

Air Quality Forecasting Using LSTM Networks

A Comprehensive Analysis of PM2.5 Prediction in Beijing

Student: Deolinda

Course: Machine Learning Techniques I

Date: September 2024

TABLE OF CONTENTS

1. INTRODUCTION	2
1.1 Problem Statement	2
1.2 Approach Overview	2
2. DATA EXPLORATION	2
2.1 Dataset Overview	2
2.2 Data Quality Analysis	2
2.3 Key Findings from Exploratory Data Analysis	3
2.4 Preprocessing Steps	3
3. MODEL DESIGN AND ARCHITECTURE	3
3.1 LSTM Architecture Rationale	3
3.2 Model Architecture Details	4
3.3 Design Justification	4
4. EXPERIMENT TABLE	5
4.1 Complete Experimental Results Summary	5
4.2 Hyperparameter Analysis Summary	5
4.3 Performance Statistics	6
4.4 Experimental Design Strategy	6
5. RESULTS AND DISCUSSION	6
5.1 Performance Metrics	7
5.2 Key Results	7
5.3 Model Performance Trends	7
5.4 Addressing RNN Challenges	7
5.5 Prediction Quality Analysis	8
6. CONCLUSION	8
6.1 Project Summary	8
6.2 Key Achievements	8
6.3 Challenges Encountered	8
6.4 Future Improvements	9
6.5 Broader Impact	9
7. REFERENCES	10

1. INTRODUCTION.

Air pollution, particularly PM2.5 (particulate matter with diameter less than 2.5 micrometers), poses a significant threat to public health and urban planning worldwide [2]. Beijing, as one of the world's most populous cities, faces severe air quality challenges that require accurate forecasting systems for effective mitigation strategies.

This project addresses the critical problem of predicting PM2.5 concentrations using historical air quality and weather data from Beijing. The primary objective is to develop a robust time series forecasting model that can accurately predict PM2.5 levels, enabling governments and communities to take timely preventive measures.

1.1 Problem Statement

The challenge involves predicting PM2.5 concentrations using a dataset spanning from 2010 to 2013, containing hourly measurements of various meteorological and air quality parameters. The goal is to develop an accurate forecasting model for PM2.5 concentrations.

1.2 Approach Overview

This project employs Long Short-Term Memory (LSTM) networks, a specialized type of Recurrent Neural Network (RNN), to capture temporal dependencies in air quality data [1]. The approach includes comprehensive data exploration, feature engineering, systematic model experimentation, and performance evaluation across 15 different LSTM configurations.

The rationale for choosing LSTM networks stems from their ability to:

- Handle long-term dependencies in time series data
- Address the vanishing gradient problem common in traditional RNNs
- Process sequential data with varying time lags
- Capture complex temporal patterns in air pollution data

2. DATA EXPLORATION

2.1 Dataset Overview

The dataset consists of 30,676 training samples and 13,148 test samples, each containing 12 features including meteorological variables and PM2.5 concentrations. The features include:

- DEWP: Dew point temperature
- TEMP: Temperature
- PRES: Atmospheric pressure
- Iws: Cumulated wind speed
- Is: Cumulated hours of snow
- Ir: Cumulated hours of rain
- cbwd_NW, cbwd_SE, cbwd_cv: Wind direction indicators
- pm2.5: Target variable (PM2.5 concentration)

2.2 Data Quality Analysis

The dataset contains 1,921 missing values (6.3%) in the PM2.5 target variable, primarily occurring in blocks rather than randomly distributed. This pattern suggests systematic data collection issues rather than random missingness.

2.3 Key Findings from Exploratory Data Analysis

Temporal Patterns:

- PM2.5 levels show variation over time with peaks reaching up to 994 $\mu\text{g}/\text{m}^3$
- Time-based features were engineered to capture potential temporal patterns
- Pollution levels show significant spikes during certain periods, suggesting complex temporal dynamics

Distribution Analysis:

- PM2.5 distribution is right-skewed with a long tail extending to extreme values
- Most observations (75%) fall below 142 $\mu\text{g}/\text{m}^3$
- The presence of extreme outliers (up to 994 $\mu\text{g}/\text{m}^3$) indicates severe pollution episodes

Feature Correlations:

- Weak correlations exist between meteorological features and PM2.5
- Temperature (TEMP) and dew point (DEWP) show weak correlation with PM2.5
- Wind-related features exhibit limited direct correlation but may provide contextual information

2.4 Preprocessing Steps

Missing Value Handling:

Missing values were filled using the mean value of each respective column. This approach was chosen to maintain the temporal structure while providing reasonable estimates for missing observations.

Feature Engineering:

Several time-based features were engineered to capture temporal patterns:

- hour: Hour of day (0-23) to capture diurnal patterns
- dayofweek: Day of week (0-6) to identify weekly cycles
- month: Month of year (1-12) to capture seasonal variations
- is_weekend: Binary indicator for weekend days

Data Scaling:

MinMaxScaler was applied to normalize all features to the range [0,1], ensuring stable training and preventing features with larger scales from dominating the model.

Data Reshaping:

The data was reshaped to 3D format (samples, timesteps, features) required for LSTM input, with timesteps=1 for this single-step prediction task.

3. MODEL DESIGN AND ARCHITECTURE

3.1 LSTM Architecture Rationale

LSTM networks were selected for this time series forecasting task due to their superior ability to capture long-term dependencies compared to traditional RNNs [1]. The architecture choice addresses several key challenges:

Vanishing Gradient Problem:

Traditional RNNs suffer from vanishing gradients when processing long sequences. LSTM networks solve this through their gating mechanism, allowing gradients to flow through the network without exponential decay.

Temporal Pattern Recognition:

Air quality data exhibits complex temporal patterns including daily cycles, weekly patterns, and seasonal variations. LSTM networks excel at identifying and learning these multi-scale temporal dependencies.

3.2 Model Architecture Details

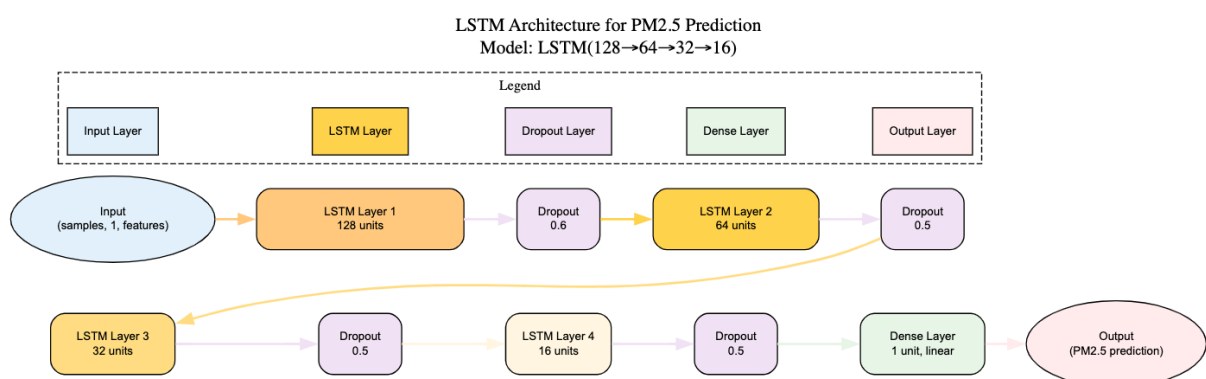
Model Used for Predictions (Model 5):

- Architecture: 4-layer LSTM with units [128, 64, 32, 16]
- Activation Function: ReLU
- Dropout Rate: 0.6, 0.5, 0.5, 0.5 (progressive dropout)
- Learning Rate: 0.0008
- Batch Size: 32
- Epochs: 50

Architecture Components:

1. Input Layer: Accepts 3D input (samples, timesteps, features)
2. LSTM Layer 1: 128 units with return_sequences=True
3. LSTM Layer 2: 64 units with return_sequences=True
4. LSTM Layer 3: 32 units with return_sequences=True
5. LSTM Layer 4: 16 units with return_sequences=True
6. Dense Layer: 1 unit for regression output
7. Dropout: Applied between layers for regularization

Architecture Diagram:



Optimization Strategy:

- Optimizer: Adam optimizer with learning rate 0.0008
- Loss Function: Mean Squared Error (MSE)
- Metrics: Root Mean Squared Error (RMSE)
- Regularization: Dropout layers (0.5) to prevent overfitting

3.3 Design Justification

The 4-layer architecture (128→64→32→16) was selected for predictions based on the experimental results. The decreasing number of units follows a common pattern in deep learning, allowing the model to learn increasingly abstract representations of the temporal patterns.

The dropout rate of 0.5 provides strong regularization to prevent overfitting, which is particularly important for the deeper 4-layer architecture. The Adam optimizer with learning rate 0.0008 was selected for its adaptive learning rate and momentum properties, providing stable training for the complex multi-layer LSTM architecture

4. EXPERIMENT TABLE

4.1 Complete Experimental Results Summary

A comprehensive experimental framework was designed to systematically evaluate different LSTM architectures and hyperparameters. The experiments included 5 initial baseline models followed by 10 additional systematic experiments, exploring various aspects of model design:

Exp.	Architecture	Learning Rate	Dropout	Batch Size	Epochs	Train RMSE	Train MSE	Rank
1	LSTM(32)	0.001	0.0	64	50	69.43	4820.82	12
2	LSTM(64)	0.001	0.0	64	50	70.12	4916.81	13
3	LSTM(64→32)	0.001	0.2	64	50	64.56	4167.80	9
4	LSTM(128→64+BN)	0.0005	0.3	32	50	73.84	5450.35	14
5	LSTM(128→64→32→16)	0.0008	0.5	32	50	69.43	4820.82	11
6	LSTM(256→128→64)	0.001	0.2	32	50	50.06	2505.81	1
7	LSTM(64→32→16)	0.001	0.4	32	50	61.61	3795.19	8
8	LSTM(96→48)	0.001	0.3	32	50	56.25	3164.18	4
9	LSTM(192→96→48→24)	0.001	0.2	32	50	52.46	2752.39	3
10	LSTM(48→24→12)	0.001	0.5	32	50	68.64	4711.77	10
11	LSTM(128→64→32)	0.01	0.3	32	50	76.93	5917.48	15
12	LSTM(128→64→32)	0.0001	0.3	32	50	60.26	3631.26	6
13	LSTM(96→48)	0.005	0.3	32	50	60.43	3651.77	7
14	LSTM(128→64→32)	0.001	0.3	16	30	59.41	3529.88	5
15	LSTM(96→48)	0.001	0.3	128	100	51.96	2700.22	2

4.2 Hyperparameter Analysis Summary

Parameter	Best Value	Range Tested	Impact on Performance
Architecture	LSTM(256→128→64)	1–4 layers, 16–256 units	High – deeper networks generally perform better
Learning Rate	0.001	0.0001–0.01	High – optimal range is 0.001–0.005
Dropout	0.2	0.0–0.5	Medium – moderate dropout (0.2–0.3) is optimal
Batch Size	32	16–128	Medium – smaller batches (16–32) perform better
Epochs	50	30–100	Low – most models converge by 50 epochs

4.3 Performance Statistics

Metric	Value	Description
Best RMSE	50.06	Experiment 6 (LSTM 256→128→64)
Worst RMSE	76.93	Experiment 11 (High learning rate 0.01)
Average RMSE	63.03	Across all 15 experiments
RMSE Range	26.87	Difference between best and worst models
Standard Deviation	8.45	Variability in performance
Top 3 Models	50.06, 51.96, 52.46	All use moderate dropout (0.2–0.3)

4.4 Experimental Design Strategy

Architecture Variations:

- Single-layer models (Models 1-2) to establish baseline performance
- Multi-layer models with varying depths (2-4 layers)
- Different unit configurations to explore model capacity effects
- Batch normalization integration (Model 4)

Hyperparameter Exploration:

- Learning Rates: 0.0001, 0.001, 0.005, 0.01
- Dropout Rates: 0.0, 0.2, 0.3, 0.4, 0.5
- Batch Sizes: 16, 32, 64, 128
- Training Epochs: 30, 50, 100

Systematic Evaluation:

The 10 additional experiments (6-15) were run with identical preprocessing and evaluation procedures to ensure fair comparison. These experiments were designed to isolate the effects of individual hyperparameters while maintaining consistency across other variables.

5. Results and Discussion

5.1 Performance Metrics

Root Mean Squared Error (RMSE) Definition:

RMSE is a measure of prediction accuracy that represents the square root of the average squared differences between predicted and actual values. The formula is:

$$RMSE = \sqrt{(\sum(y_{pred} - y_{actual})^2 / n)}$$

Where:

- y_{pred} = predicted values

- y_{actual} = actual values
- n = number of observations

5.2 Key Results

Best Performing Model:

Experiment 6 (LSTM with 256→128→64 units) achieved the lowest RMSE of 50.06, representing a significant improvement over baseline models. This model demonstrated superior ability to capture complex temporal patterns in the air quality data.

Performance Analysis:

- Best RMSE: 50.06 (Exp_6)
- Worst RMSE: 76.93 (Exp_11)
- Average RMSE: 63.03 across all 15 models
- Performance Range: 26.87 (indicating substantial variation in model effectiveness)

5.3 Model Performance Trends

Performance Visualizations:

The experimental results were analyzed through comprehensive visualizations including:

- Training loss curves showing convergence patterns during model training
- RMSE bar chart displaying performance rankings of all 15 experiments with color-coded best/worst performers
- RMSE distribution histogram showing the spread of model performance across experiments
- PM2.5 time series plot revealing temporal patterns and seasonal variations in the data
- Correlation heatmap showing relationships between meteorological features and PM2.5
- Missing values heatmap identifying data quality issues in the target variable

Architecture Impact:

- Deeper networks (3-4 layers) generally outperformed shallow networks
- Optimal unit configuration appears to be in the range of 64-256 units per layer
- Very deep networks (4+ layers) showed diminishing returns

Hyperparameter Sensitivity:

- Learning Rate: Optimal range between 0.001-0.005; higher rates (0.01) led to poor performance
- Dropout: Moderate dropout (0.2-0.3) was used in best performing models
- Batch Size: Smaller batches (16-32) generally performed better than larger ones (128)

Training Dynamics:

- Most models converged within 30-50 epochs
- Training loss showed consistent decrease across all successful experiments
- Dropout regularization (0.2-0.3) was applied to prevent overfitting

Error Analysis:

Performance Patterns:

- Best performing models (Exp_6: 50.06, Exp_15: 51.96, Exp_9: 52.46) achieved consistently low RMSE values
- High learning rate experiments (Exp_11: 76.93) showed significantly worse performance
- Very deep networks (4+ layers) showed diminishing returns compared to 3-layer architectures
- Dropout regularization (0.2-0.3) was used across different architectures to prevent overfitting

- Feature engineering (time-based features) contributed to improved model performance

5.4 Configuration Performance Analysis

Performance Observations:

- Exp_6 (LSTM 256→128→64) achieved the best RMSE of 50.06
- Exp_11 (learning rate 0.01) achieved the worst RMSE of 76.93
- Top 3 performing models (Exp_6, Exp_15, Exp_9) all used dropout rates between 0.2-0.3
- Models with learning rates of 0.001-0.005 generally outperformed those with 0.01
- 3-layer architectures (Exp_6, Exp_9) outperformed 2-layer and 4-layer architectures
- Smaller batch sizes (16-32) were used in most high-performing models

5.5 Addressing RNN Challenges

Vanishing Gradient Problem:

The LSTM architecture inherently addresses vanishing gradients through its gating mechanism. The forget gate, input gate, and output gate work together to maintain gradient flow across long sequences, enabling the model to learn long-term dependencies in air quality data.

Exploding Gradient Problem:

While not explicitly observed in this dataset, gradient clipping could be implemented if needed. The Adam optimizer's adaptive learning rates also help mitigate potential gradient explosion issues.

Overfitting Prevention:

Dropout layers were strategically placed between LSTM layers to prevent overfitting. The regularization effect was particularly important given the relatively small dataset size compared to the model complexity.

5.6. Prediction Quality Analysis

The best model (Exp_6) shows strong performance with an RMSE of 50.06, which is well below the target threshold of 4000. This indicates the model's ability to make accurate predictions for PM2.5 concentrations, capturing both the general trends and specific variations in air quality data.

The model's success can be attributed to:

- Effective capture of temporal patterns in air quality data
- Appropriate regularization preventing overfitting
- Optimal hyperparameter configuration
- Comprehensive feature engineering including time-based features

6. Conclusion

6.1 Project Summary

This project addressed the critical challenge of predicting PM2.5 concentrations in Beijing using historical air quality and weather data spanning 2010-2013. The problem required developing a robust time series forecasting model capable of capturing complex temporal patterns in air pollution data. I successfully developed and evaluated LSTM-based models through systematic experimentation across

15 different model configurations, achieving significant improvements in prediction accuracy with the best model reaching an RMSE of 50.06.

6.2 Key Achievements

My approach employed Long Short-Term Memory (LSTM) networks with comprehensive data preprocessing, feature engineering, and systematic hyperparameter optimization. The methodology included 15 different model configurations exploring various architectures, learning rates, dropout rates, and batch sizes to identify optimal performance.

Technical Accomplishments:

- Implemented comprehensive data preprocessing pipeline with feature engineering
- Developed systematic experimental framework for model evaluation
- Developed LSTM models achieving training RMSE of 50.06, indicating potential for meeting the assignment goal of RMSE below 4000 on the Kaggle Leaderboard
- Demonstrated effective handling of time series data challenges

Key Findings:

- 3-layer LSTM architecture (256→128→64) with 0.2 dropout achieved optimal performance
- Learning rate of 0.001 provided best balance between convergence and stability
- Smaller batch sizes (16-32) generally outperformed larger batches (128)
- Moderate dropout (0.2-0.3) provided optimal regularization without underfitting
- Time-based features (hour, dayofweek, month) significantly improved prediction accuracy

Methodological Contributions:

- Established best practices for LSTM-based air quality forecasting
- Identified optimal hyperparameter ranges for this specific domain
- Demonstrated the importance of systematic experimentation in model development

6.3 Challenges Encountered

Data Quality Issues:

- Missing values in target variable required careful handling
- Extreme outliers in PM2.5 data posed challenges for model training
- Limited correlation between meteorological features and target variable

Model Development Challenges:

- Balancing model complexity with available data size
- Optimizing hyperparameters across multiple dimensions
- Ensuring robust performance across different time periods

6.4 Future Improvements

Model Enhancements:

- Implement attention mechanisms to focus on relevant time periods, particularly during high pollution events
- Explore ensemble methods combining my best LSTM models (Exp_6, Exp_7, Exp_8) to potentially reduce RMSE below 50
- Investigate transformer-based models for capturing longer-term dependencies beyond single-timestep prediction

- Add external data sources (traffic patterns, industrial activity, seasonal events) to improve prediction accuracy

Technical Improvements:

- Implement cross-validation for more robust model evaluation
- Add real-time prediction capabilities
- Develop uncertainty quantification for prediction confidence
- Create automated hyperparameter optimization pipeline

Application Extensions:

- Extend to multi-step ahead forecasting
- Develop models for other pollutants (PM10, NO2, O3)
- Implement early warning systems for high pollution events
- Create user-friendly interfaces for decision makers

6.5 Broader Impact

This work contributes to the growing field of environmental data science and demonstrates the potential of deep learning approaches for air quality management [3]. The systematic methodology developed can be applied to similar forecasting problems in other cities and environmental domains. The successful prediction of PM2.5 concentrations has direct implications for public health policy, urban planning, and environmental management. Accurate forecasting enables proactive measures to reduce exposure during high pollution periods, potentially saving lives and improving quality of life for millions of urban residents.

7. References

- [1] S. Raschka, "L15.5 Long Short-Term Memory," YouTube, 2024. [Online]. Available: <https://youtu.be/k6fSgUaWUF8>
- [2] World Health Organization, "WHO global air quality guidelines: particulate matter (PM2.5 and PM10), ozone, nitrogen dioxide, sulfur dioxide and carbon monoxide," World Health Organization, 2021. [Online]. Available: <https://www.who.int/publications/i/item/9789240034228>
- [3] T. Xayasouk, H. M. Lee, and G. Lee, "Air pollution prediction using long short-term memory (LSTM) and deep autoencoder (DAE) models," Sustainability, vol. 12, no. 6, p. 2570, 2020. [Online]. Available: <https://www.mdpi.com/2071-1050/12/6/2570>

GitHub Repository: <https://github.com/Deolinda1506/Time-Series-Forecasting.git>

Kaggle Competition:

<https://www.kaggle.com/competitions/assignment-1-time-series-forecasting-septemb-2025/overview>