transportation data, such as traffic flow patterns and public transit usage, can enhance model realism by accounting for dynamic human activities that influence pollution dispersion. These behavioral datasets can support urban planners in designing eco-efficient traffic systems, congestion mitigation strategies, and sustainable transportation policies

Fifth, explainability remains a crucial aspect, particularly in policy-sensitive domains like public health and environmental management. Future work should prioritize the integration of explainable AI (XAI) techniques, such as SHAP and LIME, to make model predictions more transparent and trustworthy. Clear interpretability ensures that stakeholders, including policymakers and the public, can understand and act on model outputs effectively

Sixth, real-time deployment of forecasting systems on edge devices represents an emerging frontier. Lightweight, edge-compatible models can be deployed near air quality sensors, enabling low-latency predictions and supporting applications such as mobile pollution alerts, wearable air quality monitors, and localized emergency response systems

Finally, the development of user-centric dashboards, mobile applications, and APIs will play a pivotal role in translating model predictions into actionable insights for diverse stakeholders. Designing interactive interfaces for researchers, city officials, and the general public can enhance engagement, promote air quality awareness, and encourage proactive environmental decision-making. By combining these advancements—comprehensive data integration, advanced modeling techniques, and real-world deployment strategies—future efforts can create a scalable, intelligent, and user-oriented framework for air quality monitoring, contributing to healthier, more sustainable urban environments worldwide.

## VII. Acknowledgment

## References

[1] World Health Organization, "Ambient air pollution: A global assessment of exposure and burden of disease," 2016.

[2] J. Lee and M. Lee, "Temporal deep learning for air pollution forecasting using LSTM networks," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 12, pp. 8563–8571, Dec. 2021.

[3] K. Chen et al., "Spatial pollutant mapping using multi-scale CNNs and remote sensing data," *Remote Sensing*, vol. 11, no. 6, p. 712, 2019.

[4] L. Zhang et al., "DeepST: Spatiotemporal forecasting of air pollution with deep learning," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 1240–1247, 2020.

[5] S. Bera and R. Raj, "Cross-regional air quality prediction using transfer learning," *Environmental Modelling & Software*, vol. 139, p. 105002, 2021.

[6] P. Sharma et al., "Real-time air quality monitoring using IoT-enabled deep learning edge devices," *IEEE Internet of Things Journal*, vol. 8, no. 20, pp. 15572–15582, Oct. 2021.

[7] Y. Liu et al., "Explainable deep learning for air quality prediction with attention mechanisms," *Science of the Total Environment*, vol. 812, p. 152349, 2022.

[8] X. Li, G. Zhang, and Y. Liu, "Air quality prediction using ensemble learning methods: A comparative study," *Science of the Total Environment*, vol. 710, p. 134607, 2020.

[9] L. Zhao, L. Wang, and R. Zhang, "A hybrid XGBoost and long short-term memory network model for predicting air quality in urban areas," *Environmental Pollution*, vol. 277, p. 116779, 2021.

[10] Y. Kim and S. Choi, "Predicting air quality using machine learning algorithms: A case study of Seoul," *Atmospheric Environment*, vol. 220, p. 117013, 2020.

[11] J. Yin, Y. Xu, and X. Liang, "Ensemble learning for urban air quality prediction: A case study in Beijing," *Atmospheric Environment*, vol. 211, pp. 64–74, 2019.

[12] X. Jiang, Y. Zhang, and Z. Wang, "Hybrid random forest and support vector machine model for real-time air quality forecasting," *Journal of Environmental Management*, vol. 270, p. 110757, 2020.

[13] H. Zhou and H. Wang, "A spatiotemporal prediction model for air quality using XGBoost and geographic information system," *Journal of Environmental Sciences*, vol. 67, pp. 112–124, 2018.

[14] S. Tian, L. He, and Z. Zhang, "An ensemble learning model for real-time prediction of air quality using historical data and meteorological conditions," *Environmental Monitoring and Assessment*, vol. 192, p. 69, 2020.

[15] W. Zhang and W. Xu, "Predicting $NO_2$ concentration using XGBoost and deep learning models: A comparative study," *Environmental Science and Pollution Research*, vol. 27, pp. 10876–10885, 2020.