

# Machine Learning and Data Analytics Framework for Smart City Air Quality Monitoring

Ashish Birajdar\*, Faizur Rashid\*\*

\*(Department Name, JSPM UNIVERSITY, PUNE

Email: ashishbirajdar777@gmail.com)

\*\*\*\*\*

## Abstract:

The degradation of air quality has become a significant challenge impacting environmental sustainability and public health in regions undergoing rapid urbanization. Traditional monitoring systems, which depend on a small number of fixed regulatory stations, often miss localized changes in the levels of pollutants. To get around this problem, modern smart cities need air quality monitoring solutions that can be scaled up, work in real time, and are smart.

Recent improvements in machine learning, deep learning, the Internet of Things (IoT), and data analytics have made it possible to make air quality management systems that can predict and change. This study presents a holistic framework that amalgamates sensor calibration, spatiotemporal deep learning models, and explainable artificial intelligence (XAI) to enhance the precision of air quality forecasting. The framework uses multidimensional time-series data, such as weather, traffic, and pollution data, to make predictions more accurate.

A comparative analysis of models including ARIMA, LSTM, Transformer, and Graph Neural Networks (GNNs) illustrates that deep learning techniques substantially exceed conventional statistical methods in elucidating intricate dependencies within air pollution data. SHAP-based explainability also helps us understand model predictions and find the most important things that affect PM2.5 levels.

The proposed framework shows how AI-driven solutions could make smart cities better places to live by improving environmental monitoring, helping people make better decisions, and improving their quality of life.

**Keywords — Smart City, Air Quality Monitoring, Data Analytics, Machine Learning, PM2.5, Power BI, IoT Data**

\*\*\*\*\*

## I. INTRODUCTION

Air pollution remains a significant threat to ecological equilibrium and public health in urban areas. Reports from global health organizations indicate that prolonged exposure to polluted air significantly contributes to respiratory and cardiovascular disorders [5]. Pollution levels in big cities have gone up a lot because of the rapid growth

of the population, industry, transportation systems, and changing weather.

The idea of smart cities is to use cutting-edge digital technologies, like data analytics and communication systems, to make city life better [8], [9]. Environmental monitoring, especially checking the quality of the air, is an important part of this vision. But traditional monitoring systems depend on costly and poorly spread-out stations, which makes it hard

to get detailed information about how pollution changes in different areas [1].

Using cheap IoT-based sensors has improved coverage and made it possible to collect data in real time. Even though they have this advantage, these sensors often have problems with accuracy because of interference from the environment and problems with calibration [2]. Machine learning techniques have also shown promise in finding complicated links between environmental factors and levels of pollution [7]. Advanced deep learning models, such as LSTM networks and Transformer architectures, have exhibited robust efficacy in time-series forecasting tasks [4], [7].

There are still some problems that need to be solved, though, like combining different data sources, modeling spatial dependencies, making sure the model is easy to understand, and making it work for real-time applications. This study tackles these issues by suggesting a unified machine learning framework that merges predictive modeling, sensor calibration, and explainable AI methodologies for monitoring air quality in smart cities..

## **2. Literature Review**

### **2.1 Conventional Air Quality Monitoring Systems**

#### **Standard Air Quality Monitoring Systems**

Traditionally, high-precision regulatory stations with advanced analytical tools have been used to check air quality by measuring pollutants like PM<sub>2.5</sub>, PM<sub>10</sub>, NO<sub>2</sub>, SO<sub>2</sub>, CO, and O<sub>3</sub>. Even though these systems give accurate measurements, they are not widely used because they are expensive to install and maintain, which means they don't cover a lot of space [1].

Also, these systems mostly help with descriptive analysis and don't have strong predictive abilities. Statistical methods like ARIMA have been used for making predictions, but they don't work well when the relationships between different parts of the environment are complex and not linear.

### **2.2 IoT-Based Sensor Networks in Smart Cities**

The advent of IoT-enabled air quality sensors has revolutionized environmental monitoring by facilitating extensive and continuous data collection in urban areas [11], [12]. These sensors give us high-resolution time data that helps us understand how pollution changes over time.

But low-cost sensors are often affected by things like temperature and humidity, which can cause them to give wrong readings and drift over time [2]. To tackle these issues, researchers have looked into calibration methods that use machine learning to make data more reliable [1]. Even with these efforts, getting consistent performance across different sites is still a problem.

### **2.3 Machine Learning for Air Quality Prediction**

Machine learning methods have gotten a lot of attention for predicting air quality indices. Conventional methods like linear regression and support vector regression demonstrate moderate efficacy yet encounter difficulties with intricate temporal dependencies [19], [20].

Deep learning models, especially LSTM networks, can find patterns that happen over time and long-term dependencies in time-series data [7]. Transformer-based models have recently become strong alternatives because their attention mechanisms let them better model long-range interactions [4].

## 2.4 Spatio-Temporal Modelling Using Graph Neural Networks

There are both time and space factors that affect air pollution. Wind flow and traffic movement are two examples of things that can change the concentration of pollutants at a certain place. Graph Neural Networks effectively represent these relationships by modeling sensor locations as nodes within a graph structure [6].

Spatio-temporal models can accurately predict complex interactions between different areas by combining spatial graph representations with temporal learning mechanisms [6], [29]

## 2.5 Explainable Artificial Intelligence (XAI)

One of the biggest problems with deep learning models is that they are hard to understand. In areas like environmental monitoring, it's important to understand how models make decisions in order to make policies and gain the public's trust.

SHAP and other explainable AI methods give information about how each input variable affects the output by giving each one an importance value [3]. This helps stakeholders find the most important things that affect pollution levels and makes it easier for them to make smart choices.

## 3. Research Gap

Although substantial progress has been made, key gaps remain:

1. Limited integration of heterogeneous datasets (pollution + traffic + weather).
2. Insufficient spatio-temporal modeling in dense urban networks.

3. Lack of robust ML-based calibration frameworks for low-cost sensors.
4. Minimal application of explainable AI for policy-relevant insights.
5. Limited real-time scalable implementations.

This research proposes an integrated architecture addressing these challenges simultaneously.

## . Methodology

### 4.1 Research Design

This study adopts a **quantitative, predictive modelling research design** integrating data analytics, machine learning, and explainable AI techniques. The objective is to develop and evaluate a smart city air quality monitoring framework capable of:

1. Improving prediction accuracy of PM2.5 and AQI levels
2. Enhancing reliability of low-cost IoT sensors
3. Capturing spatial and temporal pollution dynamics
4. Providing explainable insights for urban governance

The research follows a structured pipeline:

**Data Collection → Data Preprocessing → Sensor Calibration → Model Development → Model Evaluation → Explainability Analysis**

### 4.2 Data Collection and Data Sources

To ensure comprehensive modelling, multi-source heterogeneous datasets were collected over a **24-month period** at hourly resolution.

#### 4.2.1 Air Quality Data

Data included the following pollutants:

- PM2.5 (primary dependent variable)
- PM10
- NO<sub>2</sub>
- CO
- O<sub>3</sub>

These were collected from:

1. **Regulatory Monitoring Stations** (reference-grade analyzers)
2. **Low-Cost IoT Sensors** deployed across residential, traffic, and industrial zones

Regulatory stations provided ground-truth values for calibration and benchmarking.

#### 4.2.2 Meteorological Data

Meteorological variables significantly influence pollutant dispersion and concentration. The following parameters were included:

- Temperature (°C)
- Relative humidity (%)
- Wind speed (m/s)
- Wind direction (degrees)
- Atmospheric pressure (hPa)

These variables help capture diffusion, inversion layers, and pollutant trapping mechanisms.

#### 4.2.3 Traffic and Urban Activity Data

Traffic density is a major contributor to urban PM2.5 levels. The following variables were integrated:

- Vehicle counts per hour
- Congestion index
- Peak-hour indicators
- Road type classification (arterial vs residential)

The inclusion of traffic data addresses a major limitation in many previous studies that rely solely on meteorological variables.

#### 4.3 Data Preprocessing

High-quality modelling requires robust preprocessing.

##### 4.3.1 Missing Value Treatment

Missing values were handled using:

- Linear interpolation for short gaps (< 3 hours)
- K-nearest neighbour imputation for longer gaps

This prevented bias caused by listwise deletion.

##### 4.3.2 Outlier Detection

Outliers were identified using:

- Z-score thresholding ( $|Z| > 3$ )
- Interquartile Range (IQR) method

Extreme anomalies caused by sensor malfunction were removed.

### 4.3.3 Normalization

Min-Max normalization was applied:

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

This ensured model stability and faster convergence for deep learning algorithms.

### 4.3.4 Feature Engineering

Temporal dependencies were enhanced by creating:

1. lag features (t-1, t-2, t-24)
2. Rolling averages (6-hour and 24-hour windows)
3. Time-based features (hour of day, day of week, season)

This significantly improved model performance by embedding cyclical patterns.

## 4.4 Sensor Calibration Using Machine Learning

Low-cost sensors are known to suffer from drift, humidity interference, and cross-sensitivity. To address this, a **Random Forest regression calibration model** was implemented

### 4.4.1 Calibration Model Inputs

- Raw low-cost PM2.5 reading
- Temperature
- Humidity
- Wind speed
- Time features

### 4.4.2 Calibration Target

Regulatory station PM2.5 reading (ground truth)

### 4.4.3 Why Random Forest?

Random Forest was selected because:

- It handles non-linear relationships
- It is robust to multicollinearity
- It performs well with noisy sensor data

### 4.4.4 Calibration Performance Metrics

- RMSE (Root Mean Square Error)
- MAE (Mean Absolute Error)
- Percentage error reduction

Calibration reduced average error by approximately 30%, improving low-cost sensor reliability substantially.

## 4.5 Predictive Modelling Framework

Four models were implemented for comparative analysis.

### 4.5.1 ARIMA (Baseline Statistical Model)

ARIMA (Auto Regressive Integrated Moving Average) was used as a traditional time-series baseline model.

Advantages:

- Simple and interpretable
- Effective for stationary time-series

Limitations:

- Cannot model complex non-linear dependencies
- Does not incorporate spatial features

### 4.5.2 Long Short-Term Memory (LSTM)

LSTM networks are specialized recurrent neural networks designed to capture long-term dependencies.

Architecture:

- Input layer
- Two LSTM layers (64 and 32 units)
- Dropout (0.2)
- Dense output layer

Activation: ReLU

Optimizer: Adam

Loss Function: Mean Squared Error

LSTM captures temporal dependencies such as delayed pollutant accumulation.

#### 4.5.3 Transformer-Based Time Series Model

Transformers use **self-attention mechanisms** to model long-range dependencies without recurrent connections.

Advantages:

- Handles long sequences efficiently
- Captures global dependencies
- Parallelizable during training

Multi-head attention layers were implemented to weigh pollutant interactions dynamically.

#### 4.5.4 Spatio-Temporal Graph Neural Network (ST-GNN)

To model spatial dependencies between sensor nodes:

- Sensors were represented as graph nodes
- Edges were defined based on geographic proximity and traffic connectivity
- Graph convolution layers captured spatial correlations

- LSTM layers modelled temporal dynamics
- This architecture enabled modelling of pollution dispersion across neighbourhoods.

#### 4.6 Model Training and Validation

- Dataset split: 80% training, 20% testing
- Cross-validation: 5-fold
- Early stopping used to prevent overfitting
- Hyperparameters optimized using grid search

#### 4.7 Evaluation Metrics

The following metrics were used:

**Root Mean Square Error (RMSE):**

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

Measures prediction deviation magnitude.

**Mean Absolute Error (MAE)**

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Provides average absolute prediction error.

**R-Squared (R<sup>2</sup>)**

Indicates the proportion of variance explained by the model.

### 5. Results

#### 5.1 Comparative Model Performance

The predictive models were evaluated on the test dataset.

Model	RMSE	MAE	R <sup>2</sup>
ARIMA	15.3	11.7	0.62
LSTM	10.4	7.9	0.78
Transformer	8.7	6.2	0.85
ST-GNN	7.9	5.8	0.88

### Interpretation

- ARIMA performed poorly due to inability to model non-linear patterns.
- LSTM significantly improved accuracy by capturing temporal dependencies.
- Transformer outperformed LSTM by modelling long-range interactions.
- ST-GNN achieved the highest accuracy by incorporating spatial relationships.

This confirms that pollution dynamics are both temporal and spatial.

### 5.2 Sensor Calibration Results

Before calibration:

- Average RMSE: 18.2
- High variance during high humidity conditions

After calibration using Random Forest:

- 30% average error reduction

- Improved stability during seasonal transitions
- RMSE reduced to 12.5

Calibration especially improved performance during winter inversion conditions, when low-cost sensors previously overestimated PM2.5.

### 5.3 Explainability Analysis Using SHAP

To ensure model transparency, SHAP values were computed for the best-performing ST-GNN model.

#### 5.3.1 Feature Importance Ranking

Feature	Contribution (%)
Traffic Density	35%
Temperature	16%
Humidity	12%
Wind Speed	18%
PM10	10%
Time of Day	9%

#### 5.3.2 Key Insights

1. Traffic density was the most influential predictor, validating urban emission hypotheses.
2. Wind speed negatively correlated with PM2.5 levels, confirming dispersion effects.
3. Humidity increased PM2.5 readings due to hygroscopic particle growth.
4. Peak traffic hours showed amplified pollution spikes.

#### 5.3.3 Policy-Relevant Interpretation

Explainability findings suggest:

- Traffic control policies during peak hours could significantly reduce PM2.5.
- Urban planning should consider ventilation corridors.
- Meteorological forecasts can enhance predictive warnings.

SHAP visualization confirmed that model decisions align with environmental science principles, enhancing trustworthiness.

## 6. Discussion

The findings confirm that deep learning and graph-based models outperform traditional statistical approaches in air quality prediction. The integration of traffic and meteorological data enhances prediction reliability. Furthermore, explainable AI provides actionable insights for policymakers, improving trust in AI-based systems.

This framework can be deployed in smart cities to enable proactive interventions such as traffic rerouting, public advisories, and industrial regulation adjustments.

## 7. Implications for Smart Cities

1. Real-time pollution alerts.
2. Data-driven urban planning.
3. Predictive environmental governance.
4. Reduced healthcare burden via early warnings.
5. Scalable architecture for multiple cities.

## 8. Limitations

- Model performance depends on data quality.
- Calibration may require periodic retraining.
- Transformer models demand high computational resources.

## 9. Future Research Directions

1. Cross-city transfer learning.
2. Integration with satellite remote sensing.
3. Edge computing deployment for real-time inference.
4. Reinforcement learning for adaptive pollution control strategies.

## 10. Conclusion

This study demonstrates that machine learning, deep learning, and explainable AI provide powerful tools for smart city air quality monitoring. The proposed integrated framework enhances predictive accuracy, sensor reliability, and policy interpretability. Deploying such systems can significantly improve environmental sustainability and public health outcomes in urban environments.

## References

- [1] N. Castell et al., “Can commercial low-cost sensor platforms contribute to air quality monitoring and exposure estimates?,” *Environment International*, vol. 99, pp. 293–302, 2017, doi: 10.1016/j.envint.2016.12.007.
- [2] W. Jiao et al., “Community air sensor network (CAIRSENSE) project: Evaluation of low-cost sensor performance,” *Atmospheric Measurement Techniques*, vol. 9, no. 11, pp. 5281–5292, 2016, doi: 10.5194/amt-9-5281-2016.
- [3] S. M. Lundberg and S. I. Lee, “A unified approach to interpreting model predictions,” *Advances in Neural Information Processing Systems*, vol. 30, pp. 4765–4774, 2017.
- [4] A. Vaswani et al., “Attention is all you need,” *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [5] World Health Organization, “Air pollution and health,” 2021. [Online]. Available:

<https://www.who.int>

- [6] Y. Wu et al., “Graph neural networks for spatio-temporal modelling of air quality data,” *Journal of Machine Learning Research*, vol. 21, pp. 1–30, 2020.
- [7] L. Zhang et al., “Deep learning models for air quality prediction: A systematic review,” *Environmental Science & Technology*, vol. 55, no. 7, pp. 4087–4103, 2021, doi: 10.1021/acs.est.0c06770.
- [8] H. Chourabi et al., “Understanding smart cities: An integrative framework,” in *Proc. 45th Hawaii Int. Conf. System Sciences*, 2012, pp. 2289–2297.
- [9] R. Kitchin, “The real-time city? Big data and smart urbanism,” *GeoJournal*, vol. 79, no. 1, pp. 1–14, 2014.
- [10] M. Batty et al., “Smart cities of the future,” *European Physical Journal Special Topics*, vol. 214, pp. 481–518, 2012.
- [11] S. Zanella et al., “Internet of Things for smart cities,” *IEEE Internet of Things Journal*, vol. 1, no. 1, pp. 22–32, 2014.
- [12] A. Gubbi et al., “Internet of Things (IoT): A vision, architectural elements, and future directions,” *Future Generation Computer Systems*, vol. 29, no. 7, pp. 1645–1660, 2013.
- [13] J. Manyika et al., “Big data: The next frontier for innovation,” McKinsey Global Institute, 2011.
- [14] C. Perera et al., “Context-aware computing for the Internet of Things: A survey,” *IEEE Communications Surveys & Tutorials*, vol. 16, no. 1, pp. 414–454, 2014.
- [15] T. Yigitcanlar et al., “Smart cities: An effective urban development and management model?,” *Australian Planner*, vol. 52, no. 1, pp. 27–34, 2015.
- [16] X. Zheng et al., “Urban computing: Concepts, methodologies, and applications,” *ACM Transactions on Intelligent Systems*, vol. 5, no. 3, pp. 1–55, 2014.
- [17] M. Mohammadi et al., “Deep learning for IoT big data and streaming analytics,” *IEEE Communications Magazine*, vol. 56, no. 10, pp. 28–35, 2018.
- [18] J. Lee et al., “Predictive analytics for smart cities,” *IEEE Access*, vol. 6, pp. 68763–68773, 2018.
- [19] A. Ahmed et al., “Air quality forecasting using machine learning algorithms,” *Atmospheric Environment*, vol. 225, 2020.
- [20] H. Guo et al., “A review of deep learning models for air quality prediction,” *Environmental Pollution*, vol. 265, 2020.
- [21] Y. Li et al., “Traffic prediction in smart cities using machine learning,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 10, pp. 1–12, 2020.
- [22] Z. Lv et al., “Big data analytics for smart cities,” *Future Generation Computer Systems*, vol. 81, pp. 580–591, 2018.
- [23] S. Rathore et al., “Real-time big data analytical architecture for remote sensing application,” *IEEE Journal of Selected Topics in Applied Earth Observations*, vol. 8, no. 10, pp. 4610–4621, 2015.
- [24] M. Chen et al., “Machine-to-machine communications in ultra-dense networks,” *IEEE Communications Magazine*, vol. 53, no. 1, pp. 66–72, 2015.
- [25] A. Botta et al., “Integration of cloud computing and IoT: A survey,” *Future Generation Computer Systems*, vol. 56, pp. 684–700, 2016.
- [26] J. Lin et al., “A survey on Internet of Things: Architecture, enabling technologies,” *IEEE Internet of Things Journal*, vol. 4, no. 5, pp. 1125–1142, 2017.
- [27] P. Bellavista et al., “Convergence of MANET and WSN in IoT urban scenarios,” *IEEE Sensors*

Journal, vol. 13, no. 10, pp. 3558–3567, 2013.

[28] S. K. Ghosh et al., “Air pollution monitoring using IoT and machine learning,” *IEEE Sensors Letters*, vol. 3, no. 6, pp. 1–4, 2019.

[29] Y. Ma et al., “Spatiotemporal prediction of air quality using deep learning,” *IEEE Transactions on Big Data*, vol. 6, no. 2, pp. 1–10, 2020.

[30] D. Blei et al., “Latent Dirichlet allocation,” *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.