

ASSIGNMENT TWO

PAPER NAME: Data Mining and Machine

Learning PAPER CODE: COMP809

TOTAL MARKS: 100

Students' Names:Yize(Serena) Wang.....GuangLiang(Ricky) Yang.....

Students' IDs:23198583..... 23205919.....

- Due date: 09 Jun 2024 midnight NZ time.
- **Late penalty:** maximum late submission time is 24 hours after the due date. In this case, a **5% late penalty** will be applied.
- Submit the actual code (no screenshot) separately with appropriate comments for each task.

Note: This assignment should be completed by a group of two students and both students **MUST contribute in each part**. **Submission:** a soft copy needs to be submitted through the canvas assessment link.

INSTRUCTIONS:

1. **The following actions may be deemed to constitute a breach of the General Academic Regulations Part 7: Academic Discipline,**
 - Communicating with or collaborating with another person regarding the Assignment
 - Copying from any other student work for your Assignment
 - Copying from any third-party websites unless it is an open book Assignment
 - Uses any other unfair means
2. **Please email DCT.EXAM@AUT.AC.NZ if you have any technical issues with your submission on Canvas immediately**
3. **Attach your code for all the datasets in.**

COMP809 Assignment2-Part B

Prediction of PM2.5 With Multi-layer Perceptron (MLP) and
Long Short-term Memory (LSTM)

Yize (Serena) Wang
Student ID: 23198583

GuangLiang (Ricky) Yang
Student ID: 23205919

09 June 2024

Catalogue

1. Introduction	8
1.1. Problem.....	8
1.2. Check dataset quality	8
1.3 Pre-Processing	14
1.3.1 Delete ‘Solar_Rad’ column	14
1.3.2 Converted the Data Type.....	14
1.3.3 Check and process duplicate values, missing values, errors, and outliers	15
1.3.4. Interpolation Effect Review.....	20
1.3.5. Extract lag1 and lag2 Feature.....	20
1.3.6. Extract Time Feature	21
1.3.7. Extract Wind Direction Feature.....	21
1.3.8. Final dataset	22
2. Features Selection.....	23
2.1. Extract Method and Reason	23
2.2. Influence on PM Concentration.....	28
2.3. Provide graphical visualisation of variation of PM variation.	31
2.4. Provide summary statistics of the PM concentration.....	32
2.5. Summary statistic of the highest correlation predictors.....	32
2.6. Normalization Data.....	33
3. Experimental Methods	33
3.1. Split Dataset	33
3.2. Workflow	34
4. Multilayer Perceptron (MLP)	35
4.1. MLP Description	35
4.2. Learning Rate Analysis.....	36
4.3. Neurons Distribution Analysis	38
4.4. MLP Analysis	41
4.4.1. Best Architecture	41
4.4.2. Possible Reasons for Variation.....	42
5. Long Short-Term Memory (LSTM).....	42
5.1. LSTM Introduction	42
5.1.1. LSTM.....	42
5.1.2. LSTM Architecture	43
5.1.3. How LSTM Differs from MLP	44
5.1.4. Impact of Number of Neurons and Batch Size.....	44
5.2. Epoch Optimal Size.....	45
5.2.1. Identify an appropriate cost function	45

5.2.2. Cost Function Scores	45
5.2.3 Report the Summary Statistics	49
5.2.4 Best Epoch Size	50
5.3. Batch Optimal Size.....	52
5.3.1. Report the summary statistics.....	52
5.3.2. Best Batch Size	57
5.4. Neurons Optimal Number	59
5.4.1. Report the summary statistics.....	59
5.4.2. Best Neurons Number.....	65
5.5. LSTM Conclusion	67
6. Model Comparison	67
6.1. Visually Compare Model Performance.	67
6.2. Compare the Model Performance using RMSE	71
7. Conclusion	72
Reference	73
Appendix 1. Source Code.....	73
Appendix 2. Abbreviations	73

List of Figures

FIGURE 1 DATA EXTRACTION FROM ENVIRONMENTAL AUCKLAND DATA PORTAL (PART 1)	8
FIGURE 2 DATA EXTRACTION FROM ENVIRONMENTAL AUCKLAND DATA PORTAL (PART 2).....	9
FIGURE 3 ORIGINAL DATASET SCREEN SHOT	9
FIGURE 4. DATASET OVERVIEW WITH DATASET SHAPE	10
FIGURE 5 DATASET INFORMATION	10
FIGURE 6 MISSING VALUES HEAT MAP	11
FIGURE 7 DATASET HISTOGRAM DISTRIBUTION.....	11
FIGURE 8 DATASET OUTLIERS WITH BOX PLOT.....	12
FIGURE 9 DATASET TIME SERIES PLOT	13
FIGURE 10 DATASET AFTER REMOVING SOLAR_RAD	14
FIGURE 11 DUPLICATE CHECK	15
FIGURE 12 DATASET TIME SERIES ANALYSIS	15
FIGURE 13 MISSING VALUES CHECK IN THE DATASET	16
FIGURE 14 TIME SERIES DATA OF AIR QUALITY AND METEOROLOGICAL(ORIGINAL).....	17
FIGURE 15 TIME SERIES DATA OF AIR QUALITY AND METEOROLOGICAL (LINEAR INTERPOLATION).....	18
FIGURE 16 TIME SERIES DATA OF AIR QUALITY AND METEOROLOGICAL (KNN INTERPOLATION).....	18
FIGURE 17 ERROR COMPARISON OF IMPUTATION METHODS	19
FIGURE 18 BOXPLOT OF AIR QUALITY AND METEOROLOGICAL DATA	20
FIGURE 19 DATA REVIEW WITH LAG1 AND LAG2.....	21
FIGURE 20 DATA REVIEW WITH NEW FEATURES.....	21
FIGURE 21 WIND DIRECTION FEATURES	22
FIGURE 22 FINAL DATASET REVIEW.....	22
FIGURE 23 EACH FACTOR'S RELATIONSHIP WITH PM2.5	24
FIGURE 24 PEARSON CORRELATION MATRIX (TOP 8 FEATURES)	26
FIGURE 25 SPEARMAN CORRELATION MATRIX (TOP 8 FEATURES)	27
FIGURE 26 ALL FACTORS' PEARSON CORRELATION WITH PM2.5 PLOT	27
FIGURE 27 ALL FACTORS' SPEARMAN CORRELATION WITH PM2.5 PLOT	28
FIGURE 28 MULTIPLE LINEAR REGRESSION SUMMARY	29
FIGURE 29 RESIDUAL PLOT FOR PM2.5 WITH MULTIPLE FEATURES	30
FIGURE 30 . QQ PLOT FOR PM2.5 WITH MULTIPLE FEATURES	30
FIGURE 31 PM2.5 CONCENTRATION OVER TIME	31
FIGURE 32 YEARLY PM2.5 CONCENTRATION DISTRIBUTION	31
FIGURE 33 DATASET SPLITTING FOR TRAINING AND TESTING	33
FIGURE 34 WORKFLOW OVERVIEW	34
FIGURE 35 MLP OVERVIEW	36
FIGURE 36 MLP LEARNING RATE ANALYSIS WORKFLOW.....	37
FIGURE 37 MLP LEARNING RATE COMPARISON	38
FIGURE 38 MLP NEURONS DISTRIBUTION ANALYSIS WORKFLOW.....	39
FIGURE 39 MLP NEURONS DISTRIBUTION COMPARISON	41
FIGURE 40 AN UNROLLED RECURRENT NEURAL NETWORK.....	43
FIGURE 41 LSTM ARCHITECTURE	43
FIGURE 42 LSTM TUNING WORKFLOW.....	46
FIGURE 43 LSTM MEAN TRAIN AND TEST MAE LOSS OVER 20 EPOCHS FOR EPOCHS = 20	47
FIGURE 44 LSTM MEAN TRAIN AND TEST MAE LOSS OVER 20 EPOCHS FOR EPOCHS = 60	47
FIGURE 45 LSTM MEAN TRAIN AND TEST MAE LOSS OVER 20 EPOCHS FOR EPOCHS = 100.....	48
FIGURE 46 LSTM MEAN TRAIN AND TEST MAE LOSS OVER 20 EPOCHS FOR EPOCHS = 200.....	48
FIGURE 47 LSTM MEAN TRAIN AND TEST MAE LOSS OVER 20 EPOCHS FOR EPOCHS = 500.....	49
FIGURE 48 LSTM TEST HUBER LOSS DISTRIBUTION FOR DIFFERENT EPOCHS CONFIGURATIONS.....	51
FIGURE 49 LSTM TIME DISTRIBUTION FOR DIFFERENT EPOCHS CONFIGURATIONS.....	51
FIGURE 50 LSTM MEAN TRAIN AND TEST HUBER LOSS OVER 200 EPOCHS FOR BATCH SIZE = 4.....	52
FIGURE 51 LSTM MEAN TRAIN AND TEST HUBER LOSS OVER 200 EPOCHS FOR BATCH SIZE = 8	53
FIGURE 52 LSTM MEAN TRAIN AND TEST HUBER LOSS OVER 200 EPOCHS FOR BATCH SIZE = 16.....	53
FIGURE 53 LSTM MEAN TRAIN AND TEST HUBER LOSS OVER 200 EPOCHS FOR BATCH SIZE = 32.....	54
FIGURE 54 LSTM MEAN TRAIN AND TEST HUBER LOSS OVER 200 EPOCHS FOR BATCH SIZE = 64.....	54
FIGURE 55 LSTM MEAN TRAIN AND TEST HUBER LOSS OVER 200 EPOCHS FOR BATCH SIZE = 128.....	55
FIGURE 56 LSTM MEAN TRAIN AND TEST HUBER LOSS OVER 200 EPOCHS FOR BATCH SIZE = 256.....	55
FIGURE 57 LSTM TEST HUBER LOSS DISTRIBUTION FOR DIFFERENT BATCH SIZE CONFIGURATIONS.....	57
FIGURE 58 LSTM TIME COST DISTRIBUTION FOR DIFFERENT BATCH SIZE CONFIGURATIONS.....	58
FIGURE 59 LSTM MEAN TRAIN AND TEST FOR NEURON COUNT = 1 ON HIDDEN LAYER.....	59
FIGURE 60 LSTM MEAN TRAIN AND TEST FOR NEURON COUNT = 2 ON HIDDEN LAYER.....	60

FIGURE 61 LSTM MEAN TRAIN AND TEST FOR NEURON COUNT = 4 ON HIDDEN LAYER.....	60
FIGURE 62 LSTM MEAN TRAIN AND FOR NEURON COUNT = 8 ON HIDDEN LAYER	61
FIGURE 63 LSTM MEAN TRAIN AND TEST FOR NEURON COUNT = 16 ON HIDDEN LAYER	61
FIGURE 64 LSTM MEAN TRAIN AND TEST FOR NEURON COUNT = 32 ON HIDDEN LAYER	62
FIGURE 65 LSTM MEAN TRAIN AND FOR NEURON COUNT = 64 ON HIDDEN LAYER.....	62
FIGURE 66 LSTM MEAN TRAIN AND TEST FOR NEURON COUNT = 128 ON HIDDEN LAYER.....	63
FIGURE 67 LSTM MEAN TRAIN AND TEST FOR NEURON COUNT = 256 ON HIDDEN LAYER.....	63
FIGURE 68 LSTM TEST HUBER LOSS DISTRIBUTION FOR DIFFERENT NEURON COUNT CONFIGURATIONS.....	65
FIGURE 69 LSTM TIME COST DISTRIBUTION FOR DIFFERENT NEURON COUNT CONFIGURATIONS	66
FIGURE 70 MLP AND LSTM PERFORMANCE COMPARISON ON LINEAR INTERPOLATION DATASET	68
FIGURE 71 MLP AND LSTM PERFORMANCE COMPARISON ON KNN INTERPOLATION DATASET	69
FIGURE 72 MLP AND LSTM PERFORMANCE COMPARISON ON HIGH QUALITY DATASET	69
FIGURE 73 PERFORMANCE COMPARISON OF MLP AND LSTM ON HIGH QUALITY DATASET.....	71

List of Tables

TABLE 1 ORIGINAL DATASET STRUCTURE	15
TABLE 2 DATASET ERROR SUMMARY	16
TABLE 3 DATASET OUTLIER SUMMARY	17
TABLE 4 IMPUTATION METHODS COMPARISON	19
TABLE 5 FINAL DATASET STRUCTURE.....	22
TABLE 6 COMPARISON OF STATISTICAL ANALYSIS METHODS	25
TABLE 7 CORRELATION COEFFICIENTS FOR DIFFERENT PARAMETERS	26
TABLE 8 SUMMARY STATISTICS OF PM2.5	32
TABLE 9 SUMMARY STATISTICS OF HIGHEST CORRELATION PREDICTORS	32
TABLE 10 SUMMARY STATISTICS FOR NORMALIZED DATA.....	33
TABLE 11 MLP PERFORMANCE WITH DIFFERENT LEARNING RATES.....	37
TABLE 12 MLP PERFORMANCE WITH NEURON DISTRIBUTION	40
TABLE 13 SUMMARY STATISTICS FOR EACH EPOCH CONFIGURATION	50
TABLE 14 SUMMARY STATISTICS FOR ALL BATCH SIZE CONFIGURATIONS.....	56
TABLE 15 SUMMARY STATISTICS FOR DIFFERENT NEURON COUNT CONFIGURATIONS	64
TABLE 16 MLP MODEL SUMMARY.....	67
TABLE 17 LSTM MODEL SUMMARY	68
TABLE 18 PERFORMANCE COMPARISON OF MLP AND LSTM.....	71

1. Introduction

1.1. Problem

Air pollution is a significant environmental issue that poses serious risks to public health. Among various pollutants, particulate matter (PM) smaller than 2.5 micrometres (PM_{2.5}) has been identified as having the strongest correlation with cardiovascular diseases. Therefore, accurately predicting PM_{2.5} concentrations is crucial for mitigating its adverse health effects.

In this assignment, we aim to build predictive models for PM_{2.5} concentrations using two advanced machine learning techniques: Multi-Layer Perceptron (MLP) and Long Short-Term Memory (LSTM) networks. These models will help forecast PM_{2.5} levels based on historical data and various environmental factors.

1.2. Check dataset quality

Download dataset: The dataset for this experiment comes from the Auckland Environmental Data Portal. It includes hourly measurement data collected from January 1, 2019, to December 31, 2023, from the Penrose Station in Auckland (ID:7). The data includes SO₂, NO, NO₂, Solar Radiation (W/m²), Air Temperature (°C), Relative Humidity (%), Wind Direction (°), Wind Speed (m/s) and PM_{2.5} (Figure 1 and 2). After downloading the data, manually remove the initial header information to prepare the dataset for analysis (Figure 3).

The screenshot shows the 'Export Data' section of the Auckland Environmental Data Portal. The left sidebar has a 'Data Set' tab selected. The main form is titled 'Export Data' and contains the following fields:

- Prefill from Template: dropdown menu with 'Template...' and 'Clear Form' buttons.
- Date Range: dropdown menu set to 'Custom'.
- Custom Range: input field showing '2019-05-01 00:00 to 2024-04-30 00:00'.
- Time Zone: dropdown menu set to 'First Data Set's Time Zone'.
- Interval/Points: dropdown menu set to 'Points as recorded'.
- Export Format: dropdown menu set to 'CSV'.
- Single/Multi File: radio buttons for 'Single Time-Aligned File' (selected) and 'One File Per Data Set'.
- Rounding: radio buttons for 'Full Precision' (selected) and 'Round Data to Default Specification'.
- Include Grade Codes?: radio buttons for 'Yes' (selected) and 'No'.
- Include Qualifiers?: radio buttons for 'Yes' (selected) and 'No'.
- Include Interpolation Types?: radio buttons for 'Yes' (selected) and 'No'.

Below the form, there is a 'Data Sets' section with a table showing one row of data:

Location	Data Set	Conversion Option
Entire Period of Record: 2003-04-18 00:00 (UTC+12:00) - 2024-05-01 00:00 (UTC+12:00)	Overlapping Period of Record: 2003-04-18 00:00 (UTC+12:00) - 2024-05-01 00:00 (UTC+12:00)	<input type="button" value="Hide Data Sets"/>

Figure 1 Data extraction from Environmental Auckland Data Portal (Part 1)

Location	Data Set	Conversion Option
7 - Penrose	SO2.24-Hour Aggregate ($\mu\text{g}/\text{m}^3$)@7	Average in Micrograms per cubic metre
7 - Penrose	NO2.24-Hour Aggregate ($\mu\text{g}/\text{m}^3$)@7	Average in Micrograms per cubic metre
7 - Penrose	NO2.24-Hour Aggregate ($\mu\text{g}/\text{m}^3$)@7	Average in Micrograms per cubic metre
7 - Penrose	Solar Rad.24-Hour Aggregate (W/m^2)@7	Average in Watts per square metre
7 - Penrose	Air Temp.24-Hour Aggregate ($^\circ\text{C}$)@7	Average in Celsius
7 - Penrose	Rel Humidity.24-Hour Aggregate (%)@7	Average in Percent
7 - Penrose	Wind Dir.Hourly Aggregate ($^\circ$)@7	Average in Degrees
7 - Penrose	Wind Speed m/s.24-Hour Aggregate (m/s)@7	Average in Metres per second
7 - Penrose	PM10.24-Hour Aggregate ($\mu\text{g}/\text{m}^3$)@7	Average in Micrograms per cubic metre
7 - Penrose	PM2.5.24-Hour Aggregate ($\mu\text{g}/\text{m}^3$)@7	Average in Micrograms per cubic metre

Figure 2 Data extraction from Environmental Auckland Data Portal (Part 2)

start_time	end_time	PM2.5	SO2	NO	NO2	Solar_Rad	Temp	Humidity	Wind_Dir	Wind_Speed
1/1/19 0:00	1/1/19 1:00	NaN	NaN	NaN	NaN		20	78.1	219	3.05
1/1/19 1:00	1/1/19 2:00	NaN	NaN	NaN	NaN		19.5	78.95	217	2.8
1/1/19 2:00	1/1/19 3:00	NaN	NaN	NaN	5.3		19	79.85	219	2.45
1/1/19 3:00	1/1/19 4:00	NaN	NaN	NaN	4.65		19.5	79.9	219.5	2.3
1/1/19 4:00	1/1/19 5:00	NaN	NaN	NaN	4.15		20	78.9	218.5	2.7
1/1/19 5:00	1/1/19 6:00	NaN	NaN	NaN	4.2		19.5	78.75	218	2.75
1/1/19 6:00	1/1/19 7:00	NaN	NaN	NaN	5		19.5	77.9	216.5	2.7
1/1/19 7:00	1/1/19 8:00	NaN	NaN	NaN	4.6		20.5	74.2	215.5	3.15
1/1/19 8:00	1/1/19 9:00	NaN	NaN	NaN	4.2		21	71.5	217	3.55
1/1/19 9:00	1/1/19 10:00	NaN	NaN	NaN	4.45		21	70.5	216.5	3.75
1/1/19 10:00	1/1/19 11:00	NaN	NaN	NaN	4.55		21.5	68.45	217	3.75
1/1/19 11:00	1/1/19 12:00	NaN	NaN	NaN	4.95		22	67.2	219	3.95
1/1/19 12:00	1/1/19 13:00	NaN	NaN	NaN	4.65		22	68.65	225	4.4
1/1/19 13:00	1/1/19 14:00	NaN	NaN	NaN	4.4		22	69.4	225	4.6
1/1/19 14:00	1/1/19 15:00	NaN	NaN	NaN	3.85		22.5	67.35	224	4.8
1/1/19 15:00	1/1/19 16:00	NaN	NaN	NaN	3		23	66.65	232.5	5.15
1/1/19 16:00	1/1/19 17:00	NaN	NaN	NaN	2.55		22.5	67.25	236	5.35
1/1/19 17:00	1/1/19 18:00	NaN	NaN	NaN	2.9		22	68	231.5	5.15
1/1/19 18:00	1/1/19 19:00	NaN	NaN	NaN	4.2		21.5	70.4	221	4.75
1/1/19 19:00	1/1/19 20:00	NaN	NaN	NaN	5.35		21	71.2	216	4.35
1/1/19 20:00	1/1/19 21:00	NaN	NaN	NaN	5.8		21	70.5	215	3.9
1/1/19 21:00	1/1/19 22:00	NaN	NaN	NaN	5.4		20.5	71.85	214	3.9
1/1/19 22:00	1/1/19 23:00	NaN	NaN	NaN	4.4		20	73.45	212	3.95
1/1/19 23:00	2/1/19 0:00	NaN	NaN	NaN	2.85		20	73.95	213	4

Figure 3 Original Dataset Screen Shot

Explore dataset

Dataset Overview

(43800, 11)											
	start_time	end_time	PM2.5	SO2	NO	NO2	Solar_Rad	Temp	Humidity	Wind_Dir	Wind_Speed
0	1/1/19 0:00	1/1/19 1:00	NaN	NaN	NaN	NaN	NaN	20.0	78.10	219.0	3.05
1	1/1/19 1:00	1/1/19 2:00	NaN	NaN	NaN	NaN	NaN	19.5	78.95	217.0	2.80
2	1/1/19 2:00	1/1/19 3:00	NaN	NaN	NaN	5.30	NaN	19.0	79.85	219.0	2.45
3	1/1/19 3:00	1/1/19 4:00	NaN	NaN	NaN	4.65	NaN	19.5	79.90	219.5	2.30
4	1/1/19 4:00	1/1/19 5:00	NaN	NaN	NaN	4.15	NaN	20.0	78.90	218.5	2.70

Figure 4. Dataset Overview with dataset shape

The dataset has 11 columns with 43,800 records. And through call head() to check whether the data was downloaded.

```
RangeIndex: 43800 entries, 0 to 43799
Data columns (total 11 columns):
 #   Column      Non-Null Count  Dtype  
 ---  --          --          --    
 0   start_time  43800 non-null   object 
 1   end_time    43800 non-null   object 
 2   PM2.5       35778 non-null   float64
 3   SO2         38694 non-null   float64
 4   NO          38768 non-null   float64
 5   NO2         39047 non-null   float64
 6   Solar_Rad   0 non-null     float64 
 7   Temp        40325 non-null   float64
 8   Humidity    40471 non-null   float64
 9   Wind_Dir    39162 non-null   float64
 10  Wind_Speed  40236 non-null   float64
 dtypes: float64(9), object(2)
 memory usage: 3.7+ MB
```

Figure 5 Dataset information

By calling the info() method, check data types, detect missing values, and understand memory usage. The start_time and end_time columns are of the object type, not numerical. The remaining columns are of type float64, indicating numerical data.

- There are missing values in the PM2.5, SO2, NO, NO2, Temp, Humidity, Wind_Dir, and Wind_Speed columns.
- The Solar_rad column has no data.
- This used 3.7+ MB memory

Check the data distribution

Missing values Heat Map: Use a Heat Map to see which fields have missing values and the distribution of missing values.

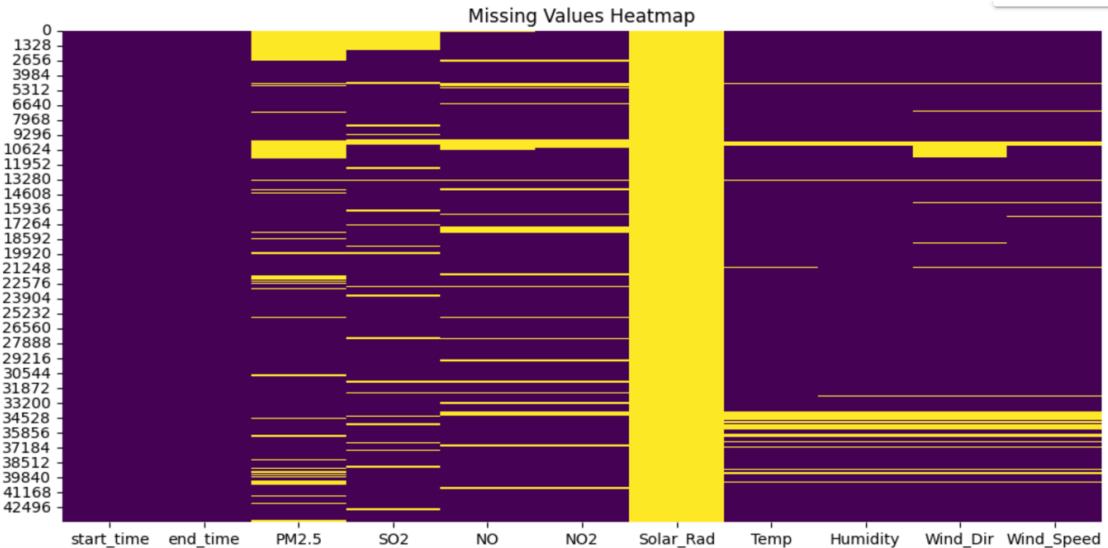


Figure 6 Missing values Heat Map

Through the analysis of the missing value heatmap, we observe that Solar_Rad is entirely missing. PM2.5, SO2, NO, and NO2 have missing values scattered across the dataset. Temp, Humidity, Wind_Dir, and Wind_Speed have a few missing values before row 21,248, no missing values between rows 21,248 and 31,872, and more severe missing values between rows 33,200 and 39,840.

Histograms: Check the distribution of each field to understand the overall shape of the data, whether there is skewness, kurtosis, etc.

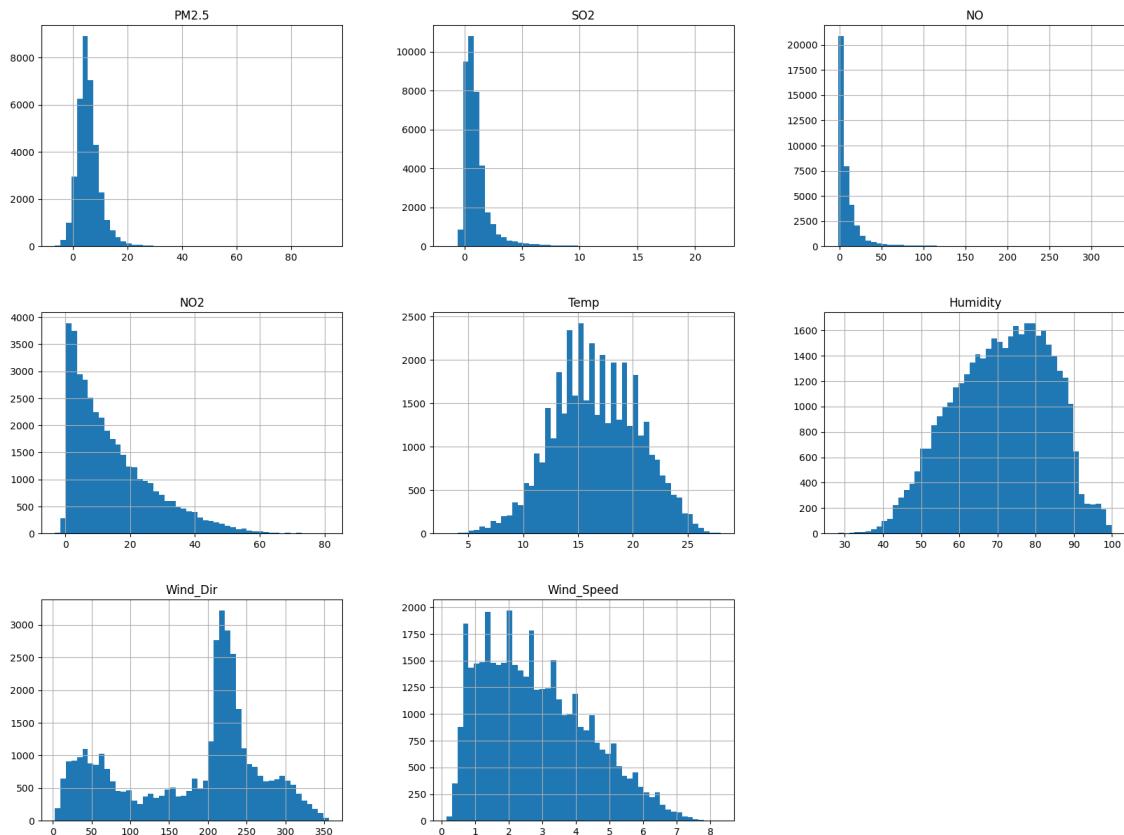


Figure 7 Dataset histogram distribution

From the histograms, we observe the presence of skewness and kurtosis but no outliers in any of the fields. PM2.5, SO2, NO, NO2, Humidity, Wind_Dir, and Wind_Speed are not symmetric, while Temp is nearly symmetric. PM2.5, SO2, NO, NO2, and Wind_Speed show positive skewness (right-skewed), while Humidity and Wind_Dir show negative skewness (left-skewed).

- PM2.5 has a bell-shaped distribution, with values mainly between -10 and 20, peaking around 5.
- SO2 has a roughly bell-shaped distribution, with values mainly between -2 and 5, peaking around 1.
- NO has a non-bell-shaped distribution, with values mainly between 0 and 50, peaking around 5.
- NO2 has a non-bell-shaped distribution, with values mainly between 0 and 60, peaking around 5.
- Temp has a bell-shaped distribution, with values mainly between 5 and 27, peaking around 78.
- Wind_Dir has a roughly bell-shaped distribution, with values mainly between 0 and 350, peaking around 225.
- Wind_Speed has a non-bell-shaped distribution, with values mainly between 0 and 7.5, peaking around 0.5.

Box plots: Detect outliers in the data. Box plots can clearly display the quartiles, maximum values, minimum values, and outliers of the data.

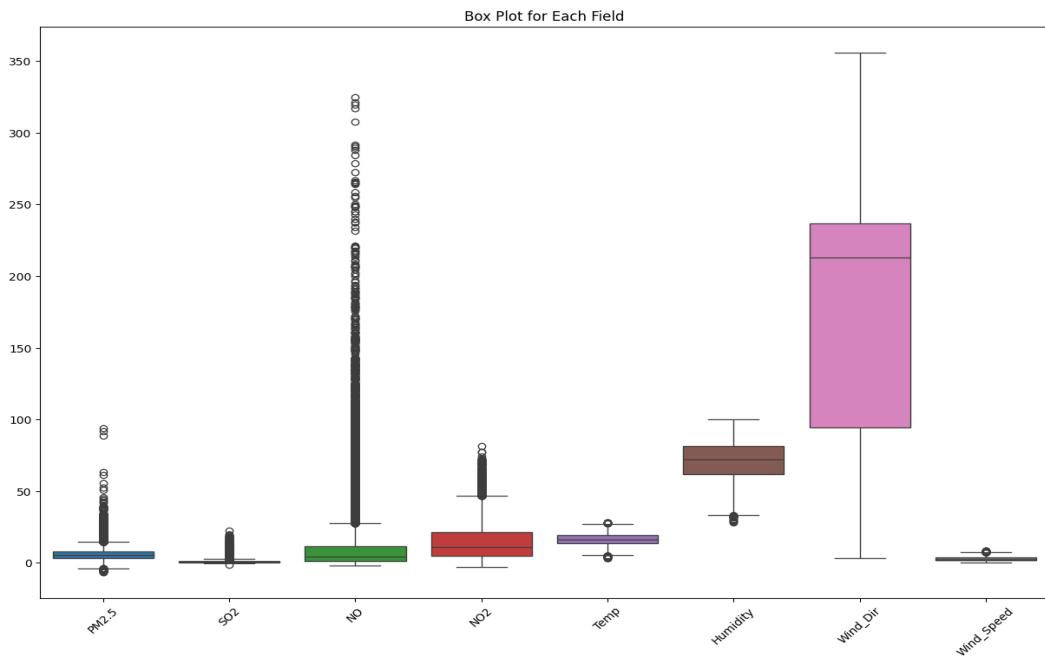


Figure 8 Dataset outliers with Box Plot

Through the observation of box plots, it was found that PM2.5, SO2, NO, NO2, Temp, Humidity, and Wind_Speed all contain outliers. Notably, NO, PM2.5, NO2, and SO2 have a significant number of outliers on the upper side. In contrast, Wind_Dir does not have any outliers.

- PM2.5: The data is symmetric and primarily distributed within a small range of values (2-10).
- SO2: The box is short and close to the bottom, indicating that most data is highly concentrated, mainly within the range of 0-2.

- NO: The box is lower, with most data values being small. The median is lower, indicating skewness, and data is mainly distributed between 0-18.
- NO₂: Similar to NO, the box is lower with most data values being small. The median is lower, indicating skewness, with data primarily distributed between 5-20.
- Temp: The data is symmetric, with 50% of the data distributed between 15-20.
- Humidity: The data is nearly symmetric, but the box is slightly higher, primarily distributed between 70-90.
- Wind_Dir: The values are quite dispersed, ranging from 0-360, with a noticeable skewness towards the higher middle range. The data is mainly distributed between 90-240.
- Wind_Speed: The box is very short and close to the bottom, indicating that most data is quite similar, primarily distributed around 0.

Time series plots: For time series data, time series plots can be drawn to view the changing trend of each variable over time.

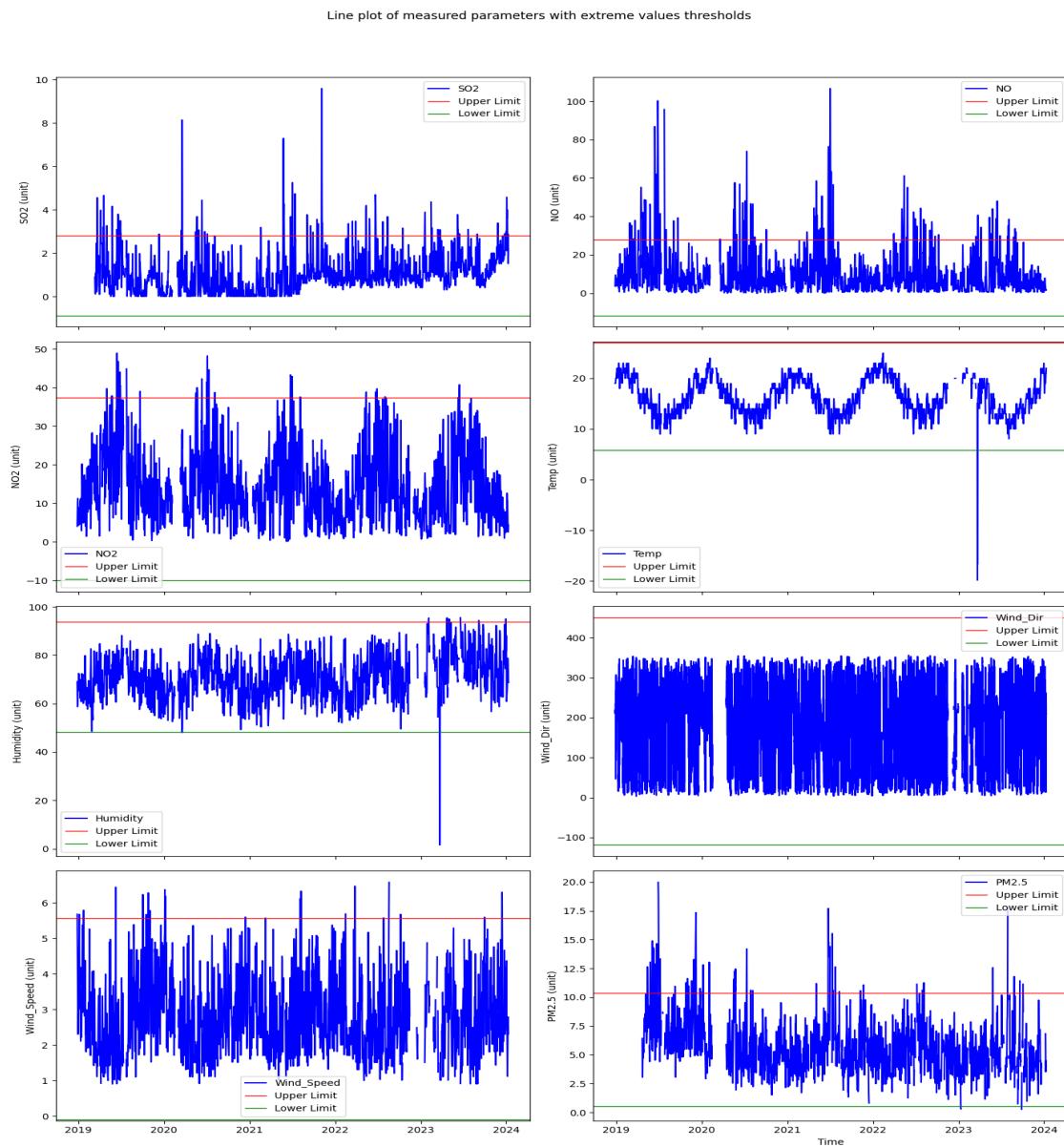


Figure 9 Dataset Time Series Plot

Through the time series plots, there is no clear trend in the data, but Wind_Dir and NO exhibit significant fluctuations, indicating weaker stationarity. Humidity and NO2 also show fluctuations, suggesting poor stationarity.

Conclusion

After examining the dataset using various visualization techniques such as Heat Map, histograms, box plots, and time series plots, several issues have been identified. These include missing values, outliers, uneven distributions, and potential data errors. It is evident that the dataset requires preprocessing to address these issues and ensure its suitability for analysis. Steps such as imputing missing values, handling outliers, and correcting data errors will be essential to prepare the dataset for further analysis. By addressing these challenges effectively, we can improve the quality and reliability of the dataset, leading to more accurate and robust model predictions.

1.3 Pre-Processing

1.3.1 Delete 'Solar_Rad' column

Based on the above analysis, since the entire 'Solar_Rad' column has no data, it will be deleted.

	start_time	end_time	PM2.5	SO2	NO	NO2	Temp	Humidity	Wind_Dir	Wind_Speed
0	1/1/19 0:00	1/1/19 1:00	NaN	NaN	NaN	NaN	20.0	78.10	219.0	3.05
1	1/1/19 1:00	1/1/19 2:00	NaN	NaN	NaN	NaN	19.5	78.95	217.0	2.80
2	1/1/19 2:00	1/1/19 3:00	NaN	NaN	NaN	5.30	19.0	79.85	219.0	2.45
3	1/1/19 3:00	1/1/19 4:00	NaN	NaN	NaN	4.65	19.5	79.90	219.5	2.30
4	1/1/19 4:00	1/1/19 5:00	NaN	NaN	NaN	4.15	20.0	78.90	218.5	2.70

Figure 10 Dataset After Removing Solar_Rad

1.3.2 Converted the Data Type

Columns are converted to appropriate data types. This preprocessing step is crucial for ensuring the data is in the correct format for further analysis and modelling. Specifically:

- Converting 'start_time' and 'end_time' to datetime: This ensures that time-based operations and analyses can be performed accurately.

- Converting other columns to float: The columns 'Temp', 'Humidity', 'Wind_Speed', 'NO', 'NO2', 'SO2', 'PM2.5', and 'Wind_Dir' are converted to float to facilitate numerical computations and statistical analyses.

Table 1 Original Dataset Structure

Column Name	Data Type	Description	Unit
start_time	datetime64[ns]	Start time of measurement	datetime
end_time	datetime64[ns]	End time of measurement	datetime
PM2.5	float64	PM2.5 concentration	µg/m³
SO2	float64	SO2 concentration	µg/m³
NO	float64	NO concentration	µg/m³
NO2	float64	NO2 concentration	µg/m³
Temp	float64	Temperature	°C
Humidity	float64	Humidity	%
Wind_Dir	float64	Wind direction	Degrees
Wind_Speed	float64	Wind speed	m/s

1.3.3 Check and process duplicate values, missing values, errors, and outliers

Check:

- Check Duplicate Values: Check if there are any duplicate records in the data, as duplicates may affect the training of the model and the accuracy of results.

Check duplicate values: 0

Figure 11 Duplicate Check

There are no duplicate values.

- Check Time: Check whether data collection is performed on an hourly basis per row of the dataset. In time series analysis, especially for Long Short-Term Memory (LSTM) networks, it is critical to maintain a continuous sequence of data points. Discontinuities disrupt the temporal relationships that LSTMs rely on, leading to unstable training and inaccurate predictions.

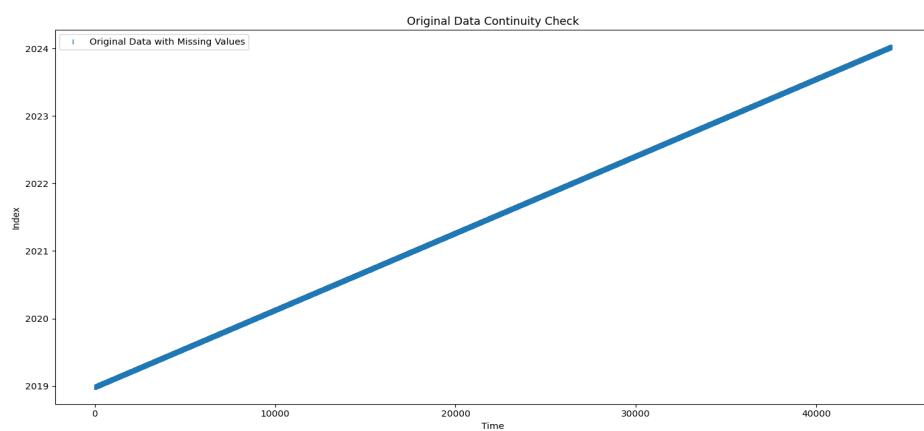


Figure 12 Dataset time series analysis

By observing the line chart, the original time series data is recorded continuously on an hourly basis.

- Check Missing Values: Examine whether there are any missing values in the data, as missing values may lead to bias and inaccuracy in the model.

```
start_time      0
end_time       0
PM2.5        8022
SO2          5106
NO           5032
NO2          4753
Temp         3475
Humidity     3329
Wind_Dir      4638
Wind_Speed    3564
dtype: int64
duplicated: 0
```

Figure 13 Missing Values Check in the Dataset

- Check Data Errors: Inspect if there are any obvious errors or anomalies in the data, such as values outside reasonable ranges or inconsistent data, which may need to be corrected or excluded.

Table 2 Dataset Error Summary

Column Name	Rule	Error Count	Percentage
Temp	Temperature should be between -20°C and 40°C	29	0.07%
Humidity	Humidity should be between 0% and 100%	13	0.03%
Wind_Speed	Wind speed should be between 0 and 60 m/s	47	0.11%
Wind_Dir	Wind direction should be between 0 and 360 degrees	0	0.00%
NO	NO should be non-negative	48	0.11%
NO2	NO2 should be non-negative	0	0.00%
SO2	SO2 should be non-negative	1025	2.34%
PM2.5	PM2.5 should be non-negative	0	0.00%

- Check Outliers: Assess whether there are any outliers in the data, as outliers may adversely impact the model and require further processing or removal.

Through the analysis of the box plot, it was found that there are numerous outliers in the data. We can identify outliers using the Interquartile Range (IQR) method. The IQR is the range between the first quartile (Q1, the 25th percentile) and the third quartile (Q3, the 75th percentile), representing the box plot's interquartile range. Outliers are defined as data points that fall below the lower bound or above the upper bound, where:

- Lower Bound = 25th percentile - 1.5 * IQR
- Upper Bound = 75th percentile + 1.5 * IQR

Any data points below the lower bound or above the upper bound are considered outliers.

Table 3 Dataset Outlier Summary

Column Name	Outlier Count	Percentage
SO2	1620	3.668478
NO	2827	6.401721
NO2	548	1.240942
Temp	13	0.029438
Humidity	100	0.226449
Wind_Dir	0	0.000000
Wind_Speed	342	0.774457
PM2.5	1172	2.653986

Processing:

The typical approach for handling these issues involves deleting outliers and interpolating missing values. Since LSTM models require continuous data, interpolation will be chosen as the preferred method of handling missing values.

- Interpolating values: We address the missing data by applying linear interpolation and KNN imputation methods. The effectiveness of these methods is compared using mean squared error (MSE) and mean absolute error (MAE).

We visualize the continuity of each parameter before and after imputation.

1. Original Data:

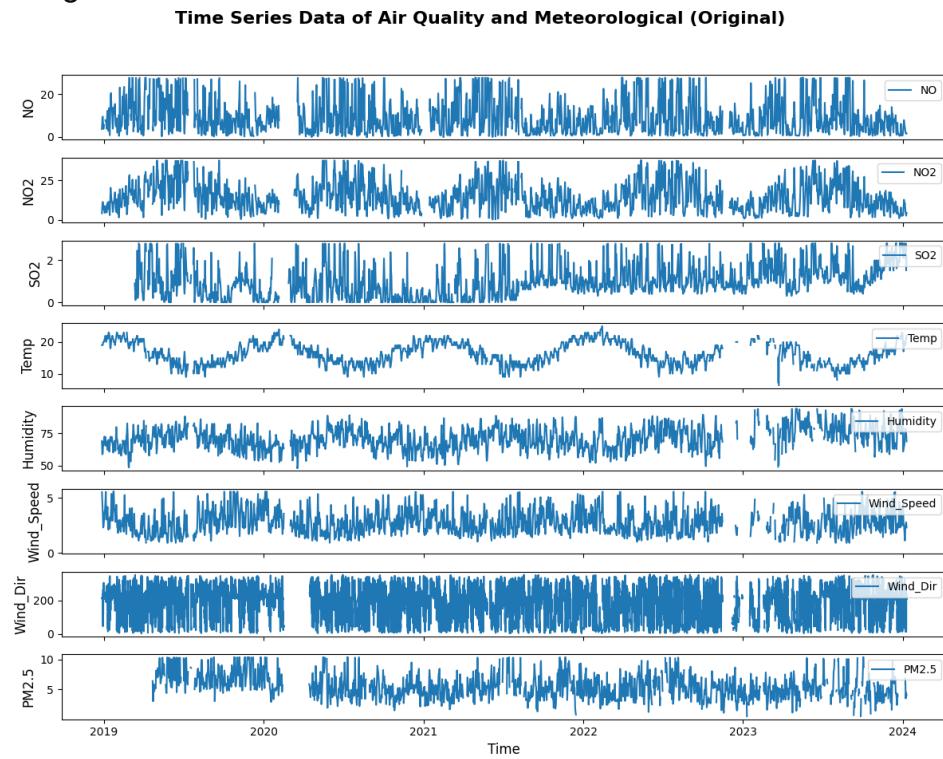


Figure 14 Time Series Data of Air Quality and Meteorological(Original)

2. Linear Interpolation: Linear interpolation is used to fill in missing values by linearly estimating values based on surrounding data points. Special consideration is given to Wind_Dir due to its circular nature.

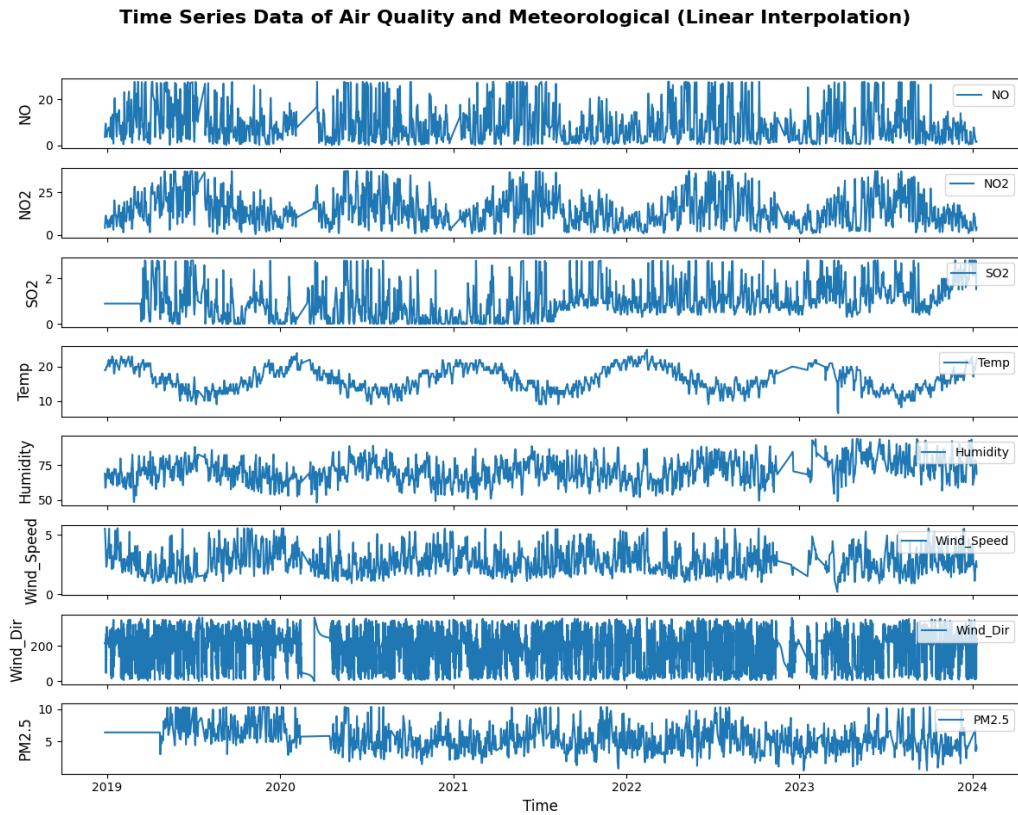


Figure 15 Time Series Data of Air Quality and Meteorological (Linear Interpolation)

3. KNN Interpolation: KNN imputation fills in missing values by considering the mean value of the k-nearest neighbors($k=5$). Wind_Dir is also handled using circular statistics.

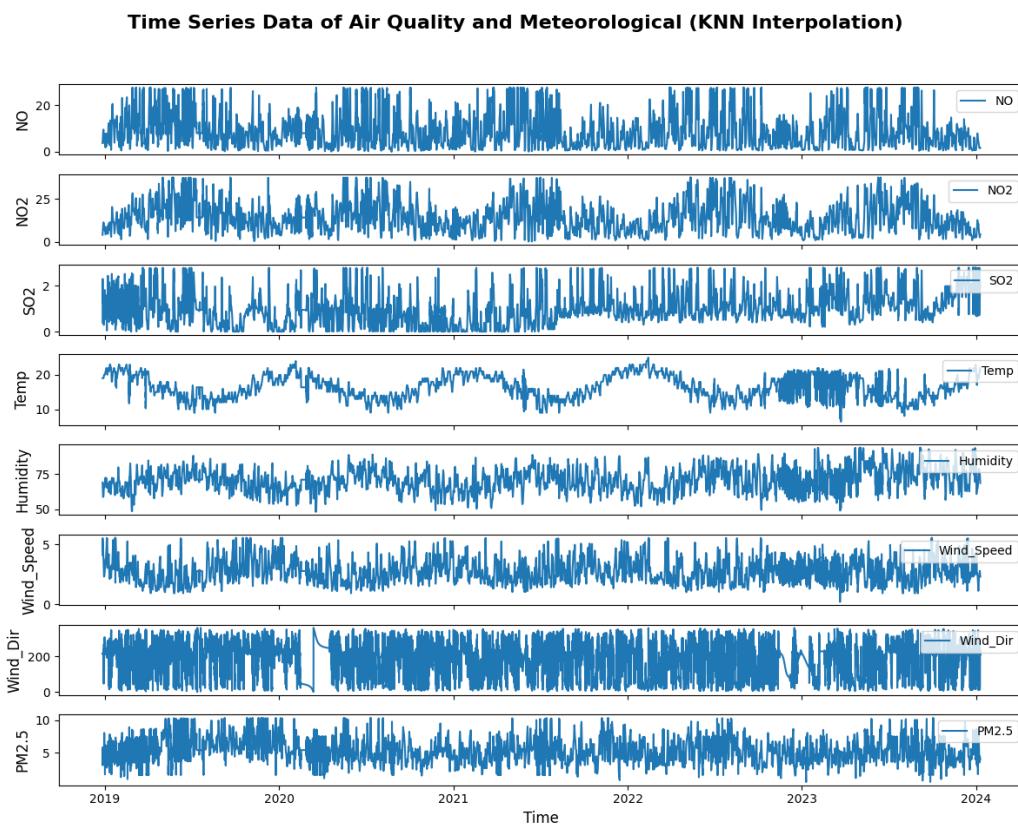


Figure 16 Time Series Data of Air Quality and Meteorological (KNN Interpolation)

Special Consideration for Wind Direction (Wind_Dir):

Wind_Dir is a circular variable, ranging from 0 to 360 degrees. This characteristic makes it different from other numeric variables. Simple numerical imputation methods like linear interpolation and KNN may not be appropriate as they do not account for the circular nature of the data. For example, both 1 degree and 359 degrees represent almost the same direction, but standard methods might treat them as far apart. We use circular statistics by converting wind direction to sine and cosine components, interpolate these components, and then convert back to angles.

Error Comparison:

We compare the effectiveness of the two imputation methods using MSE and MAE metrics.

Table 4 Imputation Methods Comparison

Column	MSE_Linear	MAE_Linear	MSE_KNN	MAE_KNN
0 Temp	0.100153	0.074096	1.811506	0.340597
1 Humidity	2.777202	0.357834	14.774094	0.955219
2 Wind_Speed	0.026439	0.035978	0.156855	0.098458
3 NO	2.799736	0.364949	14.217260	1.151493
4 NO2	1.524455	0.237843	11.171873	0.734658
5 SO2	0.009968	0.013619	0.138475	0.108127
6 PM2.5	0.352098	0.144336	1.658225	0.482346
7 Wind_Dir	979.095523	6.294835	979.095523	6.294835

Visualization of Errors:

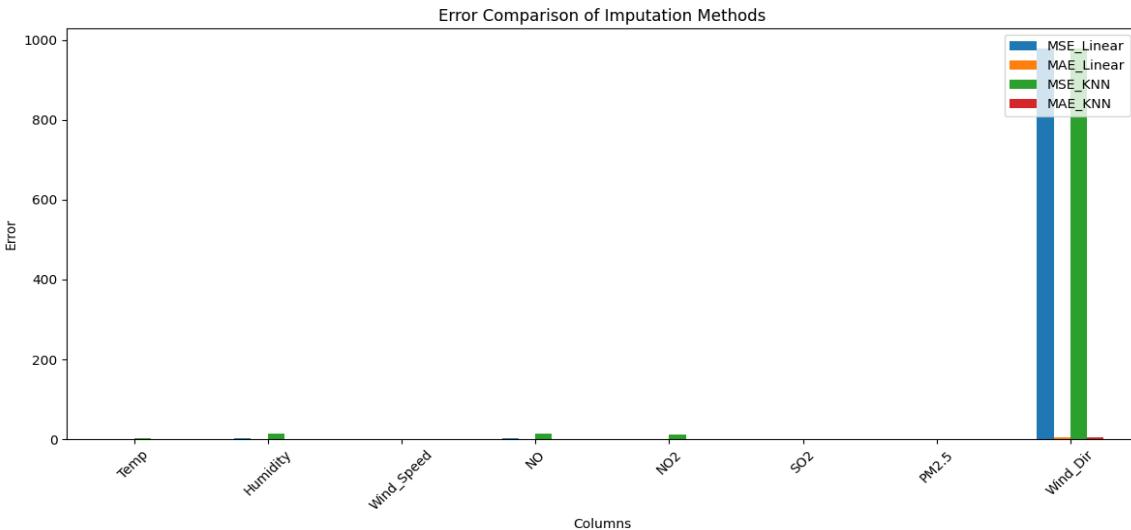


Figure 17 Error Comparison of Imputation Methods

Conclusion

Based on the error metrics (MSE and MAE), linear interpolation generally performed better than KNN interpolation for most columns. The KNN method showed significantly higher errors, especially for Wind_Dir, NO, and NO2. Therefore, linear interpolation is the preferred method for imputing missing values in this dataset. However, for Wind_Dir, we used a circular statistics approach which is more suitable

for its nature. Both methods yielded similar results for Wind_Dir, indicating the necessity of specialized techniques for circular data.

1.3.4. Interpolation Effect Review

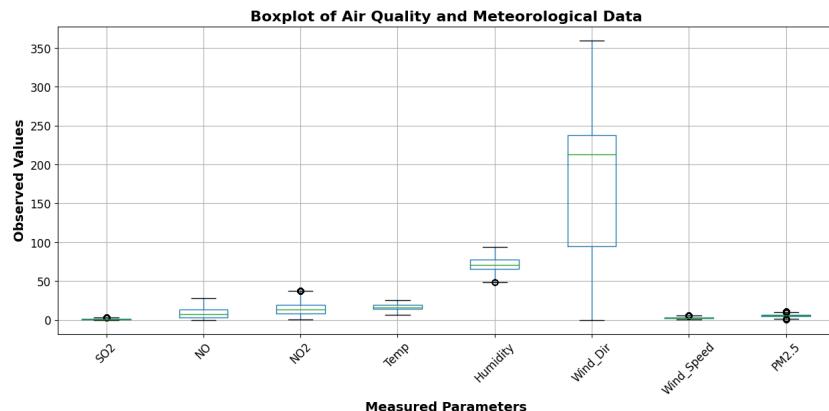


Figure 18 Boxplot of Air Quality and Meteorological Data

The box plot illustrates the air quality and meteorological data after invalid data and outliers were replaced with null values. Post-interpolation for outlier handling, the data has become smoother, mitigating some extreme values, but outliers are still observable in the box plot. Due to the inherent variability or natural outliers in the data, interpolation does not eliminate these outliers as they are part of the true data distribution. The median is positioned near the center of the box, indicating effective data processing. Notably, "Wind_Dir" shows a wide range of values, whereas other values like "PM2.5" and "NO2" remain within a narrower range.

1.3.5. Extract lag1 and lag2 Feature

We will add two columns, Two PM_lag measurements, lag1 and lag2, to help the model understand how past values influence future PM2.5 levels. This temporal dependency is crucial for time series forecasting tasks.

- lag1: PM2.5 value from 1 hour ago. This feature helps capture the short-term temporal dependency of PM2.5 levels.
- lag2: PM2.5 value from 2 hours ago. This feature helps capture the medium-term temporal dependency of PM2.5 levels.

We are using data from 2019-01-01 00:00 to 2023-12-31 00:00. When creating the lag features, we realized that the lag requires prior data, so we downloaded data **from 2018-12-27 00:00 to 2024-01-19 23:00 and applied all the above preprocessing steps**, which are **not repeated here**.

Handling Missing Values for lag1, and lag2: To ensure the completeness of the dataset, we deleted rows where the future lag1, or lag2 was NaN. These missing values occur when there is not enough data to calculate the future or lagged PM2.5 values. Deleting these rows ensures that the model has a complete dataset for training and prediction.

	start_time	end_time	SO2	NO	NO2	Temp		
2	2018-12-27 02:00:00	2018-12-27 03:00:00	0.89375	4.010417	4.452083	19.0		
3	2018-12-27 03:00:00	2018-12-27 04:00:00	0.89375	4.214583	4.672917	19.0		
4	2018-12-27 04:00:00	2018-12-27 05:00:00	0.89375	4.418750	4.893750	19.0		
5	2018-12-27 05:00:00	2018-12-27 06:00:00	0.89375	4.622917	5.114583	19.0		
6	2018-12-27 06:00:00	2018-12-27 07:00:00	0.89375	4.827083	5.335417	19.0		
	Humidity	Wind_Dir	Wind_Speed	PM2.5	lag1	lag2	Month	\
2	68.195833	213.5	5.53125	6.427083	6.427083	6.427083	12	
3	67.754167	213.5	5.53125	6.427083	6.427083	6.427083	12	
4	67.312500	212.0	5.53125	6.427083	6.427083	6.427083	12	
5	66.870833	211.5	5.49375	6.427083	6.427083	6.427083	12	
6	66.429167	212.0	5.45625	6.427083	6.427083	6.427083	12	

Figure 19 Data Review with Lag1 and Lag2

1.3.6. Extract Time Feature

Through analysis, it has been observed that PM2.5 levels often exhibit seasonal and weekly patterns. For instance, due to seasonal factors, higher pollution levels may be observed in certain months, or due to traffic patterns, on specific days of the week. Therefore, we will extract Month and Day_Of_Week from the start_time column to capture these seasonal and weekly patterns.

- Month: The month of the year (1-12).
- Day_Of_Week: The day of the week (0-6, where 0 represents Monday and 6 represents Sunday).

We applied one-hot encoding to the Month and Day_Of_Week columns and converted the one-hot encoded columns to float64 type for consistency and to avoid potential issues with model training.

Month_10	Month_11	Month_12	Day_0	Day_1	Day_2	Day_3	Day_4	Day_5	Day_6
0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0

Figure 20 Data Review with New Features

1.3.7. Extract Wind Direction Feature

Wind direction is represented by degrees from 0 to 360. Due to its circular nature, it is not a standard numerical value. Given the 360-degree range of wind direction, we simplify the data into 16 categories based on the standard 16-point compass rose, like N, NNE, NE, ENE, E, ESE, SE, SSE, S, SSW, SW, WSW, W, WNW, NW, NNW. We processed this into 16 columns and used one-hot encoding to process and convert the one-hot encoded columns to float64 type for consistency and to avoid potential issues with model training. For example: N (North wind). If there is a north wind, it is marked as 1; if there is no north wind, it is marked as 0. By discretizing the wind direction in this manner, we avoid the discontinuity between 0 and 360 degrees, making the data more manageable for algorithms. This enhances the model's predictive capability, interpretability, and efficiency.

Wind_Type_E	44158	non-null	float64
Wind_Type_ENE	44158	non-null	float64
Wind_Type_ESE	44158	non-null	float64
Wind_Type_N	44158	non-null	float64
Wind_Type_NE	44158	non-null	float64
Wind_Type_NNE	44158	non-null	float64
Wind_Type_NNW	44158	non-null	float64
Wind_Type_NW	44158	non-null	float64
Wind_Type_S	44158	non-null	float64
Wind_Type_SE	44158	non-null	float64
Wind_Type_SSE	44158	non-null	float64
Wind_Type_SSW	44158	non-null	float64
Wind_Type_SW	44158	non-null	float64
Wind_Type_W	44158	non-null	float64
Wind_Type_WNW	44158	non-null	float64
Wind_Type_WSW	44158	non-null	float64

Figure 21 Wind Direction Features

1.3.8. Final dataset

	start_time	end_time	SO2	NO	NO2	Temp	Humidity	Wind_Speed	PM2.5	lag1	...	Month_10	Month_11	Month_12	Day_0	Day_1	Day_2	Day_3	Day_4	Day_5	Day_6
2	2018-12-27 02:00:00	2018-12-27 03:00:00	0.89375	4.010417	4.452083	19.0	68.195833	5.53125	6.427083	6.427083	...	0.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0
3	2018-12-27 03:00:00	2018-12-27 04:00:00	0.89375	4.214583	4.672917	19.0	67.754167	5.53125	6.427083	6.427083	...	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
4	2018-12-27 04:00:00	2018-12-27 05:00:00	0.89375	4.418750	4.893750	19.0	67.312500	5.53125	6.427083	6.427083	...	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
5	2018-12-27 05:00:00	2018-12-27 06:00:00	0.89375	4.622917	5.114583	19.0	66.870833	5.49375	6.427083	6.427083	...	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
6	2018-12-27 06:00:00	2018-12-27 07:00:00	0.89375	4.827083	5.335417	19.0	66.429167	5.45625	6.427083	6.427083	...	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0

5 rows x 46 columns

Figure 22 Final Dataset Review

Table 5 Final Dataset Structure

Column Name	Data Type	Meaning	Contains Null Values
PM2.5	float64	PM2.5 concentration	No
SO2	float64	SO2 concentration	No
NO	float64	NO concentration	No
NO2	float64	NO2 concentration	No
Temp	float64	Temperature	No
Humidity	float64	Humidity	No
Wind_Speed	float64	Wind Speed	No
lag1	float64	PM2.5 concentration 24 hours before	No
lag2	float64	PM2.5 concentration 48 hours before	No
Wind_Type_E	float64	Wind direction: East	No
Wind_Type_ENE	float64	Wind direction: East-Northeast	No
Wind_Type_ESE	float64	Wind direction: East-Southeast	No
Wind_Type_N	float64	Wind direction: North	No
Wind_Type_NE	float64	Wind direction: Northeast	No
Wind_Type_NNE	float64	Wind direction: North-Northeast	No
Wind_Type_NNW	float64	Wind direction: North-Northwest	No
Wind_Type_NW	float64	Wind direction: Northwest	No

Wind_Type_S	float64	Wind direction: South	No
Wind_Type_SE	float64	Wind direction: Southeast	No
Wind_Type_SSE	float64	Wind direction: South-Southeast	No
Wind_Type_SSW	float64	Wind direction: South-Southwest	No
Wind_Type_SW	float64	Wind direction: Southwest	No
Wind_Type_W	float64	Wind direction: West	No
Wind_Type_WNW	float64	Wind direction: West-Northwest	No
Wind_Type_WSW	float64	Wind direction: West-Southwest	No
Month_1	float64	Month: January	No
Month_2	float64	Month: February	No
Month_3	float64	Month: March	No
Month_4	float64	Month: April	No
Month_5	float64	Month: May	No
Month_6	float64	Month: June	No
Month_7	float64	Month: July	No
Month_8	float64	Month: August	No
Month_9	float64	Month: September	No
Month_10	float64	Month: October	No
Month_11	float64	Month: November	No
Month_12	float64	Month: December	No
Day_0	float64	Day of the week: Monday	No
Day_1	float64	Day of the week: Tuesday	No
Day_2	float64	Day of the week: Wednesday	No
Day_3	float64	Day of the week: Thursday	No
Day_4	float64	Day of the week: Friday	No
Day_5	float64	Day of the week: Saturday	No
Day_6	float64	Day of the week: Sunday	No

The above table is the final dataset and dataset structure after preprocessing.

2. Features Selection

2.1. Extract Method and Reason

We first visually inspected scatter plots for each predictor factor against PM2.5 concentration. Due to the massive dataset, it wasn't feasible to observe each factor's relationship with PM2.5 directly. Therefore, we randomly selected 500 data points and plotted scatter plots.

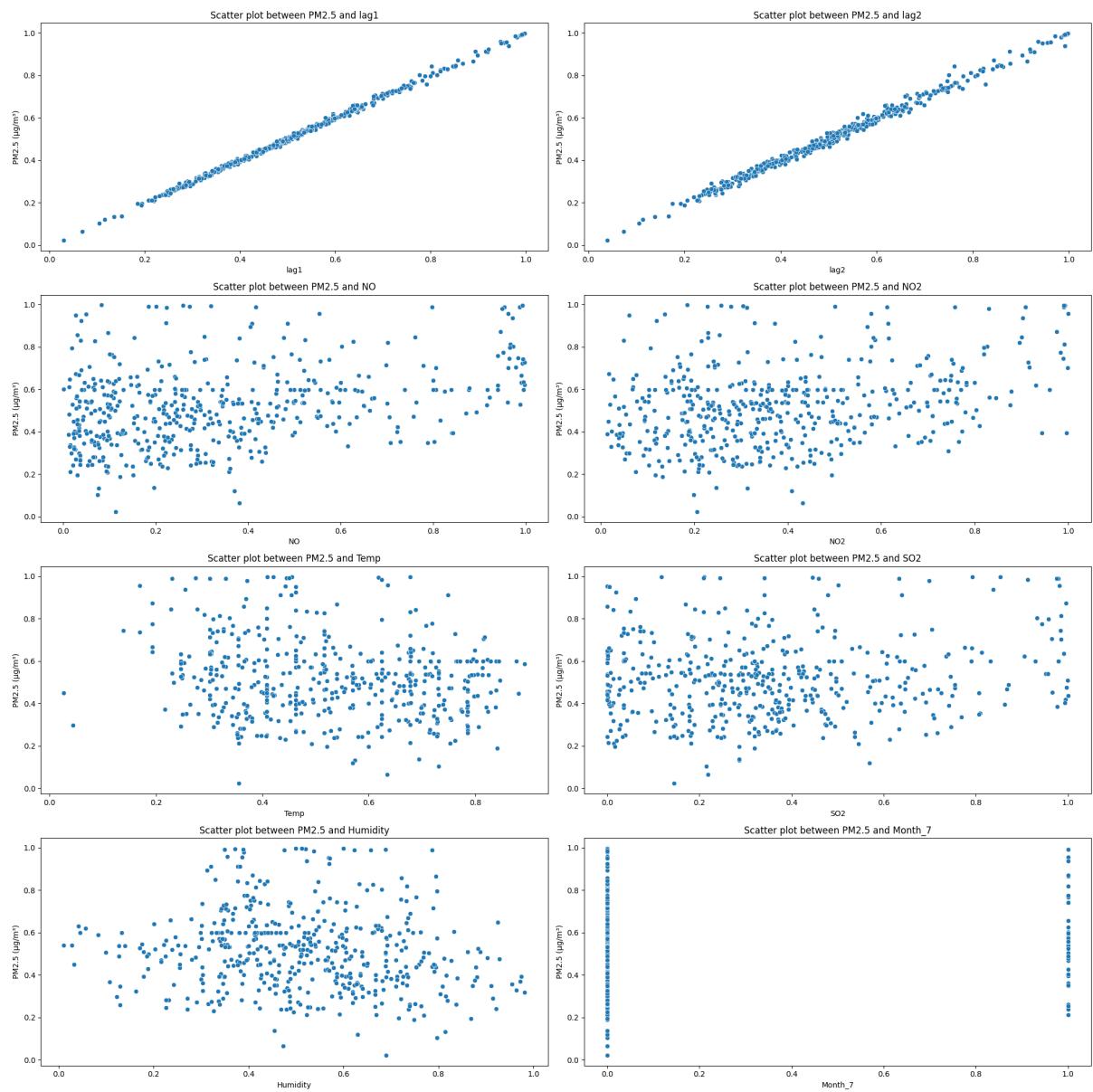


Figure 23 Each Factor's Relationship with PM2.5

After observing the scatter plots, we found that Log1 and Log2 exhibit a linear distribution, indicating a linear relationship with PM2.5 concentration. However, the scatter plots of other variables show scattered distributions, indicating non-linear relationships with PM2.5 concentration.

Table 6 Comparison of Statistical Analysis Methods

Method	Advantages	Disadvantages
Pearson Correlation	Simple and intuitive to compute, suitable for linear relationships	Only applicable to linear relationships, sensitive to outliers, assumes normal distribution
Spearman Rank Correlation	Suitable for nonlinear relationships and non-normal data	Only captures monotonic relationships, cannot handle complex nonlinear relationships
Kendall's Tau	Suitable for small samples and ordered data	High computational complexity, inefficient for large datasets
Chi-Squared Test	Suitable for categorical data	Only applicable to categorical data, cannot handle continuous data
Mutual Information	Can capture any type of relationship, including nonlinear	High computational complexity, requires large amounts of data for accurate estimation
ANOVA (Analysis of Variance)	Suitable for continuous dependent variables and categorical independent variables	Assumes normal distribution and homogeneity of variances
Lasso Regression	Performs both feature selection and regression, suitable for high-dimensional data	High computational complexity, requires parameter tuning
Decision Trees	Can capture complex nonlinear relationships, intuitive and easy to interpret	Prone to overfitting, needs to be combined with other methods like Random Forest for stability

There are several feature selection methods available, each with its own advantages and disadvantages. Pearson correlation is suitable for linear relationships, while Spearman rank correlation and decision trees are suitable for non-linear relationships. However, since we will be using MLP and LSTM models for training, and decision trees may suffer from overfitting issues, we will initially use Pearson correlation for screening, followed by comparison using Spearman rank correlation.

Table 7 Correlation Coefficients for Different Parameters

Parameters	Pearson_Correlation	Spearman_Correlation
Lag1	0.999043	0.998773
Lag2	0.996238	0.995631
NO	0.370475	0.350023
NO2	0.366291	0.330139
Temp	-0.19689	-0.16745
SO2	0.176059	0.149126
Humidity	-0.148571	-0.158478
Month_7	0.143165	0.129658
Month_10	-0.111916	-0.112441
Month_5	0.0866676	0.071302
Month_3	-0.0861609	-0.0693056
Wind_Type_SW	0.084419	0.0846687
Month_6	0.0758995	0.069287
Wind_Type_NE	-0.0755627	-0.0735857
Wind_Type_NNE	-0.0723409	-0.0634552
Wind_Type_WSW	0.0587131	0.0658568
Wind_Type_ENE	-0.0581009	-0.0637431
Day_0	-0.0539499	-0.0541181
Month_1	-0.0519881	-0.037144

Pearson Correlation Matrix (Top 8 Features)

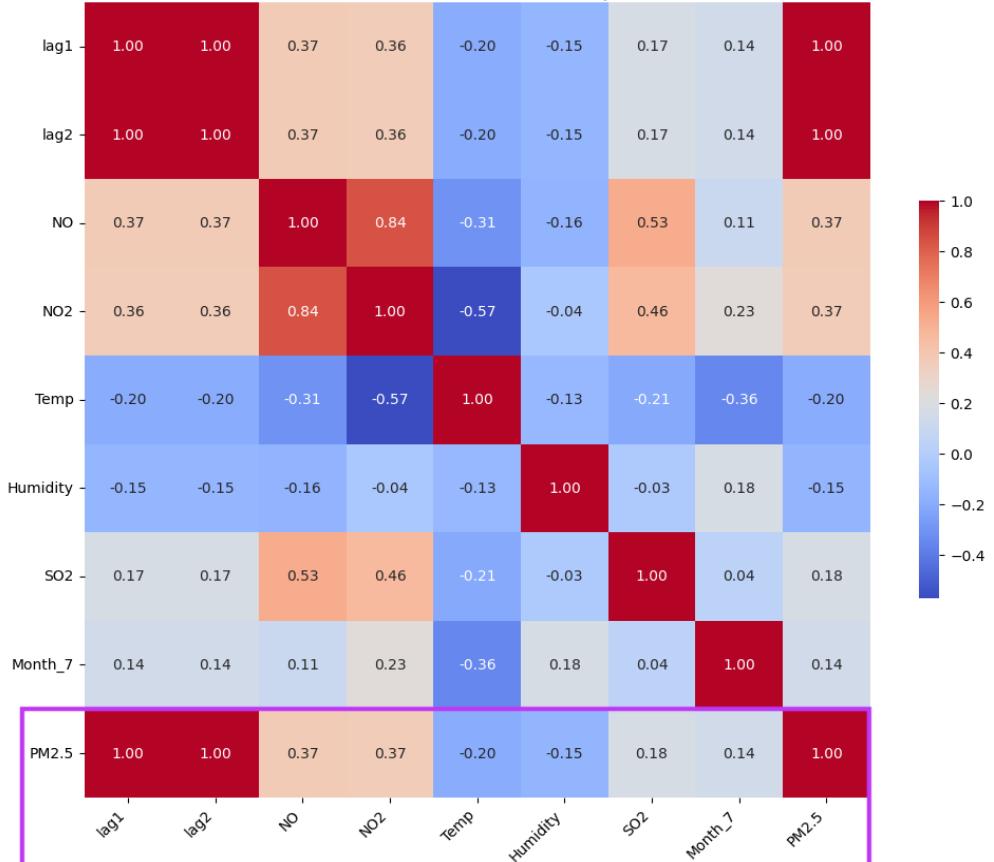


Figure 24 Pearson Correlation Matrix (Top 8 Features)

Spearman Correlation Matrix (Top 8 Features)

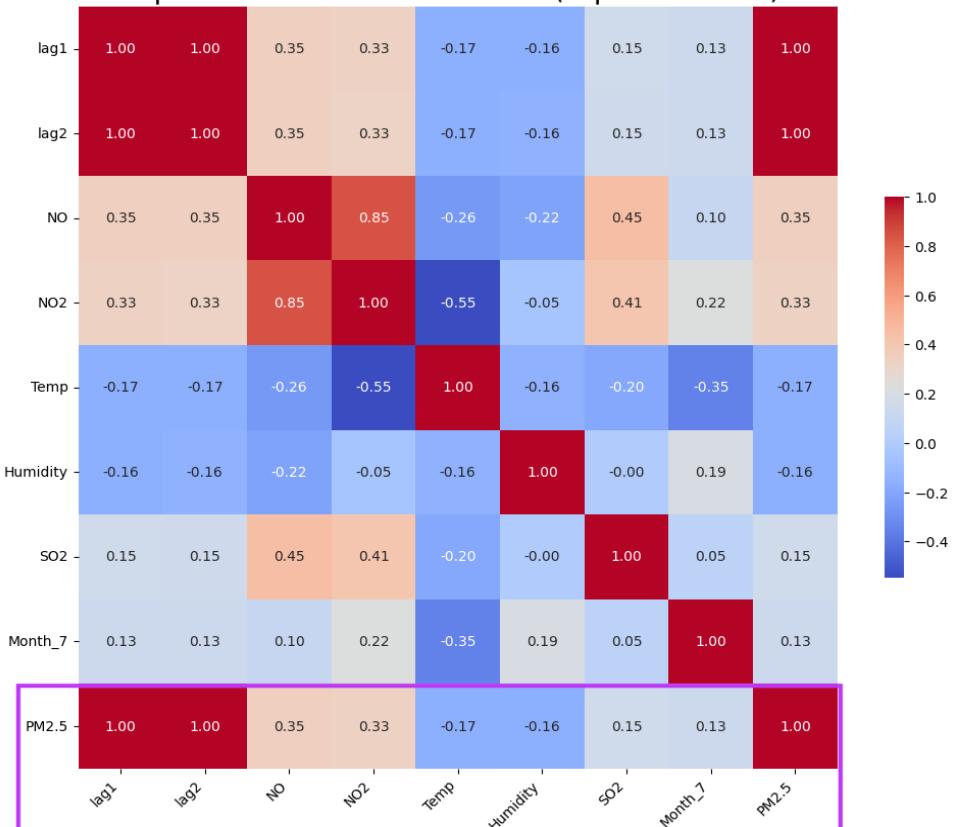


Figure 25 Spearman Correlation Matrix (Top 8 Features)

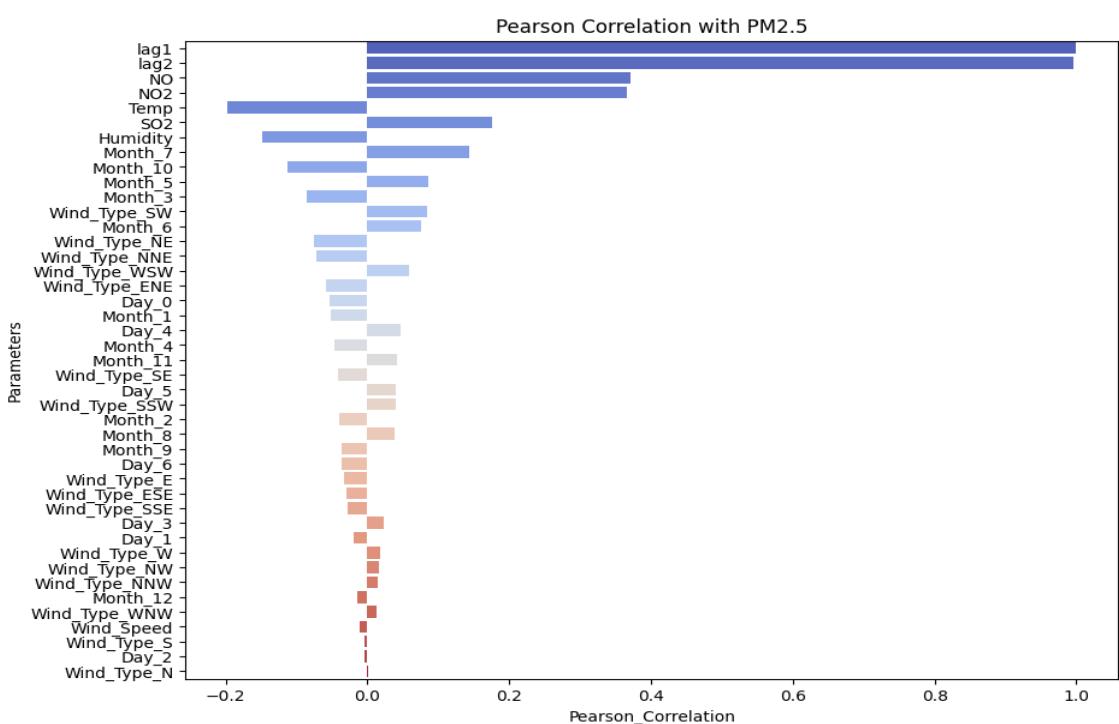


Figure 26 All factors' Pearson Correlation with PM2.5 Plot

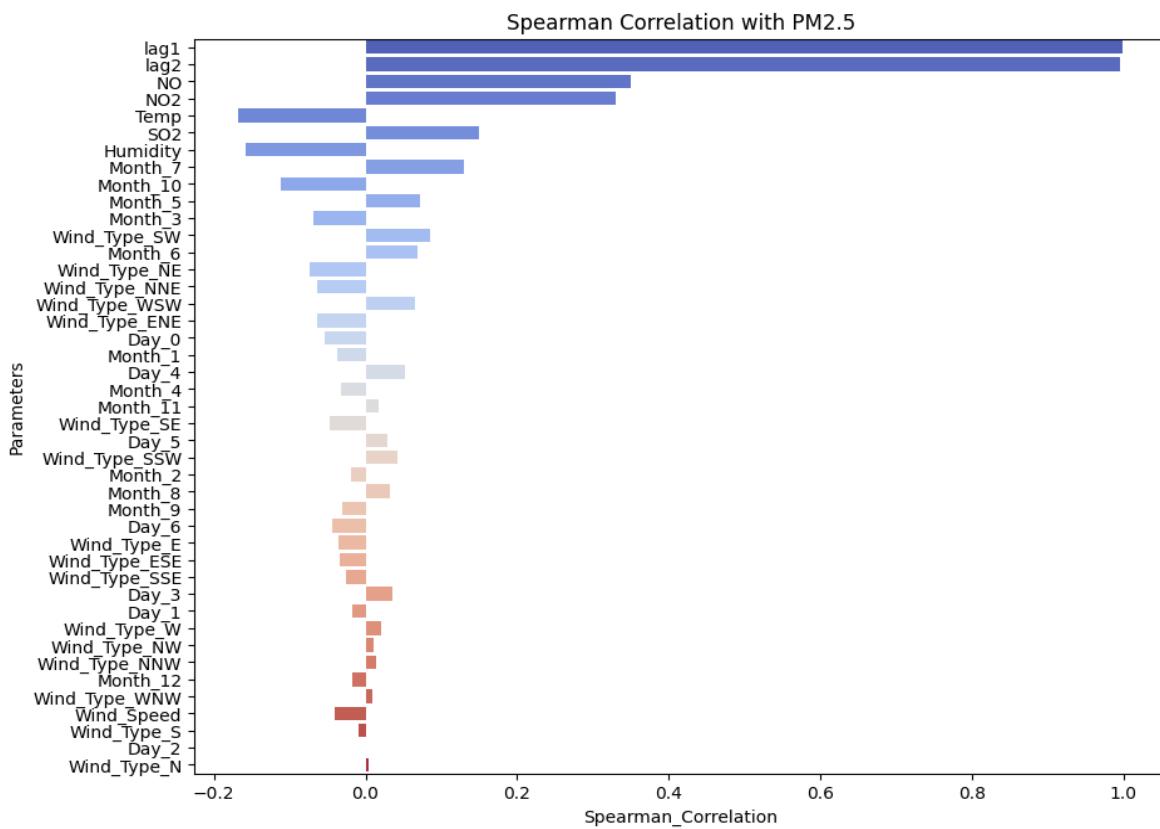


Figure 27 All factors' Spearman Correlation with PM2.5 Plot

Through Pearson and Spearman rank correlation, the closer the absolute value is to 1, the higher the correlation with PM2.5. Through the Heat Map, we can observe the correlation of each feature with PM 2.5. The closer the colour is to red, the positive correlation, and the closer the colour is to blue, the negative correlation. The table, Heat Map and table visualisation above show the reverse order sorting of Pearson and Spearman rank correlation. We found that the top 5 features for both Pearson and Spearman rank correlation are the same, namely: Log1, Log2, NO, NO2, and Temp. Therefore, we will select these 5 features. To ensure the correctness of feature selection, we will further analyse the top 8 features, Log1, Log2, NO, NO2, Temp, SO2, Humidity and Month_7.

2.2. Influence on PM Concentration

Using linear regression, we analyse the impact of these 8 features on PM2.5. Linear regression quantifies the linear influence of each predictor on PM2.5 and tests the significance of the regression coefficients.

Hypothesis tests H0 and H1 are as follows:

- H0: There is no significant linear relationship between predictors and PM2.5 (that is, the regression coefficient is zero).
- H1: The predictor has a significant linear relationship with PM2.5 (that is, the regression coefficient is not zero).

Regression Analysis for PM2.5 with Multiple Features
OLS Regression Results

Dep. Variable:	PM2.5	R-squared:	1.000			
Model:	OLS	Adj. R-squared:	1.000			
Method:	Least Squares	F-statistic:	4.420e+07			
Date:	Sat, 08 Jun 2024	Prob (F-statistic):	0.00			
Time:	06:36:52	Log-Likelihood:	2.1075e+05			
No. Observations:	44158	AIC:	-4.215e+05			
Df Residuals:	44149	BIC:	-4.214e+05			
Df Model:	8					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	0.0012	7.25e-05	16.207	0.000	0.001	0.001
lag1	1.9642	0.001	1608.196	0.000	1.962	1.967
lag2	-0.9662	0.001	-791.795	0.000	-0.969	-0.964
NO	-0.0003	7.64e-05	-3.436	0.001	-0.000	-0.000
NO2	0.0006	9.88e-05	5.785	0.000	0.000	0.001
Temp	-7.672e-05	7.24e-05	-1.060	0.289	-0.000	6.52e-05
SO2	-0.0001	4.51e-05	-2.678	0.007	-0.000	-3.24e-05
Humidity	-0.0003	5.77e-05	-5.601	0.000	-0.000	-0.000
Month_7	0.0001	3.82e-05	3.034	0.002	4.11e-05	0.000
Omnibus:	29733.047	Durbin-Watson:	1.337			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	15966049.945			
Skew:	-1.968	Prob(JB):	0.00			
Kurtosis:	96.070	Cond. No.	282.			

Figure 28 Multiple Linear Regression Summary

By observing the linear regression model results, the R-squared value indicates the proportion of total variance explained by the model. A value of 1 indicates a perfect fit.

- Lag1: Lag1 is positively correlated with PM2.5. When Lag1 increases by 1 unit, PM2.5 increases by 1.9642 units, and the p-value is 0, indicating that Lag1 has a highly significant impact on PM2.5.
- Lag2: Lag2 is negatively correlated with PM2.5. When Lag2 decreases by 1 unit, PM2.5 decreases by 0.9662 units, and the p-value is 0, indicating that Lag2 has a highly significant impact on PM2.5.
- NO: NO is negatively correlated with PM2.5. When NO concentration decreases by 1 unit, PM2.5 concentration decreases by an average of 0.0003 units, and the p-value is 0.001, indicating that NO has a significant impact on PM2.5.
- NO2: NO2 is positively correlated with PM2.5. When NO2 increases by 1 unit, PM2.5 concentration increases by an average of 0.0006 units, and the p-value is 0, indicating that NO2 has a significant impact on PM2.5.
- Temp: Temp is negatively correlated with PM2.5. When Temp decreases by 1 unit, PM2.5 concentration decreases by an average of 7.672e-05 units, and the p-value is 0.289, which is above the 0.05 threshold, indicating that Temp does not have a significant impact on PM2.5.
- SO2: SO2 is negatively correlated with PM2.5. When SO2 decreases by 1 unit, PM2.5 concentration decreases by an average of 0.0001 units, and the p-

value is 0.007, indicating that SO2 has a significant impact on PM2.5.

- Humidity: Humidity is negatively correlated with PM2.5. When Humidity decreases by 1 unit, PM2.5 concentration decreases by an average of 0.0003 units, and the p-value is 0, indicating that Humidity has a significant impact on PM2.5.
- Month_7: Month_7 is positively correlated with PM2.5. When Month_7 increases by 1 unit, PM2.5 increases by 0.0001 units, and the p-value is 0.002, indicating that Month_7 has a significant impact on PM2.5.

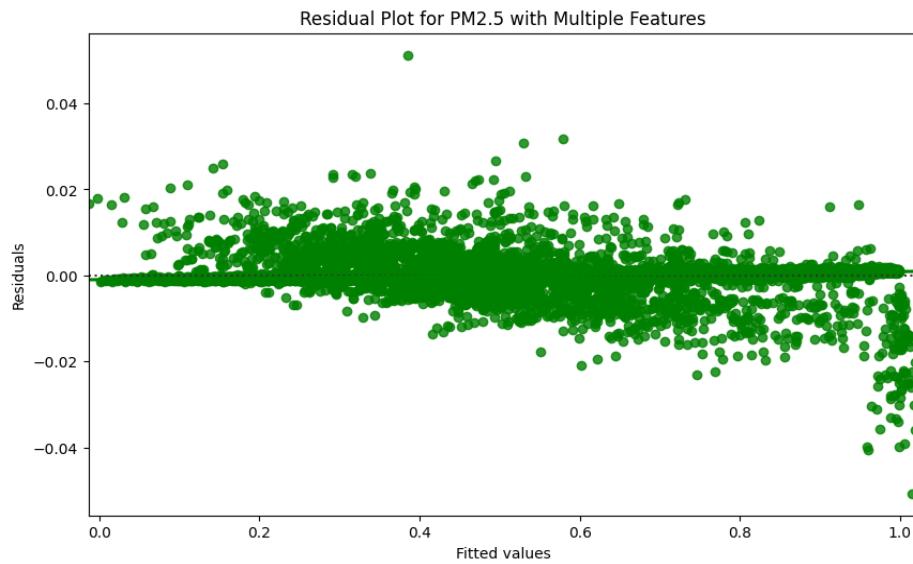


Figure 29 Residual Plot for PM2.5 with Multiple Features

The data in the residual plot are randomly distributed near the zero line with no obvious pattern, indicating that the linear regression model fits the relationship between PM2.5 and the selected features well.

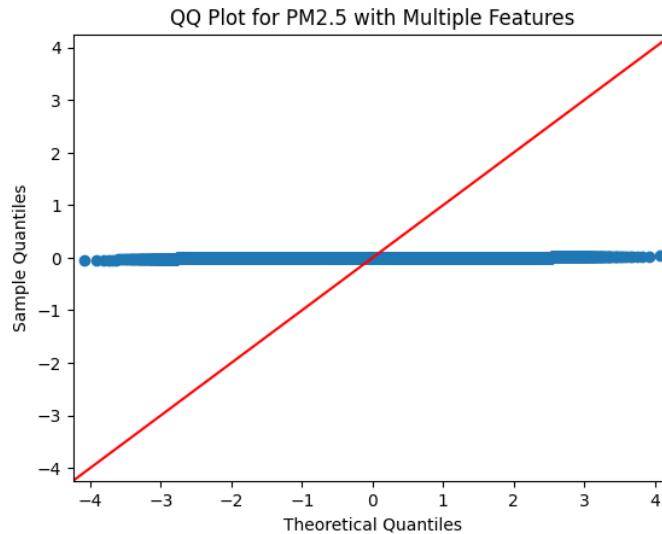


Figure 30 . QQ Plot for PM2.5 with Multiple Features

The data points in the QQ plot are not distributed along the 45-degree straight line, indicating that the residuals are not normally distributed.

In summary, Log1, Log2, NO, NO2, Temp, SO2 and Month_7's p value is less than

the significance level (0.05), then the null hypothesis is rejected and it is considered that the feature has a significant linear relationship with PM2.5. The temp's p-value is bigger than the significant level, the null hypothesis is not rejected and it is considered that PM2.5 does not have a significant linear relationship. So Log1, Log2, NO, NO2, Temp, SO2 and Month_7 are useful for prediction. And we decided to remove Temp as a feature. Combining the results from Pearson and Spearman rank correlations, **we finally decide on the following 5 variables: Lag1, Lag2, NO, NO2, and SO2.**

2.3. Provide graphical visualisation of variation of PM variation.

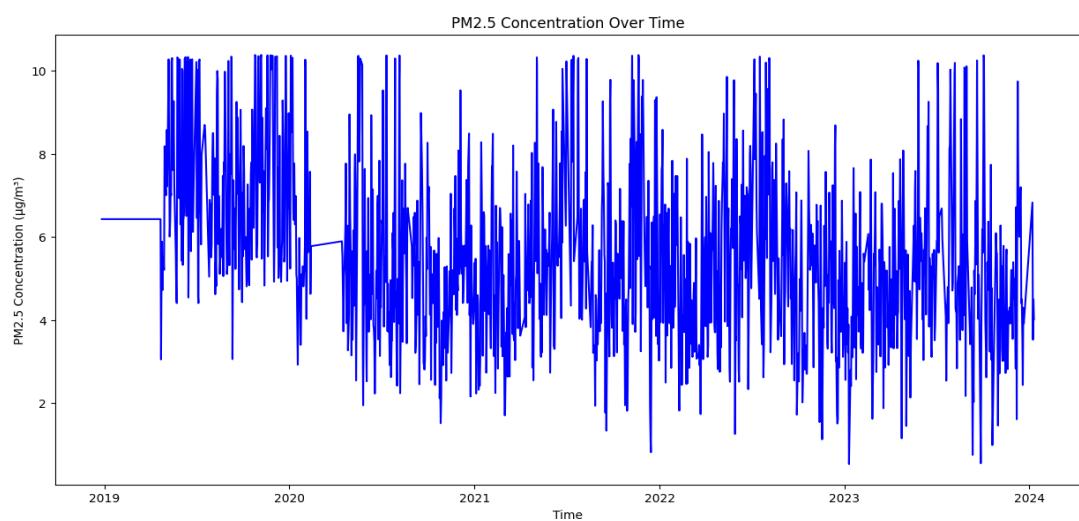


Figure 31 PM2.5 Concentration Over Time

This line plot shows the PM2.5 concentration over time from 2019 to 2024. The concentrations fluctuate significantly, with values ranging from 2 to 10 $\mu\text{g}/\text{m}^3$. Notable peaks are observed throughout the years, indicating periods of high pollution. For instance, early 2019 shows concentrations consistently above 6 $\mu\text{g}/\text{m}^3$, while fluctuations are more pronounced in later years.

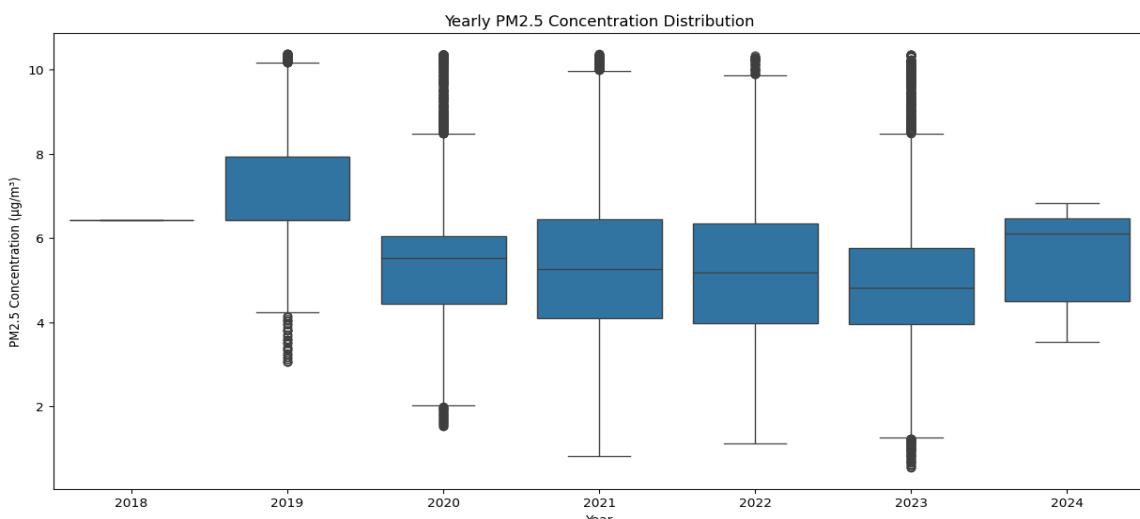


Figure 32 Yearly PM2.5 Concentration Distribution

This box plot shows yearly PM2.5 concentration distributions from 2018 to 2024. The median concentration fluctuates around 5-6 µg/m³. In 2019, the median is 6.5 µg/m³, with outliers reaching up to 10 µg/m³. In 2024, the median is approximately 6 µg/m³, with less variability.

2.4. Provide summary statistics of the PM concentration.

Table 8 Summary Statistics of PM2.5

Statistic	Value
Count	44158.000000
Mean	5.649257
Std	1.805623
Min	0.525000
25%	4.372917
50%	5.597917
75%	6.493750
Max	10.383333

This table summarizes the key statistics of PM2.5 concentration. The mean concentration is 5.65 µg/m³ with a standard deviation of 1.81. The data ranges from a minimum of 0.53 to a maximum of 10.38. The 25th, 50th (median), and 75th percentiles are 4.37, 5.60, and 6.49, respectively, indicating a moderately spread distribution.

2.5. Summary statistic of the highest correlation predictors

Table 9 Summary Statistics of Highest Correlation Predictors

Parameters	Mean	Median	Std	Min	Max	25%	75%
Lag1	5.64931	5.59792	1.80561	0.525	10.3833	4.37292	6.49375
Lag2	5.64937	5.59792	1.8056	0.525	10.3833	4.37292	6.49375
NO	9.38783	7.53125	7.58012	0.0042	27.725	3.35625	13.1042
NO2	14.5415	13.0563	8.60759	0.1063	37.425	8	19.6391
SO2	0.9936	0.89375	0.71734	0	2.7958	0.49375	1.36875

This table provides summary statistics for key predictors of PM2.5 concentration.

Lag1 and Lag2 show similar distributions with means around 5.65 $\mu\text{g}/\text{m}^3$ and standard deviations of approximately 1.81. NO has a mean of 9.39 $\mu\text{g}/\text{m}^3$, with a wide range (0.0042 to 27.73). NO2 has the highest mean at 14.54 $\mu\text{g}/\text{m}^3$, and SO2 shows a lower mean of 0.99 $\mu\text{g}/\text{m}^3$.

2.6. Normalization Data

Normalized data, the maximum and minimum values are easily affected by outliers. In order to ensure the training effect of the model MLP (multi-layer perceptron) and LSTM (long short-term memory), we use MinMaxScaler to normalize the data and convert it to [0,1] to ensure uniform scaling across all features. Using normalization after selecting features can reduce the processing of unused features and increase work efficiency.

Table 10 Summary Statistics for Normalized Data

Statistic	lag1	lag2	NO	NO2	Temp	PM2.5
count	43800.000	43800.000	43800.000	43800.000	43800.000	43800.000
mean	0.508	0.508	0.339	0.383	0.541	0.508
std	0.189	0.189	0.273	0.232	0.182	0.189
min	0.000	0.000	0.000	0.000	0.000	0.000
25%	0.374	0.374	0.120	0.207	0.396	0.374
50%	0.494	0.494	0.274	0.340	0.529	0.494
75%	0.610	0.610	0.473	0.519	0.692	0.609
max	1.000	1.000	1.000	1.000	1.000	1.000

3. Experimental Methods

3.1. Split Dataset

Training Set: The training set is used to train the machine learning models, allowing them to learn patterns and relationships within the data.

Testing Set: The testing set is used to evaluate the performance of the trained models, ensuring their ability to generalize to unseen data.

Out of a total of 44,158 records, 70% of the data will be used for training, and 30% for testing.

- Training Set Size (70%): $44,158 \times 0.70 = 30,910$
- Testing Set Size (30%): $44,158 \times 0.30 = 13,248$

```
X_train = (30910, 5)
X_test = (13248, 5)
y_train = (30910,)
y_test = (13248,)
```

Figure 33 Dataset Splitting for Training and Testing

3.2. Workflow

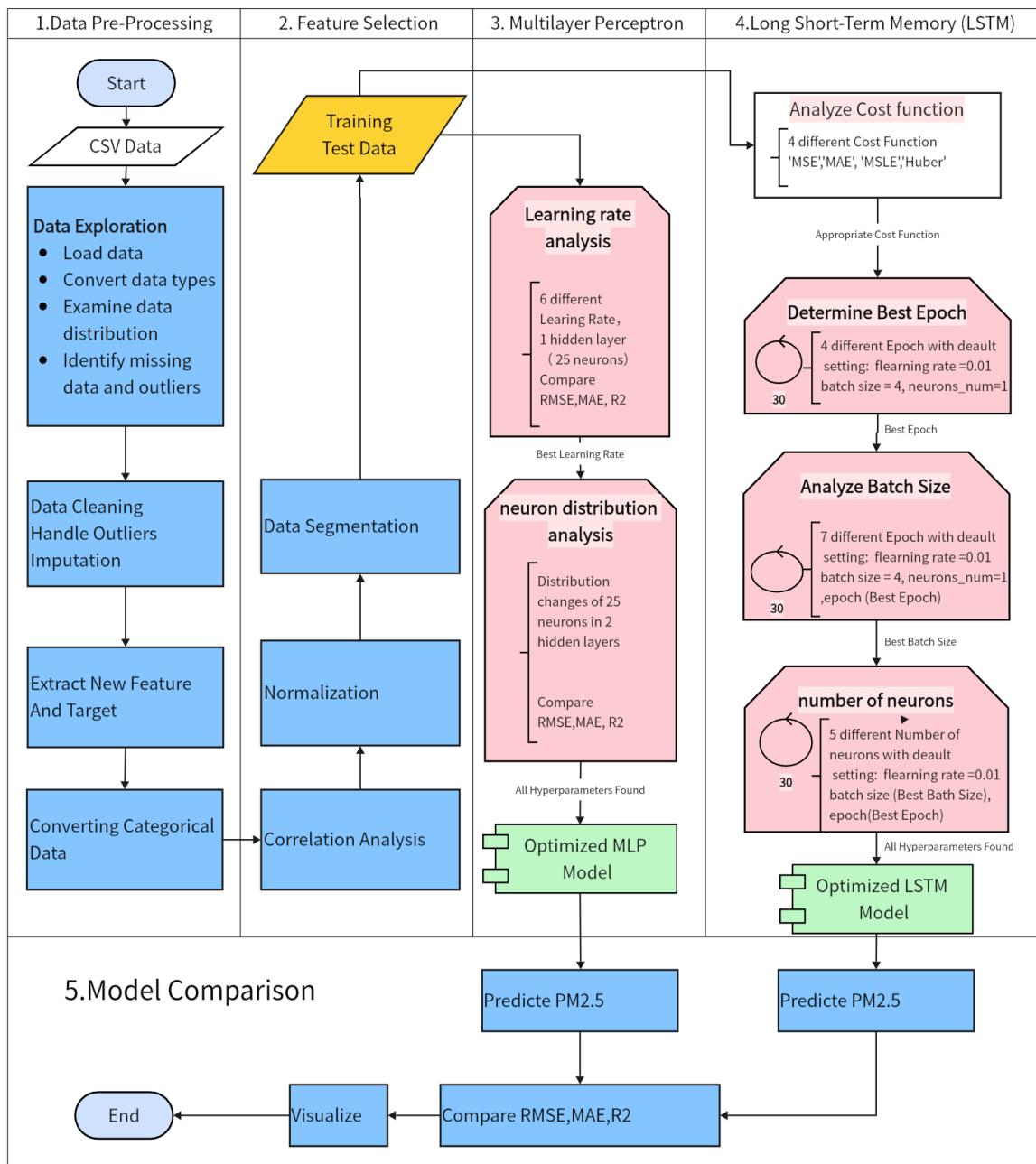


Figure 34 Workflow Overview

Workflow Process:

1. **Data Preprocessing:** Ensure that all data has the same temporal resolution, identify and handle missing data and outliers, and normalize the data as needed.
2. **Feature Selection:** Identify the five attributes with the highest correlation to PM2.5 concentration using methods such as Pearson correlation to ensure that the most relevant features are used in the predictive model.
3. **Model Training and Evaluation:**
 - Use the splitted Training and Test dataset

- MLP Model: Determine the optimal learning rate and evaluate the impact of varying the number of neurons across layers to select the best architecture.
 - LSTM Model: Train the LSTM model using the Adaptive Moment Estimation (ADAM) optimizer. Employ a two-step Rolling Window to generate and evaluate 3D data sequences. Determine the optimal number of epochs, batch size, and neurons in the hidden layers to find the best settings.
4. Performance Metrics: Evaluate models using performance metrics such as root mean square error (RMSE), mean absolute error (MAE), and correlation coefficient (R^2). Perform a visual and statistical comparison of actual and predicted PM2.5 values to determine the better-performing model between MLP and LSTM.

The project aims to find robust and accurate prediction models for PM2.5 concentrations by taking an integrated approach from data preprocessing to detailed model evaluation. The insights gained can inform public health strategies and improve air quality management, contributing to improved health outcomes and environmental policies.

4. Multilayer Perceptron (MLP)

4.1. MLP Description

A Multilayer Perceptron (MLP) is a fundamental feedforward neural network model consisting of one or more hidden layers (intermediate layers) and an output layer, with neurons fully connected between layers. MLP can be used to solve classification and regression problems, and its hidden layers enable the model to learn nonlinear patterns and relationships in the data.

Advantages: powerful modelling capability, versatility, scalability, automatic feature extraction.

Disadvantages: long training time, high computational resource requirement, prone to overfitting, complex hyper-parameter tuning

Application Scenarios: Image Recognition, Natural Language Processing, Financial Prediction, Medical Diagnosis

MLP mainly consists of three components:

- Input Layer: The layer that receives input data, where each input feature corresponds to a node in the input layer. The number of nodes equals the dimension of the input features.
- Hidden Layers: One or more layers located between the input layer and the output layer. Each hidden layer consists of multiple neurons, fully connected to neurons in the adjacent layers. The role of the hidden layers is to perform nonlinear transformations and feature extraction on the input features, thereby learning complex patterns and relationships in the data.

- Output Layer: The final layer that outputs the model's prediction results. For classification problems, the output layer usually uses a softmax activation function to convert the outputs into a probability distribution over classes. For regression problems, the output layer typically does not use an activation function or uses a linear activation function.

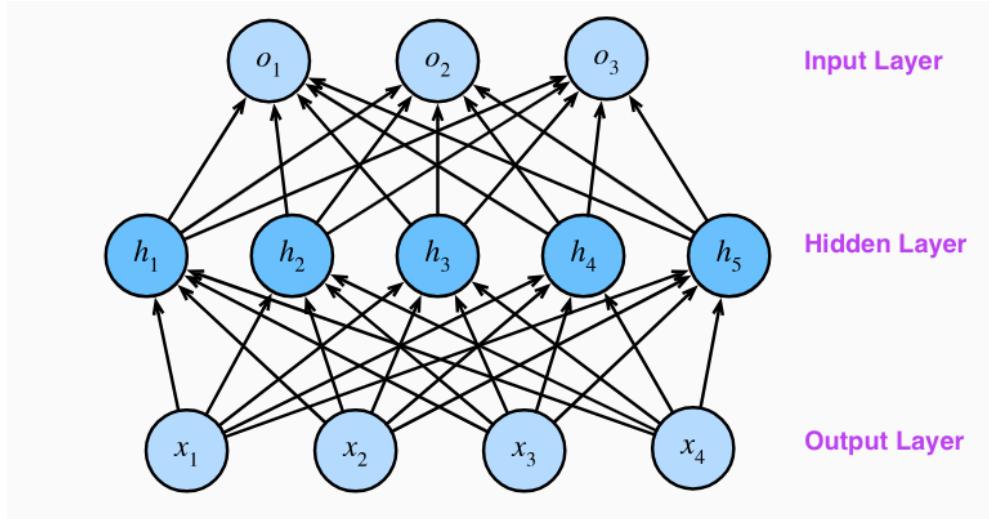


Figure 35 MLP Overview

This multi-layer perceptron has 4 inputs, 3 outputs, and its hidden layer contains 5 hidden units. The input layer does not involve any computation, so using this network to produce output only requires implementing the computations of the hidden and output layers. Therefore, the number of layers in this multilayer perceptron is 2. Both layers are fully connected. Every input affects every neuron in the hidden layer, and every neuron in the hidden layer affects every neuron in the output layer.

4.2. Learning Rate Analysis

In this section, we utilize the MLP regressor from the `sklearn.neural_network` module to create a baseline model for predicting PM2.5 levels.

We use a single hidden layer with 25 neurons and the default values for all other parameters. To identify the optimal learning rate that maximizes model performance on the test dataset, we test multiple values: [0.0001, 0.001, 0.01, 0.1, 0.5, 1.0]. These values cover a broad range of learning rates from very small to relatively large. For each learning rate, we train the MLP regressor model on the training data and evaluate its performance using Test data. The evaluation metrics used are:

- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- Mean Absolute Error (MAE)
- R² (Coefficient of Determination)

The best learning rate is the one that results in the lowest MSE.

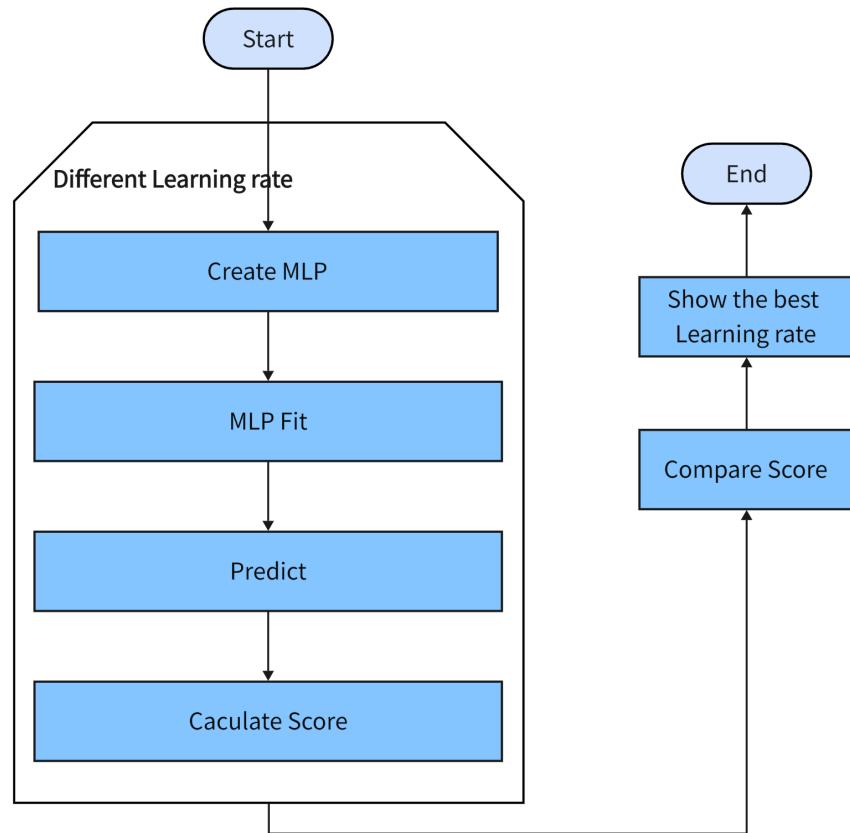


Figure 36 MLP Learning Rate Analysis Workflow

Experimental Result

After testing different learning rates, the best learning rate was identified based on the lowest MSE. The following table summarizes the results of the learning rate experimentation:

Table 11 MLP Performance with Different Learning Rates

Learning Rate	MSE	RMSE	MAE	R ²
0.0001	0.024147	0.155386	0.120110	0.319865
0.001	0.023626	0.153702	0.118947	0.334522
0.01	0.023467	0.153183	0.118599	0.338998
0.1	0.023963	0.154794	0.120160	0.325026
0.5	0.026101	0.161464	0.124994	0.265015
1.0	0.029039	0.169608	0.131417	0.183542



Figure 37 MLP Learning Rate Comparison

Analysis

The best learning rate is determined to be 0.01 based on the lowest MSE. The corresponding performance metrics for this learning rate are as follows:

- MSE: 0.023467
- RMSE: 0.153183
- MAE: 0.118599
- R²: 0.338998

Among the tested learning rates, 0.01 achieved the lowest MSE and RMSE, indicating the best performance. Higher learning rates like 0.1, 0.5, and 1.0 resulted in significantly higher error metrics, indicating poorer model performance. Lower learning rates like 0.0001 and 0.001 also did not perform as well as 0.01, confirming that 0.01 is the optimal learning rate for this MLP model.

Conclusion

The best learning rate is determined to be 0.01 with a single hidden layer with 25 neurons

4.3. Neurons Distribution Analysis

In this section, we experiment with a two-layer MLP model to determine the optimal distribution of neurons across the two hidden layers that give the highest accuracy with a learning rate of 0.01. Unlike the previous single-layer model, we now distribute the total 25 neurons between the two layers in various configurations. The neurons are transferred from the first hidden layer to the second iteratively, in steps of 1 neuron. For instance, in the first iteration, the first hidden layer will have 24 neurons and the second will have 1; in the second iteration, the first layer will have 23 neurons and the second will have 2, and so on.

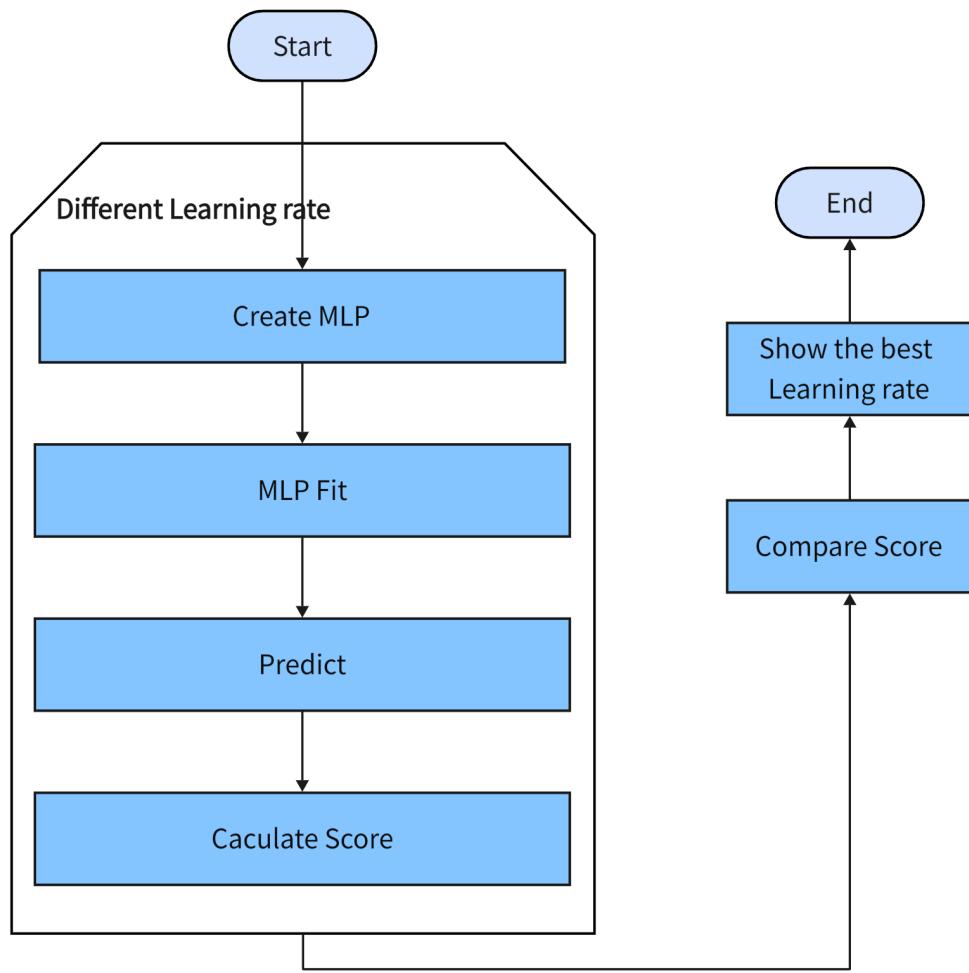


Figure 38 MLP Neurons Distribution Analysis Workflow

Experimental Result

After testing different configurations, the best configuration of neurons across the two layers was identified based on the lowest MSE. The following table summarizes the results of the neuron distribution experimentation:

Table 12 MLP Performance with Neuron Distribution

Neurons in Layer 1	Neurons in Layer 2	MSE	RMSE	MAE	R ²
24	1	0.035875	0.189408	0.147522	-0.000000
23	2	0.017146	0.130943	0.098805	0.522066
22	3	0.017147	0.130946	0.098988	0.522044
21	4	0.017245	0.131320	0.098821	0.519306
20	5	0.017290	0.131493	0.099792	0.518044
19	6	0.017092	0.130738	0.098767	0.523504
18	7	0.017128	0.130875	0.098778	0.522565
17	8	0.017296	0.131512	0.100239	0.517899
16	9	0.017098	0.130761	0.098864	0.523395
15	10	0.016687	0.129180	0.097766	0.534848
14	11	0.017315	0.131586	0.099289	0.517356
13	12	0.016969	0.130265	0.098204	0.527005
12	13	0.016927	0.130102	0.098077	0.528183
11	14	0.017272	0.131423	0.099161	0.518570
10	15	0.016905	0.130018	0.098705	0.528790
9	16	0.017319	0.131601	0.100276	0.517253
8	17	0.017308	0.131561	0.099204	0.517541
7	18	0.017035	0.130517	0.099602	0.524880
6	19	0.017045	0.130557	0.098956	0.524294
5	20	0.017066	0.130637	0.099230	0.524295
4	21	0.017489	0.132244	0.100297	0.511995
3	22	0.017221	0.131228	0.098931	0.519982
2	23	0.017923	0.133875	0.101590	0.500419
1	24	0.035910	0.189499	0.147217	-0.000967

Analysis

The best configuration with a learning rate of 0.01 is Layer 1 neurons: 15 and Layer 2 neurons: 10.

The plot below shows the error metrics (MSE, RMSE, MAE) as the number of neurons in the first layer decreases and the number of neurons in the second layer increases.

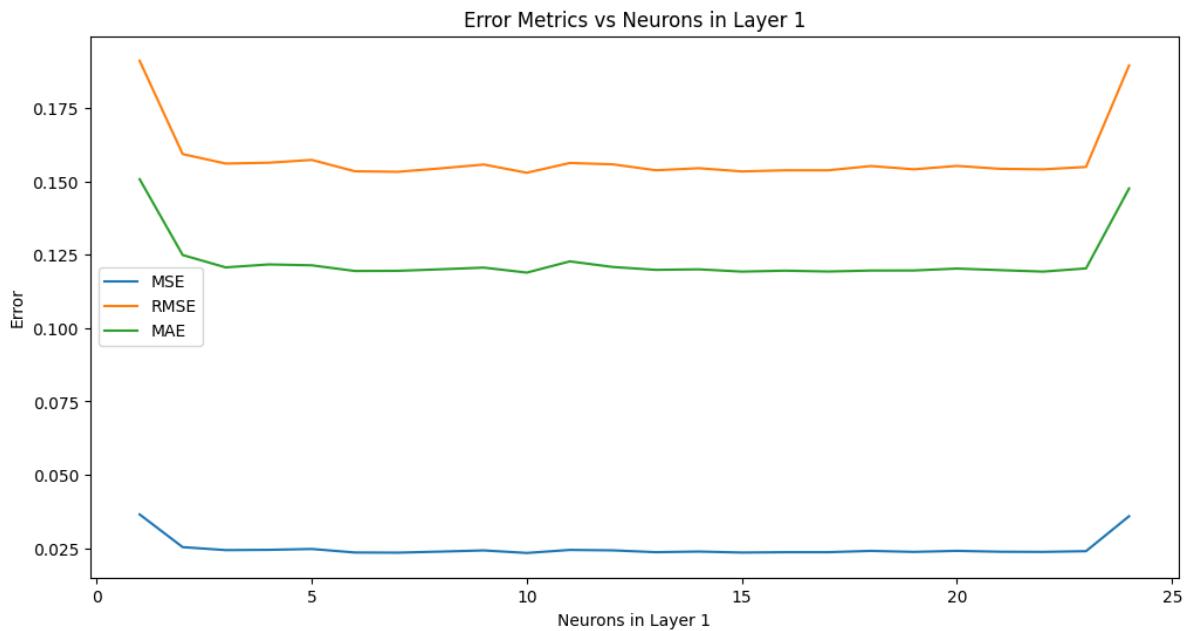


Figure 39 MLP Neurons Distribution Comparison

Conclusion

The best configuration of neurons is found to be 15 neurons in the first hidden layer and 10 neurons in the second hidden layer when learning rate is 0.01. This configuration achieved the lowest MSE and RMSE, indicating the best performance. The error metrics (MSE, RMSE, MAE) as the number of neurons in the first layer decreases and the number of neurons in the second layer increases are shown in the plot.

4.4. MLP Analysis

4.4.1. Best Architecture

We conducted a series of experiments with a single hidden layer and two hidden layers, distributing 25 neurons to find the optimal neuron distribution structure. The single hidden layer experiment serves as a baseline for comparison. In searching for the best neuron distribution structure, we used the following performance metrics: Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and the coefficient of determination (R^2).

- **Single-layer with 25 neurons (baseline):**
 - MSE: 0.023467
 - RMSE: 0.153183
 - MAE: 0.118599
 - R^2 : 0.338998
 - Learning Rate: 0.01

Configurations with the first layer having 24 neurons and the second layer having 1 neuron, and the first layer having 1 neuron and the second layer having 24 neurons, both resulted in MSE values higher than the baseline and were not considered further. Among the others with lower MSE than the baseline, the one with the first layer having 15 neurons and the second layer having 10 neurons achieved the lowest MSE (0.016687).

- **Two-layer structure:** 15 neurons in the first layer, 10 neurons in the second layer
 - MSE: 0.016687
 - RMSE: 0.129208
 - MAE: 0.100523
 - R²: 0.495781
 - Learning Rate: 0.01

Conclusion

Through our experiments, we determined that the optimal two-layer structure consists of 15 neurons in the first layer and 10 neurons in the second layer, with a learning rate of 0.01.

4.4.2. Possible Reasons for Variation

- Complexity of Features: Some features may benefit more from deeper representations which two layers can provide.
- Learning Dynamics: Different splits may affect the learning dynamics, allowing the network to capture different aspects of the data.
- Regularisation Effect: Spreading neurons across layers can act as a form of regularisation, potentially reducing overfitting.

5. Long Short-Term Memory (LSTM)

5.1. LSTM Introduction

5.1.1. LSTM

Long Short-Term Memory (LSTM) networks are a type of recurrent neural network (RNN) designed to learn long-term dependencies. They are particularly effective at handling sequential data and overcoming the vanishing gradient problem.

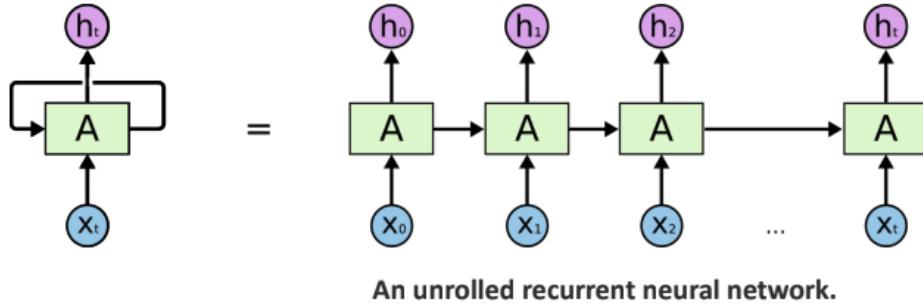


Figure 40 An unrolled recurrent neural network

5.1.2. LSTM Architecture

LSTM consists of several components: the state function (cell state and hidden state) and three gates (input gate, forget gate, and output gate).

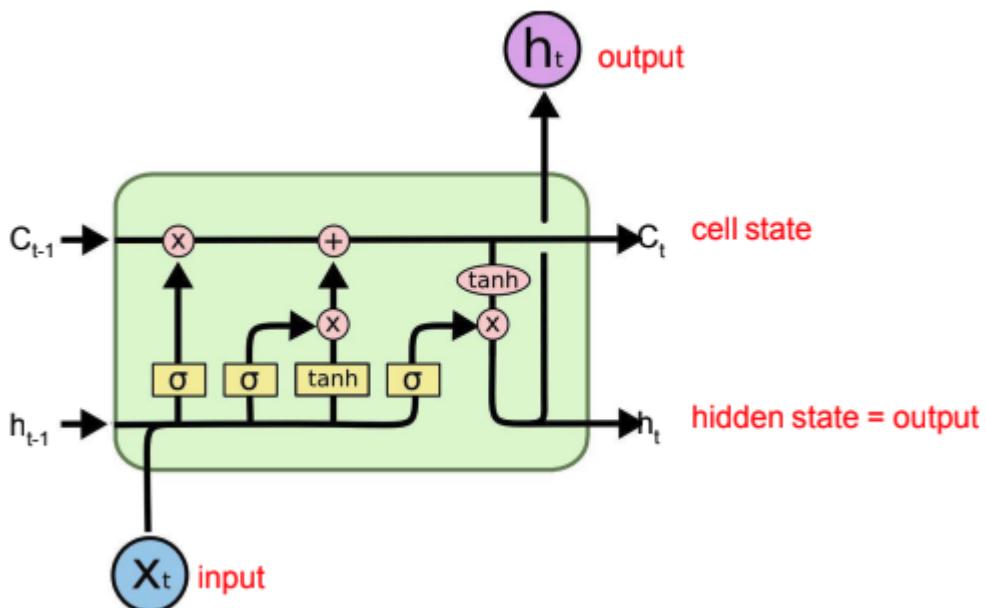


Figure 41 LSTM Architecture

State Function: Memory cell stores long-term information across time steps. It manages and updates this information through forget, input, and output gates, capturing long and short-term dependencies effectively.

- Cell State: This is the long-term memory of the LSTM cell. It runs through the entire sequence with minor linear interactions, making it easier to remember information across long sequences.
- Hidden State: This is the short-term memory that gets passed to the next cell. It's a filtered version of the cell state.
- Updating Cell State: The new cell state is calculated by combining the old cell state (modulated by the forget gate) and the candidate cell state (modulated by the input gate).

Gates: Gates are mechanisms that regulate the flow of information into and out of the memory cell. There are three primary gates in an LSTM:

- Forget Gate: This gate decides what information to discard from the cell state. It takes the previous hidden state and the current input and passes them through a sigmoid activation function.
- Input Gate: This gate decides what new information to store in the cell state. It consists of two parts:
 - A sigmoid layer that decides which values to update.
 - A tanh layer that creates a vector of new candidate values that could be added to the cell state.
- Output Gate: This gate decides the output of the LSTM cell. It first passes the previous hidden state and the current input through a sigmoid function. Then, it multiplies the output by the tanh of the updated cell state to decide the final output.

5.1.3. How LSTM Differs from MLP

Sequence Handling:

- LSTM: Designed to handle sequential data by maintaining long-term dependencies through its cell state. It processes data one step at a time, maintaining a memory of previous inputs.
- MLP: A feedforward neural network that processes all inputs simultaneously without any inherent memory of previous inputs.

Architecture:

- LSTM: Includes gates (forget, input, output) and cell state to manage memory and control the flow of information.
- MLP: Consists of multiple fully connected layers with neurons that transform inputs into outputs through weights and biases.

Training Challenges:

- LSTM: Designed to overcome the vanishing gradient problem commonly faced by traditional RNNs.
- MLP: While not specifically designed for sequence data, it can suffer from vanishing gradients in deep architectures.

5.1.4. Impact of Number of Neurons and Batch Size

Impact of Number of Neurons

- Increased Neurons: Adding more neurons can increase the model's capacity to learn complex patterns, potentially improving performance. However, it can also lead to overfitting, especially with limited data.
- Decreased Neurons: Fewer neurons reduce the model's complexity, which can be beneficial for preventing overfitting but may lead to underfitting if the model is too simple to capture the data's patterns.

Impact of Batch Size

- Large Batch Size: This leads to more stable gradient estimates and can speed up training due to parallel processing advantages. However, it may result in poorer generalisation.
- Small Batch Size: Provides noisier gradient estimates, which can lead to better generalisation but slower training times. It also increases the chance of finding better local minima in the loss landscape.

Conclusion

Choosing the right number of neurons and batch size is crucial for the performance of both LSTMs and MLPs. Too many neurons can cause overfitting, while too few may lead to underfitting. The batch size affects the stability and speed of training, with larger batches providing more stable updates and smaller batches potentially leading to faster convergence but requiring careful tuning.

5.2. Epoch Optimal Size

5.2.1. Identify an appropriate cost function

The cost function, also known as the loss function, is crucial for measuring the model's performance during training. It quantifies the difference between the predicted outputs and the actual values. For time series forecasting and regression tasks, we will compare the following cost functions:

- MSE (Mean Squared Error): If outliers are present, MSE gives higher weight to outliers, sacrificing the predictive accuracy for normal data points, which can affect the overall performance of the model.
- MAE (Mean Absolute Error): Less sensitive to outliers, treating all errors equally, making it more robust. More inclusive of variations without disproportionately penalising outliers.
- Huber Loss: Combines the best of both MSE and MAE, offering good resistance to outliers while maintaining a balanced weight distribution. The rate of decrease in Huber loss lies between that of MSE and MAE, addressing the slow decrease issue of MAE.

Conclusion

Given that our dataset, despite outlier handling, still contains some unavoidable outliers (as observed in the box plot), Comparing different cost functions, Huber loss stands out by balancing the advantages of MSE and MAE. It considers outliers while maintaining a balanced weight distribution. Therefore, the most appropriate loss function for our data is Huber loss.

5.2.2. Cost Function Scores

We conducted 30 training runs, keeping the learning rate at 0.01, the batch size at 4 and the number of hidden layer neurons at 1. Train the model with different epoch settings: 20, 60, 100, 200, and 500. Record the training and test Huber loss for each epoch to evaluate the model's performance over time. The model was trained using

the Huber loss as the cost function. For each run, we randomly selected continuous sub-data segments from the training and test sets to ensure variability and robustness in the evaluation.

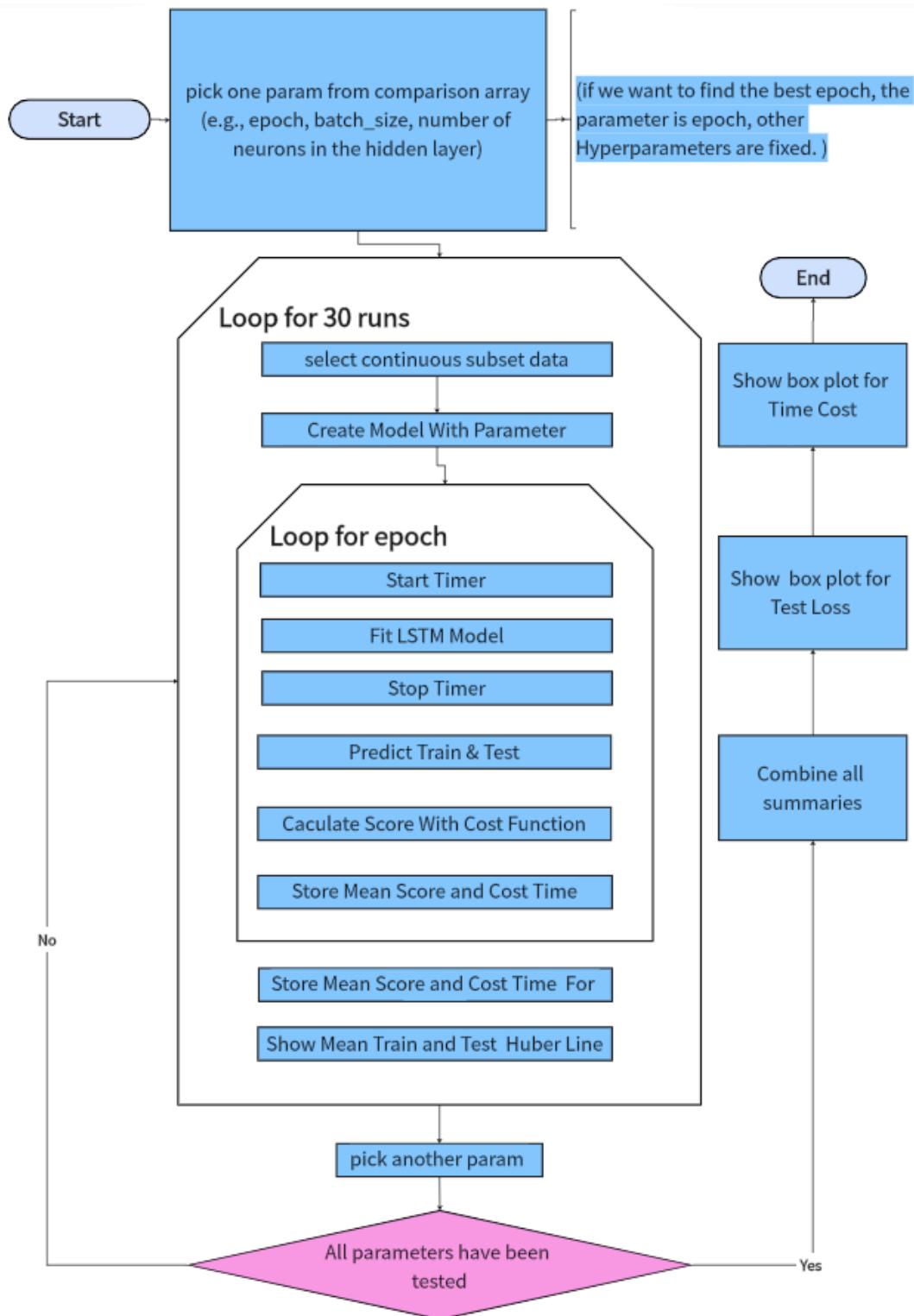


Figure 42 LSTM Tuning Workflow

The following line plots illustrate the train and test Huber loss values over different epochs for analysis of overfitting.

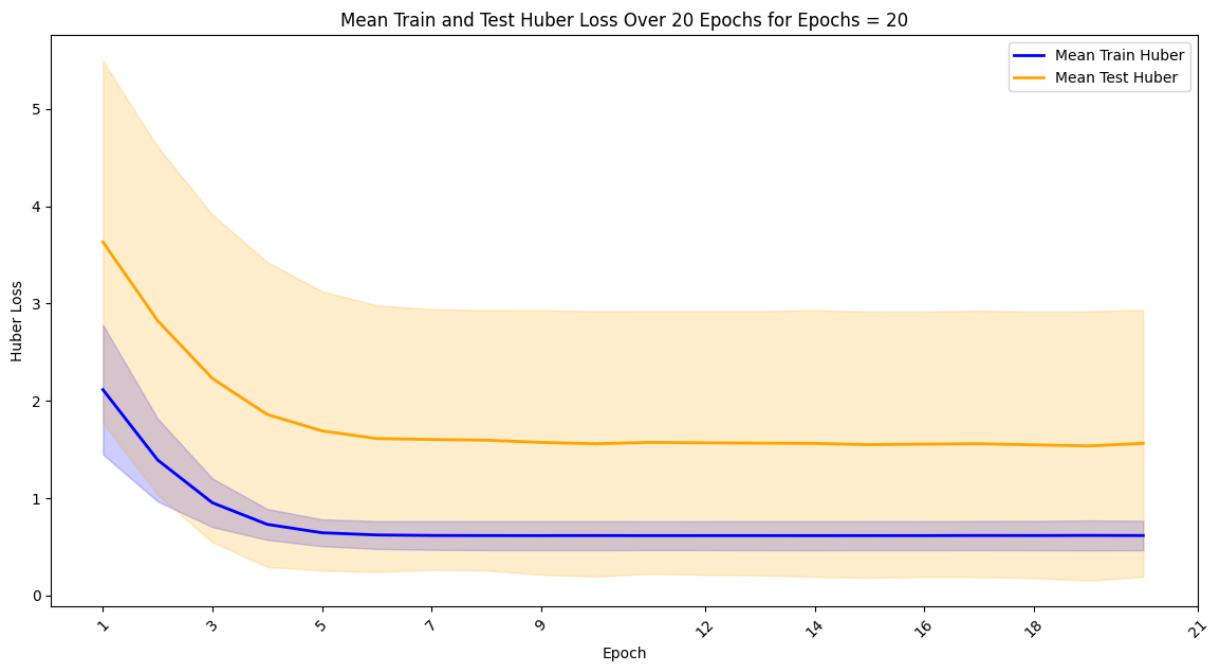


Figure 43 LSTM Mean Train and Test MAE Loss Over 20 Epochs for Epochs = 20

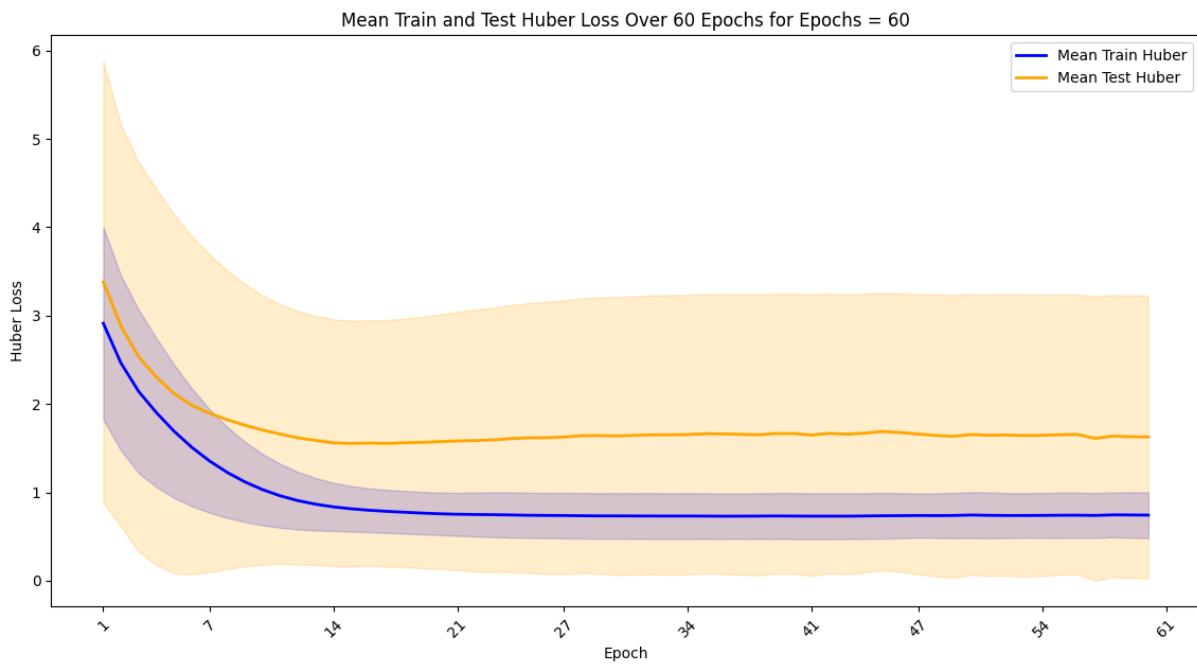


Figure 44 LSTM Mean Train and Test MAE Loss Over 20 Epochs for Epochs = 60

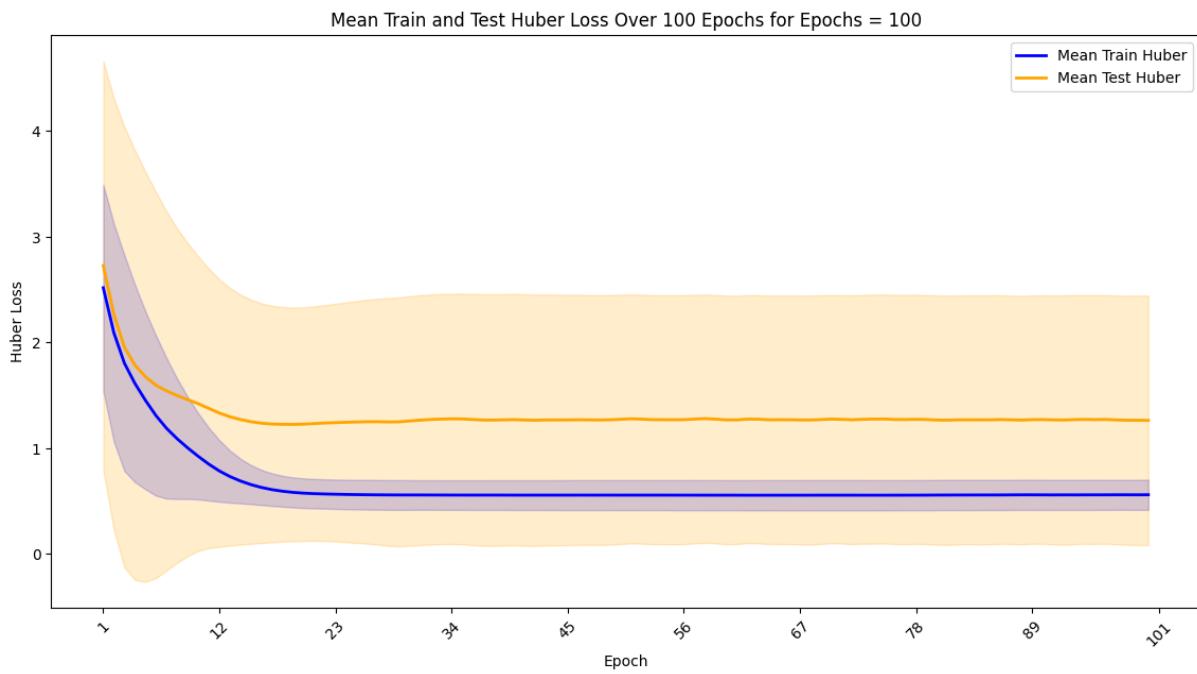


Figure 45 LSTM Mean Train and Test MAE Loss Over 20 Epochs for Epochs = 100

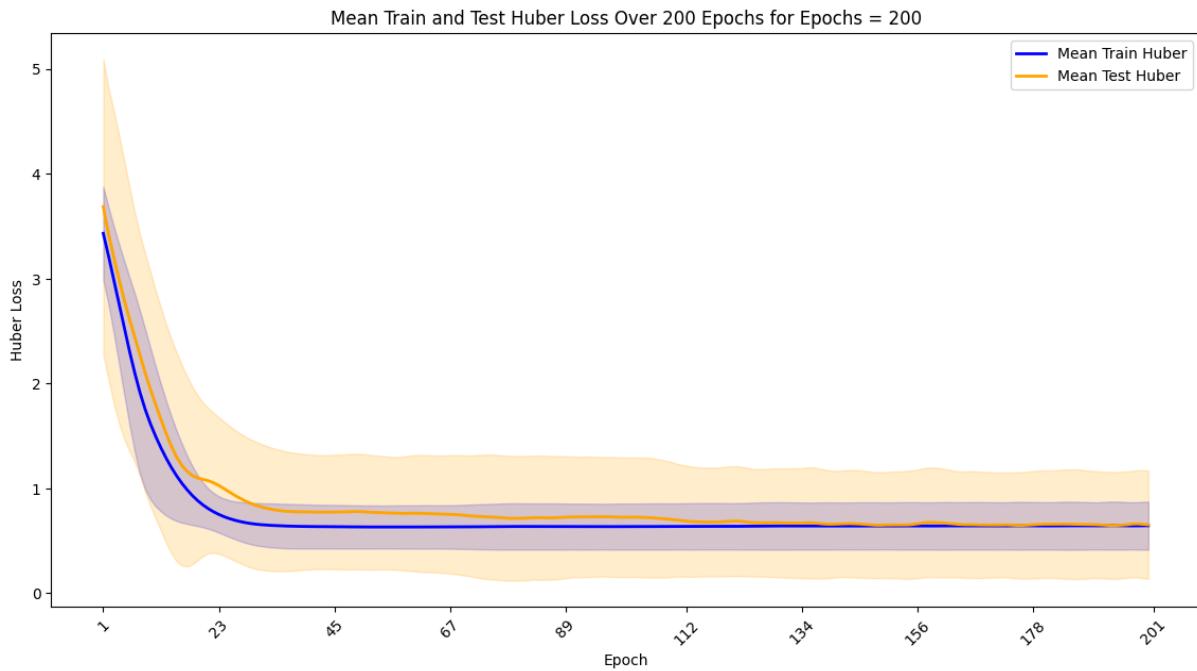


Figure 46 LSTM Mean Train and Test MAE Loss Over 20 Epochs for Epochs = 200

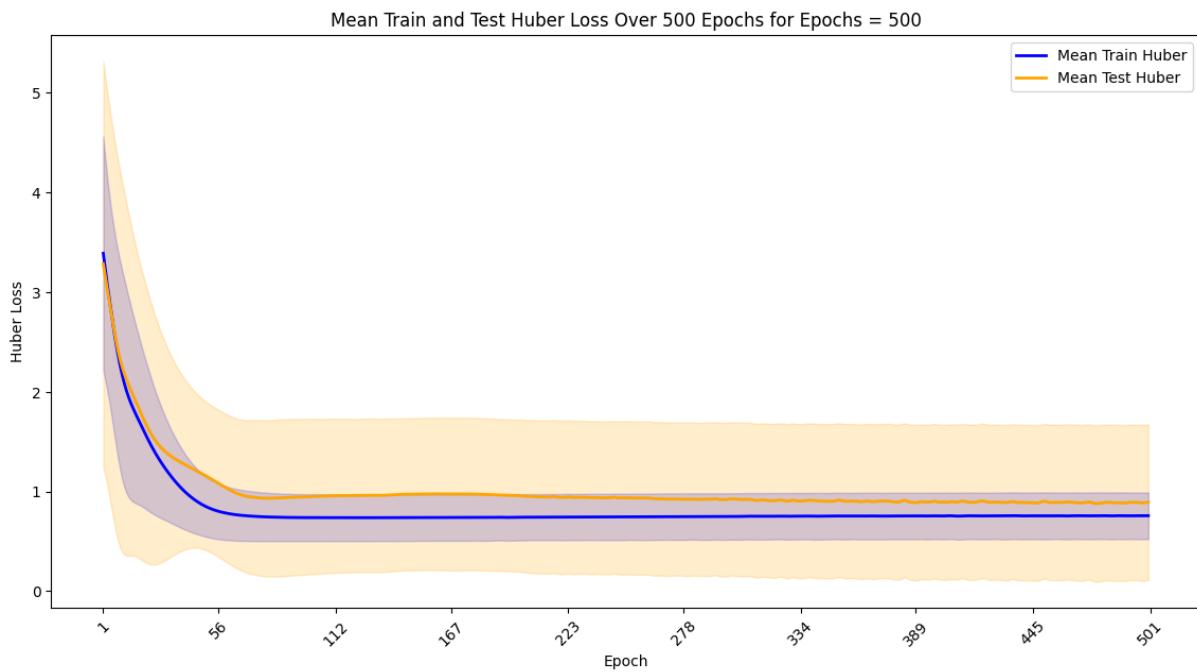


Figure 47 LSTM Mean Train and Test MAE Loss Over 20 Epochs for Epochs = 500

Both the train and test losses show a general decreasing trend initially, indicating that the model is learning and improving. However, after a certain point, fluctuations suggest overfitting or instability. Specific analysis (Mean Train and Test Huber Loss) is as follows:

- Epochs = 20: The Huber loss decreases rapidly within the first few epochs and then stabilises. However, the test loss shows higher variability, indicating potential underfitting.
- Epochs = 60: Similar to the 20 epochs configuration, but with more training iterations, the Huber loss shows a more consistent downward trend. The test loss still exhibits some fluctuations.
- Epochs = 100: With 100 epochs, the training loss continues to decrease, and the test loss shows a more stable pattern. This configuration seems to balance between underfitting and overfitting.
- Epochs = 200: The 200 epochs configuration achieves the lowest and most stable test loss, indicating a good fit for the data without significant overfitting.
- Epochs = 500: Although the training loss continues to decrease, the test loss shows more variability and outliers, suggesting that the model may be overfitting the training data.

5.2.3 Report the Summary Statistics

For each epoch configuration, the summary statistics are calculated as follows.

Table 13 Summary Statistics for Each Epoch Configuration

Epochs	Test Huber Mean	Test Huber Std Dev	Test Huber Min	Test Huber Max	Time Mean (s)
20	1.562835	1.372978	0.090328	4.629751	3.689531
60	1.627785	1.600514	0.259126	4.795913	6.054313
100	1.263646	1.178932	0.016680	3.323727	9.385436
200	0.654080	0.514997	0.026291	1.742128	16.444361
500	0.892497	0.779357	0.081133	2.242917	36.588192

- Test Huber Mean: The mean test Huber loss decreases consistently with an increase in epochs, indicating improved model performance. At 20 epochs, the mean test Huber loss is 1.562835, while at 500 epochs, it reduces significantly to 0.892497. This trend demonstrates that more training epochs contribute to better fitting of the model to the data.
- Test Huber Std Dev: The standard deviation of test Huber loss also decreases with more epochs, from 1.372978 at 20 epochs to 0.779357 at 500 epochs. This reduction suggests that the model's performance becomes more consistent as it trains longer.
- Test Huber Min and Max: The minimum and maximum test Huber loss values generally decrease as epochs increase, highlighting fewer outliers and better overall performance. For example, the maximum test Huber loss drops from 4.629751 at 20 epochs to 2.242917 at 500 epochs.
- Time Mean (s): However, the mean training time per epoch increases significantly with more epochs, from 3.689531 seconds at 20 epochs to 36.588192 seconds at 500 epochs, indicating a substantial increase in computational cost.

Conclusion

Increasing the number of epochs improves model performance but also increases computational time. A balance must be found, with 200 epochs offering a good trade-off, achieving a mean test Huber loss of 0.654080 with a mean training time of 16.444361 seconds.

5.2.4 Best Epoch Size

These box plots provide a clear visualisation of the Test Huber Loss distribution and time distribution for different epochs. The box plot helps in understanding the spread and skewness of the loss values, indicating the median, quartiles, and potential outliers. It is particularly useful for comparing the central tendency and variability of loss values for each epoch configuration, so we use it to help us analyse.

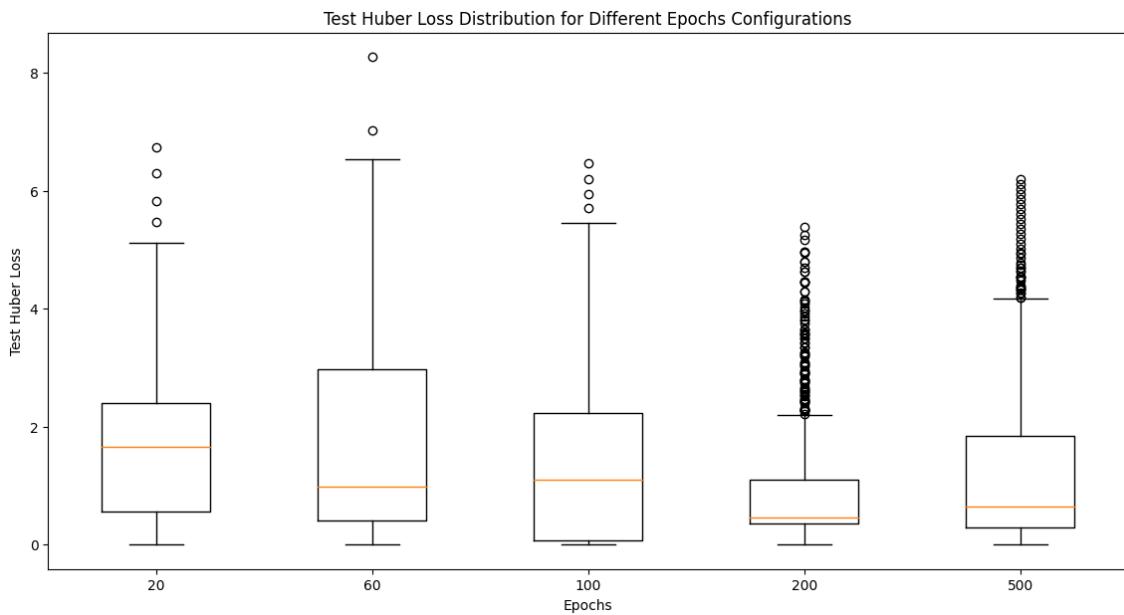


Figure 48 LSTM Test Huber Loss Distribution for Different Epochs Configurations

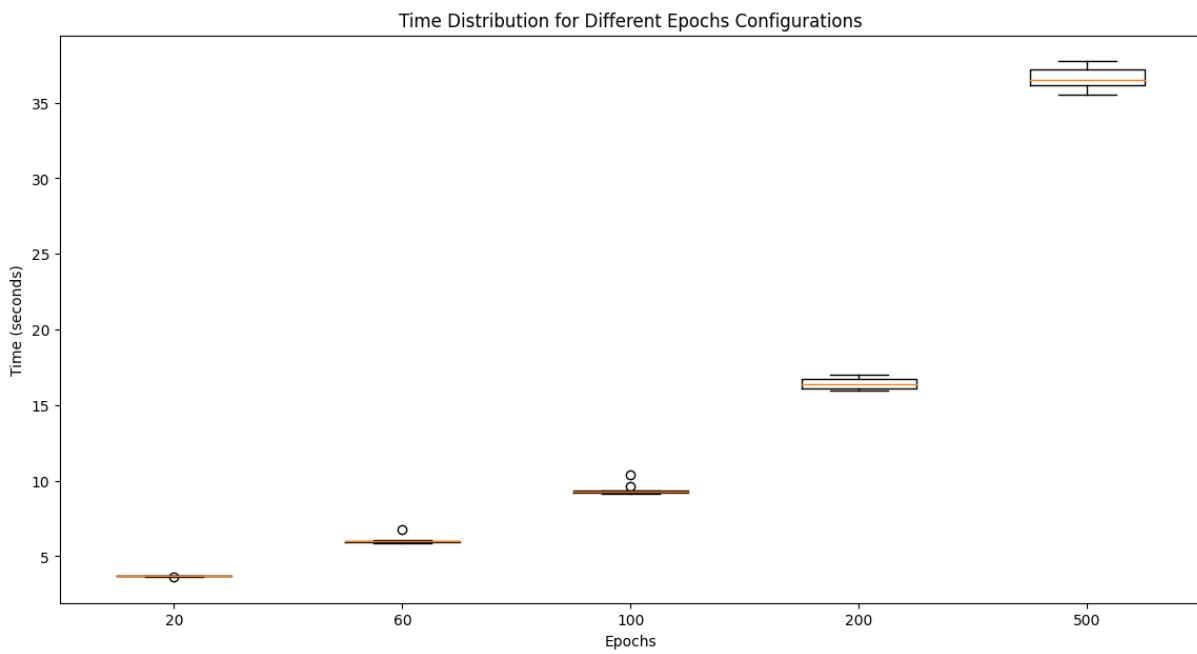


Figure 49 LSTM Time Distribution for Different Epochs Configurations

- Huber Loss Distribution: From the box plot, we can see that the 200 epochs configuration has the most concentrated Huber loss distribution with the lowest median and relatively small standard deviation, indicating stable performance. The 500 epochs configuration, although having a low median, shows a large number of outliers, suggesting potential overfitting issues.
- Performance Comparison: Comparing the model performance across different epoch configurations, we find that the 200 epochs configuration performs the best on the test set. Although the 500 epochs configuration also shows good performance on the training set, it does not perform as well on the test set and shows signs of overfitting.

- Training Time Distribution: The plot shows that the training time increases linearly with the number of epochs. The 500 epochs configuration takes the longest time to train, while the 20 epochs configuration is the fastest.

Conclusion

Considering the training and test Huber loss, training time, and performance stability, the 200 epochs configuration is determined to be the best choice. This configuration yields the best performance on the test set with moderate training time and avoids overfitting issues.

5.3. Batch Optimal Size

5.3.1. Report the summary statistics

We conducted 30 training runs, setting the learning rate to 0.01, epochs to 200 based on optimal epoch analysis, and using Huber loss as the loss function. Different batch sizes (4, 8, 16, 32, 64, 128, 256) were used to evaluate the performance of the model. Summary statistics of the cost function mean, standard deviation, minimum and maximum values and running time for each batch size.

The following line plots illustrate the train and test Huber loss values over different batch sizes for analysis of overfitting.

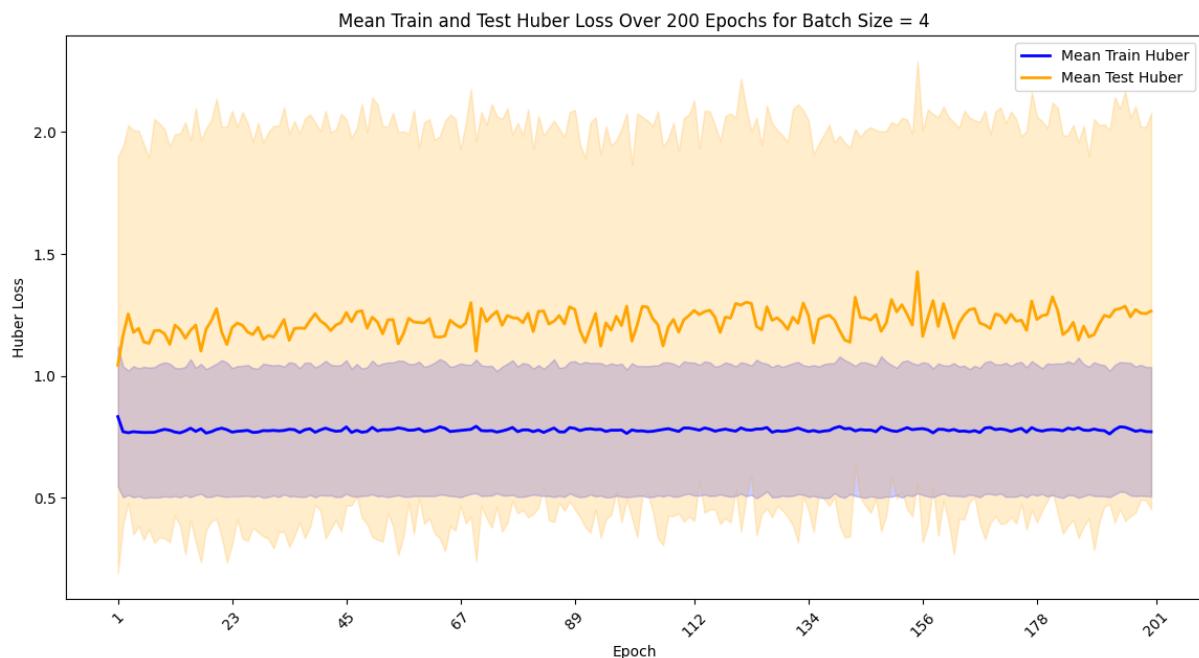


Figure 50 LSTM Mean Train and Test Huber Loss Over 200 Epochs for Batch Size = 4



Figure 51 LSTM Mean Train and Test Huber Loss Over 200 Epochs for Batch Size = 8



Figure 52 LSTM Mean Train and Test Huber Loss Over 200 Epochs for Batch Size = 16

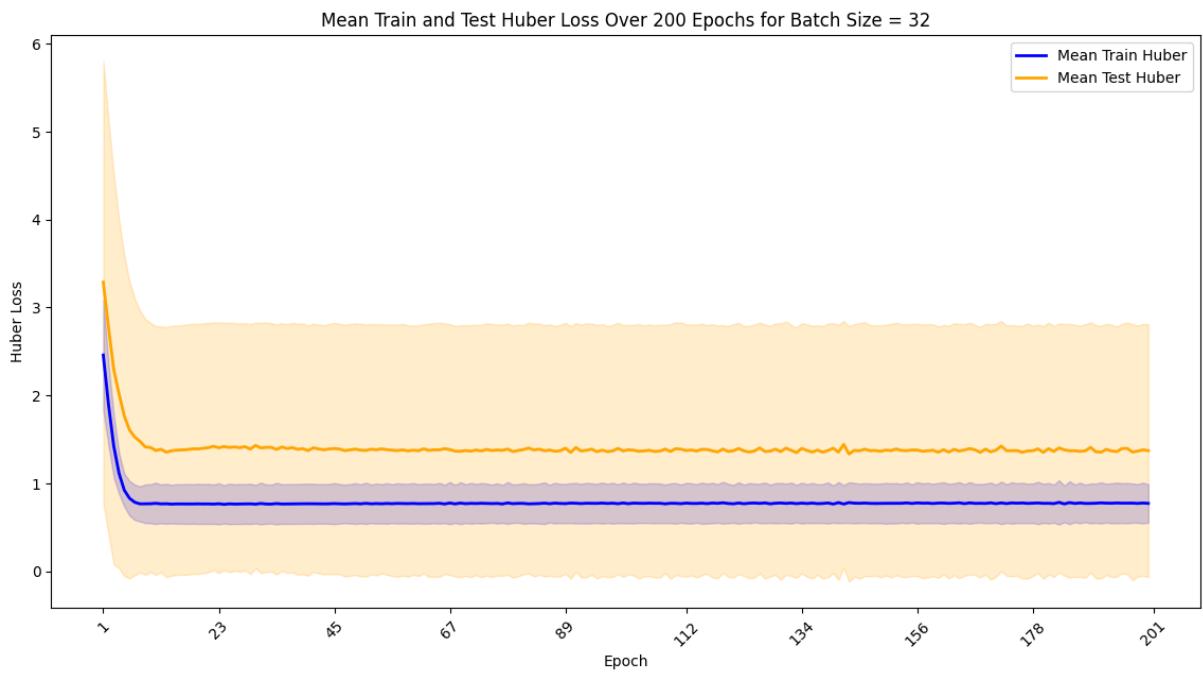


Figure 53 LSTM Mean Train and Test Huber Loss Over 200 Epochs for Batch Size = 32

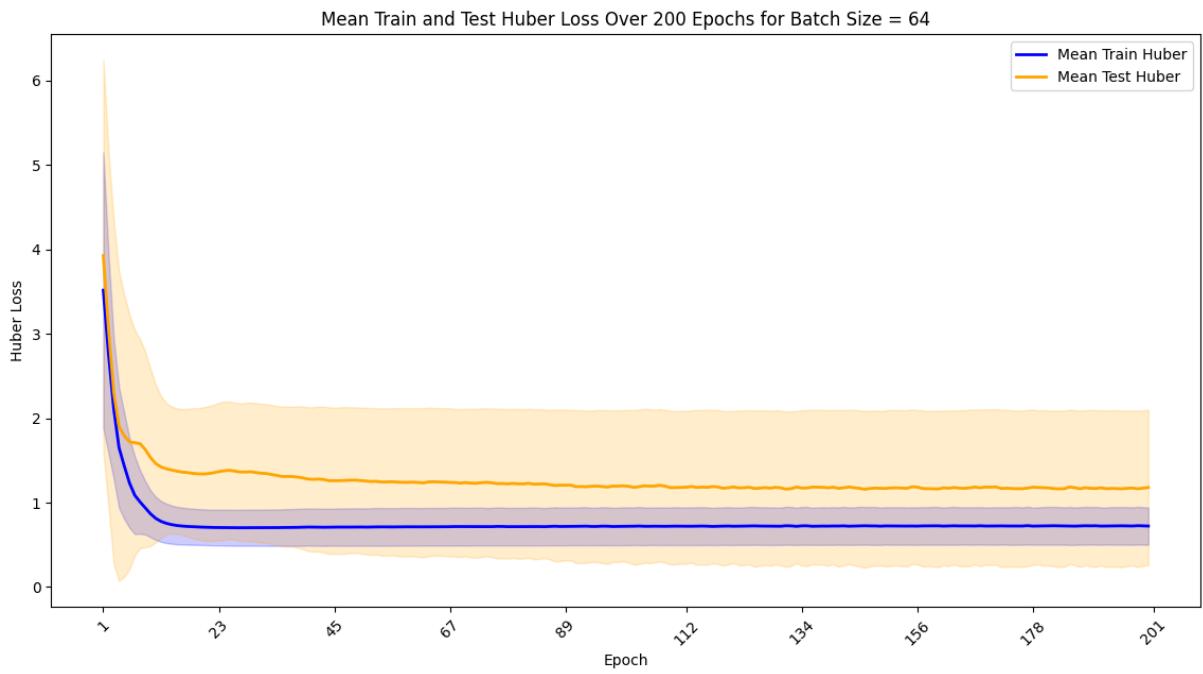


Figure 54 LSTM Mean Train and Test Huber Loss Over 200 Epochs for Batch Size = 64

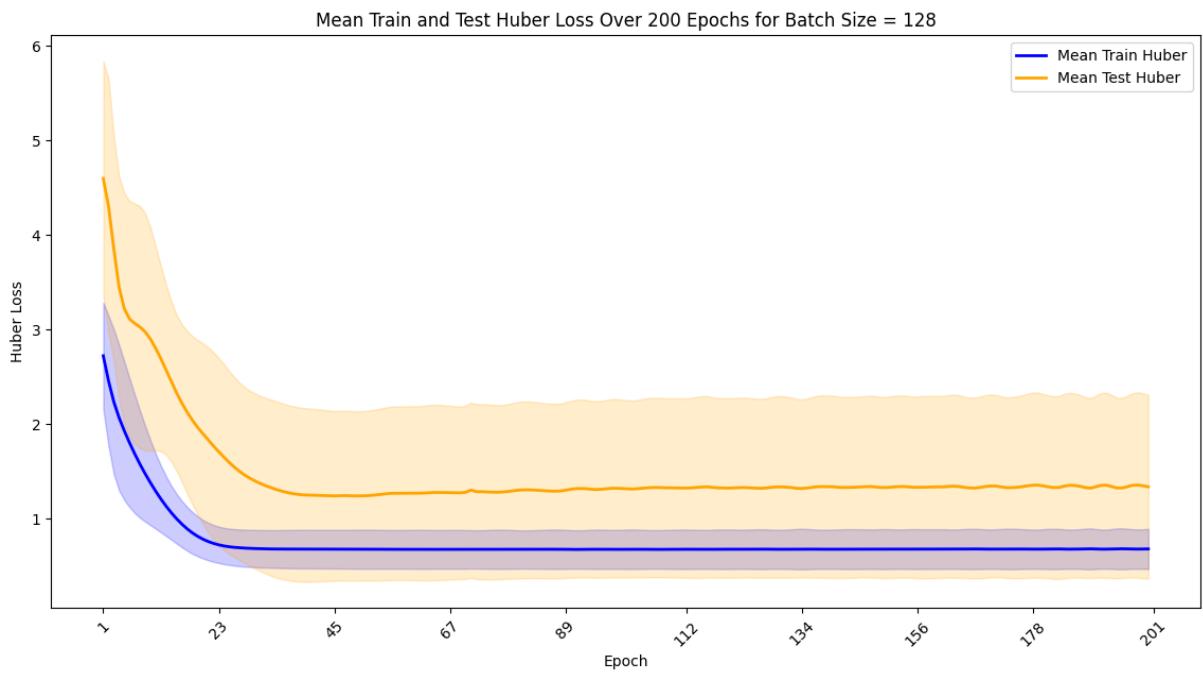


Figure 55 LSTM Mean Train and Test Huber Loss Over 200 Epochs for Batch Size = 128

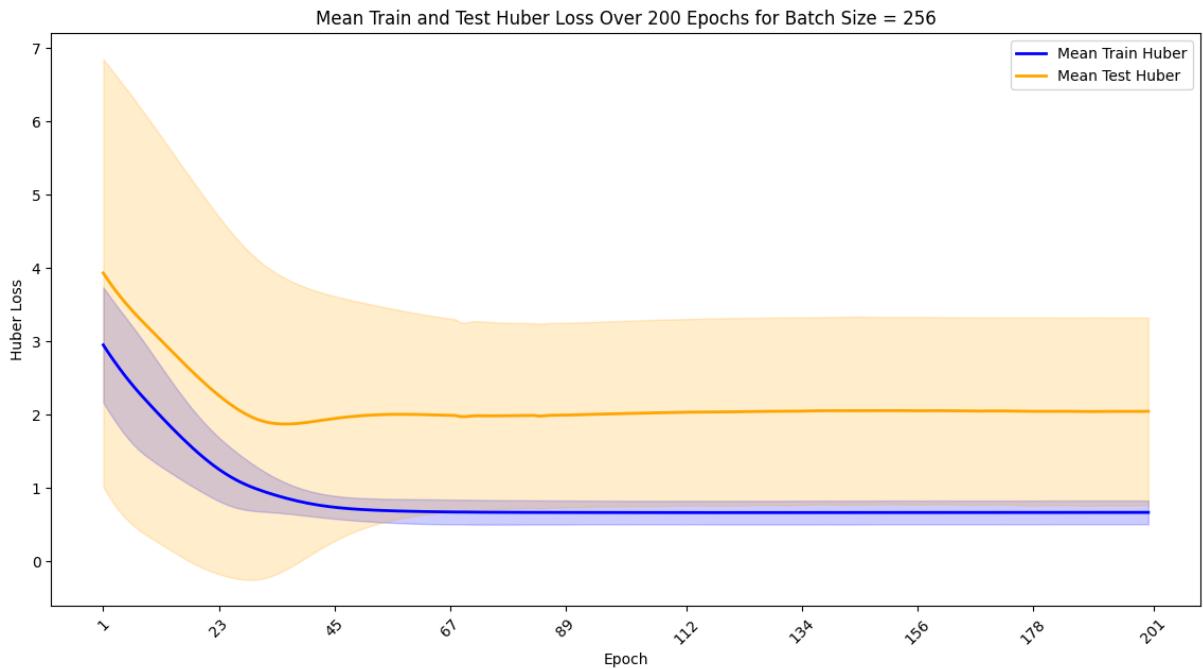


Figure 56 LSTM Mean Train and Test Huber Loss Over 200 Epochs for Batch Size = 256

- Batch Size = 4: The Huber loss shows significant variability, indicating instability in training.
- Batch Size = 8: The Huber loss remains relatively high and variable, suggesting underfitting.
- Batch Size = 16: The Huber loss is lower and more stable, indicating a good balance between learning and generalisation.

- Batch Size = 32: The Huber loss starts to show more variability, with some signs of overfitting.
- Batch Size = 64: The Huber loss is low and stable, indicating good performance without significant overfitting.
- Batch Size = 128: The Huber loss increases again, with more variability and potential overfitting.
- Batch Size = 256: The Huber loss is high, with significant variability and many outliers, indicating overfitting.

Summary Statistics:

The summary statistics for the test Huber loss and time for all batch size configurations are as follows

Table 14 Summary Statistics for All Batch Size Configurations

Batch Size	Test Huber Mean	Test Huber Std Dev	Test Huber Min	Test Huber Max	Time Mean (s)
4	1.265298	0.811906	0.283959	2.833586	50.246278
8	1.615765	1.056485	0.434209	4.311317	32.066006
16	1.167724	0.896273	0.284821	3.112693	22.863513
32	1.370608	1.438615	0.002350	3.794633	20.857284
64	1.179884	0.920753	0.214852	2.964329	17.825495
128	1.337338	0.971068	0.182127	3.297163	16.079806
256	2.046280	1.277735	0.093482	3.988119	16.116098

- Batch Size and Mean Test Huber Loss:
 - Smaller batch sizes, such as 4, show a higher mean test Huber loss (1.265298) compared to larger batch sizes. The standard deviation is relatively lower, indicating more consistent performance.
 - Batch sizes like 16 and 64 exhibit lower mean test Huber loss (1.167724 and 1.179884, respectively), suggesting they are more effective in reducing loss.
 - The batch size of 256 has the highest mean test Huber loss (2.046280), indicating poorer performance.
- Minimum and Maximum Loss:
 - Minimum test Huber loss values remain low across all batch sizes, but the maximum test Huber loss varies significantly.
 - Smaller batch sizes (4, 8) show a wider range between the minimum and maximum test Huber loss, while larger batch sizes (128, 256) exhibit a narrower range.
- Time Mean:

- Smaller batch sizes require more time for training, with batch size 4 taking the longest (50.246278 seconds).
- Larger batch sizes, like 128 and 256, show significantly reduced mean training times (16.079806 and 16.116098 seconds, respectively), highlighting efficiency improvements.

Conclusion

Batch sizes of 16 and 64 provide a good balance between low test Huber loss and reasonable training time, making them optimal choices for this LSTM model. Batch sizes that are too small or too large negatively impact performance, either through increased loss or inefficiencies in training time.

5.3.2. Best Batch Size

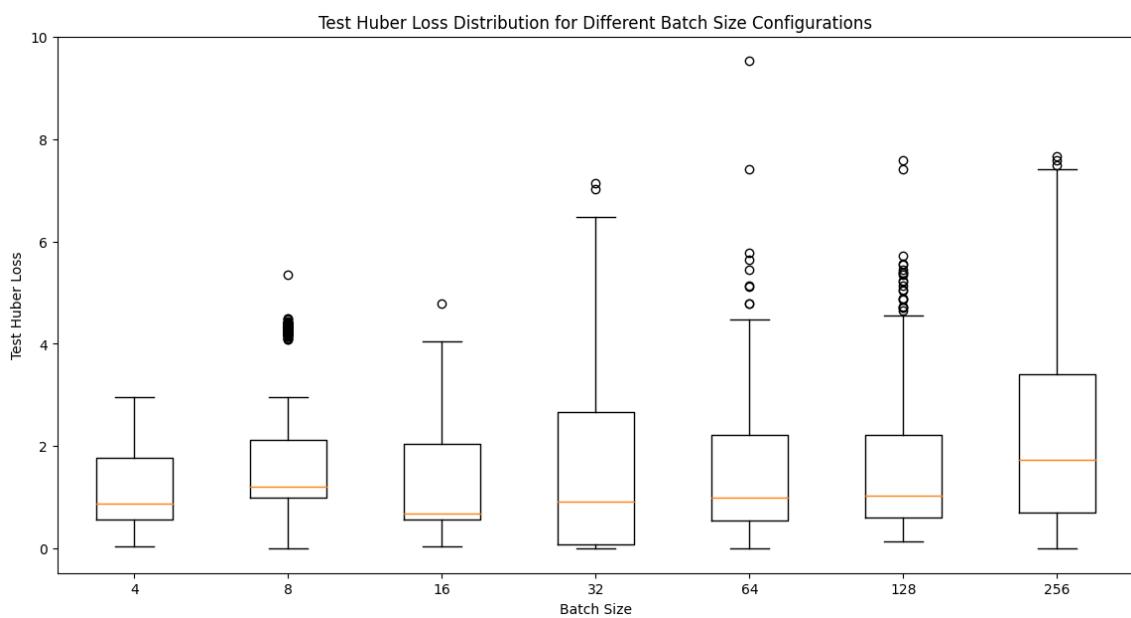


Figure 57 LSTM Test Huber Loss Distribution for Different Batch Size Configurations

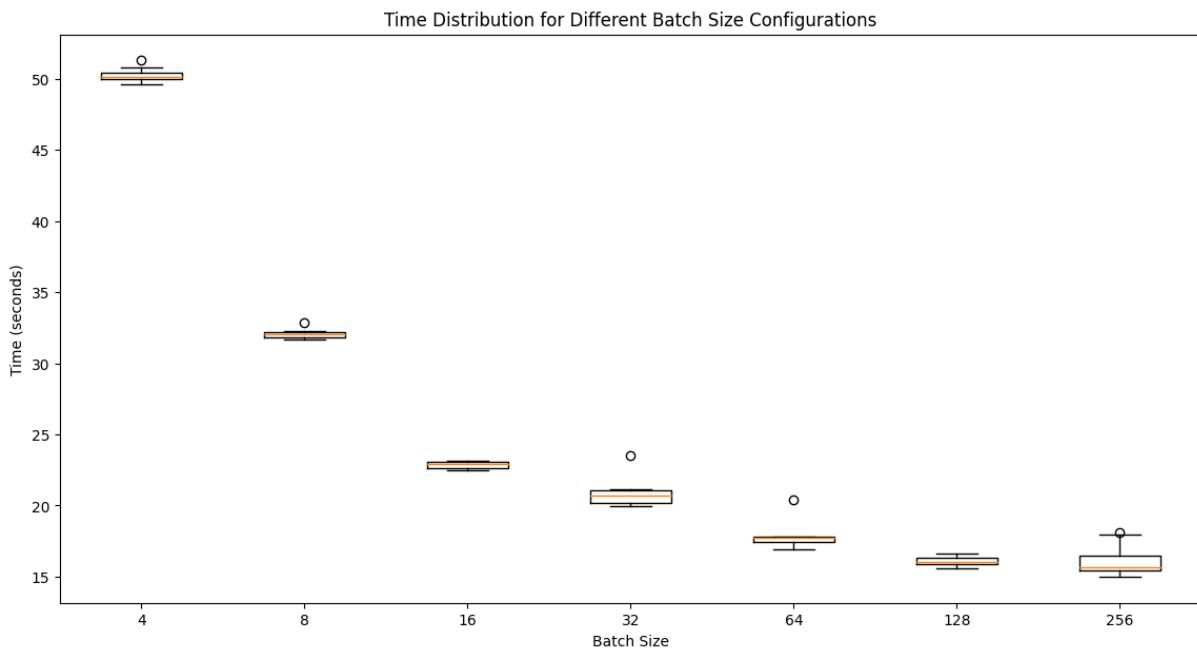


Figure 58 LSTM Time Cost Distribution for Different Batch Size Configurations

- Huber Loss Distribution: From the box plot, we can see that the batch size 16 configuration has a relatively concentrated Huber loss distribution with a low median and small standard deviation, indicating stable performance. Batch sizes 4 and 64 also show good performance, with low variability in test Huber loss. In contrast, the batch size 256 configuration shows a larger number of outliers, suggesting potential overfitting or instability.
- Performance Comparison: Comparing the model performance across different batch size configurations, we find that batch sizes 16 and 64 perform the best on the test set. The batch size 256 configuration, while reducing training time, shows higher test Huber loss and greater variability, indicating overfitting.
- Training Time Distribution: The plot shows that the training time decreases as the batch size increases. Smaller batch sizes result in longer training times due to more frequent updates to the model parameters. Larger batch sizes reduce the number of updates, thus decreasing the overall training time.

Conclusion:

Considering the training and test Huber loss, training time, and performance stability, the batch size 64 configuration is determined to be the best choice. This configuration yields the best performance on the test set with moderate training time and avoids significant overfitting issues. Batch size 16 also performs well but with longer training time. Larger batch sizes (128, 256) reduce training time but result in higher test Huber loss and greater variability in performance.

5.4. Neurons Optimal Number

5.4.1. Report the summary statistics

We conducted 30 training runs with different neuron configurations, keeping the learning rate at 0.01 and the batch size at 64, with epochs set to 200. The neuron configurations tested were 1, 2, 4, 8, 16, 32, 64, 128, and 256. The primary objective was to evaluate how varying the number of neurons affects the model's performance in terms of Huber loss and training time.

The following line plots illustrate the train and test Huber loss values over different numbers of neurons for analysis of overfitting.

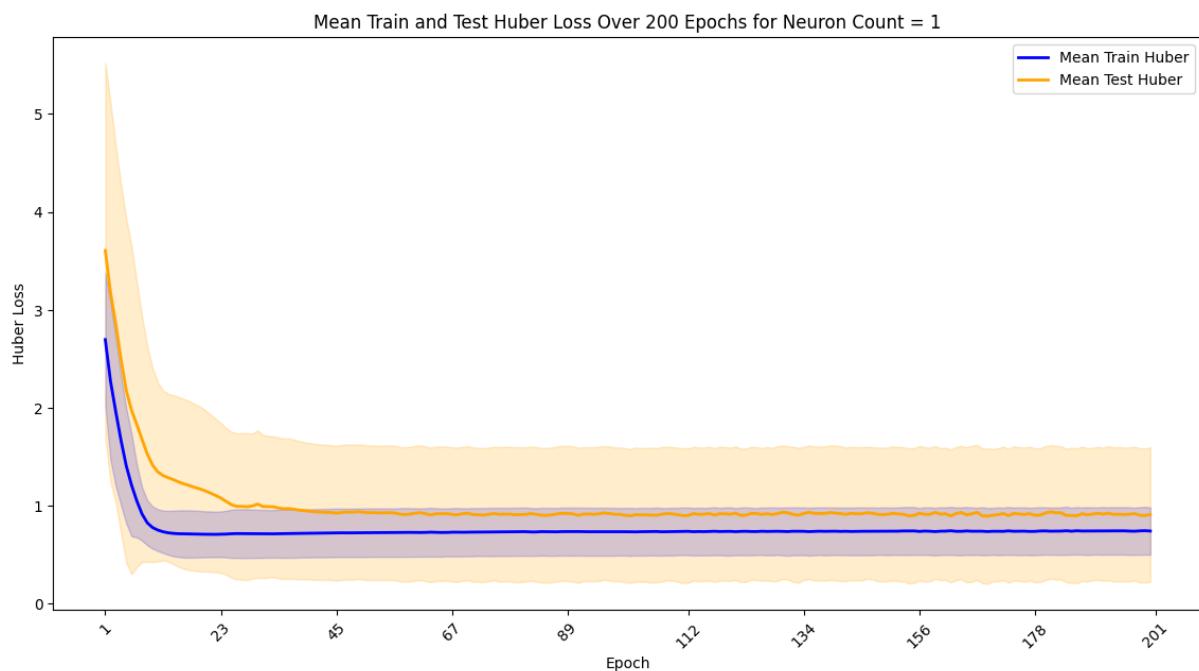


Figure 59 LSTM Mean Train and Test for Neuron Count = 1 on Hidden Layer

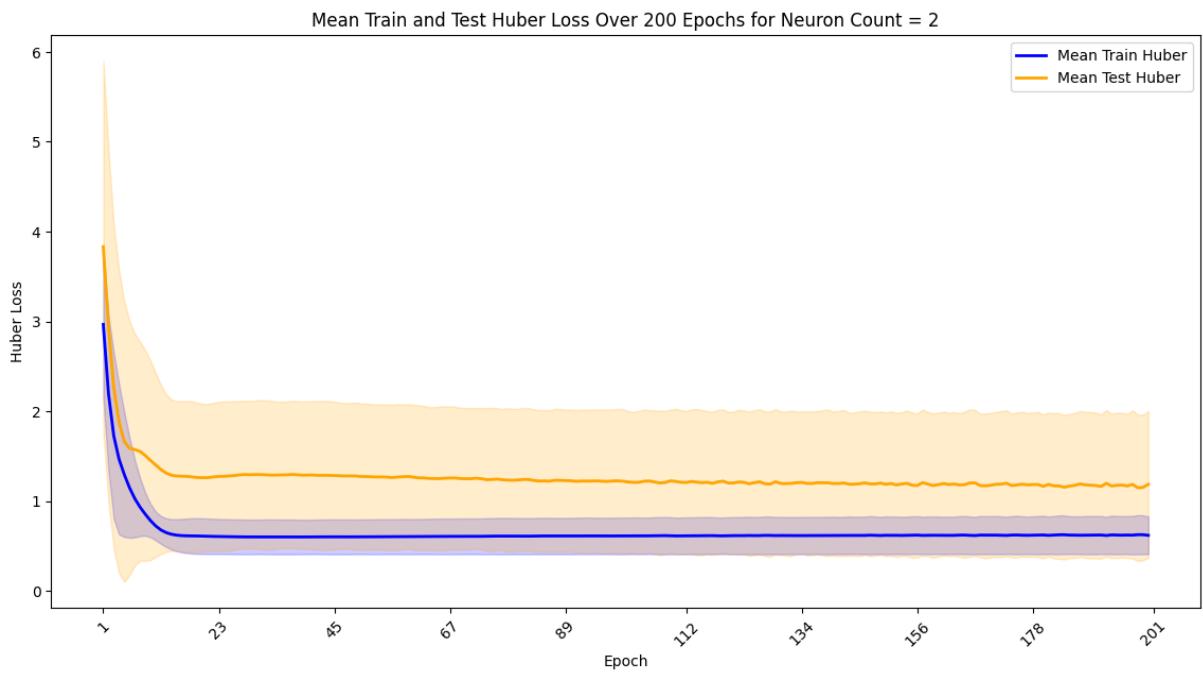


Figure 60 LSTM Mean Train and Test for Neuron Count = 2 on Hidden Layer

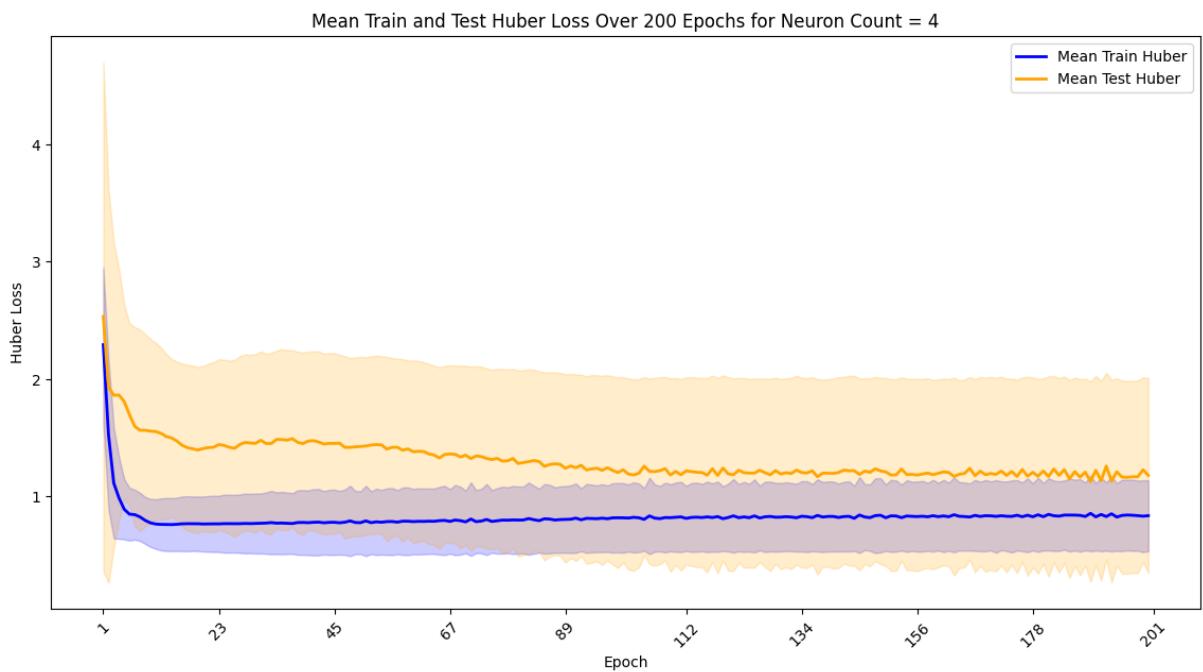


Figure 61 LSTM Mean Train and Test for Neuron Count = 4 on Hidden Layer

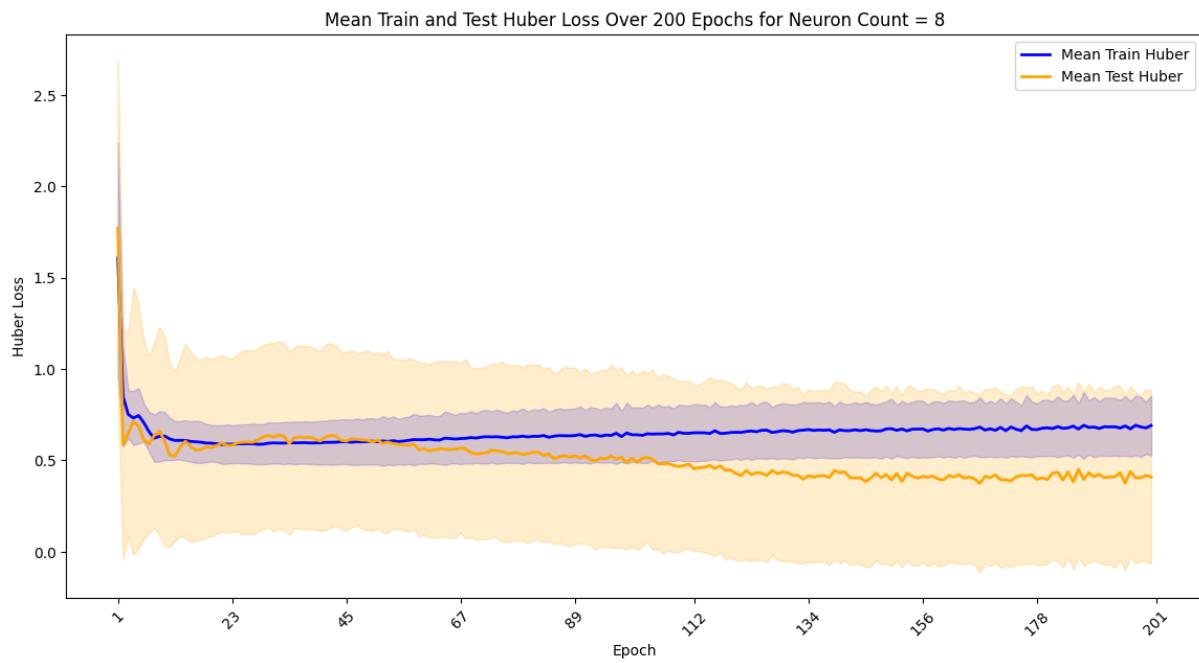


Figure 62 LSTM Mean Train and Test for Neuron Count = 8 on Hidden Layer

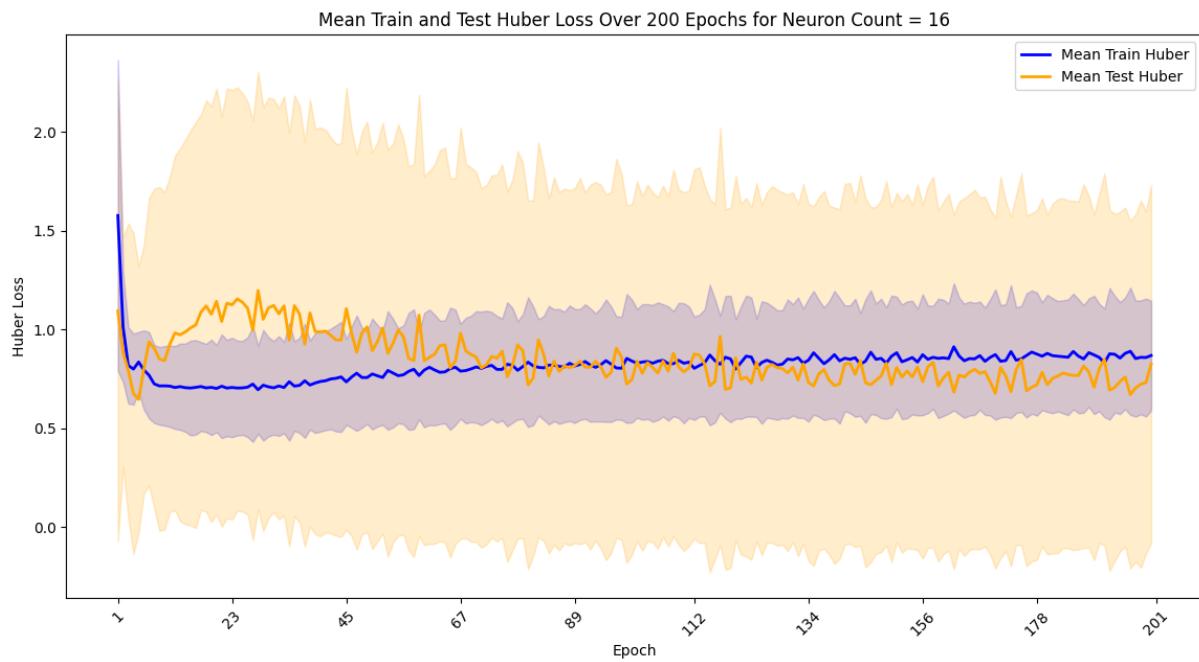


Figure 63 LSTM Mean Train and Test for Neuron Count = 16 on Hidden Layer

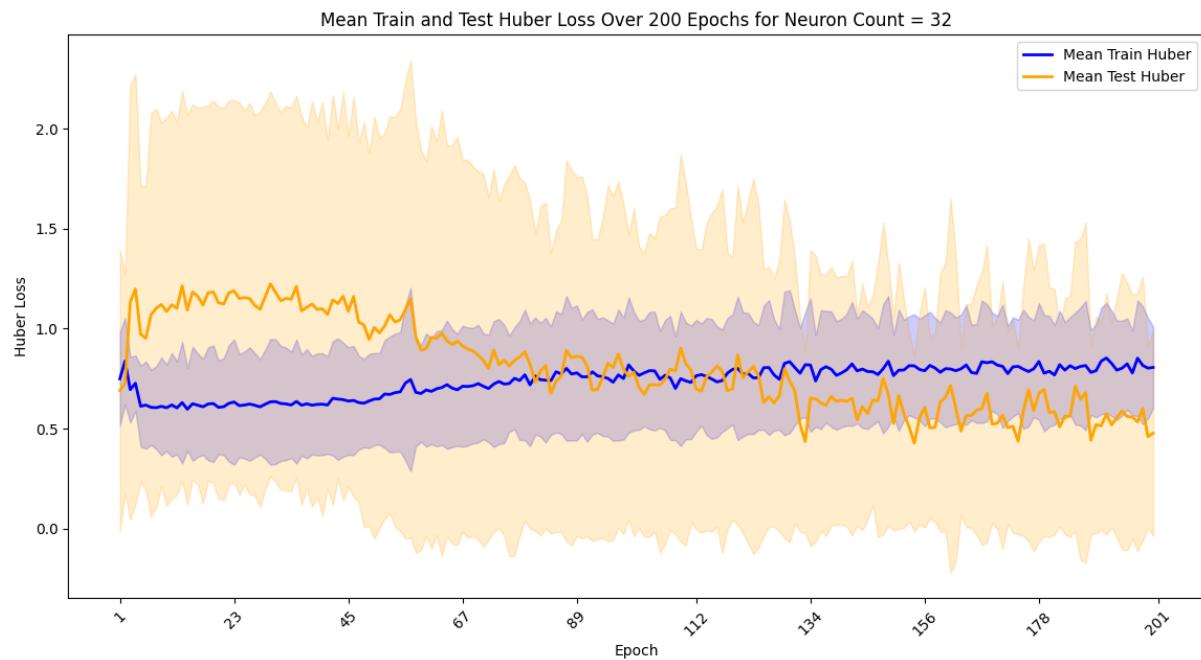


Figure 64 LSTM Mean Train and Test for Neuron Count = 32 on Hidden Layer

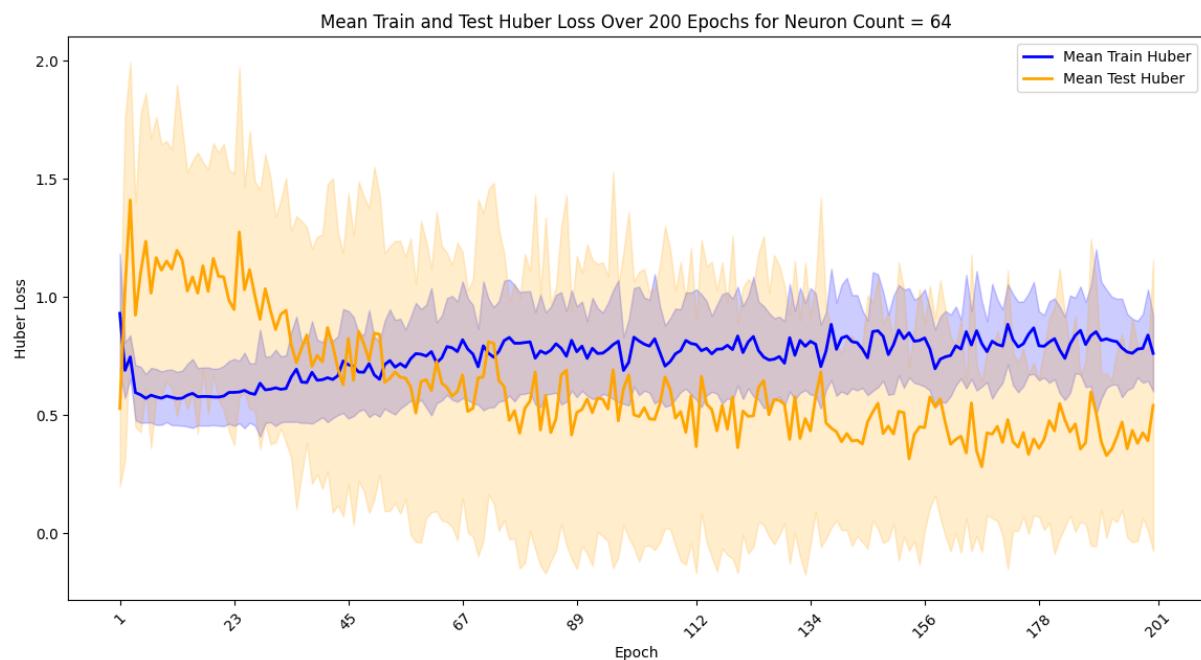


Figure 65 LSTM Mean Train and for Neuron Count = 64 on Hidden Layer

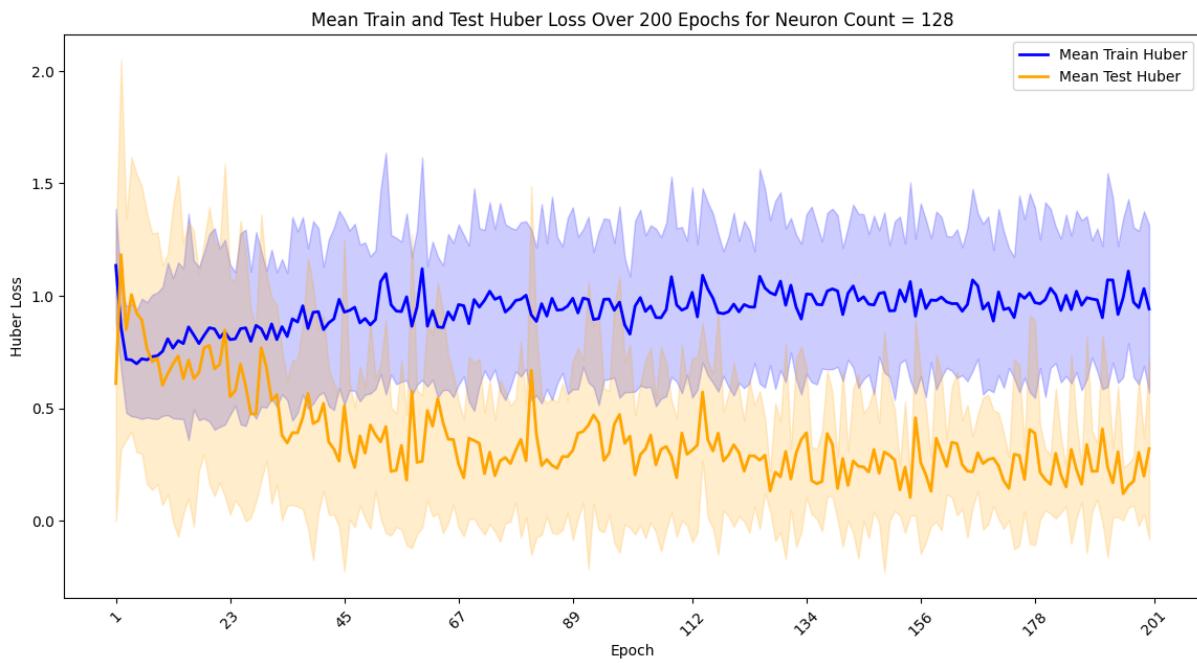


Figure 66 LSTM Mean Train and Test for Neuron Count = 128 on Hidden Layer

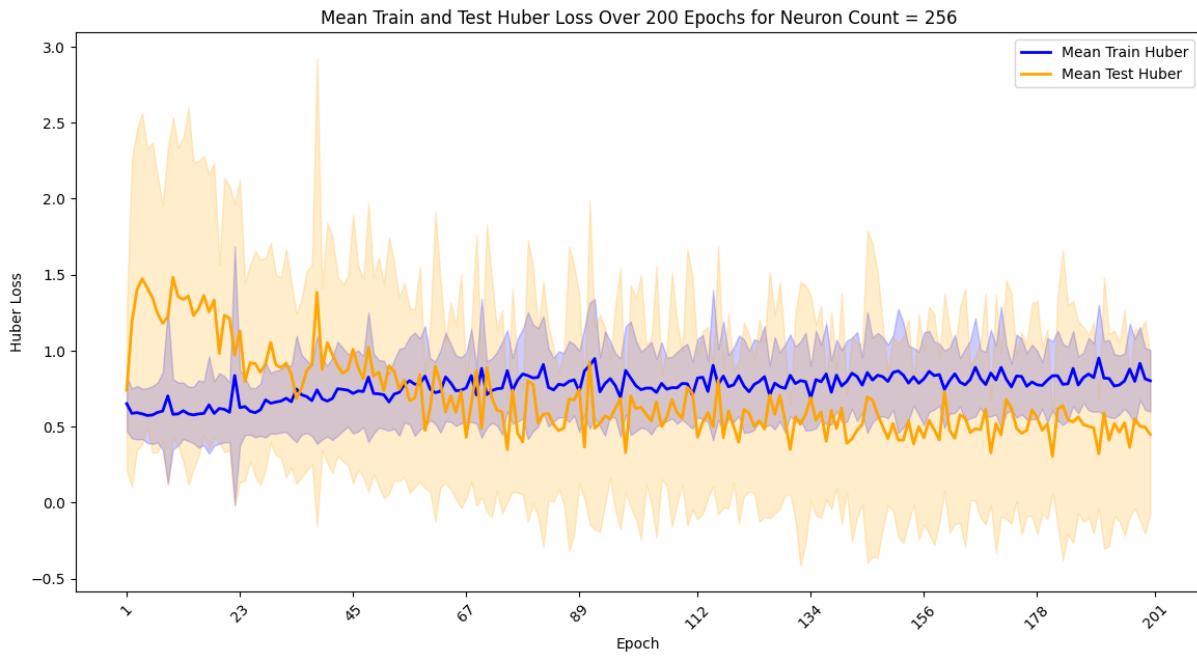


Figure 67 LSTM Mean Train and Test for Neuron Count = 256 on Hidden Layer

- Neuron Count = 1: The Huber loss decreases rapidly within the first few epochs and then stabilises. However, the test loss shows higher variability, indicating potential underfitting.
- Neuron Count = 2: Similar to the 1 neuron configuration, but with more neurons, the Huber loss shows a more consistent downward trend. The test loss still exhibits some fluctuations.

- Neuron Count = 4: With 4 neurons, the training loss continues to decrease, and the test loss shows a more stable pattern. This configuration seems to balance between underfitting and overfitting.
- Neuron Count = 8: The 8 neurons configuration achieves the lowest and most stable test loss, indicating a good fit for the data without significant overfitting.
- Neuron Count = 16: Although the training loss continues to decrease, the test loss shows more variability and outliers, suggesting that the model may be overfitting the training data.
- Neuron Count = 32, 64, 128, 256: These configurations show increasing overfitting with the number of neurons, as indicated by the higher test loss and variability.

Summary Statistics

The table below summarises the test Huber loss and training time for different neuron count configurations.

Table 15 Summary Statistics for Different Neuron Count Configurations

Neuron Count	Test Huber Mean	Test Huber Std Dev	Test Huber Min	Test Huber Max	Time Mean
1	0.912088	0.686218	0.091705	1.736611	17.591369
2	1.186427	0.817511	0.067574	2.460284	17.505485
4	1.178749	0.829548	0.185987	3.262839	17.532434
8	0.409530	0.475699	0.026940	1.569849	17.670490
16	0.825020	0.903861	0.027515	2.631970	17.632412
32	0.476491	0.512857	0.013359	1.341394	17.977479
64	0.539497	0.617093	0.023242	1.947430	18.349120
128	0.320828	0.403075	0.045227	1.189516	18.702594
256	0.448271	0.519563	0.000160	1.626913	18.194403

- Neuron Count and Mean Test Huber Loss:
 - Lower neuron counts, such as 1 and 2, show a moderate test Huber loss (0.912088 and 1.186427, respectively), but the loss is relatively high compared to other configurations.
 - Neuron counts like 8 and 128 exhibit significantly lower mean test Huber loss (0.409530 and 0.320828, respectively), suggesting better performance in reducing loss.
 - A very high neuron count, like 256, results in a higher mean test Huber loss (0.448271) compared to 128, indicating that extremely large neuron counts might not always be beneficial.
- Minimum and Maximum Loss:

- Minimum test Huber loss values remain low across all neuron counts, but the maximum test Huber loss varies significantly.
- Configurations like 1, 2, and 4 show higher maximum test Huber loss values (1.736611, 2.460284, and 3.262839, respectively), indicating less consistent performance.
- Neuron counts of 32 and 64 exhibit lower maximum test Huber loss values (1.341394 and 1.947430, respectively), suggesting more stable performance.
- Time Mean:
 - Training time increases slightly with the increase in neuron count. Smaller neuron counts, such as 1, 2, and 4, require less time (around 17.5 seconds), whereas larger neuron counts, such as 128 and 256, take slightly more time (around 18.2 seconds).

Conclusion

Neuron counts of 8 and 128 provide a good balance between low test Huber loss and reasonable training time, making them optimal choices for this LSTM model. Extremely low or very high neuron counts negatively impact performance, either through increased loss or inefficiencies in training time.

5.4.2. Best Neurons Number

The following box plots illustrate the distribution of test Huber loss and training time for each neuron count configuration.

Test Huber Loss Distribution

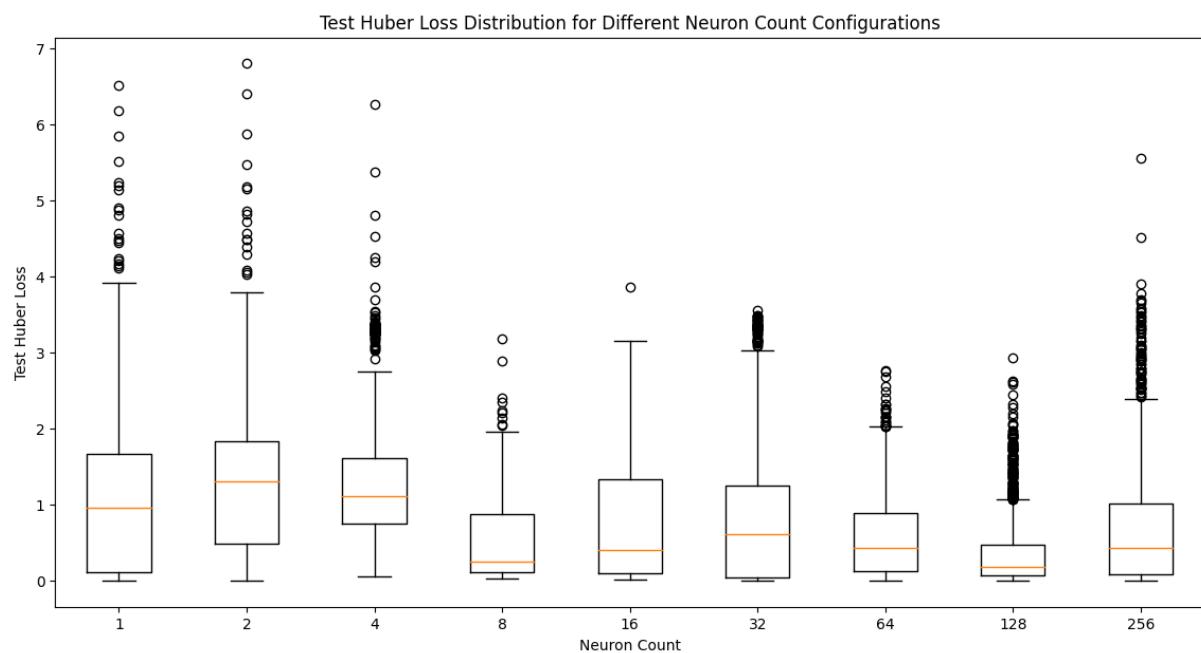


Figure 68 LSTM Test Huber Loss Distribution for Different Neuron Count Configurations

Training Time Distribution

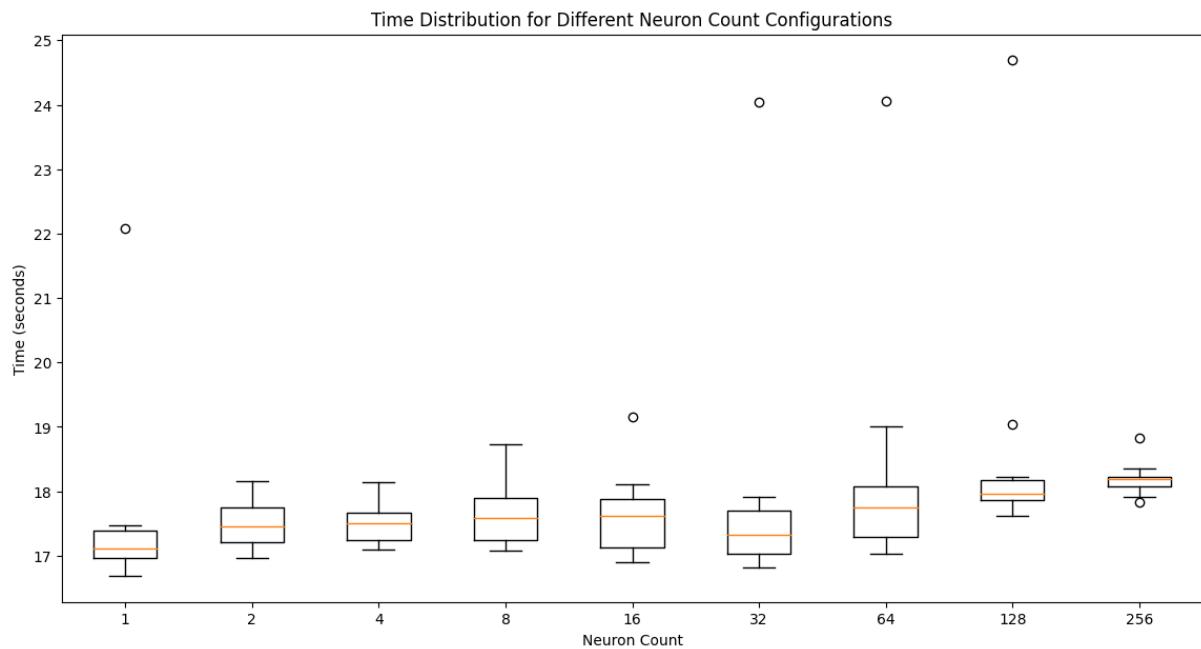


Figure 69 LSTM Time Cost Distribution for Different Neuron Count Configurations

- Huber Loss Distribution: From the box plot, we can see that the neuron count of 8 has the most concentrated Huber loss distribution with the lowest median and relatively small standard deviation, indicating stable performance. Configurations with higher neuron counts, while showing lower test loss in some cases, also exhibit a larger number of outliers, suggesting potential overfitting issues.
- Performance Comparison: Comparing the model performance across different neuron count configurations, we find that the 8 neurons configuration performs the best on the test set. Although configurations with higher neuron counts also show good performance on the training set, they do not perform as well on the test set and show signs of overfitting.
- Training Time Distribution: The plot shows that the training time slightly increases with the number of neurons, but not significantly. The configurations with larger neuron counts take slightly longer to train, but the increase in time is relatively minor.

Conclusion

Considering the training and test Huber loss, training time, and performance stability, the 8 neurons configuration is determined to be the best choice. This configuration yields the best performance on the test set with moderate training time and avoids overfitting issues.

5.5. LSTM Conclusion

Based on the analysis conducted, with fixed learning rate = 0.1 and loss function set to Huber, the following hyper-parameter configurations are recommended for the LSTM model to achieve the best performance in predicting PM2.5 levels:

- Number of Epochs: 200
- Batch Size: 64
- Number of Hidden Layer Neurons: 8

Together, these configurations provide the best balance between training time, model performance, and generalisation ability.

6. Model Comparison

6.1. Visually Compare Model Performance.

To compare the performance of the MLP and LSTM models in predicting PM2.5 levels, we followed these steps:

1. Train both models using the best hyper-parameters found predictions on the test set.

Table 16 MLP Model Summary

Attribute	Value
Hidden Layers	(15, 10)
Number of iterations	13
Learning rate	0.01

Table 17 LSTM Model Summary

Layer (type)	Output Shape	Param #
Istm_271 (LSTM)	(None, 8)	448
dropout_542 (Dropout)	(None, 8)	0
dense_542 (Dense)	(None, 8)	72
dropout_543 (Dropout)	(None, 8)	0
dense_543 (Dense)	(None, 1)	9
Total params		529
Trainable params		529
Non-trainable params		0
Learning Rate		0.01
Best Epoch		200
Best Batch Size		64
Best Neuron Count		8

2. Plot the actual vs. predicted PM2.5 values for both models.

Visual Prediction Comparison:

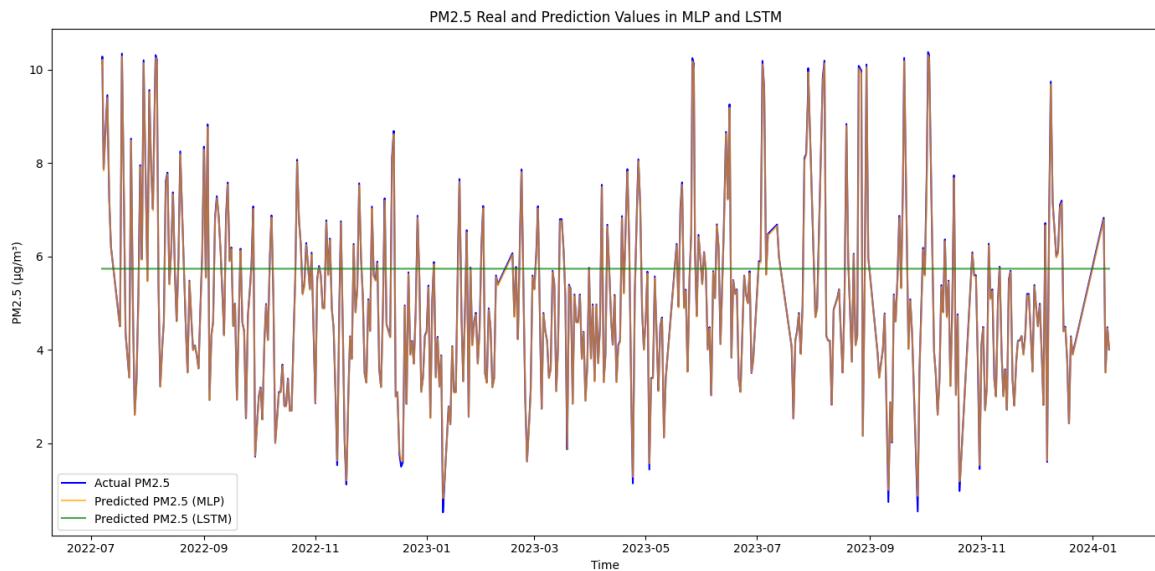


Figure 70 MLP and LSTM Performance Comparison on Linear Interpolation Dataset

In the initial results (as shown above), we observed that the LSTM predictions formed a straight line, indicating very poor performance. We consider this issue is due to data problems, as we initially used linear interpolation to handle outliers.

Data Processing Improvement: KNN Interpolation

To improve data quality, we decided to use KNN interpolation instead. The data processing steps, feature selection, and model configuration remain unchanged, so we won't elaborate on these details again. The improved results are shown in the figure below:

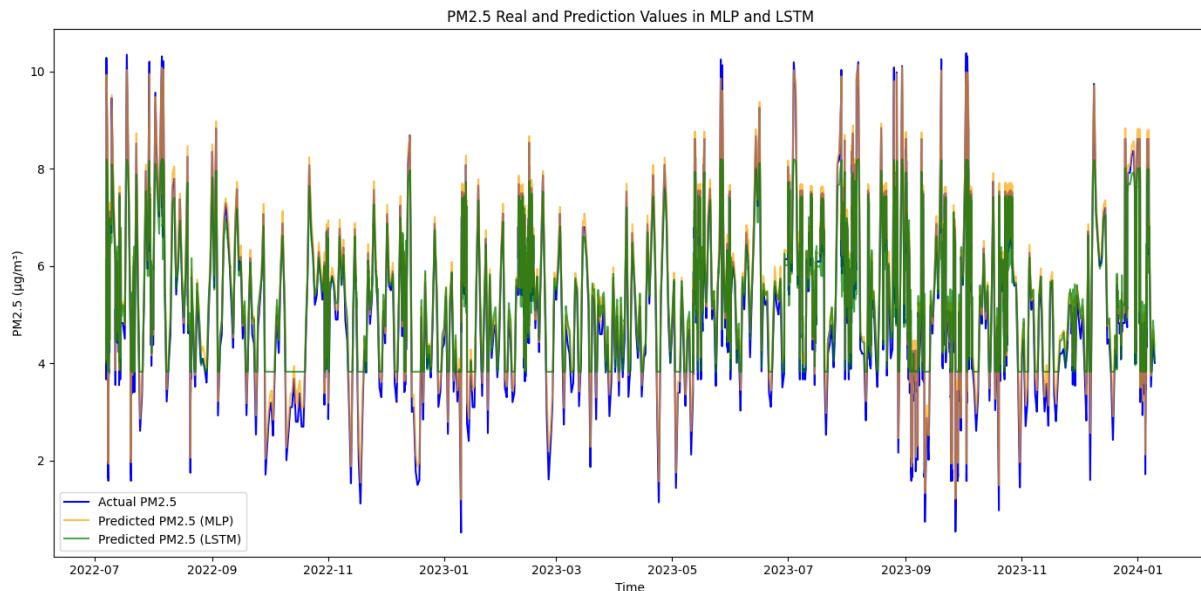


Figure 71 MLP and LSTM Performance Comparison on KNN Interpolation Dataset

In the improved results, although the LSTM predictions no longer form a straight line, there are still noticeable errors. This indicates that while data processing has improved, it is still insufficient to significantly enhance model performance.

Data Selection: High-Quality Data Segment

To further improve model performance, we selected data from March 18, 2022, to June 6, 2023. This segment has higher data quality and is more suitable for model training. The processing steps, feature selection, and model configuration remain the same, and the specifics are not reiterated here. The results after training are shown in the figure below:

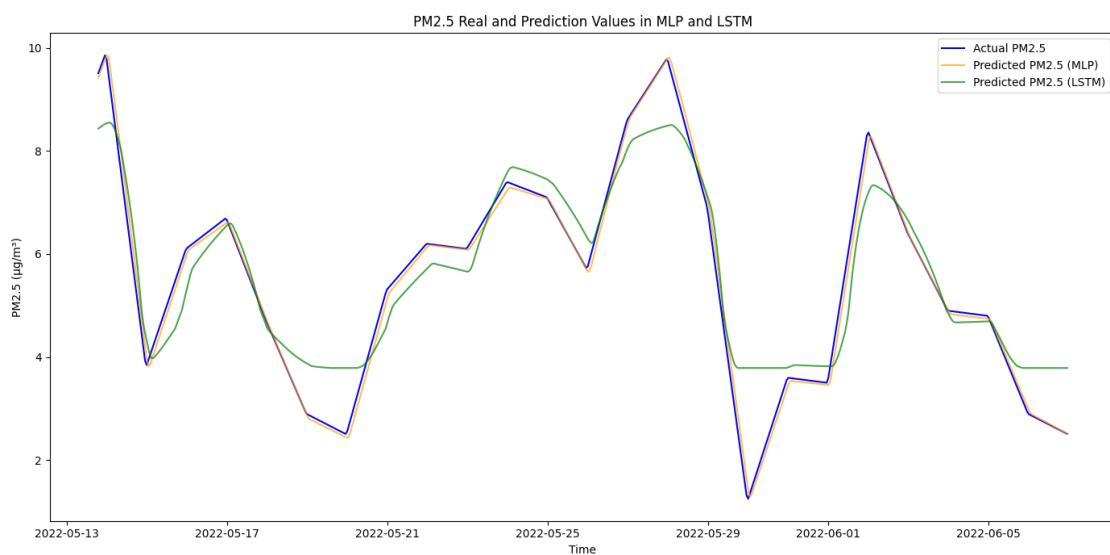


Figure 72 MLP and LSTM Performance Comparison on High Quality Dataset

From the final results, we observe a significant improvement in model performance. Therefore, **we will use this high-quality data segment for training and use the**

trained model for prediction and subsequent analysis.

The line plot above shows the actual PM2.5 values alongside the predicted values from both the MLP and LSTM models over a specified time period. The blue line represents the actual PM2.5 values, while the orange and green lines represent the predicted values from the MLP and LSTM models, respectively. Here's a detailed analysis of the plot:

- Trend Matching:
 - Both the MLP and LSTM models capture the overall trend of the actual PM2.5 values. The peaks and troughs in the actual PM2.5 levels are reflected in the predictions of both models.
- Prediction Accuracy:
 - The MLP model's predictions (orange line) are generally closer to the actual PM2.5 values (blue line) compared to the LSTM model's predictions (green line). This is particularly evident in the time periods where there are sharp changes in PM2.5 levels, such as the peaks around May 25 and May 29.
 - The LSTM model's predictions tend to deviate more from the actual values, especially during periods of rapid increase or decrease in PM2.5 levels.
- Overlapping Predictions:
 - In some periods, the MLP and LSTM predictions overlap significantly with each other and the actual values, indicating that both models are performing well in those specific intervals. For example, around May 17, both models are able to closely follow the actual PM2.5 values.
- Model Differences:
 - The MLP model appears to be more robust in capturing sudden changes in PM2.5 levels, whereas the LSTM model shows a lag or smoother transition which may be due to its sequential nature. This results in the LSTM model sometimes missing the sharp peaks and dips in the actual data.
 - There are instances where the LSTM model underestimates the PM2.5 values, such as around May 29, where the green line falls below the blue line significantly.

Conclusion

- The MLP model demonstrates better performance in terms of prediction accuracy, as indicated by its closer alignment with the actual PM2.5 values. This aligns with the performance metrics observed earlier (lower RMSE and MAE, higher R²).
- The LSTM model, while capable of capturing overall trends, struggles more with the precise prediction of rapid changes in PM2.5 levels, leading to higher prediction errors in such periods.

6.2. Compare the Model Performance using RMSE

[5 marks]

For both models, provide root mean square error (RMSE), Mean Absolute Error (MAE), and correlation coefficient (R^2) to quantify the prediction performance of each model.

Compare the performance of both MLP and LSTM using RMSE. Which model performed better? Justify your finding. [2.5 marks]

Based on this requirement of the assignment and this part requirement, we compare the performance of the MLP and LSTM models using the RMSE (Root Mean Squared Error), MAE (Mean Absolute Error), and R^2 (Coefficient of Determination) metrics.

Table 18 Performance Comparison of MLP and LSTM

Metric	MLP	LSTM
RMSE	0.1523230754119499	0.6834049238839086
MAE	0.11225667685215686	0.531406333841138
R^2	0.9940461421020389	0.8789903142856332

Using the bar chart makes it visually observe the performance of the MLP and LSTM models using the RMSE, MAE, and R^2 metrics.

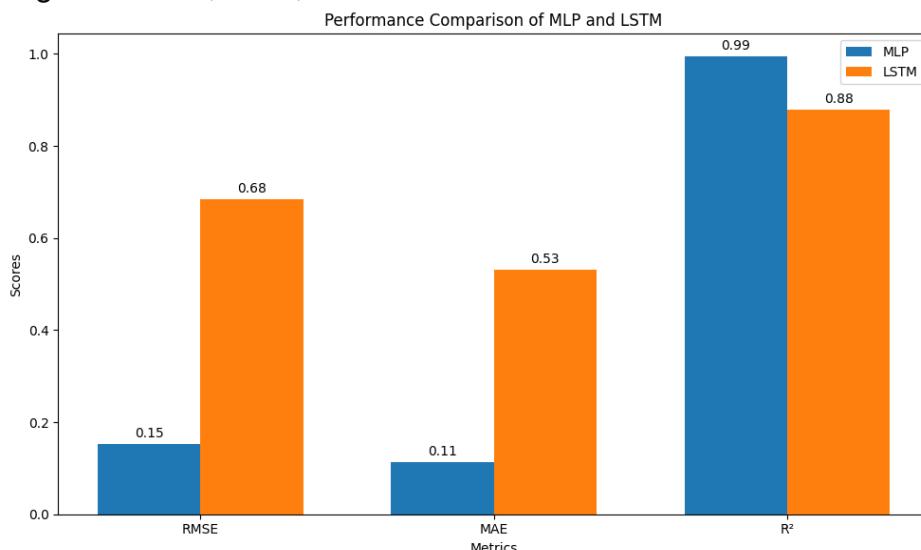


Figure 73 Performance Comparison of MLP and LSTM on High Quality Dataset

- RMSE: MLP: 0.1523 and LSTM: 0.6834
MLP performs better than LSTM in terms of RMSE because a lower The MLP model has a significantly lower RMSE compared to the LSTM model, indicating that the MLP model has smaller prediction errors on average. RMSE is sensitive to large errors, and a lower value means better performance in terms of error magnitude.
- MAE: MLP: 0.1123 and LSTM: 0.5314

The MLP model also has a lower MAE compared to the LSTM model, showing that it has more accurate predictions on average. MAE is less sensitive to outliers than RMSE, but in this case, it still shows the MLP model as superior.

- R²: MLP: 0.9940 and LSTM: 0.8789
The MLP model has a higher R² value, indicating that it explains a greater proportion of the variance in the target variable. An R² close to 1 means that the model fits the data very well.

Justify My Finding:

- Performance Metrics: All three performances metrics (RMSE, MAE, and R²) consistently show that MLP has lower prediction errors and a better fit to the data compared to LSTM.
- Error Measurement: RMSE is particularly important as it penalizes larger errors more than MAE, and the significantly lower RMSE for MLP suggests that it handles outliers or large errors better than LSTM for this dataset.
- Variance Explanation: The high R² value for MLP indicates that it captures the variance in the data much better than LSTM, suggesting a superior understanding of the underlying data patterns.

Conclusion

The MLP model performs better across all three metrics (RMSE, MAE, and R²) compared to the LSTM model. The significantly lower RMSE and MAE values suggest that the MLP model makes more accurate predictions with fewer and smaller errors. The higher R² value for the MLP model indicates a better fit to the data, meaning the MLP model's predicted values closely follow the actual PM2.5 values, whereas the LSTM model's predictions exhibit higher variability and less accuracy.

7. Conclusion

In this project, we aimed to predict PM2.5 levels using two different machine learning models: A Multilayer Perceptron (MLP) and a Long Short-Term Memory (LSTM) network. We conducted extensive experiments to optimise the hyper-parameters for both models and compared their performance in terms of root mean square error (RMSE), mean absolute error (MAE), and R-squared (R²).

- Best Epoch Analysis: We found that 200 epochs yielded the best balance between training time and model performance. This configuration minimised the test MAE loss while avoiding significant overfitting.
- Batch Size Analysis: A batch size of 64 was identified as optimal, providing a good trade-off between training stability and computational efficiency. This batch size minimised test MAE loss and maintained consistent training times across runs.
- Neuron Count Analysis: Our experiments indicated that a neuron count of 8 in the LSTM's hidden layer was optimal. This configuration produced the lowest test MAE loss, demonstrating the model's capability to capture the temporal dependencies in the data effectively.
- Model Comparison: The MLP outperformed the LSTM in this particular task. The MLP achieved an RMSE of 0.15, an MAE of 0.11, and an R² of 0.99, whereas the LSTM had an RMSE of 0.68, an MAE of 0.53, and an R² of 0.88.

Despite the LSTM's ability to model sequential data, the MLP's simpler architecture provided more accurate predictions for this dataset.

Reference

Brownlee, J. (2023). "How to Tune LSTM Hyper-parameters with Keras for Time Series Forecasting". *Machine Learning Mastery*. Retrieved June 7, 2024, from <https://machinelearningmastery.com/tune-lstm-hyperparameters-keras-time-series-forecasting/> [https://machinelearningmastery.com/tune-lstm-hyperparameters-keras-time-series-forecasting/]

Kristiani, E., Lin, H., Lin, J. R., Chuang, Y. H., Huang, C. Y., & Yang, C. T. (2022). "Short-term prediction of PM2. 5 using LSTM deep learning methods". *Sustainability*, 14(4), 2068

Appendix 1. Source Code

https://drive.google.com/file/d/1E8XmWFPwJCr3y00tGUCkNEaACQP-BN8e/view?usp=drive_link

Appendix 2. Abbreviations

RMSE: Root Mean Square Error
MAPE: Mean Absolute Percentage Error
MAE: Mean Absolute Error
MSE: Mean Square Error
LSTM: Long Short-Term Memory
GRU: Gated Recurrent Unit
RNN: Recurrent Neural Network
CNN: Convolutional Neural Network
PM: Particulate Matter
TEMP: Temperature
CO: Carbon Monoxide
NO: Nitrogen Monoxide
NOx: Nitrogen Oxide
SO: Sulphur Oxide
KNN: K-Nearest Neighbours
CSV: Comma-Separated Values