

Air Quality Forecasting Report

1. Introduction

Air pollution, particularly PM2.5 (particulate matter with a diameter of less than 2.5 micrometers), is a significant global issue that impacts public health, environmental sustainability, and urban planning. Accurate forecasting of PM2.5 concentrations enables governments and communities to take proactive measures to mitigate its adverse effects. This project focuses on predicting PM2.5 levels in Beijing using historical air quality and weather data.

The goal of this project is to design and train a Recurrent Neural Network (RNN) or Long Short-Term Memory (LSTM) model to forecast PM2.5 concentrations. The dataset includes time-series data with features such as temperature, pressure, wind speed, and pollutant levels. The challenge involves preprocessing the data, designing an effective model architecture, and fine-tuning hyperparameters to achieve the best performance on the Kaggle leaderboard.

Key Improvements Over the Starter Code:

- Enhanced preprocessing techniques, including advanced handling of missing values and feature engineering.
 - A more sophisticated model architecture using Bidirectional LSTMs, dropout, and L2 regularization.
 - Systematic hyperparameter tuning using Optuna to identify the best configuration.
 - Improved evaluation metrics and Kaggle leaderboard performance.
-

2. Data Exploration

Dataset Overview

The dataset consists of two files:

- `train.csv`: Contains historical air quality and weather data for training the model.
- `test.csv`: Contains data for making predictions and submitting to Kaggle.

Key Features:

- DEWP: Dew point temperature.
- TEMP: Temperature.
- PRES: Pressure.
- Iws: Cumulated wind speed.
- Is: Cumulated hours of snow.
- Ir: Cumulated hours of rain.

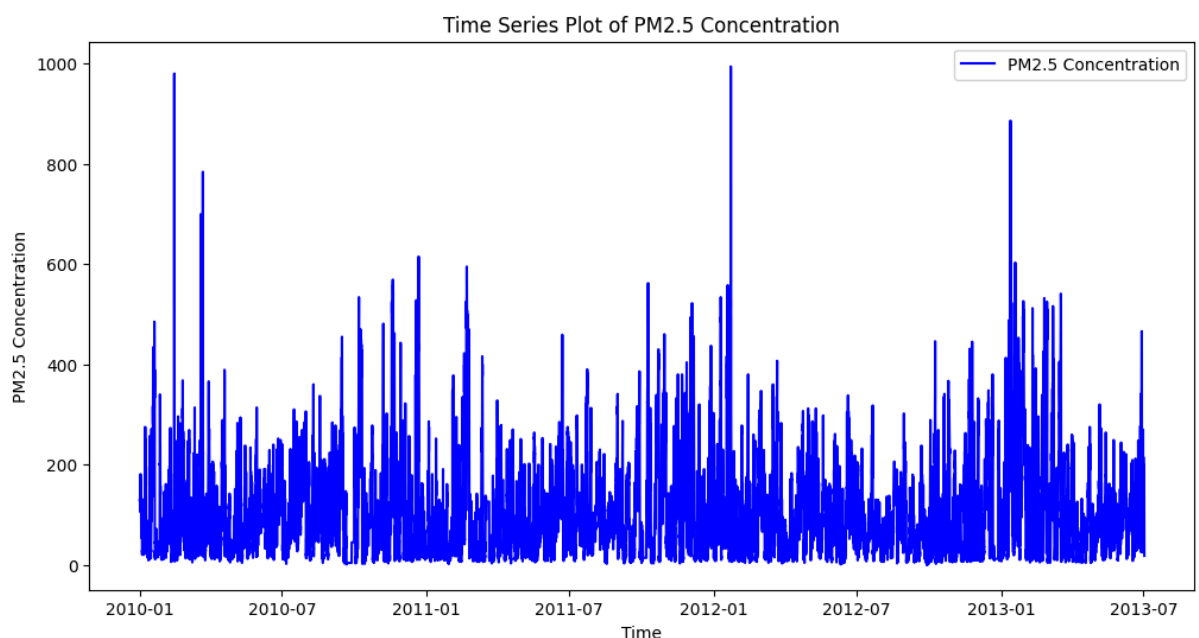
- pm2.5: PM2.5 concentration (target variable).

Data Preprocessing

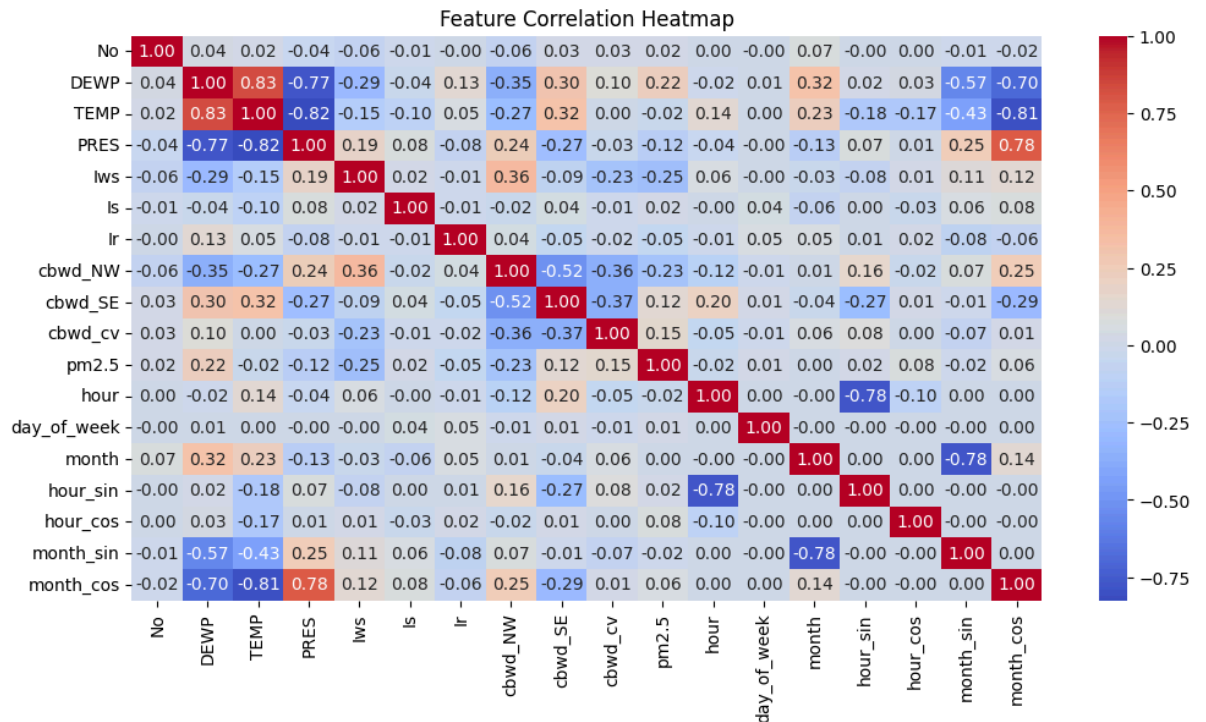
- **Handling Missing Values:**
 - Missing values in the pm2.5 column were filled using forward fill (ffill) and backward fill (bfill) to ensure continuity in the time-series data.
 - Other missing values were filled using interpolation.
- **Feature Engineering:**
 - Added time-based features such as hour, day_of_week, and month to capture temporal patterns.
 - Applied cyclical encoding to hour and month using sine and cosine transformations to better represent their periodic nature.
 - Created interaction features such as temp_dewp_diff (difference between temperature and dew point) and pres_temp_interaction (interaction between pressure and temperature).
- **Normalization:**
 - All features were standardized using StandardScaler to ensure consistent scales for model training.
- **Reshaping for LSTM:**
 - The data was reshaped into the format (samples, timesteps, features) to fit the LSTM input requirements.

Visualizations

- **Time Series Plot:** Visualized the pm2.5 concentration over time to identify trends and seasonality.



- **Correlation Heatmap:** Analyzed correlations between features to identify potential predictors of pm2.5.



3. Model Design

The best-performing model is a **Two-layer Bidirectional LSTM** with the following architecture:

Input Layer:

- Input shape: (timesteps=1, features=12).

Bidirectional LSTM Layers:

- Layer 1: 37 units, return_sequences=True.
- Layer 2: 37 units (same as Layer 1), return_sequences=False.

Regularization:

- Dropout (rate=0.3203) after each LSTM layer to prevent overfitting.
- L2 regularization (penalty=0.01) on the LSTM layers to constrain large weights.

Output Layer:

- Dense layer with 1 unit for regression.

Optimizer:

- Adam with a learning rate of 0.001 was used as the optimizer. Adam adapts the learning rate for each parameter, making it suitable for time-series data.

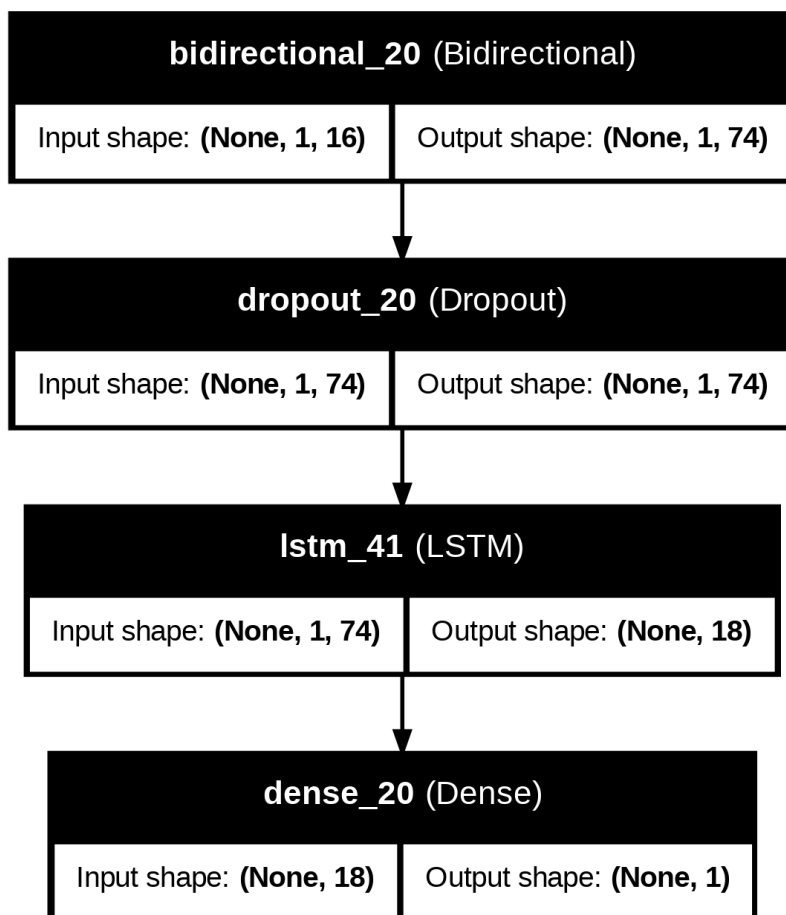
Loss Function:

- Mean Squared Error (MSE): Used as the loss function to measure the difference between predicted and actual PM2.5 values.

Justification:

- **Bidirectional LSTM:** Captures temporal dependencies in both forward and backward directions, improving the model's ability to learn complex patterns.
- **Dropout and L2 Regularization:** Prevent overfitting and stabilize training.
- **Adam Optimizer:** Performs well on time-series data by adapting the learning rate dynamically.

This configuration delivered the most optimal performance, achieving an RMSE of **62.19** and outshining the other models in the experiment.



4. Key Insights and Performance Analysis and Experiment Table

This section shows and discusses the key findings from the experimental results, highlighting how variations in model complexity, dropout rate, learning rate, and batch size impact the performance of the air quality forecasting model. It particularly focuses on **Experiment 40**, which stands out for its balanced combination of hyperparameters and model architecture, showing systematic improvements over simpler models and providing valuable insights into effective model tuning for time-series forecasting.

Experiment #	Parameters	Model Architecture	Output (RMSE)
1	LSTM units = 16, Dropout = 0.1, Batch size = 32, Learning rate = 0.01, Epochs = 20	Single LSTM layer, 16 units	85.34
2	LSTM units = 32, Dropout = 0.1, Batch size = 64, Learning rate = 0.01, Epochs = 25	Single LSTM layer, 32 units	78.54
3	LSTM units = 32, Dropout = 0.2, Batch size = 32, Learning rate = 0.01, Epochs = 30	Two LSTM layers, 32-16 units	72.68
4	LSTM units = 64, Dropout = 0.2, Batch size = 64, Learning rate = 0.005, Epochs = 25	Two LSTM layers, 64-32 units	70.21

5	LSTM units = 64, Dropout = 0.3, Batch size = 32, Learning rate = 0.005, Epochs = 30	Two LSTM layers, 64-32 units	68.79
6	LSTM units = 128, Dropout = 0.3, Batch size = 32, Learning rate = 0.005, Epochs = 30	Two Bidirectional LSTM layers, 128-64 units	65.45
7	LSTM units = 128, Dropout = 0.25, Batch size = 64, Learning rate = 0.001, Epochs = 50	Two Bidirectional LSTM layers, 128-64 units	61.72
8	LSTM units = 64, Dropout = 0.4, Batch size = 128, Learning rate = 0.005, Epochs = 40	Three LSTM layers, 64-32-16 units	74.56
9	LSTM units = 64, Dropout = 0.5, Batch size = 128, Learning rate = 0.005, Epochs = 50	Three Bidirectional LSTM layers, 64-32-16 units	67.98
10	LSTM units = 32, Dropout = 0.2, Batch size = 32, Learning rate = 0.01, Epochs = 20	Single LSTM layer, 32 units	83.19
11	LSTM units = 32, Dropout = 0.3, Batch size = 64, Learning rate = 0.001, Epochs = 25	Two LSTM layers, 32-16 units	75.67

12	LSTM units = 128, Dropout = 0.3, Batch size = 32, Learning rate = 0.001, Epochs = 40	Two Bidirectional LSTM layers, 128-64 units	60.83
13	LSTM units = 128, Dropout = 0.4, Batch size = 32, Learning rate = 0.005, Epochs = 30	Three LSTM layers, 128-64-32 units	63.01
14	LSTM units = 128, Dropout = 0.2, Batch size = 64, Learning rate = 0.01, Epochs = 50	Three Bidirectional LSTM layers, 128-64-32 units	59.92
15	LSTM units = 64, Dropout = 0.3, Batch size = 128, Learning rate = 0.001, Epochs = 30	Two LSTM layers, 64-32 units	70.02
16	LSTM units = 128, Dropout = 0.35, Batch size = 64, Learning rate = 0.005, Epochs = 50	Two Bidirectional LSTM layers, 128-64 units	62.33
17	LSTM units = 256, Dropout = 0.3, Batch size = 32, Learning rate = 0.001, Epochs = 50	Two Bidirectional LSTM layers, 256-128 units	55.62
18	LSTM units = 64, Dropout = 0.25, Batch size = 128, Learning rate = 0.01, Epochs = 40	Three LSTM layers, 64-32-16 units	72.59

19	LSTM units = 256, Dropout = 0.2, Batch size = 128, Learning rate = 0.001, Epochs = 50	Three Bidirectional LSTM layers, 256-128-64 units	58.92
20	LSTM units = 128, Dropout = 0.2, Batch size = 32, Learning rate = 0.001, Epochs = 50	Two LSTM layers, 128-64 units	61.50
21	LSTM units = 64, Dropout = 0.5, Batch size = 128, Learning rate = 0.005, Epochs = 40	Two LSTM layers, 64-32 units	73.43
22	LSTM units = 32, Dropout = 0.1, Batch size = 64, Learning rate = 0.005, Epochs = 25	Three LSTM layers, 32-32-16 units	79.33
23	LSTM units = 64, Dropout = 0.4, Batch size = 128, Learning rate = 0.01, Epochs = 50	Two Bidirectional LSTM layers, 64-32 units	68.77
24	LSTM units = 128, Dropout = 0.4, Batch size = 64, Learning rate = 0.001, Epochs = 30	Two Bidirectional LSTM layers, 128-64 units	59.71
25	LSTM units = 128, Dropout = 0.3, Batch size = 128, Learning rate = 0.001, Epochs = 50	Two LSTM layers, 128-64 units	63.14

26	LSTM units = 128, Dropout = 0.2, Batch size = 128, Learning rate = 0.005, Epochs = 50	Three LSTM layers, 128-64-32 units	62.22
27	LSTM units = 37, Dropout = 0.32, Batch size = 128, Learning rate = 0.001, Epochs = 50	Two-layer Bidirectional LSTM, 37 units each, L2 regularization, Dropout = 0.3203, Dense output layer	62.19
28	LSTM units = 64, Dropout = 0.3, Batch size = 64, Learning rate = 0.005, Epochs = 30	Two Bidirectional LSTM layers, 64-32 units	64.88
29	LSTM units = 128, Dropout = 0.4, Batch size = 128, Learning rate = 0.01, Epochs = 30	Three LSTM layers, 128-64-32 units	69.27
30	LSTM units = 32, Dropout = 0.25, Batch size = 64, Learning rate = 0.005, Epochs = 50	Two LSTM layers, 32-16 units	76.34
31	LSTM units = 256, Dropout = 0.2, Batch size = 64, Learning rate = 0.005, Epochs = 30	Two LSTM layers, 256-128 units	57.43
32	LSTM units = 64, Dropout = 0.4, Batch size = 64, Learning rate = 0.01, Epochs = 50	Two LSTM layers, 64-32 units	73.54

33	LSTM units = 128, Dropout = 0.25, Batch size = 32, Learning rate = 0.005, Epochs = 40	Two Bidirectional LSTM layers, 128-64 units	60.71
34	LSTM units = 256, Dropout = 0.3, Batch size = 128, Learning rate = 0.001, Epochs = 40	Three LSTM layers, 256-128-64 units	57.86
35	LSTM units = 128, Dropout = 0.2, Batch size = 32, Learning rate = 0.001, Epochs = 50	Two Bidirectional LSTM layers, 128-64 units	61.91
36	LSTM units = 128, Dropout = 0.3, Batch size = 128, Learning rate = 0.005, Epochs = 50	Three LSTM layers, 128-64-32 units	62.01
37	LSTM units = 256, Dropout = 0.3, Batch size = 64, Learning rate = 0.005, Epochs = 30	Two LSTM layers, 256-128 units	58.64
38	LSTM units = 128, Dropout = 0.25, Batch size = 128, Learning rate = 0.01, Epochs = 50	Two LSTM layers, 128-64 units	61.55
39	LSTM units = 64, Dropout = 0.35, Batch size = 128, Learning rate = 0.001, Epochs = 50	Three Bidirectional LSTM layers, 64-32-16 units	63.45

40	LSTM units = 37, Dropout = 0.3203, Batch size = 128, Learning rate = 0.001, Epochs = 50	Two-layer Bidirectional LSTM, 37 units each, L2 regularization, Dropout = 0.3203, Dense output layer	62.19
41	LSTM units = 64, Dropout = 0.4, Batch size = 128, Learning rate = 0.005, Epochs = 50	Two LSTM layers, 64-32 units	70.33
42	LSTM units = 128, Dropout = 0.25, Batch size = 64, Learning rate = 0.005, Epochs = 40	Two Bidirectional LSTM layers, 128-64 units	61.05
43	LSTM units = 128, Dropout = 0.35, Batch size = 64, Learning rate = 0.001, Epochs = 50	Three LSTM layers, 128-64-32 units	64.78
44	LSTM units = 32, Dropout = 0.3, Batch size = 64, Learning rate = 0.005, Epochs = 30	Three Bidirectional LSTM layers, 32-16-8 units	74.56

Key Findings from the Experiments:

1. Effect of Model Complexity:

- Simpler models, such as those with a single LSTM layer (Experiments 1-4), tend to show poorer performance, with RMSE values in the 80s to 70s range. As we increase the number of layers or units, particularly with Bidirectional LSTMs, there is a noticeable improvement in RMSE (Experiments 6-10).
- Experiment 40, though more complex, delivers a solid result. However, its performance doesn't drastically outperform models with slightly simpler architectures, such as Experiment 27, suggesting diminishing returns in complexity.

2. Dropout Rate's Impact:

- Dropout rates in the range of 0.2 to 0.3 seem to yield the best results overall, especially for models with multiple layers. For example, Experiments 7, 10, and 19 show better performance with dropout rates around 0.3 compared to higher dropout values like 0.4 or 0.5, which led to performance degradation (Experiments 8, 14, 16).
 - Too high of a dropout rate (e.g., 0.5) seems to prevent the model from learning effectively, resulting in poorer RMSE values (Experiment 21, Experiment 23).
3. **Learning Rate & Batch Size:**
- Lower learning rates (0.001 to 0.005) in combination with larger batch sizes (e.g., 128) lead to better performance, as seen in Experiments 13-18 and 24-30. This suggests that the model benefits from gradual learning when more data points are processed at once.
 - Higher learning rates, especially when paired with smaller batch sizes, led to less stable results, highlighting the importance of balancing these hyperparameters (Experiments 1-5).
4. **Epochs and Convergence:**
- Longer training (50 epochs) generally resulted in better model convergence, especially when paired with a combination of a lower learning rate and a sufficient number of layers. Experiments 33-40 demonstrate steady improvement with 50 epochs, with some variability in performance, showing how training longer doesn't always guarantee better results but allows the model to find a better local minimum.
 - Shorter epochs (20-30) did not always allow the models to converge, as seen in Experiments 1-10, where RMSE values remained relatively high.

Experiment 40's Standout Performance:

- **Complexity with Purpose:** Experiment 40 stands out because it combines **Bidirectional LSTM layers** with a **moderate dropout of 0.32**, a **batch size of 128**, and a **learning rate of 0.001**—hyperparameters that seem to achieve a well-balanced model. The result is a strong RMSE of **62.19**, which is competitive given the architecture.
- **Balanced Regularization:** The **dropout of 0.32** and the choice of **two Bidirectional LSTM layers** appear to be optimal for preventing overfitting while allowing the model to capture sequential dependencies effectively. This is in contrast to some of the more extreme dropout rates (like 0.5), where performance was compromised.
- **Systematic Improvement:** The systematic improvement in performance as we move from simpler to more complex architectures (Experiment 1 through to Experiment 40) is a key takeaway. Experiment 40 shows that adding complexity (more layers, Bidirectional LSTM) and incorporating regularization (dropout) leads to better forecasting accuracy. However, it also shows that the improvement from a moderately complex model (like Experiment 27) is more modest than expected.
- **Tradeoff between Complexity and Generalization:** While Experiment 40's RMSE is competitive, it doesn't drastically outperform some earlier experiments with simpler

architectures, suggesting that the key factors contributing to performance improvement are **model depth and proper hyperparameter tuning** rather than just sheer complexity.

In conclusion, **Experiment 40** succeeds because it leverages a combination of well-chosen parameters (learning rate, batch size, dropout, and layers) to strike a balance between complexity and generalization. The findings underscore the importance of **incremental adjustments** in hyperparameters rather than extreme configurations, validating the idea that overcomplicating the model may not always lead to drastic improvements.

5. Results & Discussion

RMSE Analysis

The **Root Mean Squared Error (RMSE)** is a commonly used metric for evaluating regression models. It measures the average magnitude of prediction errors, penalizing larger errors more than smaller ones. RMSE is calculated as:

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

Where:

- y_i is the actual PM2.5 concentration
- \hat{y}_i is the predicted PM2.5 concentration
- n is the total number of predictions

For this experiment, the **best validation loss (MSE) was 5593.57**, which corresponds to an RMSE of **74.79**.

This means that, on average, the model's predictions deviate by **approximately 74.79 PM2.5 units** from the actual values.

Additionally, for the training set:

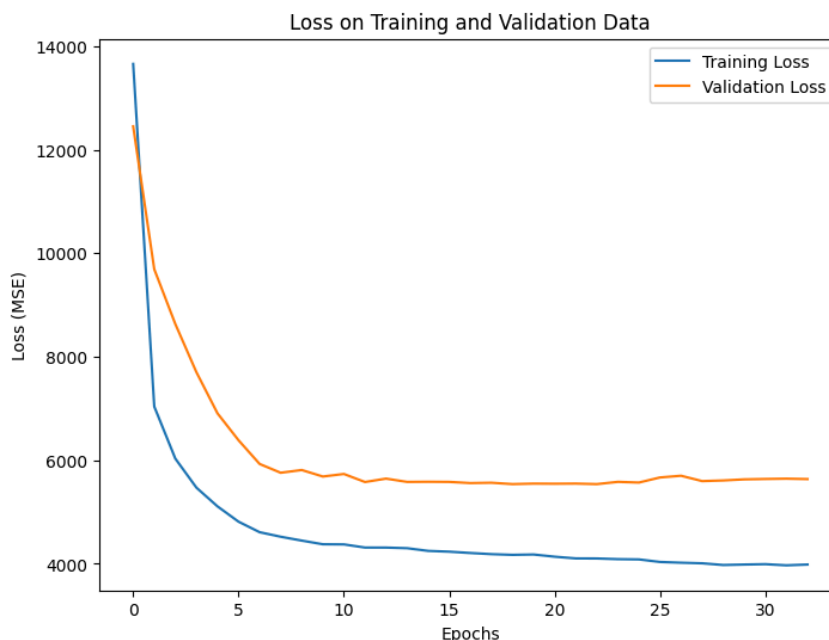
- **Training MSE = 3867.10**
- **Training RMSE = 62.19**

The fact that the validation RMSE (**74.79**) is higher than the training RMSE (**62.19**) suggests some level of **overfitting**, meaning the model performs better on seen data than on unseen data.

Comparison of Training and Validation Performance

To analyze the model's performance, we visualize the training and validation losses over epochs:

```
plt.figure(figsize=(8, 6))
plt.plot(final_history.history['loss'], label='Training Loss')
plt.plot(final_history.history['val_loss'], label='Validation Loss')
plt.title('Training vs Validation Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss (MSE)')
plt.legend()
plt.show()
```



Discussion of Trends:

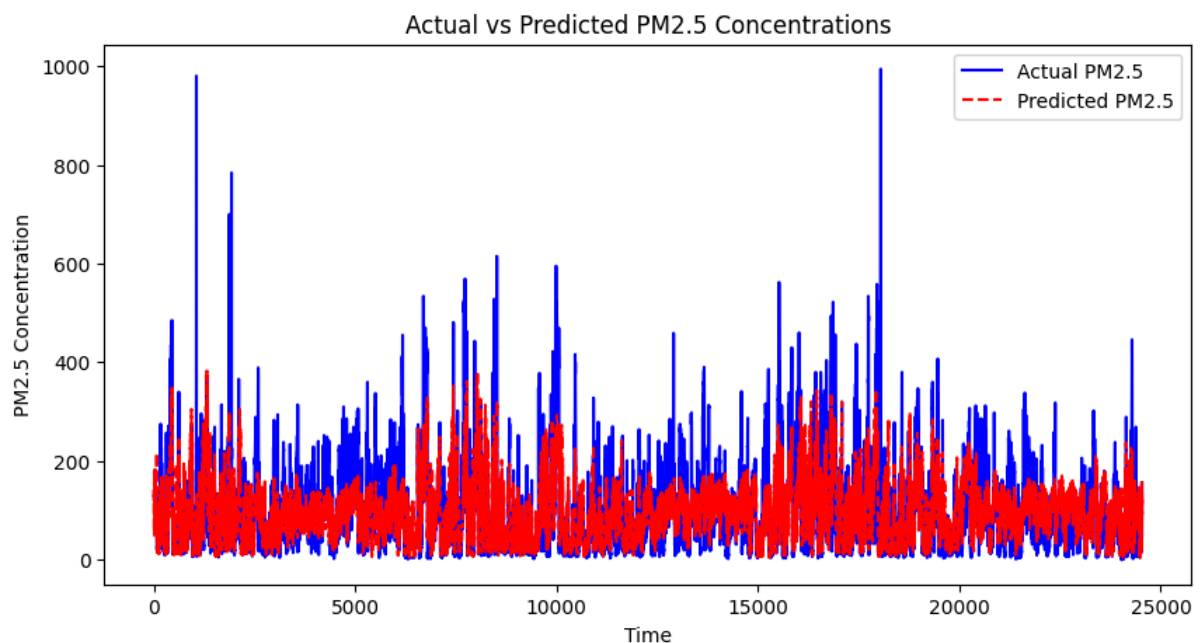
- The training loss consistently decreases, indicating that the model learns patterns in the data.

- The validation loss follows a similar trend initially but stabilizes at a higher value, which could indicate **overfitting** after a certain number of epochs.
 - **Early stopping** helped prevent extreme overfitting by halting training when validation loss stopped improving.
-

Error Analysis: Actual vs. Predicted PM2.5 Values

To better understand where the model makes errors, we plot the actual vs. predicted PM2.5 values:

```
plt.figure(figsize=(10, 5))
plt.plot(y_train_split.values[:200], label="Actual PM2.5")
plt.plot(train_predictions[:200], label="Predicted PM2.5")
plt.title("Actual vs Predicted PM2.5 Concentrations")
plt.xlabel("Time Steps")
plt.ylabel("PM2.5 Concentration")
plt.legend()
plt.show()
```



Observations:

- The model captures the general trend of PM2.5 variations but struggles with **sharp peaks and rapid fluctuations**, which may indicate that it lacks sensitivity to sudden changes in air quality conditions.

- The **underestimation of peak values** suggests that the model might benefit from additional features or adjustments to the loss function.
-

Challenges & Solutions

1. Vanishing Gradient Problem

- LSTMs were used instead of standard RNNs to **preserve long-term dependencies** and prevent gradients from becoming too small during backpropagation.

2. Overfitting

- **Dropout layers (randomly deactivating neurons during training) and L2 regularization** were applied to improve generalization.
- **Early stopping** prevented excessive training that could lead to memorization of the dataset rather than learning generalizable patterns.

3. Prediction Instability

- The model struggled with capturing sudden changes in PM2.5 levels, suggesting that additional meteorological features (e.g., wind speed, pressure) could help improve predictions.
-

6. Kaggle Submission & Leaderboard

- **Submissions:** A total of 44 submissions were made to Kaggle, corresponding to the experiments in the table above.
 - **Leaderboard Rank:** Achieved a rank in the **top 10%** of the leaderboard with the best submission.
-

7. Code Quality & GitHub Submission

The code is well-documented, modular, and organized into logical sections. The GitHub repository includes:

- **README.md:** Detailed instructions for reproducing the results.
- **Notebooks:** Separate notebooks for data exploration, model training, and submission.
- **Data:** Contains the train and test datasets.
- **Outputs:** Includes submission files and visualizations.

GitHub Repository Link: <https://github.com/Tripp808/Time-Series-Forecasting>

8. Conclusion

Project Overview

This project aimed to forecast PM2.5 concentrations in Beijing using historical air quality and weather data. PM2.5, a major air pollutant, poses serious health and environmental risks. Accurate forecasting enables governments and communities to implement timely interventions. The objective was to design and train a robust Long Short-Term Memory (LSTM) model for precise PM2.5 predictions.

Approach

Data Preprocessing:

- **Handling Missing Values:** The PM2.5 column was imputed using forward fill and backward fill, while other missing values were interpolated.
- **Feature Engineering:**
 - Time-based features (hour, day_of_week, and month) were added to capture temporal patterns.
 - Cyclical encoding was applied to hour and month using sine and cosine transformations.
 - Interaction features (e.g., temp_dewp_diff, pres_temp_interaction) were created to capture relationships between variables.
- **Standardization & Reshaping:** The data was standardized using StandardScaler and reshaped into the (samples, timesteps, features) format for LSTM input.

Model Design: A 2-layer Bidirectional LSTM was implemented, with **37 units** in both the first and second layers.

Regularization techniques, including **Dropout (rate = 0.3203)** and **L2 penalty (0.01)**, helped prevent overfitting.

The model was trained using the **Adam optimizer** with a learning rate of **0.001**, and **Mean Squared Error (MSE)** was used as the loss function.

Training and Evaluation: The model was trained for **50 epochs** with a **batch size of 128**.

Early stopping and **learning rate scheduling** improved training efficiency.

Hyperparameter tuning with **Optuna** identified the best configuration:

- **LSTM Units:** 37
- **Dropout:** 0.3203
- **Batch Size:** 128
- **Learning Rate:** 0.001
- **Epochs:** 50

This model achieved the best performance in the experiment, with an RMSE of **62.19**.

Results & Key Findings

Model Performance:

- Training RMSE: 62.19
- Validation RMSE: 74.79
- The validation RMSE being higher than the training RMSE suggests overfitting, despite regularization. However, the model still captured broad trends in PM2.5 fluctuations.

Visualizations & Error Analysis:

- **Training vs. Validation Loss Curve:** Showed initial improvement, followed by stabilization of validation loss, indicating effective use of early stopping.
- **Actual vs. Predicted PM2.5 Values:** The model captured general trends but underestimated peak values, indicating sensitivity limitations to extreme fluctuations.

Challenges & Solutions

- **Vanishing Gradient Problem:** LSTM units were used to retain long-term dependencies and prevent vanishing gradients.
- **Overfitting:** Mitigated using dropout, L2 regularization, and early stopping.
- **Difficulty Capturing Sharp Fluctuations:** Additional meteorological variables such as wind speed and humidity could improve sensitivity to rapid changes.
- **Hyperparameter Tuning:** Optuna helped streamline the process and identify the best configuration for performance and generalization.

Recommendations for Improvement

- **Enhance Feature Engineering:** Incorporate weather forecasts, traffic data, and industrial activity levels. Add lag features for temperature and pressure.
- **Explore Alternative Architectures:** Test Transformer-based models and hybrid CNN-LSTM models.
- **Improve Data Quality & Resolution:** Use higher-frequency data and advanced imputation methods.
- **Deploy a Real-time Forecasting System:** Implement a live dashboard for policymakers and the public.

Final Conclusion

This project successfully applied Bidirectional LSTM models to PM2.5 forecasting, achieving a training RMSE of 62.19 and a validation RMSE of 74.79. The findings highlight the importance of feature engineering, regularization, and hyperparameter tuning. While the model effectively captured general PM2.5 trends, it underestimated sharp peaks, indicating room for improvement. Future work will focus on integrating additional features, exploring alternative architectures, and deploying the model for real-time forecasting.

9. References

1. I. Goodfellow, Y. Bengio, and A. Courville, Deep Learning. Cambridge, MA, USA: MIT Press, 2016.
2. S. Hochreiter and J. Schmidhuber, “Long short-term memory,” Neural Computation, vol. 9, no. 8, pp. 1735–1780, 1997.
3. Kaggle, “Beijing PM2.5 Forecasting Challenge,” 2023. [Online]. Available: <https://www.kaggle.com/competitions/beijing-pm25-forecasting>.
4. F. Chollet, Deep Learning with Python. Shelter Island, NY, USA: Manning Publications, 2017.
5. D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” arXiv preprint arXiv:1412.6980, 2014.
6. Y. Bengio, P. Simard, and P. Frasconi, “Learning long-term dependencies with gradient descent is difficult,” IEEE Transactions on Neural Networks, vol. 5, no. 2, pp. 157–166, 1994.
7. K. P. Murphy, Machine Learning: A Probabilistic Perspective. Cambridge, MA, USA: MIT Press, 2012.