



Intelligent Air Quality Index Prediction System Using Machine Learning and Deep Learning Techniques

UPPADA LAKSHMI

PG Scholar, Department of MCA, DNR College, Bhimavaram, Andhra Pradesh

A. Naga Raju

(Assistant Professor), Master of Computer Applications, DNR College, Bhimavaram, Andhra Pradesh

ABSTRACT

Air pollution has emerged as one of the most critical environmental challenges worldwide, significantly affecting human health, ecosystems, and climate. Accurate prediction of the Air Quality Index (AQI) plays a vital role in mitigating risks, enabling authorities to take proactive measures and helping individuals make informed decisions. This project presents an intelligent air quality prediction system that leverages machine learning and deep learning techniques to estimate AQI based on multiple environmental pollutant parameters. The system is designed using a user-friendly graphical interface developed with Python's Tkinter library, allowing users to upload datasets, train multiple models, and perform AQI predictions seamlessly. The dataset consists of various pollutant concentrations such as PM_{2.5}, PM₁₀, NO, NO₂, NO_x, NH₃, CO, SO₂, O₃, Benzene, Toluene, and Xylene, which are widely recognized as key contributors to air pollution. After preprocessing, including handling missing values and normalization using StandardScaler, the dataset is split into training and testing sets. The proposed system incorporates multiple regression algorithms, including Linear Regression, Random Forest Regressor, Support Vector Regression (SVR), and a Deep Neural Network model implemented using TensorFlow. Each model is trained independently, and their performance is evaluated using Mean Squared Error (MSE). The system automatically selects the best-performing model based on the lowest MSE, ensuring optimal prediction accuracy.

A significant feature of the system is its ability to categorize AQI values into standard air quality levels such as Good, Satisfactory, Moderate, Poor, Very Poor, and Severe. This classification helps users easily interpret the results without requiring technical knowledge. Additionally, the system provides a graphical comparison of model performances, enabling users to understand the effectiveness of different algorithms. The integration of traditional machine learning models with deep learning techniques enhances the robustness and flexibility of the system. While simpler models like Linear Regression offer interpretability, advanced models like Random Forest and Neural Networks capture complex nonlinear relationships in the data, improving prediction accuracy. Overall, this project demonstrates a scalable, efficient, and intelligent approach to air quality prediction. It can be extended for real-time monitoring by integrating IoT sensors and deployed as a web or mobile application for broader accessibility. The

system contributes to environmental sustainability by promoting awareness and enabling data-driven decision-making for pollution control.

KEYWORDS: Air Quality Index (AQI), Machine Learning, Neural Networks, Random Forest, Support Vector Regression, Environmental Monitoring, Data Analytics, Pollution Prediction, Smart Systems, AI-based Prediction

I. INTRODUCTION

Air pollution has become a global concern due to rapid industrialization, urbanization, and increased vehicular emissions. Poor air quality has been directly linked to severe health issues such as respiratory diseases, cardiovascular problems, and premature deaths. Monitoring and predicting air quality is therefore essential for safeguarding public health and ensuring environmental sustainability. The Air Quality Index (AQI) is a standardized metric used to communicate how polluted the air currently is or how polluted it is forecasted to become. It converts complex air pollution data into a single numerical value, making it easier for the general public to understand. However, accurately predicting AQI is a challenging task due to the complex interactions between various pollutants and environmental factors. Traditional statistical methods often fail to capture the nonlinear relationships present in air pollution data. With the advancement of Artificial Intelligence (AI) and Machine Learning (ML), more sophisticated models have been developed to improve prediction accuracy. These models can learn patterns from historical data and provide reliable forecasts.

This project aims to develop an intelligent AQI prediction system using multiple machine learning and deep learning techniques. The system allows users to upload datasets, train different models, and select the best-performing model based on evaluation metrics. By incorporating algorithms such as Linear Regression, Random Forest, Support Vector Regression, and Neural Networks, the system ensures both accuracy and flexibility. The use of a graphical user interface (GUI) enhances usability, making the system accessible even to non-technical users. Users can input pollutant values and instantly obtain AQI predictions along with corresponding air quality categories. This feature makes the system practical for real-world applications, including environmental monitoring agencies and smart city initiatives. Another important aspect of the system is data preprocessing. Real-world datasets often contain missing or inconsistent values, which can negatively impact model performance. The system addresses this issue by cleaning the dataset and applying feature scaling to ensure consistency. In addition, the system provides a visual comparison of model performances, enabling users to understand which algorithm performs best under given conditions. This comparative approach not only improves prediction accuracy but also offers insights into the strengths and weaknesses of different models. In conclusion, this project combines the power of machine learning and deep learning to create a robust AQI prediction system. It serves as a foundation for future enhancements such as real-time prediction, integration with IoT devices, and

deployment on cloud platforms. By providing accurate and timely air quality predictions, the system contributes to better environmental management and public health awareness.

II. LITERATURE SURVEY (WITH EXISTING METHODS)

Numerous studies have been conducted in the field of air quality prediction, employing various statistical and machine learning techniques. Traditional approaches primarily relied on linear regression and time-series models such as ARIMA (AutoRegressive Integrated Moving Average). While these methods provided a basic understanding of air pollution trends, they often struggled with nonlinear relationships and complex interactions between pollutants. With the advancement of machine learning, researchers began exploring algorithms such as Decision Trees, Random Forests, and Support Vector Machines (SVM). Random Forest, an ensemble learning method, has been widely used due to its ability to handle high-dimensional data and reduce over fitting. Studies have shown that Random Forest models outperform traditional regression techniques in predicting AQI due to their robustness and accuracy. Support Vector Regression (SVR) is another widely used method for AQI prediction. It is particularly effective in handling nonlinear data by using kernel functions. Several research works have demonstrated that SVR provides high accuracy in predicting pollutant concentrations, especially when combined with feature engineering techniques. In recent years, deep learning models such as Artificial Neural Networks (ANN), Convolutional Neural Networks (CNN), and Long Short-Term Memory (LSTM) networks have gained popularity. Neural Networks are capable of capturing complex nonlinear patterns in large datasets, making them suitable for air quality prediction. LSTM models, in particular, have been used for time-series forecasting of AQI due to their ability to retain long-term dependencies.

Hybrid models combining machine learning and deep learning techniques have also been proposed. For example, integrating Random Forest with Neural Networks has shown improved prediction performance by leveraging the strengths of both approaches. Additionally, feature selection techniques such as Principal Component Analysis (PCA) have been used to reduce dimensionality and improve model efficiency. Recent studies have also focused on real-time air quality monitoring using IoT-based sensors. These systems collect real-time pollutant data and use cloud-based machine learning models for prediction. Such approaches enable continuous monitoring and early warning systems for air pollution. Despite these advancements, challenges remain in terms of data quality, model generalization, and computational complexity. Many models require large datasets and significant computational resources, which may not be feasible in all scenarios. The proposed system in this project addresses these challenges by combining multiple machine learning models with a deep learning approach and selecting the best-performing model based on evaluation metrics. This ensures both accuracy and efficiency, making the system suitable for practical applications.

III. EXISTING SYSTEM

Existing air quality prediction systems primarily rely on traditional statistical methods and basic machine learning techniques. These systems often use linear regression models or time-series forecasting methods such as ARIMA to predict AQI values. While these approaches are simple and easy to implement, they have several limitations. One major drawback of traditional systems is their inability to handle nonlinear relationships between different air pollutants. Air quality is influenced by multiple factors, including meteorological conditions and chemical interactions between pollutants. Linear models fail to capture these complex relationships, resulting in lower prediction accuracy. Another limitation is the reliance on a single model for prediction. Most existing systems do not provide a comparative analysis of multiple algorithms, which can lead to suboptimal performance. Additionally, these systems often lack automated model selection mechanisms, requiring manual intervention to choose the best model.

Many existing systems also do not include user-friendly interfaces, making them less accessible to non-technical users. The absence of visualization tools further limits the interpretability of results. Users are often required to analyze raw numerical outputs without any clear understanding of air quality categories. Furthermore, traditional systems may not handle missing or inconsistent data effectively, leading to unreliable predictions. Data preprocessing is often overlooked, which negatively impacts model performance. In contrast, the proposed system overcomes these limitations by integrating multiple machine learning and deep learning models, performing proper data preprocessing, and providing an intuitive GUI. It also includes automated model selection and visualization features, making it more efficient, accurate, and user-friendly compared to existing systems.

IV. PROPOSED METHOD

The proposed system is an intelligent Air Quality Index (AQI) prediction framework that integrates multiple machine learning and deep learning models to provide accurate and reliable air quality forecasts. Unlike traditional systems that rely on a single algorithm, this system adopts a multi-model approach and automatically selects the best-performing model based on evaluation metrics. The system accepts historical air pollution data consisting of key pollutant parameters such as PM_{2.5}, PM₁₀, NO, NO₂, NO_x, NH₃, CO, SO₂, O₃, Benzene, Toluene, and Xylene. These pollutants are known to significantly influence AQI levels and are widely used in environmental monitoring studies. The dataset is preprocessed by removing missing values and normalizing features using standard scaling techniques to ensure consistency and improve model performance. The system implements multiple regression algorithms including Linear Regression, Random Forest Regressor, Support Vector Regression (SVR), and a Deep Neural Network. Each model is trained independently, and its performance is evaluated using Mean Squared Error (MSE). The system then automatically selects the best model with the lowest error, ensuring optimal prediction accuracy. A user-friendly graphical interface developed using Tkinter allows users to upload datasets, train models, compare performance, and predict AQI values in real-time. The predicted AQI is further categorized into standard air

quality levels such as Good, Moderate, and Severe, enhancing interpretability for non-technical users. The proposed system is scalable and can be extended for real-time air quality monitoring by integrating IoT sensors and cloud-based data processing. Recent research highlights the importance of combining machine learning and deep learning models for improved AQI prediction accuracy and real-time decision-making. Thus, the system provides a robust, efficient, and user-centric solution for environmental monitoring and public health awareness.

V. IMPLEMENTATION

The implementation of the Air Quality Prediction System is carried out using Python, leveraging libraries such as Tkinter for GUI development, Pandas and NumPy for data handling, Scikit-learn for machine learning models, and TensorFlow for deep learning. The first step in the implementation is dataset acquisition and preprocessing. The system allows users to upload a CSV dataset containing pollutant values and corresponding AQI. Missing values are removed using data cleaning techniques, and feature selection is performed based on relevant pollutant parameters. The dataset is then divided into training and testing sets using an 80:20 split. Feature scaling is applied using StandardScaler to normalize the input values. This step is crucial for algorithms such as Support Vector Regression and Neural Networks, which are sensitive to feature magnitudes.

The system implements four models:

1. Linear Regression – a simple baseline model.
2. Random Forest Regressor – an ensemble model that improves accuracy by combining multiple decision trees.
3. Support Vector Regression – effective for nonlinear relationships using kernel functions.
4. Neural Network – a deep learning model with multiple dense layers for capturing complex patterns.

Each model is trained using the training dataset and evaluated on the testing dataset. The performance of each model is measured using Mean Squared Error (MSE), which quantifies the difference between predicted and actual AQI values. A key feature of the implementation is automated model selection. The system compares the MSE values of all trained models and selects the one with the lowest error as the best model. This eliminates the need for manual selection and ensures optimal performance. The graphical user interface (GUI) is designed using Tkinter, providing buttons for dataset upload, model training, model selection, and prediction. Users can input pollutant values through text fields, and the system displays the predicted AQI along with its category. Additionally, a visualization feature is implemented using Matplotlib to display a bar chart comparing the performance of different models. This helps users understand which algorithm performs best. The system is modular, allowing easy integration of additional



International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

models or features. It can be extended to support real-time data streaming, cloud deployment, and mobile applications. Overall, the implementation demonstrates a practical and efficient approach to AQI prediction, combining usability with advanced machine learning techniques.

ALGORITHMS

The system utilizes multiple machine learning and deep learning algorithms to predict AQI, ensuring robustness and accuracy.

1. Linear Regression

Linear Regression is a statistical method that models the relationship between independent variables and the dependent variable (AQI). It assumes a linear relationship and is used as a baseline model due to its simplicity and interpretability.

2. Random Forest Regressor

Random Forest is an ensemble learning algorithm that constructs multiple decision trees during training and outputs the average prediction. It reduces overfitting and improves accuracy. Studies show that ensemble models like Random Forest perform well on structured air quality datasets .

3. Support Vector Regression (SVR)

SVR is an extension of Support Vector Machines for regression problems. It uses kernel functions to map data into higher-dimensional space, enabling it to handle nonlinear relationships effectively.

4. Neural Network (Deep Learning)

The system uses a feedforward neural network with multiple hidden layers. Neural networks are capable of learning complex nonlinear patterns in data. Recent studies indicate that deep learning models such as LSTM and ANN outperform traditional models in AQI prediction tasks .

5. Model Evaluation (MSE)

Mean Squared Error (MSE) is used to evaluate model performance. It measures the average squared difference between predicted and actual values. Lower MSE indicates better model performance.

6. Best Model Selection Algorithm

The system compares MSE values of all trained models and selects the model with the minimum error. This automated approach ensures optimal prediction accuracy without manual intervention.

These algorithms collectively enhance the system's ability to handle both linear and nonlinear patterns, making it reliable for real-world air quality prediction.

SYSTEM DESIGN

The system architecture is designed as a modular and user-friendly framework that integrates data processing, machine learning, and user interaction components.

1. Input Layer (Data Acquisition)

The system accepts input in the form of a CSV dataset containing pollutant concentrations and AQI values. Users upload the dataset through the GUI interface.

2. Data Preprocessing Module

This module handles data cleaning, missing value removal, and feature selection. It also performs normalization using StandardScaler to ensure uniform data distribution.

3. Training Module

The training module consists of multiple machine learning and deep learning algorithms. Each model is trained independently using the processed dataset. The system supports:

- Linear Regression
- Random Forest
- Support Vector Regression
- Neural Network

4. Evaluation Module

After training, each model is evaluated using Mean Squared Error. The evaluation results are stored and used for comparison.

5. Model Selection Module

This module automatically selects the best-performing model based on the lowest MSE. Automated model selection improves efficiency and ensures optimal predictions.

6. Prediction Module

Users can input pollutant values through the GUI. The selected model processes the input data and predicts AQI. The system also categorizes the AQI into predefined levels.

7. Visualization Module

A graphical representation of model performance is provided using bar charts. This helps users compare different algorithms visually.

8. User Interface (GUI)

The GUI is developed using Tkinter and includes:

- Dataset upload button
- Model training buttons
- Model selection option



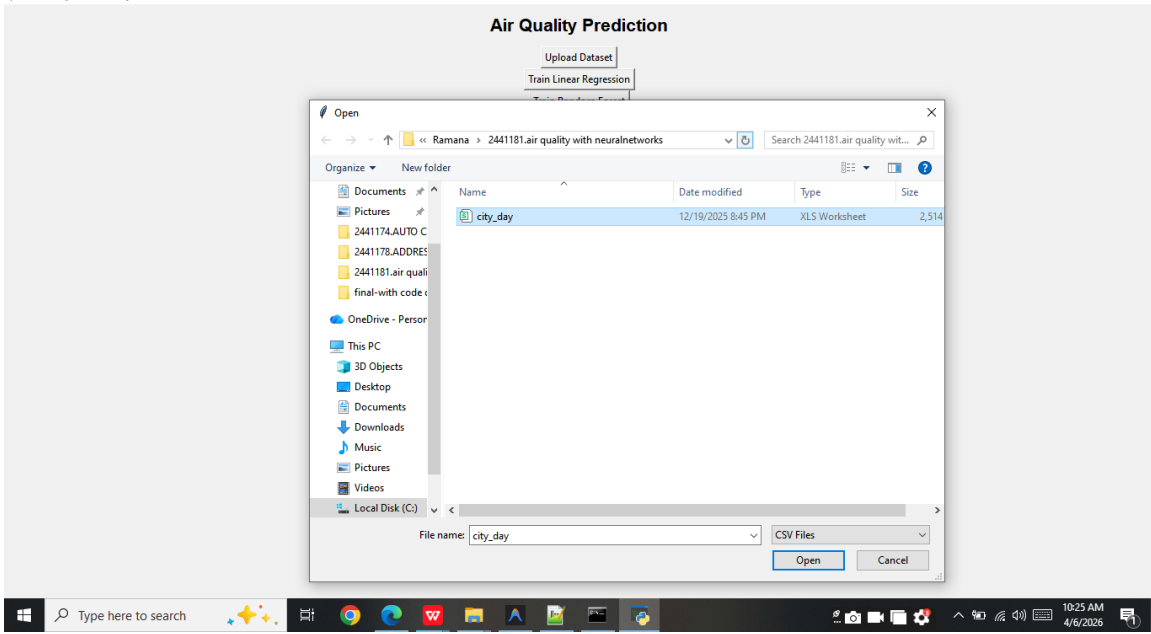
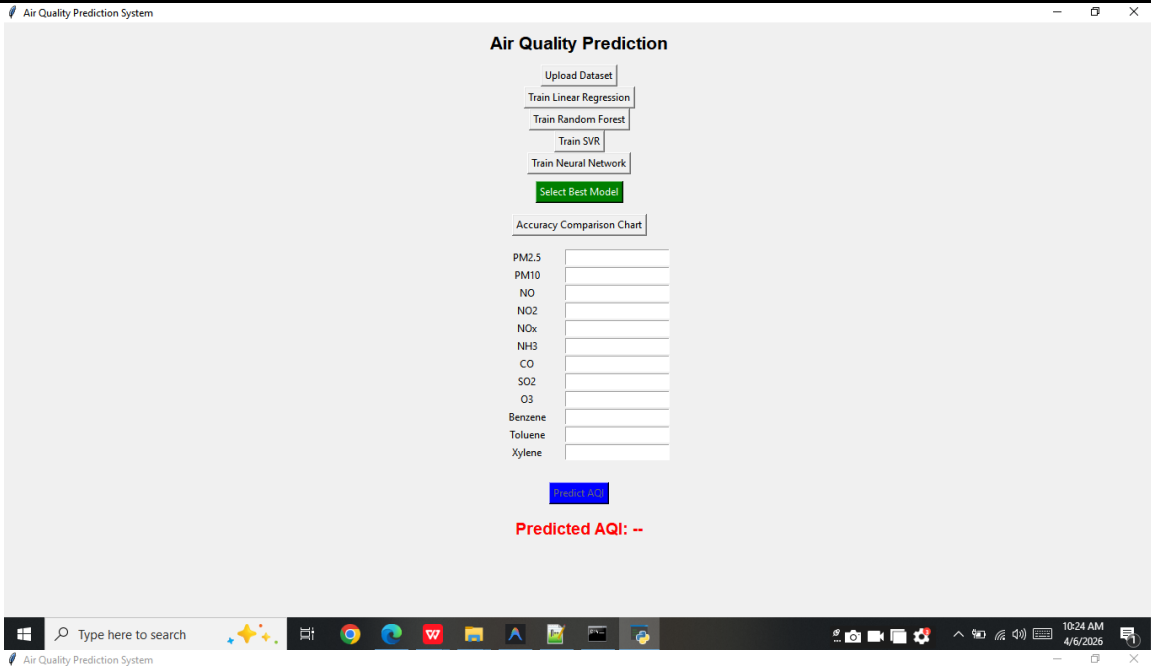
- Input fields for pollutants
- Prediction display

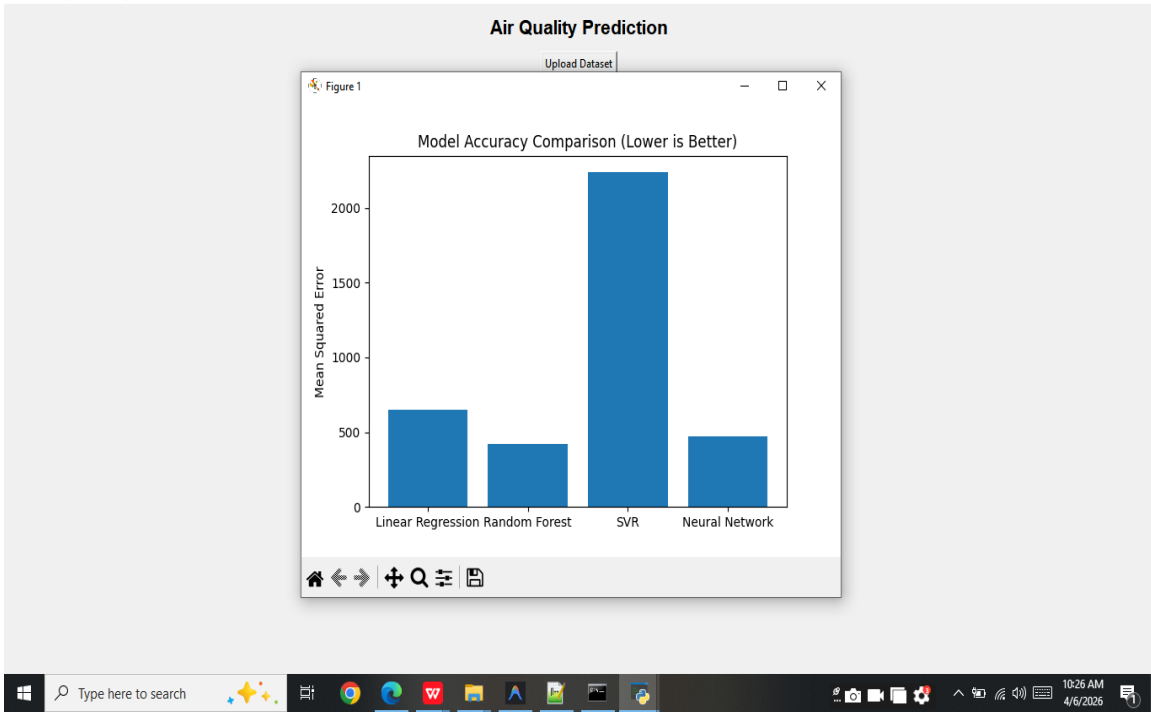
9. System Workflow

1. Upload dataset
2. Preprocess data
3. Train models
4. Evaluate models
5. Select best model
6. Input pollutant values
7. Predict AQI

Recent system designs emphasize integrating multiple models and real-time data for improved accuracy and scalability . The proposed design follows this approach, ensuring flexibility and future scalability.

SYSTEM DESIGN IMAGES





Air Quality Prediction

Upload Dataset

Train Linear Regression

Train Random Forest

Train SVR

Train Neural Network

Select Best Model

Accuracy Comparison Chart

PM2.5	2
PM10	2
NO	34
NO2	44
NOx	54
NH3	23
CO	78
SO2	89
O3	50
Benzene	37
Toluene	56
Xylene	61

Predict AQI

Predicted AQI: 86.35
Category: Satisfactory
Model: Random Forest

Air Quality Prediction System

VI. CONCLUSION

The Air Quality Prediction System developed in this project provides an efficient and intelligent solution for predicting AQI using machine learning and deep learning techniques. By integrating multiple models and selecting the best-performing one, the system ensures high accuracy and reliability. The use of various algorithms such as Linear Regression, Random Forest, Support Vector Regression, and Neural Networks allows the system to capture both linear and complex nonlinear relationships in air pollution data. Automated model selection based on Mean Squared Error enhances performance without requiring manual intervention. The graphical user interface makes the system accessible to users with minimal technical knowledge, enabling easy dataset upload, model training, and AQI prediction. The inclusion of AQI categorization further improves usability by presenting results in an understandable format. The system also demonstrates the importance of data preprocessing and feature scaling in improving model accuracy. Visualization features provide valuable insights into model performance, aiding in better decision-making. Recent advancements in air quality prediction highlight the growing importance of machine learning and deep learning techniques for environmental monitoring and public health protection. The proposed system aligns with these advancements and offers a scalable solution that can be extended to real-time applications. Future enhancements may include integration with IoT sensors for real-time data collection, deployment on cloud platforms, and the use of advanced models such as LSTM for time-series forecasting. Overall, the project contributes to environmental sustainability by providing a practical tool for air quality prediction, helping individuals and authorities take proactive measures to reduce pollution and protect public health.

REFERENCES

1. Liu, Q., et al. (2024). *Air Quality Prediction Using Machine Learning*. Atmosphere.
2. Méndez, M., et al. (2023). *Machine Learning Algorithms for Air Quality Forecasting*. Artificial Intelligence Review.
3. Agbehadji, I., et al. (2024). *Systematic Review of ML & DL for Air Quality Prediction*.
4. Study on Visakhapatnam AQI Prediction (2023). Chemosphere.



5. Rajesh, M., et al. (2025). *Real-time Air Quality Prediction Framework*. Scientific Reports.
6. Tang, D., et al. (2024). *Review of ML for Air Quality Modeling*.
7. Anggraini, T., et al. (2024). *Global AQI Prediction using ML*.
8. Singh, S. (2025). *Deep Learning for PM2.5 Prediction*.
9. Alawi, O. A., et al. (2024). *AI Models for Air Pollution Prediction*.
10. ScienceDirect (2024). *Deep Learning for AQI Prediction Survey*.
11. Science of Total Environment (2025). *ML Advances in AQI Prediction*.
12. Sidhu, K. K., et al. (2024). *AQI Prediction using ML in India*.
13. Pahari, S., et al. (2025). *Hybrid CNN-BiLSTM AQI Prediction Model*.
14. Chen, J., et al. (2025). *Deep Classifier Kriging for AQI*.
15. Masud, K. I., et al. (2026). *Benchmarking ML Models for AQI Prediction*.