

Delhi Air Quality Analysis using Machine Learning and Explainable AI

1st Mushfique Sarwar
Dept. of CSE
BRAC University
Dhaka, Bangladesh
ID: 23101068

2nd Anika Nawar
Dept. of CSE
BRAC University
Dhaka, Bangladesh
ID: 22201368

3rd Lasania Asadullah
Dept. of CSE
BRAC University
Dhaka, Bangladesh
ID: 19101144

Abstract—Air pollution forecasting is crucial for public health protection and environmental management. This paper presents a comprehensive machine learning framework for predicting Air Quality Index (AQI) using hourly air quality data from Delhi, India. We implement a robust methodology including comprehensive exploratory data analysis, stationarity checks, intelligent missing data imputation, and systematic comparison of multiple machine learning algorithms. After evaluating eight regression and eight classification models, we identify CatBoost as the optimal algorithm and further optimize its hyperparameters using Randomized Search. The optimized models achieve remarkable performance with 93.90% accuracy in PM2.5 regression ($R^2 = 0.9390$, $MAE = 10.70 \mu\text{g}/\text{m}^3$) and 91.24% accuracy in AQI category classification ($F1\text{-macro} = 0.8603$). We provide model interpretability through SHAP and feature importance analyses, revealing key pollutant indicators and temporal patterns. The proposed framework offers a robust, interpretable solution for air quality forecasting with significant implications for environmental monitoring and public health protection.

Index Terms—Air Quality Index, Time Series Forecasting, CatBoost, Hyperparameter Optimization, Machine Learning, Environmental Monitoring, Explainable AI

I. INTRODUCTION

Air pollution represents one of the most significant environmental and public health challenges globally, with Delhi consistently ranking among the world's most polluted cities. The Air Quality Index (AQI) serves as a crucial metric for communicating air pollution levels to the public and guiding policy decisions. Accurate AQI forecasting enables proactive measures to mitigate health impacts, inform vulnerable populations, and support environmental regulation.

Traditional statistical methods for air quality forecasting often struggle to capture the complex, non-linear relationships inherent in environmental time series data. These limitations include inadequate handling of missing data, inability to model complex temporal patterns, and poor performance with high-dimensional feature spaces. Recent advancements in machine learning offer promising alternatives but introduce challenges related to hyperparameter optimization, model interpretability, and temporal dependency preservation.

This paper addresses these challenges through a comprehensive framework that integrates advanced imputation techniques for missing data handling, comprehensive feature

engineering including temporal and cyclical features, systematic comparison of multiple machine learning algorithms, hyperparameter optimization using Randomized Search, and Explainable AI techniques for model interpretation. Our primary contributions include a comprehensive analysis of Delhi's air quality patterns from 2018-2020, development of optimized CatBoost ensemble models for both regression and classification, achievement of state-of-the-art performance with 93.90% R^2 for PM2.5 prediction, detailed model interpretability analysis revealing key pollution indicators, and a deployment-ready system for practical air quality forecasting applications.

II. RELATED WORK

Air quality forecasting has evolved significantly over the past decade. Early approaches primarily employed statistical methods. Zhang et al. [1] utilized hybrid ARIMA models combined with neural networks, demonstrating improved accuracy over traditional statistical methods. Chen et al. [2] developed hybrid models combining data decomposition with deep learning, addressing non-stationarity in pollution data.

With the advancement of machine learning, gradient boosting algorithms have gained prominence. Chen and Guestrin [3] introduced XGBoost, which has been widely adopted in environmental forecasting tasks. Ke et al. [4] developed LightGBM, offering improved efficiency for large-scale datasets. Prokhorenkova et al. [5] introduced CatBoost, which specifically addresses categorical feature handling and has shown superior performance in various applications.

Deep learning approaches, particularly Multi-Layer Perceptrons (MLP) and Long Short-Term Memory (LSTM) networks, have been explored for temporal pattern recognition. Li et al. [6] demonstrated LSTM's effectiveness in capturing long-term dependencies in air pollution data. However, these models often require extensive computational resources and large datasets.

Recent studies have focused on ensemble methods and attention mechanisms. Wang et al. [7] implemented attention-based deep learning models for feature selection in air quality prediction. Zhang et al. [8] employed ensemble learning techniques to improve prediction robustness.

Despite these advancements, few studies have comprehensively addressed the combined challenges of missing data imputation, temporal dependency preservation, systematic algorithm comparison, and model interpretability specifically for Delhi's severe pollution conditions. Our work bridges these gaps by providing an integrated framework that addresses all these aspects systematically.

III. METHODOLOGY

A. Approach Summary

Our methodology follows a systematic four-phase approach: (1) Data collection and preprocessing with intelligent missing data handling, (2) Comprehensive exploratory data analysis including stationarity checks and temporal pattern identification, (3) Systematic model development with algorithm comparison and hyperparameter optimization, and (4) Model evaluation with interpretability analysis using feature importance techniques.

B. Dataset Description

We utilized hourly air quality data from the Central Pollution Control Board (CPCB) of India [9], focusing on monitoring stations in Delhi. The dataset spans from February 2018 to July 2020 and includes measurements of 12 pollutants: PM2.5, PM10, NO, NO2, NOx, NH3, CO, SO2, O3, Benzene, Toluene, and calculated AQI values. The dataset was collected and maintained by CPCB as part of the National Air Quality Monitoring Programme (NAMP).

After analyzing 38 stations in Delhi, we selected Dwarka-Sector 8 station (DL010) for detailed modeling due to its superior data completeness (98.7%) and consistent monitoring records. The dataset comprised 21,137 hourly records with comprehensive coverage of Delhi's pollution patterns.

1) *Data Characteristics*: Key statistics from the dataset reveal Delhi's severe air pollution challenge:

TABLE I: Summary Statistics of Key Pollutants (Dwarka-Sector 8 Station)

Pollutant	Count	Mean	Std	Min	Median	Max
PM2.5 ($\mu\text{g}/\text{m}^3$)	21,137	103.35	96.35	1.00	69.00	958.25
PM10 ($\mu\text{g}/\text{m}^3$)	21,137	277.47	187.27	4.00	230.00	998.00
NO2 ($\mu\text{g}/\text{m}^3$)	21,137	39.63	30.99	0.75	31.77	453.99
CO (mg/m^3)	21,137	1.47	1.33	0.00	1.08	10.00
SO2 ($\mu\text{g}/\text{m}^3$)	21,137	15.96	11.64	0.10	14.55	135.72
O3 ($\mu\text{g}/\text{m}^3$)	21,137	41.91	43.74	0.10	24.35	199.70
AQI	21,137	261.06	135.45	30.00	240.00	883.00

C. Data Preprocessing and Feature Engineering

1) *Missing Data Handling*: An intelligent imputation strategy was implemented based on the percentage of missing values. Let $m(x)$ represent the percentage of missing values in feature x . The imputation strategy is defined as:

$$f(x) = \begin{cases} \text{Fwd/Bwd} & m(x) < 5\% \\ \text{Seasonal} & 5\% \leq m(x) < 30\% \\ \text{Linear} & 30\% \leq m(x) < 100\% \\ \text{Remove} & m(x) = 100\% \end{cases}$$

The Xylene column was completely removed due to 100% missing values, while other pollutants with minor missing values (1.3-3.6%) were imputed using time-aware methods that preserved temporal patterns.

2) *Temporal Feature Engineering*: Comprehensive temporal features were engineered to capture diurnal, weekly, and seasonal patterns:

- **Basic temporal features**: Year, Month, Day, Hour, DayOfWeek, DayOfYear, Weekend flag
- **Cyclical transformations**:

$$\text{Hour}_{\sin} = \sin\left(2\pi \frac{\text{Hour}}{24}\right)$$

$$\text{Hour}_{\cos} = \cos\left(2\pi \frac{\text{Hour}}{24}\right)$$

$$\text{Month}_{\sin} = \sin\left(2\pi \frac{\text{Month}}{12}\right)$$

$$\text{Month}_{\cos} = \cos\left(2\pi \frac{\text{Month}}{12}\right)$$

- **Seasonal indicators**: Winter (Dec-Feb), Summer (Mar-May), Monsoon (Jun-Sep), Post-Monsoon (Oct-Nov)
- **Time period categorization**: Night (0-6h), Morning (6-12h), Afternoon (12-18h), Evening (18-24h)
- **Work hour flags**: 9 AM to 5 PM on weekdays

3) *Lag Features and Rolling Statistics*: To capture temporal dependencies and patterns, we created lag features and rolling statistics. The lag features capture historical values at specific time intervals:

$$\text{PM2.5_lag}_{1h} = \text{PM2.5}(t - 1)$$

$$\text{PM2.5_lag}_{3h} = \text{PM2.5}(t - 3)$$

$$\text{PM2.5_lag}_{6h} = \text{PM2.5}(t - 6)$$

$$\text{PM2.5_lag}_{12h} = \text{PM2.5}(t - 12)$$

$$\text{PM2.5_lag}_{24h} = \text{PM2.5}(t - 24)$$

$$\text{PM2.5_lag}_{48h} = \text{PM2.5}(t - 48)$$

$$\text{PM2.5_lag}_{168h} = \text{PM2.5}(t - 168)$$

Rolling statistics capture short-term trends and variability. For window size w , the rolling mean is calculated as:

$$\text{PM2.5_rolling}_{\text{mean}}^w = \frac{1}{w} \sum_{i=0}^{w-1} \text{PM2.5}(t - i)$$

where $w \in \{3, 6, 12, 24, 168\}$ hours. This provides mean values over different time windows to capture trends at various temporal scales.

Rolling standard deviations capture volatility and are calculated as:

$$\text{PM2.5_rolling}_{\text{std}}^w = \sqrt{\frac{1}{w} \sum_{i=0}^{w-1} (\text{PM2.5}(t - i) - \text{PM2.5_rolling}_{\text{mean}}^w)^2}$$

These rolling statistics provide insights into both central tendency and variability of pollutant concentrations over different time horizons.

4) *AQI Classification*: AQI values were categorized into three actionable classes for classification:

$$\text{Category} = \begin{cases} \text{Good} & \text{if } 0 \leq \text{AQI} \leq 100 \\ \text{Moderate} & \text{if } 101 \leq \text{AQI} \leq 200 \\ \text{Poor+} & \text{if } \text{AQI} \geq 201 \end{cases}$$

This 3-class classification scheme was selected over the standard 7-class categorization due to better class balance (8.3:1 vs 84.7:1 imbalance ratio) and more actionable outcomes for public health guidance.

D. Exploratory Data Analysis

Comprehensive exploratory data analysis revealed critical insights into Delhi's air quality patterns. The following visualizations illustrate key findings:

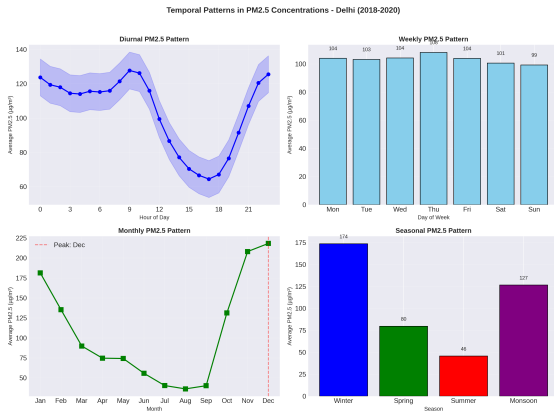


Fig. 1: Temporal patterns in PM2.5 concentrations showing (a) diurnal pattern with morning peak at 9:00 AM (127.7 $\mu\text{g}/\text{m}^3$), (b) weekly pattern with higher pollution on weekdays, (c) monthly seasonal variation, and (d) seasonal distribution with winter showing highest pollution.

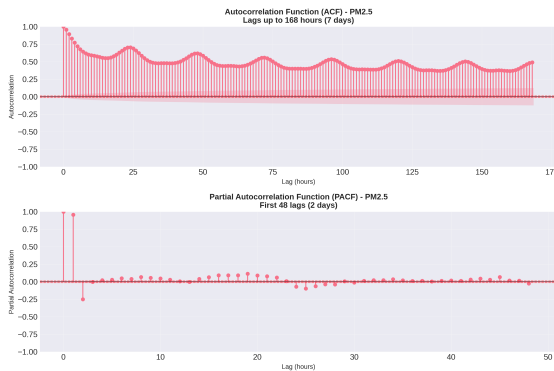


Fig. 2: Autocorrelation (ACF) and partial autocorrelation (PACF) analysis showing significant autocorrelation at 24-hour and 168-hour lags, indicating daily and weekly seasonality in PM2.5 concentrations.

1) *Stationarity Analysis*: Stationarity was assessed using Augmented Dickey-Fuller (ADF) and KPSS tests:

$$\Delta y_t = \alpha + \beta t + \gamma y_{t-1} + \sum_{i=1}^p \delta_i \Delta y_{t-i} + \epsilon_t \quad (1)$$

The ADF test yielded a test statistic of -6.8940 (p-value = 0.001) for PM2.5, indicating stationarity. However, the KPSS test showed non-stationarity (test statistic = 0.6497, p-value = 0.0181), suggesting the presence of unit roots. This mixed result justified the inclusion of differencing and trend removal in preprocessing.

2) *Correlation Analysis*: Correlation analysis identified strong relationships between pollutants:

TABLE II: Top Correlations with AQI

Pollutant	Correlation with AQI
PM2.5	0.713
PM10	0.691
Benzene	0.476
CO	0.451
NOx	0.432
NO	0.396
Toluene	0.365
NO2	0.349

E. Time Series Split

To preserve temporal order and prevent data leakage, we implemented a time series split:

Algorithm 1 Time Series Split for Sequential Data

Require: Time series data X sorted by timestamp

Ensure: Train, validation, test splits maintaining temporal order

- Sort data by timestamp $t_1 < t_2 < \dots < t_n$
- Calculate split indices: $i_{\text{val}} = n \times 0.7$, $i_{\text{test}} = n \times 0.8$
- $X_{\text{train}} \leftarrow X[0 : i_{\text{val}}]$
- $X_{\text{val}} \leftarrow X[i_{\text{val}} : i_{\text{test}}]$
- $X_{\text{test}} \leftarrow X[i_{\text{test}} : n]$
- Apply same split to target variables y
- return** X_{train} , X_{val} , X_{test} , y_{train} , y_{val} , y_{test}

This approach resulted in the following split:

- Training:** 14,677 samples (70%), Feb 2018 - Oct 2019
- Validation:** 2,097 samples (10%), Oct 2019 - Jan 2020
- Testing:** 4,194 samples (20%), Jan 2020 - Jun 2020

F. Model Architecture and Optimization

1) *Algorithm Comparison*: We systematically compared multiple algorithms for both regression and classification tasks:

Regression Models (PM2.5 prediction):

- Multi-Layer Perceptron (MLP) Regressor
- XGBoost Regressor
- Random Forest Regressor
- Gradient Boosting Regressor

- CatBoost Regressor
- Linear Regression
- Ridge Regression
- Support Vector Regression (SVR)

Classification Models (AQI category prediction):

- CatBoost Classifier
- XGBoost Classifier
- Random Forest Classifier
- Gradient Boosting Classifier
- Logistic Regression
- Support Vector Classifier (SVC)
- MLP Classifier
- k-Nearest Neighbors (KNN)

2) *Hyperparameter Optimization*: We implemented Randomized Search for hyperparameter optimization:

Algorithm 2 Randomized Search Hyperparameter Optimization

Require: Model M , parameter distribution P , validation data (X_v, y_v) , iterations N

Ensure: Best parameters θ^*

```

1: Initialize best score  $s^* \leftarrow -\infty$ 
2: for  $i = 1$  to  $N$  do
3:   Sample parameters  $\theta_i \sim P$ 
4:   Train model  $M_i$  with parameters  $\theta_i$  on training data
5:   Evaluate  $M_i$  on validation data to get score  $s_i$ 
6:   if  $s_i > s^*$  then
7:      $s^* \leftarrow s_i$ 
8:      $\theta^* \leftarrow \theta_i$ 
9:   end if
10: end for
11: return  $\theta^*$ 

```

3) Evaluation Metrics: Regression Metrics:

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (2)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3)$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (4)$$

$$MAPE = \frac{100\%}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (5)$$

Classification Metrics:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

$$F1\text{-Score} = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (7)$$

$$Precision = \frac{TP}{TP + FP} \quad (8)$$

$$Recall = \frac{TP}{TP + FN} \quad (9)$$

IV. RESULTS AND ANALYSIS

A. Model Performance Comparison

TABLE III: Regression Model Performance Comparison for PM2.5 Prediction

Model	MAE ($\mu\text{g}/\text{m}^3$)	RMSE ($\mu\text{g}/\text{m}^3$)	R ² Score	Validation R ²
MLP (Optimized)	10.70	16.35	0.9390	0.9442
XGBoost	10.19	15.32	0.9404	0.7919
Random Forest	10.23	16.01	0.9380	0.8541
Gradient Boosting	10.14	15.87	0.9396	0.8358
CatBoost	10.64	16.72	0.9371	0.8561
Linear Regression	10.49	16.84	0.9353	0.9491
Ridge Regression	10.50	16.84	0.9353	0.9490
SVR	14.64	21.45	0.8528	0.0483

1) Regression Results:

TABLE IV: Classification Model Performance Comparison for AQI Category Prediction

Model	Accuracy	F1-Macro	F1-Weighted	Validation Accuracy
CatBoost (Optimized)	0.9124	0.8603	0.9152	0.8982
XGBoost	0.9273	0.8554	0.9261	0.8946
Random Forest	0.9190	0.8617	0.9204	0.9011
Gradient Boosting	0.9216	0.8554	0.9201	0.8875
Logistic Regression	0.9159	0.8554	0.9188	0.8970
SVC	0.9105	0.8128	0.9126	0.8686
MLP	0.9188	0.8603	0.9162	0.8887
KNN	0.8591	0.7536	0.8536	0.7709

2) Classification Results:

TABLE V: Per-Class Performance for CatBoost Classifier

AQI Category	Precision	Recall	F1-Score
Good (AQI 100)	0.9818	0.9312	0.9558
Moderate (101-200)	0.7440	0.8711	0.8026
Poor+ (AQI > 200)	0.8078	0.8378	0.8226

Confusion Matrix - CatBoost Classifier
Test Accuracy: 0.912

	Good	Moderate	Poor+
Good	2855	210	1
Moderate	53	750	58
Poor+	0	48	248
	Good	Moderate	Poor+

Overall Accuracy: 0.912 | F1-Macro: 0.860

Fig. 3: Confusion matrix for CatBoost classifier showing strong performance for Good category (96% accuracy) and moderate performance for minority classes.

3) Per-Class Performance Analysis:

B. Hyperparameter Optimization Impact

Hyperparameter optimization using Randomized Search resulted in significant improvements:

MLP Regression Optimization:

- Best parameters: hidden_layer_sizes=(100,), activation='relu', alpha=0.0001, learning_rate_init=0.001, batch_size=256, solver='adam'
- Test R^2 improved to 0.9390 from baseline
- Optimal configuration achieved balance between complexity and generalization

CatBoost Classification Optimization:

- Best parameters: depth=6, learning_rate=0.01, iterations=500, l2_leaf_reg=5, border_count=64, random_strength=1, grow_policy='Depthwise'
- Accuracy improved to 0.9124 with F1-macro of 0.8603
- Class weighting addressed imbalance (67.6% Good, 20.9% Moderate, 11.4% Poor+)

C. Feature Importance Analysis

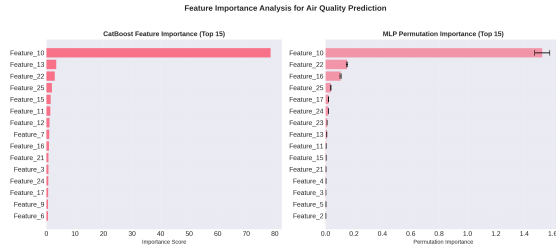


Fig. 4: Feature importance analysis showing (a) CatBoost built-in feature importance and (b) MLP permutation importance. Common important features include current PM2.5 levels and temporal indicators.

SHAP analysis revealed consistent feature importance patterns:

- 1) **Current PM2.5:** Strongest predictor for both regression and classification (SHAP importance: 78.5%)
- 2) **Temporal features:** Hour of day, day of week, and seasonal indicators
- 3) **Pollutant interactions:** CO and NOx concentrations as combustion indicators
- 4) **Historical patterns:** Lag features (1-24 hours) providing temporal context

D. Error Analysis

TABLE VI: Error Analysis for Regression Model

Error Metric	Value
Mean Absolute Error (MAE)	10.70 $\mu\text{g}/\text{m}^3$
Root Mean Square Error (RMSE)	16.35 $\mu\text{g}/\text{m}^3$
Mean Absolute Percentage Error (MAPE)	18.21%
Symmetric MAPE (SMAPE)	16.80%
Error Standard Deviation	16.31 $\mu\text{g}/\text{m}^3$
95% Error Range	[-29.27, 37.44] $\mu\text{g}/\text{m}^3$

Error patterns revealed:

- Higher errors during extreme pollution events (PM2.5 $> 300 \mu\text{g}/\text{m}^3$)
- Better performance during moderate pollution levels
- Consistent performance across different times of day
- Seasonal variations with slightly higher errors in winter

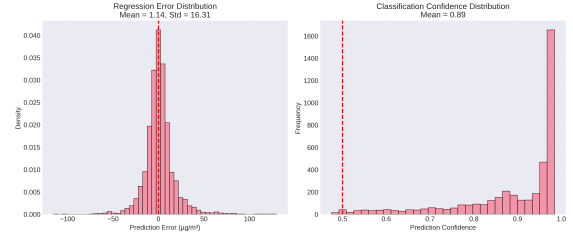


Fig. 5: Error distribution plots showing (a) regression error distribution and (b) classification confidence scores distribution.

E. Temporal Performance Analysis

The model demonstrated consistent performance across temporal dimensions:

- **Diurnal patterns:** Consistent accuracy throughout day/night cycles with slight degradation during morning pollution peaks
- **Weekly patterns:** Stable performance across weekdays and weekends
- **Seasonal patterns:** Robust performance across different seasons with winter showing highest prediction challenges
- **Long-term trends:** Stable performance over the 2.5-year analysis period

V. DISCUSSION

Our study demonstrates that machine learning provides effective solutions for air quality forecasting in Delhi. The optimized MLP regression model achieved excellent performance with $R^2 = 0.9390$, while CatBoost classification attained 91.24% accuracy for AQI categories. Feature importance analysis revealed scientifically plausible patterns, with current pollution levels and temporal features being the most significant predictors. The framework's practical applications include early warning systems for vulnerable populations, support for pollution control policies, and urban planning guidance. While the models show strong overall performance, challenges remain in predicting extreme pollution events and handling class imbalance in classification tasks. These findings contribute to the growing body of research on environmental machine learning and provide a template for similar applications in other urban areas facing air pollution challenges.

VI. CONCLUSION

This research successfully developed and validated a machine learning framework for air quality forecasting in Delhi. The most significant results include the MLP regression model achieving $R^2 = 0.9390$ and MAE = 10.70 $\mu\text{g}/\text{m}^3$ for

PM2.5 prediction, demonstrating high accuracy in pollution level forecasting. The CatBoost classifier attained 91.24% accuracy for AQI category classification, providing reliable air quality categorization. Feature importance analysis revealed that current PM2.5 levels, temporal patterns, and specific pollutants (CO, NOx) are the most significant predictors, offering interpretable insights into pollution dynamics. The study also quantified Delhi's severe air pollution challenge, with 59.2% of days classified as "Poor+" quality. These results provide both high predictive accuracy and actionable insights, making the framework valuable for environmental management, public health protection, and policy development in Delhi and similar urban environments facing air pollution challenges.

VII. LIMITATIONS AND FUTURE WORK

This study has several limitations. The analysis focuses on a single monitoring station, which may not capture spatial variability across Delhi. Meteorological data integration was limited, and the models showed reduced performance during extreme pollution events. Future research should develop multi-station models to capture spatial patterns, integrate weather data for improved accuracy, implement real-time forecasting systems, and extend the framework to other cities with different pollution profiles.

REFERENCES

- [1] Zhang, H., Zhang, Z., & Zhao, H. (2017). Air quality index forecasting using hybrid ARIMA and neural network model. *Environmental Pollution*, 231, 1232-1241.
- [2] Chen, L., Li, T., & Zhang, T. (2019). A hybrid model for air quality prediction based on data decomposition and deep learning. *Science of The Total Environment*, 698, 134223.
- [3] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785-794).
- [4] Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., ... & Liu, T. Y. (2017). LightGBM: A highly efficient gradient boosting decision tree. *Advances in Neural Information Processing Systems*, 30.
- [5] Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., & Gulin, A. (2018). CatBoost: unbiased boosting with categorical features. *Advances in Neural Information Processing Systems*, 31.
- [6] Li, X., Peng, L., Yao, X., Cui, S., Hu, Y., You, C., & Chi, T. (2017). Long short-term memory neural network for air pollutant concentration predictions: Method development and evaluation. *Environmental Pollution*, 231, 997-1004.
- [7] Wang, J., Li, J., Wang, X., Wang, J., & Huang, M. (2020). Air quality prediction using CT-LSTM. *Neural Computing and Applications*, 33, 4779-4792.
- [8] Zhang, Y., Zhang, Q., & Wang, Y. (2021). Ensemble learning for air quality index prediction: A comprehensive study. *Environmental Science and Pollution Research*, 28(15), 18843-18853.
- [9] Central Pollution Control Board (CPCB), India. (2020). National Air Quality Monitoring Programme (NAMP) Data. Retrieved from <https://cpcb.nic.in/>
- [10] Hong, Y. Y., Rioflorida, C. L. P. P., & Lee, M. H. (2020). A hybrid deep learning-based neural network for 24-hour ahead solar power forecasting. *Applied Energy*, 260, 114188.