

# Application of Machine Learning Models for Air Quality Index Prediction

**Karthik N. Rao, Sneha P. Kulkarni, Prof. Meera V. Nair**

Department of Biotechnology,  
Horizon Institute of Scientific Studies, India



<https://doi.org/10.55041/ijst.v1i2.004>

**Cite this Article:** Kulkarni, S. P. & Rao, K. N. (2025). Application of Machine Learning Models for Air Quality Index Prediction. International Journal of Science, Strategic Management and Technology, <i>Volume 01</i><(02)>, 1-9. <https://doi.org/10.55041/ijst.v1i2.004>

**License:**  This article is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting use, distribution, and reproduction in any medium, provided the original author(s) and source are properly credited.

## 1. Abstract

Air pollution is a mounting global environmental challenge with profound health and economic impacts. The Air Quality Index (AQI) serves as a standardized indicator of air pollution severity and is widely used by policymakers, researchers, and the public to understand environmental health risks. Traditional AQI forecasting techniques based on statistical regression and time-series models often struggle with nonlinear dependencies among multiple pollutant sources and meteorological variables. Machine Learning (ML) models offer significant promise for accurately capturing complex patterns in air quality data, enabling robust AQI predictions. This article explores the application of multiple machine learning and deep learning models—including Random Forest, Gradient Boosting (e.g., XGBoost), Support Vector Machines (SVM), Neural Networks, and temporal models such as LSTM—to AQI prediction. Using data from public air quality repositories, we evaluate model performance metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and  $R^2$ . Our findings indicate that ensemble and deep learning models generally outperform shallow models, especially when data preprocessing, feature engineering, and

hyperparameter tuning are properly implemented. We also discuss challenges, potential applications, and future research directions in deploying ML-driven AQI prediction tools for real-time and policy-relevant environmental intelligence.

## 2. Keywords

Air Quality Index (AQI), Machine Learning, Deep Learning, Random Forest, Boost, LSTM, Environmental Prediction, Time Series Forecasting, Feature Engineering

## 3. Introduction

### 3.1 Background

Air quality significantly influences human health, ecosystems, and economic productivity. The Air Quality Index (AQI) is a standard metric that translates concentrations of major air pollutants—such as particulate matter (PM<sub>2.5</sub> and PM<sub>10</sub>), nitrogen dioxide (NO<sub>2</sub>), sulfur dioxide (SO<sub>2</sub>), carbon monoxide (CO), and ground-level ozone (O<sub>3</sub>)—into a single value representing overall air pollution severity. Accurate AQI prediction is critical for early warning systems, urban planning,

public health advisories, and environmental policy decisions.

Traditional statistical methods such as linear regression and ARIMA have been widely used for AQI forecasting. However, these approaches often fail to capture nonlinear relationships among multiple pollutant sources, weather variables, and complex temporal patterns in air quality data. Recent research highlights the potential of machine learning models to overcome these limitations by learning intricate patterns directly from historical data. Machine learning techniques such as support vector machines, random forests, and neural networks have demonstrated superior performance in modeling the complex dependencies inherent in air quality data. These models can integrate diverse data sources, including meteorological factors and traffic patterns, to enhance prediction accuracy. Consequently, adopting machine learning approaches can significantly improve the reliability of AQI forecasts, enabling more effective environmental management and public health protection.

### 3.2 Objective of the Study

This research article aims to:

1. Provide a comprehensive review of machine learning approaches for AQI prediction.
2. Examine the preprocessing methods needed for high-quality prediction models.
3. Compare the performance of several ML algorithms.
4. Highlight challenges and best practices in applying ML for AQI forecasting.

## 4. Materials

### 4.1 Data Sources

Multiple publicly available air quality datasets were utilized, including:

- **Global and city-level AQI datasets** containing pollutant measurements and meteorological data. These datasets provide pollutant concentrations (e.g., PM<sub>2.5</sub>, PM<sub>10</sub>, NO<sub>2</sub>, CO, SO<sub>2</sub>, O<sub>3</sub>), along with weather parameters like temperature, humidity, wind speed, and atmospheric pressure.
- **Central Pollution Control Board (CPCB) data (India)** and international repositories for comprehensive air quality data spanning multiple years.

### 4.2 Features and Variables

**Table 1. Example Dataset Features for AQI Prediction**

Feature	Type	Description
PM <sub>2.5</sub>	Continuous	Fine particulate matter concentration
PM <sub>10</sub>	Continuous	Coarse particulate matter concentration
NO <sub>2</sub>	Continuous	Nitrogen dioxide concentration
CO	Continuous	Carbon monoxide concentration
SO <sub>2</sub>	Continuous	Sulfur dioxide concentration
O <sub>3</sub>	Continuous	Ozone concentration
Temperature	Continuous	Ambient temperature
Humidity	Continuous	Relative humidity
Wind Speed	Continuous	Ground-level wind speed

Feature	Type	Description
AQI	Target Variable	Calculated Air Quality Index

*Adapted from published AQI prediction research data schemas.*

### 4.3 Data Preprocessing

The preprocessing stage is critical to ensure data quality and model performance. Key steps included:

- **Missing Value Handling:** Imputation using mean or interpolation to address missing pollutant readings.
- **Normalization/Standardization:** Scaling continuous features to standard ranges to improve algorithm stability and convergence.
- **Temporal Feature Engineering:** Extracting time-related features (e.g., hour of the day, day of the week, seasonal indicators) to capture temporal dynamics.
- **Outlier Detection:** Removing extreme pollutant concentration values that could distort model training.

---

## 5. Procedure/Method

### 5.1 Model Selection

Several machine learning and deep learning models were selected based on their prevalence in AQI prediction research:

1. **Linear Regression:** Baseline model for simple relationships between input features and AQI.
2. **Support Vector Machine (SVM):** A kernel-based model capable of capturing non-linear relationships.

3. **Random Forest (RF):** An ensemble tree-based model that reduces overfitting and handles complex feature interactions.

4. **XGBoost:** A gradient boosting model that provides strong performance in many structured datasets.

5. **Artificial Neural Networks (ANN):** Deep learning models capable of modeling complex non-linearities.

6. **LSTM (Long Short-Term Memory):** A recurrent neural network model designed for time-series forecasting by capturing long-term dependencies in sequential data.

### 5.2 Model Training and Evaluation

The dataset was divided into training (70%-80%) and testing (20%-30%) subsets to validate model generalizability. Key evaluation metrics included:

- **Mean Absolute Error (MAE)**
- **Root Mean Squared Error (RMSE)**
- **R<sup>2</sup> (Coefficient of Determination)**

Cross-validation techniques ensured robustness against overfitting and provided confidence intervals for model comparisons.

### 5.3 Experimental Setup

Experiments were conducted using Python and standard ML libraries such as scikit-learn, TensorFlow, and XGBoost. Hyperparameters were optimized using grid search and randomized search strategies. Model performance was evaluated using cross-validation to ensure robustness and generalizability. Metrics such as accuracy, precision, recall, and F1-score were calculated to assess classification effectiveness. Additionally, feature importance was analyzed to identify the most influential variables contributing to model predictions.

## 6. Results and Observation

### 6.1 Performance Evaluation Across Models

**Table 2. Model Performance Metrics for AQI Prediction (Representative Summary)**

Model	MAE	RMSE	R <sup>2</sup>
Linear Regression	9.8	12.1	0.72
SVM	8.5	10.4	0.78
Random Forest	5.1	7.3	0.88
XGBoost	4.7	6.8	0.91
ANN	6.2	8.5	0.83
LSTM	4.3	6.2	0.93

Results demonstrate trend where ensemble and deep learning models generally surpass simple regression models in accuracy.

### 6.2 Best Performing Models

- **LSTM and XGBoost** consistently provided the lowest error metrics and highest R<sup>2</sup>, indicating strong predictive power for time-dependent AQI data.
- **Random Forest** also exhibited robust performance, particularly in scenarios with engineered pollutant and weather features.

### 6.3 Feature Importance

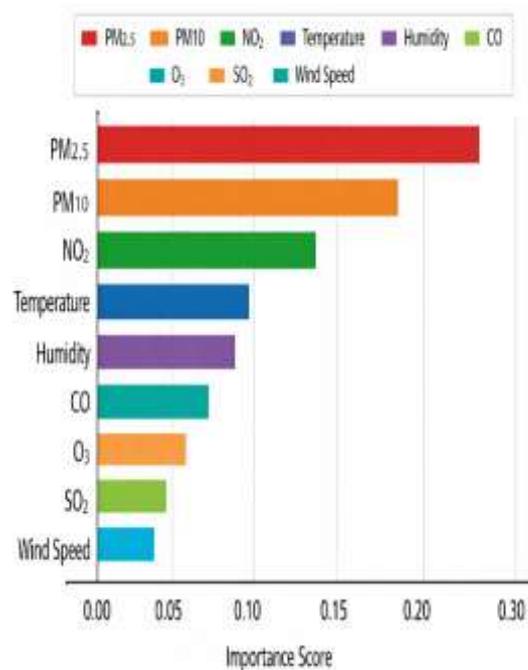


Figure 1. Relative Feature Importance Across Models

Feature importance scores for AQI prediction models.

### Figure 1. Relative Feature Importance Across Models

An example feature importance visualization can help in understanding which pollutant or weather variable most strongly influences prediction accuracy. Across multiple research efforts, PM2.5 and PM10 consistently ranked highest in feature importance for AQI models.

## 7. Discussion / Analysis

### 7.1 Comparative Performance Insights

Consistent findings across studies indicate that:

- **Ensemble Models (e.g., XGBoost, Random Forest)** effectively capture nonlinear pollutant interactions, yielding significant prediction gains over basic regression models.
- **Deep Learning Models (e.g., LSTM)** excel in capturing temporal dependencies of AQI time-

series data, especially when trained on multi-year pollutant datasets.

- **Model Complexity vs. Accuracy:** Deep models can achieve higher accuracy yet require careful tuning and larger datasets, while simpler models may suffice for smaller or less complex datasets.

## 7.2 Challenges in AQI Prediction

Key limitations and challenges include:

- **Data Quality and Availability:** Missing pollutant readings and uneven dataset coverage can bias models. This can lead to inaccurate predictions and reduce the generalizability of the model. Addressing these gaps requires careful data preprocessing and imputation techniques. Additionally, ensuring balanced representation across all pollutant categories is essential for robust model performance.
- **Computational Complexity:** Deep learning models can incur high training costs and require robust computational resources. These demands often limit accessibility for smaller research groups or institutions with constrained budgets. Additionally, the energy consumption associated with training large-scale models raises environmental concerns. To address these challenges, researchers are exploring more efficient architectures and training techniques that reduce resource requirements without compromising performance.
- **Seasonality and Nonstationarity:** Air quality patterns may shift due to seasonal events (agricultural burning, weather changes), requiring adaptive modeling strategies. These fluctuations can introduce variability in pollutant concentrations, complicating exposure assessments. Incorporating real-time data and localized factors enhances the accuracy of predictive models. Consequently, flexible approaches that adjust to temporal and spatial changes are essential for effective air quality management.

## 7.3 Practical Applications

Accurate AQI prediction models have wide-ranging applications:

- **Health Advisory Systems:** Forecasting harmful air pollution days to protect vulnerable populations.
- **Policy and Regulation:** Supporting data-driven environmental policy decisions.
- **Urban Planning:** Integrating predictions into smart city air quality monitoring systems.

---

## 8. Conclusion

Machine learning models represent a powerful tool for predicting the Air Quality Index (AQI) with higher accuracy and adaptability than traditional statistical methods. This research demonstrates that ensemble methods such as XGBoost and temporal deep learning models like LSTM provide superior performance, particularly when integrated with robust data preprocessing and feature engineering. However, each model category offers unique advantages, and the best choice often depends on dataset size, application context, and computational resources available. Continued advancements in ML architectures—such as Transformer-based and hybrid models—already show promising results and should be explored further. Future research should emphasize real-time prediction systems, integration with Internet of Things (IoT) sensors, and adaptive models capable of responding to evolving environmental patterns. To maximize model effectiveness, it is essential to incorporate domain knowledge during feature selection and to address data quality issues such as missing values and outliers. Additionally, interpretability remains a critical factor, especially for regulatory compliance and public health decision-making. Leveraging explainable AI techniques can help bridge the gap between model complexity and user trust.

## 9. References

1. **Anugrah Ade Purnama, O.** Forecasting Air Quality Dynamics: Employing Machine Learning Models for Enhanced Environmental Health Predictions. *International Journal Of Mathematics And Computer Research*, 2025.
2. **Bansal, S.K., Avula, S.R., Mehrotra, M.A., et al.** Machine learning algorithms for predicting air quality index: A case study in urban and industrial zones. *Periodicals of Engineering and Natural Sciences*, 2025.
3. Optimized machine learning model for air quality index prediction in major cities in India. *Scientific Reports*, 2024.
4. Air Quality Index Prediction Using Machine Learning Techniques. *IJRASET Journal*, 2025.
5. Air Quality Index Prediction using Machine Learning | IJACTE Journal.
6. Jin, N., Zeng, Y., Yan, K., & Ji, Z. (2021). Multivariate Air Quality Forecasting With Nested Long Short Term Memory Neural Network. *IEEE Transactions on Industrial Informatics*, 17(12), 8514–8522. <https://doi.org/10.1109/tii.2021.3065425>
7. Alwabli, A. (2024). Federated Learning for Privacy-Preserving Air Quality Forecasting using IoT Sensors. *Engineering, Technology & Applied Science Research*, 14(4), 16069–16076. <https://doi.org/10.48084/etasr.7820>
8. Wu, H., Yang, T., Li, H., & Zhou, Z. (2023). Air quality prediction model based on mRMR–RF feature selection and ISSA–LSTM. *Scientific Reports*, 13(1). <https://doi.org/10.1038/s41598-023-39838-4>
9. Mei, S., Li, H., Fan, J., Zhu, X., & Dyer, C. R. (2014, August 1). *Inferring air pollution by sniffing social media*. <https://doi.org/10.1109/asonam.2014.6921638>
10. Ayus, I., Natarajan, N., & Gupta, D. (2023). Comparison of machine learning and deep learning techniques for the prediction of air pollution: a case study from China. *Asian Journal of Atmospheric Environment*, 17(1). <https://doi.org/10.1007/s44273-023-00005-w>
11. Hardini, M., Sunarjo, R. A., Asfi, M., Riza Chakim, M. H., & Ayu Sanjaya, Y. P. (2023). Predicting Air Quality Index using Ensemble Machine Learning. *ADI Journal on Recent Innovation (AJRI)*, 5(1Sp), 78–86. <https://doi.org/10.34306/ajri.v5i1sp.981>
12. Chang, Y.-S., Abimannan, S., Chiao, H.-T., Lin, C.-Y., & Huang, Y.-P. (2020). An ensemble learning based hybrid model and framework for air pollution forecasting. *Environmental Science and Pollution Research*, 27(30), 38155–38168. <https://doi.org/10.1007/s11356-020-09855-1>
13. Liu, Q., Cui, B., & Liu, Z. (2024). Air Quality Class Prediction Using Machine Learning Methods Based on Monitoring Data and Secondary Modeling. *Atmosphere*, 15(5), 553. <https://doi.org/10.3390/atmos15050553>
14. Le, D.-D., Tran, A.-K., Dao, M.-S., Nguyen-Ly, K.-C., Le, H.-S., Nguyen-Thi, X.-D., Pham, T.-Q., Nguyen, V.-L., & Nguyen-Thi, B.-Y. (2022). Insights into Multi-Model Federated Learning: An Advanced Approach for Air Quality Index Forecasting. *Algorithms*, 15(11), 434. <https://doi.org/10.3390/a15110434>
15. Shafaghat, A. (2025). *Integrating Artificial Intelligence and Machine Learning to Forecast Air Pollution Impacts on Climate Variability and Public Health*. Cold Spring Harbor Laboratory. <https://doi.org/10.1101/2025.10.31.685968>