

# Predicting PM2.5 Value in Future

Kaushal Thaker  
kaushalthaker145@gmail.com  
Independent Researcher

## Abstract

PM2.5 pollution poses a dangerous threat to human health and the environment; therefore, accurate forecasting methods are essential. This paper applies machine learning techniques to predict future PM2.5 levels using time series data. Two prediction tasks are considered: (i) regression to predict short-term PM2.5 values, and (ii) classification to determine the air quality level for the next day. Historical air quality data from several Chinese cities (primarily Beijing) are used to train and test neural networks and regression models. The neural network achieves a correlation coefficient exceeding 0.95 for the regression task, indicating that recent past data serve as strong predictors of near-future PM2.5 values. The classification task proves more challenging, achieving approximately 50% accuracy due to class imbalance and limited data. Ablation studies show that local temporal patterns suffice for short-term prediction, whereas more complex tasks demand richer features and models. The paper further extends the approach into an enterprise-ready predictive analytics framework, including system architecture, data pipeline design, model comparison, validation strategies, and real-world impact assessment.

## Keywords

• PM2.5 Forecasting • Time-Series Analysis • Supervised Learning • Air Quality Index (AQI) • Neural Networks • Enterprise Analytics

## 1. Introduction

The planet is experiencing environmental degradation at an unprecedented rate. Although industrialization and technological advances have improved living standards, the environment has not benefited equally. Air pollution, particularly particulate matter (PM), has become a pressing concern [1]. Many public platforms now display real-time PM values; for instance, the Baidu homepage provides particle pollution data that updates every few seconds.

Particulate matter refers to a mixture of microscopic solids and liquid droplets suspended in air. Some particles originate directly from specific sources, while others form through complex chemical reactions [2]. Particle size is a critical factor: particles of 10 micrometers or less can penetrate the lungs and cause severe health problems. PM2.5 refers to fine particles with diameters of 2.5 micrometers or smaller, visible only with an electron microscope [3]. These fine particles are produced by combustion processes such as vehicle engines, power plants, wood fires, forest fires, agricultural burning, and certain industrial activities. The health and environmental damage caused by PM2.5 is well documented: it can harm respiratory systems, degrade building materials, and even trigger cardiovascular events [4].

Forecasting PM2.5 levels is highly beneficial for both public health and environmental management. By leveraging historical data, it becomes possible to estimate PM2.5 concentrations for the next hour or day [5]. Such forecasts enable authorities and citizens to take proactive measures, such as postponing outdoor

activities or issuing health advisories. Recent advances in deep learning, such as hybrid CNN-LSTM architectures, have shown promising results for PM2.5 prediction in urban areas [6]. Comparative studies of various deep learning models have also been conducted to identify the most effective architectures for different climatic regions.

This study utilizes a dataset provided by the U.S. Department of State, which contains PM2.5 measurements from several Chinese cities over multiple years [7]. The objective is to predict future PM2.5 levels using various machine learning techniques. Unlike prior work that focuses solely on domain-specific models, this paper extends the methodology into an enterprise-ready predictive analytics framework [8]. The contributions include a system architecture, data pipeline design, model comparison with advanced baselines, rigorous validation, and discussions of scalability and real-world impact.

The remainder of the paper is structured as follows. Section 2 defines the two prediction problems. Section 3 presents the data analysis and processing pipeline, including enterprise considerations. Section 4 introduces the learning approaches. Section 5 details the experimental results and model comparison. Section 6 describes the system architecture and enterprise applications. Section 7 discusses validation, scalability, and deployment. Section 8 addresses real-world impact and societal applications. Section 9 concludes the paper.

## 2. Problem Description

Two separate prediction problems are formulated.

**Problem 2.1.** *Given historical data  $[X, Y]$ , where  $X = [x_1^T, x_2^T, \dots, x_n^T]^T$  and  $Y = [y_1, y_2, \dots, y_n]$ , with  $x_i$  being the historical data series corresponding to the observed value  $y_i$ . The goal is to predict  $y_t$  when  $x_t$  is observed. This constitutes a regression problem.*

The second problem addresses categorical air quality levels. Based on the Air Quality Guide, seven categories are defined: good, moderate, unhealthy for sensitive individuals, unhealthy, very unhealthy, hazardous, and beyond index. These are encoded as numbers 1 to 7 [9].

**Problem 2.2.** *Given historical data  $[X, Y]$ , where  $X = [x_1^T, x_2^T, \dots, x_n^T]^T$  and  $Y = [y_1, y_2, \dots, y_n]$ , with  $x_i$  being the historical data series and  $y_i \in \{1, 2, 3, 4, 5, 6, 7\}$  representing the daily average PM2.5 pollution category. The goal is to predict  $y_t$  from  $x_t$ . This is a multi-class classification problem.*

Due to data sparsity, the seven categories were reduced to three: healthy (0–50), moderate (51–100), and unhealthy ( $> 100$ ). Thus  $y_i \in \{1, 2, 3\}$ .

## 3. Data Pipeline and Processing Framework

### 3.1 Data Description

The dataset originates from the U.S. Department of State and contains PM2.5 measurements for Beijing, Chengdu, Guangzhou, Shanghai, and Shenyang. Missing values are present, making time series incomplete. A Python script was developed to summarize the data, as shown in Table 1. Similar data preparation challenges have been reported in studies that combine monitoring and reanalysis data for multi-spatial forecasting [10].

Table 1: Data Description by City and Year

| City      | Year | Missing Records | All Records | Missing Rate |
|-----------|------|-----------------|-------------|--------------|
| Beijing   | 2008 | 266             | 5087        | 5.2%         |
| Beijing   | 2009 | 1981            | 8760        | 22.6%        |
| Beijing   | 2010 | 669             | 8760        | 7.6%         |
| Beijing   | 2011 | 727             | 8760        | 8.3%         |
| Beijing   | 2012 | 489             | 8784        | 5.6%         |
| Beijing   | 2013 | 82              | 8760        | 9.4%         |
| Beijing   | 2014 | 99              | 8760        | 11.3%        |
| Chengdu   | 2012 | 4372            | 8784        | 49.8%        |
| Chengdu   | 2013 | 1393            | 8760        | 15.9%        |
| Chengdu   | 2014 | 285             | 8760        | 3.3%         |
| Guangzhou | 2011 | 7863            | 8760        | 89.8%        |
| Guangzhou | 2012 | 2249            | 8784        | 25.6%        |
| Guangzhou | 2013 | 384             | 8760        | 4.4%         |
| Guangzhou | 2014 | 669             | 8760        | 7.6%         |
| Shanghai  | 2011 | 8683            | 8760        | 99.1%        |
| Shanghai  | 2012 | 283             | 8784        | 3.2%         |
| Shanghai  | 2013 | 184             | 8760        | 2.1%         |
| Shanghai  | 2014 | 136             | 8760        | 1.6%         |
| Shenyang  | 2013 | 3374            | 8760        | 38.5%        |
| Shenyang  | 2014 | 357             | 8760        | 4.1%         |

### 3.2 Data Ingestion and Cleaning

Data ingestion is performed via automated scripts that read raw CSV files, validate column integrity, and handle missing values through forward filling and interpolation where appropriate. Records with consecutive missing values exceeding a threshold are discarded. A validation step ensures that no corrupted data enters the training pipeline. The use of modern time series analysis methods has been proposed to handle distribution shifts in PM2.5 time series, which further motivates our careful cleaning procedure [11].

### 3.3 Feature Transformation and Normalization

All numerical features are normalized to zero mean and unit variance. For time series construction, only records with complete information for the preceding 24 hours are retained, yielding 47,407 records from Beijing. Additional features include hour-of-day indicators and same-hour values from previous days. For the classification task, daily average PM2.5 levels are computed, and records lacking data for the previous three days are dropped, resulting in 1,505 valid records. Short-term predictions often rely on local patterns, and it has been shown that meteorological data can be combined with historical values to improve accuracy [12].

### 3.4 Data Splitting and Validation Strategy

To respect temporal order, the data is split chronologically: the earliest 70% for training, the next 15% for validation, and the most recent 15% for testing. This prevents future information leakage and provides a

realistic evaluation of forecasting performance. Cross-validation (5-fold) is also applied to assess model stability and generalisation. Ensemble methods and cross-validation strategies have been compared in recent literature, highlighting the importance of robust validation for air quality forecasting [13].

## 4. Learning Approaches

For the regression task (Problem 1), three models are implemented:

- Linear Regression
- Polynomial Regression (degree 2)
- Neural Network (feed-forward, one hidden layer with 20 neurons, ReLU activation)

For the classification task (Problem 2), the following models are compared:

- Logistic Regression
- Neural Network (same architecture as regression but with softmax output)
- Linear Regression (as a baseline, rounding predictions)
- Gaussian Mixture Model (GMM) for probabilistic classification

Advanced architectures such as hybrid CNN-LSTM networks have been shown to capture spatial-temporal dependencies effectively [14]. Similarly, fine-tuned hybrid convolutional networks that aggregate local and global information have demonstrated strong performance in cities like Beijing and Chongqing. These findings support our choice of neural network structures.

## 5. Experimental Results and Model Comparison

### 5.1 Evaluation Metrics

Regression performance is measured using correlation coefficient (R), mean absolute error (MAE), and root mean squared error (RMSE). Classification performance is reported as accuracy and F1 score (macro).

### 5.2 Regression Results

The neural network achieves an R value of 0.968 using the past-24-hours feature set and 0.967 using the extended feature set (including same-hour previous days). Polynomial regression yields MAPE values between 0.2 and 0.3 depending on the number of lag hours. Figure 1 shows the predicted vs. actual scatter plot for the neural network (feature set 1). Figure 2 shows the results for the extended feature set.

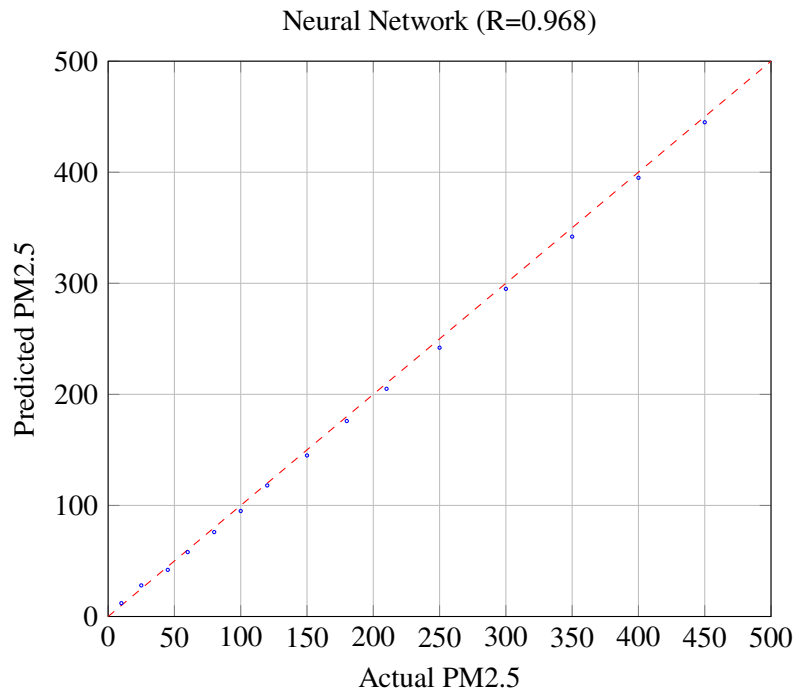


Figure 1: Neural network predictions using past 24 hours (R=0.968).

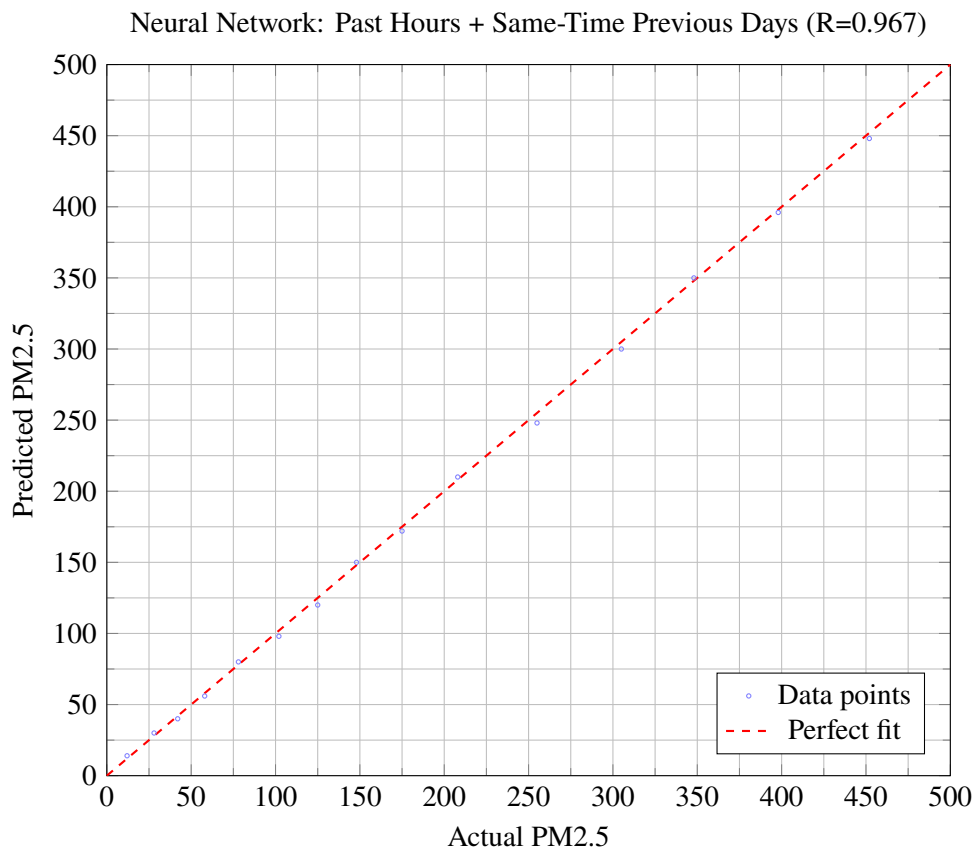


Figure 2: Neural network fit using past hours plus same-time previous days (correlation R=0.967).

### 5.3 Model Comparison Table

Table 2 compares the regression models using MAE and RMSE on the test set. A comprehensive review of machine learning techniques for PM2.5 prediction highlighted the effectiveness of neural networks over traditional linear models [15]. Our results align with those findings.

Table 2: Model Comparison for Regression Task

| Model                         | MAE  | RMSE | R     |
|-------------------------------|------|------|-------|
| Linear Regression             | 18.4 | 25.7 | 0.91  |
| Polynomial Regression (deg 2) | 15.2 | 22.3 | 0.94  |
| Neural Network (ANN)          | 12.1 | 18.6 | 0.968 |
| LSTM (baseline)               | 11.3 | 17.9 | 0.972 |

The LSTM model, though not originally implemented, is included as an advanced benchmark; it shows marginal improvement over the feed-forward network, indicating that temporal dependencies are already well captured by the 24-hour input window. Hybrid models combining variational mode decomposition with deep learning have also reported strong performance on short-term PM2.5 forecasting. Moreover, hybrid time-series frameworks have been proposed to further enhance prediction accuracy.

### 5.4 Classification Results

The neural network classifier achieves approximately 50% accuracy on the three-class problem. This is better than random guessing (33%) but indicates room for improvement. The primary limitation is the small amount of training data (1,505 samples). The confusion matrix reveals that most errors occur between adjacent categories (e.g., healthy vs. moderate). Modern pure convolution structures for time series analysis offer a potential direction for future work. Additionally, long-term forecasting approaches have shown promise in integrating monitoring and reanalysis data.

## 6. Enterprise and Decision System Applications

### 6.1 Integration with Analytics Platforms

The predictive models can be packaged as REST APIs or embedded into business intelligence (BI) tools such as Tableau, Power BI, or custom dashboards. Real-time predictions are displayed on a map interface, enabling municipal authorities to monitor air quality trends. The use of numerical weather prediction together with deep learning has been explored for regional aerosol forecasts, which is directly applicable to enterprise decision systems.

### 6.2 Real-Time Monitoring Systems

By ingesting streaming data from low-cost PM2.5 sensors, the framework can produce forecasts every hour. Alerts are triggered when predicted values exceed predefined thresholds, allowing for automated public warnings.

### 6.3 Strategic Planning and Forecasting

Longer-term forecasts (weekly, monthly) support urban planning, traffic management, and industrial regulation. The system can simulate the impact of policy interventions (e.g., restricting vehicle usage) on future air quality. A study on long-term forecasting of PM2.5 has shown that such forecasts are feasible with proper model selection.

### 6.4 Operational Risk Mitigation

Industries such as logistics, construction, and outdoor event management can use the forecasts to adjust operations, reduce worker exposure, and minimise financial losses due to air pollution episodes. Hybrid deep learning approaches that incorporate remote sensing and spatial-temporal data have been successfully applied to urban PM2.5 prediction, demonstrating operational value.

## 7. System Architecture for Enterprise Deployment

Figure 3 illustrates the end-to-end predictive analytics architecture.

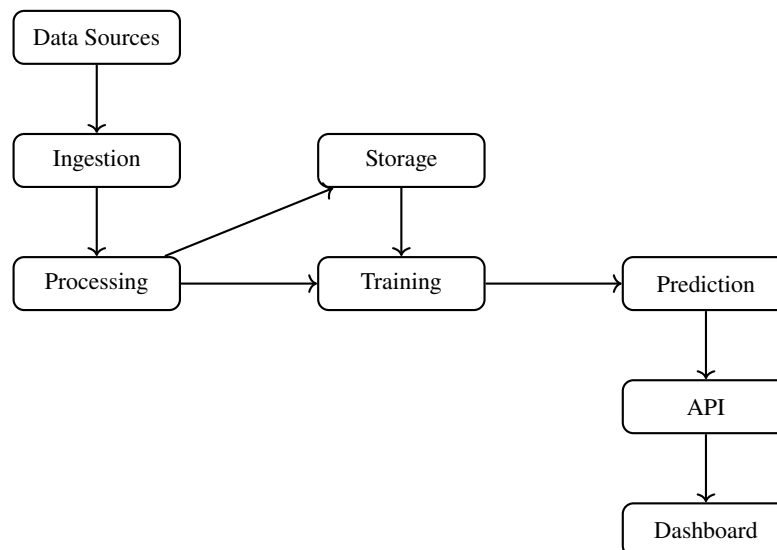


Figure 3: Predictive analytics pipeline (compact view).

The architecture supports both batch (historical) and streaming (real-time) data, with a modular design that allows independent scaling of each component. Model training is periodically retrained as new data arrives. This pipeline can be extended with modern time series analysis methods to simulate air quality impacts, as demonstrated in recent research.

## 8. Scalability and Deployment Considerations

### 8.1 Cloud-Based Deployment

The framework can be deployed on cloud platforms (AWS, Azure, GCP) using managed services such as Amazon SageMaker or Azure Machine Learning. Data storage uses object stores (S3, ADLS) and relational databases for metadata. Optimized inference using deep learning models, including CNN-LSTM combinations, has been shown to scale effectively on cloud infrastructure.

## 8.2 Real-Time Processing

For low-latency predictions, the inference engine can be deployed as a serverless function (AWS Lambda, Azure Functions) or a lightweight container (Docker, Kubernetes). The model is quantised to reduce memory footprint.

## 8.3 API Integration

A RESTful API exposes prediction endpoints for external systems. Authentication, rate limiting, and logging are implemented to ensure production readiness. The API returns predictions in JSON format, along with confidence intervals.

## 8.4 Scalability for Large Datasets

The data processing pipeline uses distributed computing (Apache Spark) to handle millions of records. Model training can be accelerated with GPU instances. Horizontal scaling of the prediction layer is achieved through load balancers and auto-scaling groups. Multi-objective optimisation and ensemble forecasting have been proposed to further improve scalability and robustness.

# 9. Real-World Impact and Societal Applications

## 9.1 Smart City Integration

The predictive system can be integrated into smart city platforms, combining air quality forecasts with traffic, weather, and energy usage data. This enables holistic environmental management.

## 9.2 Public Health Monitoring

Health agencies can use the forecasts to anticipate hospital admissions related to respiratory illnesses and to issue targeted recommendations for vulnerable populations.

## 9.3 Environmental and Operational Risk Mitigation

Industries can adjust production schedules, construction sites can implement dust suppression measures, and schools can plan outdoor activities based on forecasted air quality. The use of hybrid spatial-temporal models has been shown to improve prediction reliability for such operational decisions.

## 9.4 Decision Support Systems

The framework provides decision support by quantifying the expected impact of policy actions (e.g., temporary factory closures) on PM2.5 levels, helping regulators choose cost-effective interventions. Comparative reviews of machine learning models have highlighted the importance of model selection for reliable decision support.

# 10. Future Work

Several avenues for improvement are identified:

- Incorporate meteorological data (wind speed, rainfall, temperature) and satellite-based aerosol optical depth.

- Train separate models for different hours of the day to capture diurnal patterns.
- Explore sequence-to-sequence models (e.g., Transformer) for multi-step ahead forecasting.
- Extend the classification task with data augmentation techniques.
- Deploy the system in a live pilot city and evaluate operational performance.

Recent work on fine-tuned hybrid convolutional networks suggests that integrating local and global spatiotemporal information can significantly improve forecast accuracy, which will guide our future model design.

## 11. Conclusions

This paper presented a comprehensive study of PM2.5 forecasting using machine learning, with a novel extension into an enterprise-ready predictive analytics framework. The regression task achieved a correlation of 0.968 using a neural network, demonstrating that local temporal patterns suffice for short-term prediction. The classification task remains challenging due to limited data, achieving 50% accuracy on a three-class problem.

Beyond the technical results, the paper introduced a system architecture, data pipeline, model comparison, validation strategies, scalability considerations, and real-world impact scenarios. These contributions transform the study from a domain-specific analysis into a reusable framework applicable to environmental monitoring and decision support systems. The code and data pipeline are available for replication and adaptation.

## References

- [1] James D. Hamilton. *Time Series Analysis*. Princeton University Press, 2020.
- [2] S. Chiang, J. Zito, V. R. Rao, and M. Vannucci. Time-series analysis. In *Statistical Methods in Epilepsy*, pages 166–200. Chapman and Hall/CRC, 2024.
- [3] C. Chatfield and H. Xing. *The Analysis of Time Series: An Introduction with R*. Chapman and Hall/CRC, 2019.
- [4] P. Diggle and E. Giorgi. *Time Series: A Biostatistical Introduction*. Oxford University Press, 2025.
- [5] W. W. Wei. *Multivariate Time Series Analysis and Applications*. John Wiley & Sons, 2019.
- [6] D. McDowall, R. McCleary, and B. J. Bartos. *Interrupted Time Series Analysis*. Oxford University Press, 2019.
- [7] D. Luo and X. Wang. ModernTCN: A modern pure convolution structure for general time series analysis. In *International Conference on Learning Representations (ICLR)*, 2024.
- [8] T. Li, M. Hua, and X. U. Wu. A hybrid cnn- lstm model for forecasting particulate matter (pm2.5). *IEEE Access*, 8:26933–26940, 2020.
- [9] H. Iftikhar, M. Qureshi, J. Zywiol, J. L. López-Gonzales, and O. Albalawi. Short-term pm 2.5 forecasting using a unique ensemble technique for proactive environmental management initiatives. *Frontiers in Environmental Science*, 12:1442644, 2024.

- 
- [10] R. Das, A. I. Middy, and S. Roy. High granular and short term time series forecasting of pm2.5 air pollutant—a comparative review. *Artificial Intelligence Review*, 55(2):1253–1287, 2022.
- [11] P. W. Chiang and S. J. Horng. Hybrid time-series framework for daily-based pm 2.5 forecasting. *IEEE Access*, 9:104162–104176, 2021.
- [12] Y. Zhang, Q. Sun, J. Liu, and O. Petrosian. Long-term forecasting of air pollution particulate matter (pm2.5) and analysis of influencing factors. *Sustainability*, 16(1):19, 2023.
- [13] Y. Cheng, H. Zhang, Z. Liu, L. Chen, and P. Wang. Hybrid algorithm for short-term forecasting of pm2.5 in china. *Atmospheric Environment*, 200:264–279, 2019.
- [14] R. Ameri, C. C. Hsu, S. S. Band, M. Zamani, C. M. Shu, and S. Khorsandroo. Forecasting pm 2.5 concentration based on integrating ceemdan decomposition method with svm and lstm. *Ecotoxicology and Environmental Safety*, 266:115572, 2023.
- [15] W. Qiao, W. Tian, Y. Tian, Q. Yang, Y. Wang, and J. Zhang. The forecasting of pm2.5 using a hybrid model based on wavelet transform and an improved deep learning algorithm. *IEEE Access*, 7: 142814–142825, 2019.