# Huy_Hung_Pham_MIS780A2_Task3

September 6, 2025

# 1 MIS780 Advanced AI For Business - Assessment 2 - T2 2025

## 1.1 Task Number: Forecasting Air Quality Pollution (PM2.5 Concentration) with Time-series Data

**Student Name:** *Huy Hung Pham*

**Student ID:** *s224212292*

## 1.2 Table of Content

1. Executive Summary

2. Data Preprocessing

3. Predictive Modeling

4. Experiments Report

5. Role of GenAI

## 1. Executive Summary

Air pollution, particularly PM2.5 (particulate matter with a diameter of 2.5 micrometers or less), poses significant risks to public health, environmental sustainability, and urban planning. To address this challenge, a comprehensive multivariate time-series dataset was utilized, consisting of hourly PM2.5 concentrations along with related meteorological and pollutant indicators from March 2013 to February 2017. Data from January 2016 was reserved for testing purposes, while the earlier data served as the training set. The forecasting task was approached using Long Short-Term Memory (LSTM) neural networks, which are particularly effective at capturing complex temporal dependencies in sequential data. A series of ten experimental models were developed, varying the number of LSTM layers, hidden units, optimizers, regularization strategies, and batch sizes. The performance of each model was assessed using the Mean Absolute Error (MAE) across multiple forecasting horizons. The experiments revealed that while simpler models performed reasonably well, deeper architectures with regularization and adaptive optimizers produced more stable forecasts. Among all the configurations tested, Model 11 which comprising a 3-layer LSTM with 256, 128, and 64 units, using the Adam optimizer, kernel regularization, and early stopping achieved the most consistent error (MAE= 0.04), making it the optimal solution.

This model provides accurate hourly forecasts of PM2.5 concentrations, facilitating data-driven decision-making for pollution control, resource allocation, and public health interventions (Istiana et al., 2024). For real-world applications, continuous retraining on recent data and integration into haze and fog alert systems is recommended(Qiang Liu et al., 2019). Additionally, further

enhancements, such as incorporating attention mechanisms or hybrid ensemble approaches, could improve predictive accuracy and robustness in dynamic urban environments.

## 2. Data Preprocessing

**2.1. Initialisation**

```
[2]:  # Importing the libraries
      import numpy as np
      import pandas as pd
      import matplotlib.pyplot as plt
      import seaborn as sns
      from sklearn.preprocessing import MinMaxScaler
      from sklearn.metrics import mean_absolute_error
      from tensorflow.keras.models import Sequential
      from tensorflow.keras.layers import Dense, LSTM
      from tensorflow.keras.optimizers import SGD
      from tensorflow.random import set_seed


      set_seed(224212292)
      np.random.seed(224212292)
```

### 1.2.1  2.2 Importing dataset

```
[3]:  dataset = pd.read_csv('/content/Task 3 Dataset Air Quality.csv'
          )
      print(dataset.tail())
```

```
              datetime  PM2.5  SO2   NO2   CO    O3  TEMP    PRES  DEWP  WSPM
35059  28/02/2017 19:00   14.0  3.0  27.0  400  72.0  12.5  1013.5 -16.2   2.4
35060  28/02/2017 20:00   18.0  3.0  37.0  400  59.0  11.6  1013.6 -15.1   0.9
35061  28/02/2017 21:00   15.0  5.0  50.0  600  41.0  10.8  1014.2 -13.3   1.1
35062  28/02/2017 22:00   11.0  6.0  49.0  500  41.0  10.5  1014.4 -12.9   1.2
35063  28/02/2017 23:00   10.0  7.0  48.0  600  39.0   8.6  1014.1 -15.9   1.3
```

```
[4]:  dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 35064 entries, 0 to 35063
Data columns (total 10 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   datetime  35064 non-null  object
 1   PM2.5     35064 non-null  float64
 2   SO2       35064 non-null  float64
 3   NO2       34372 non-null  float64
 4   CO        35064 non-null  int64
 5   O3        34558 non-null  float64
 6   TEMP      35064 non-null  float64
```

2

```
7    PRES       35064 non-null   float64
8    DEWP       35064 non-null   float64
9    WSPM       35064 non-null   float64
dtypes: float64(8), int64(1), object(1)
memory usage: 2.7+ MB
```

[5]: `print(dataset.describe())`

```
              PM2.5           SO2           NO2            CO            O3  \
count  35064.000000  35064.000000  34372.000000  35064.000000  34558.000000
mean      84.518994     18.679584     58.097172   1320.073266     58.534682
std       85.482426     24.125895     36.297740   1223.774709     58.401448
min        2.000000      0.570000      2.000000    100.000000      0.214200
25%       23.000000      3.000000     29.000000    500.000000     10.000000
50%       60.000000      9.000000     51.000000   1000.000000     45.000000
75%      114.000000     23.000000     80.000000   1600.000000     84.000000
max      844.000000    257.000000    273.000000  10000.000000    390.000000

               TEMP          PRES          DEWP          WSPM
count  35064.000000  35064.000000  35064.000000  35064.000000
mean      13.672020   1012.547222      2.447850      1.860641
std       11.455204     10.263134     13.806763      1.280133
min      -16.800000    987.100000    -35.300000      0.000000
25%        3.100000   1004.000000     -8.800000      1.000000
50%       14.600000   1012.200000      3.000000      1.500000
75%       23.500000   1020.900000     15.000000      2.400000
max       41.100000   1042.000000     28.800000     10.500000
```
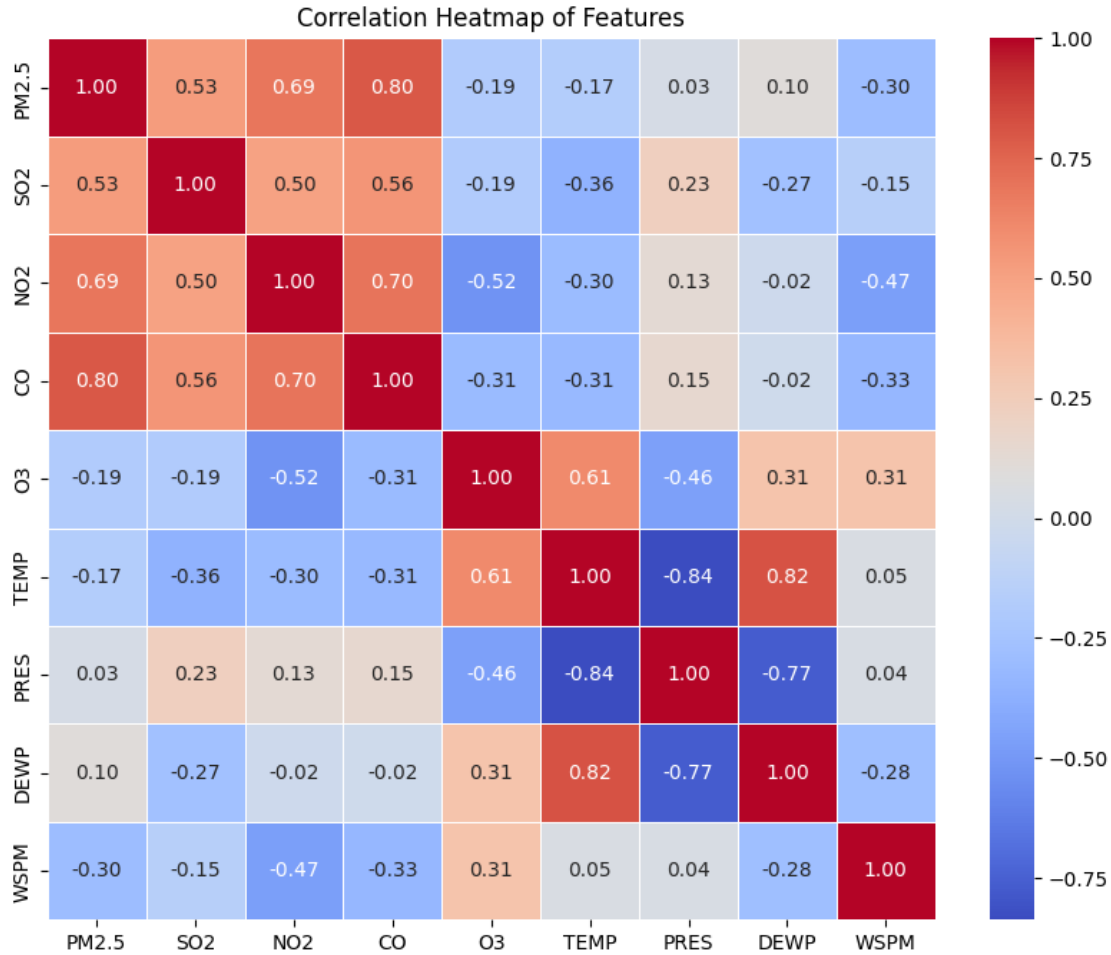
[6]: 
```python
# Convert 'datetime' column to datetime objects and set as index
dataset['datetime'] = pd.to_datetime(dataset['datetime'], format='%d/%m/%Y %H:
  ↪%M', errors='coerce')
dataset = dataset.set_index('datetime')

# Calculate the correlation matrix
correlation_matrix = dataset.corr()

# Create a heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f",␣
  ↪linewidths=.5)
plt.title('Correlation Heatmap of Features')
plt.show()
```

## Correlation Heatmap of Features



```
[7]: columns_to_drop = ['O3', 'TEMP', 'PRES', 'DEWP', 'WSPM']
     dataset = dataset.drop(columns=columns_to_drop)

     display(dataset.head())
```

|                     | PM2.5 | SO2 | NO2  | CO  |
|---------------------|-------|-----|------|-----|
| datetime            |       |     |      |     |
| 2013-03-01 00:00:00 | 5.0   | 4.0 | 12.0 | 200 |
| 2013-03-01 01:00:00 | 8.0   | 6.0 | 14.0 | 200 |
| 2013-03-01 02:00:00 | 3.0   | 5.0 | 14.0 | 200 |
| 2013-03-01 03:00:00 | 5.0   | 5.0 | 14.0 | 200 |
| 2013-03-01 04:00:00 | 5.0   | 6.0 | 21.0 | 200 |

The decision to exclude the columns 'O3', 'TEMP', 'PRES', 'DEWP', and 'WSPM' is informed by the analysis of the correlation heatmap, which indicates weak correlations between these features and 'PM2.5'—specifically, O3 (-0.19), TEMP (-0.17), PRES (0.03), DEWP (0.10), and WSPM (-0.30). In contrast, stronger correlations have been observed with 'SO2' (0.53), 'NO2' (0.68), and 'CO' (0.80). Features exhibiting weak correlations are less impactful in predictive modeling;
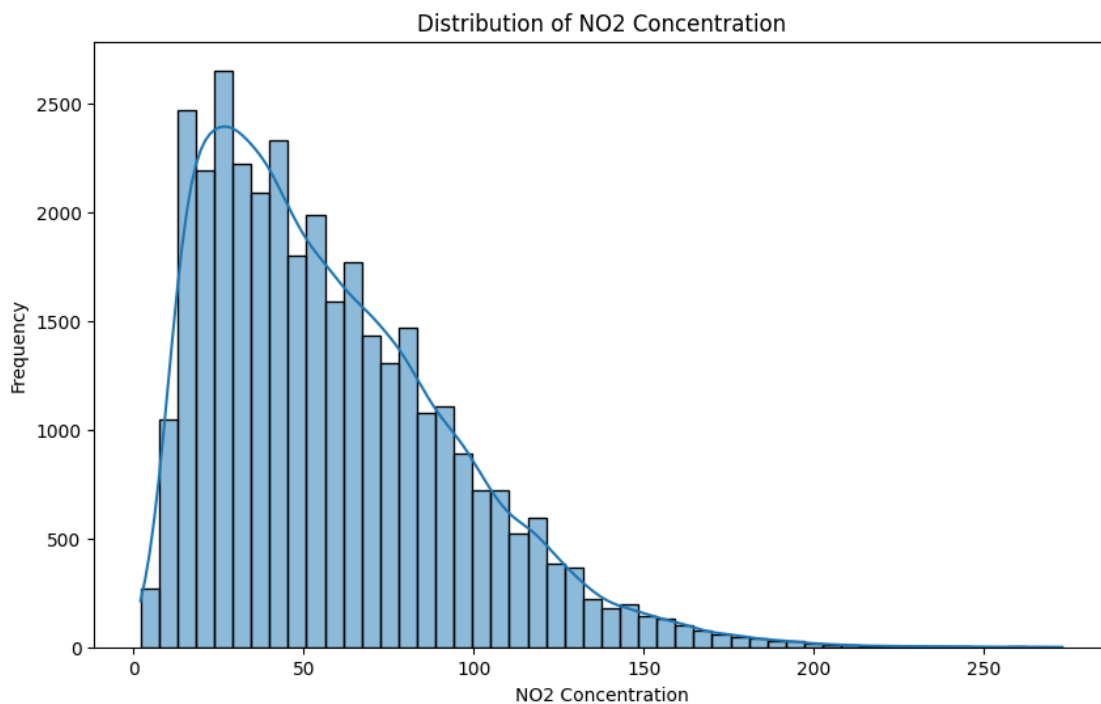
4

therefore, their exclusion can streamline the model, diminish noise, and enhance overall performance by concentrating on the more pertinent features.

### 1.2.2  2.3 Determine missing values

```
[8]: dataset.isna().sum()
```

```
[8]: PM2.5      0
     SO2        0
     NO2      692
     CO         0
     dtype: int64
```

```
[9]: # Create a histogram for the 'NO2' column to visualize its distribution
     plt.figure(figsize=(10, 6))
     sns.histplot(dataset['NO2'], kde=True, bins=50)
     plt.title('Distribution of NO2 Concentration')
     plt.xlabel('NO2 Concentration')
     plt.ylabel('Frequency')
     plt.show()
```



```
[10]: # Handle missing values in the 'NO2' column by filling with the mean
      dataset['NO2'] = dataset['NO2'].fillna(dataset['NO2'].mean())
```

5

```
# Verify that there are no more missing values in the relevant columns
print(dataset[['PM2.5', 'SO2', 'NO2', 'CO']].isna().sum())
```

```
PM2.5    0
SO2      0
NO2      0
CO       0
dtype: int64
```
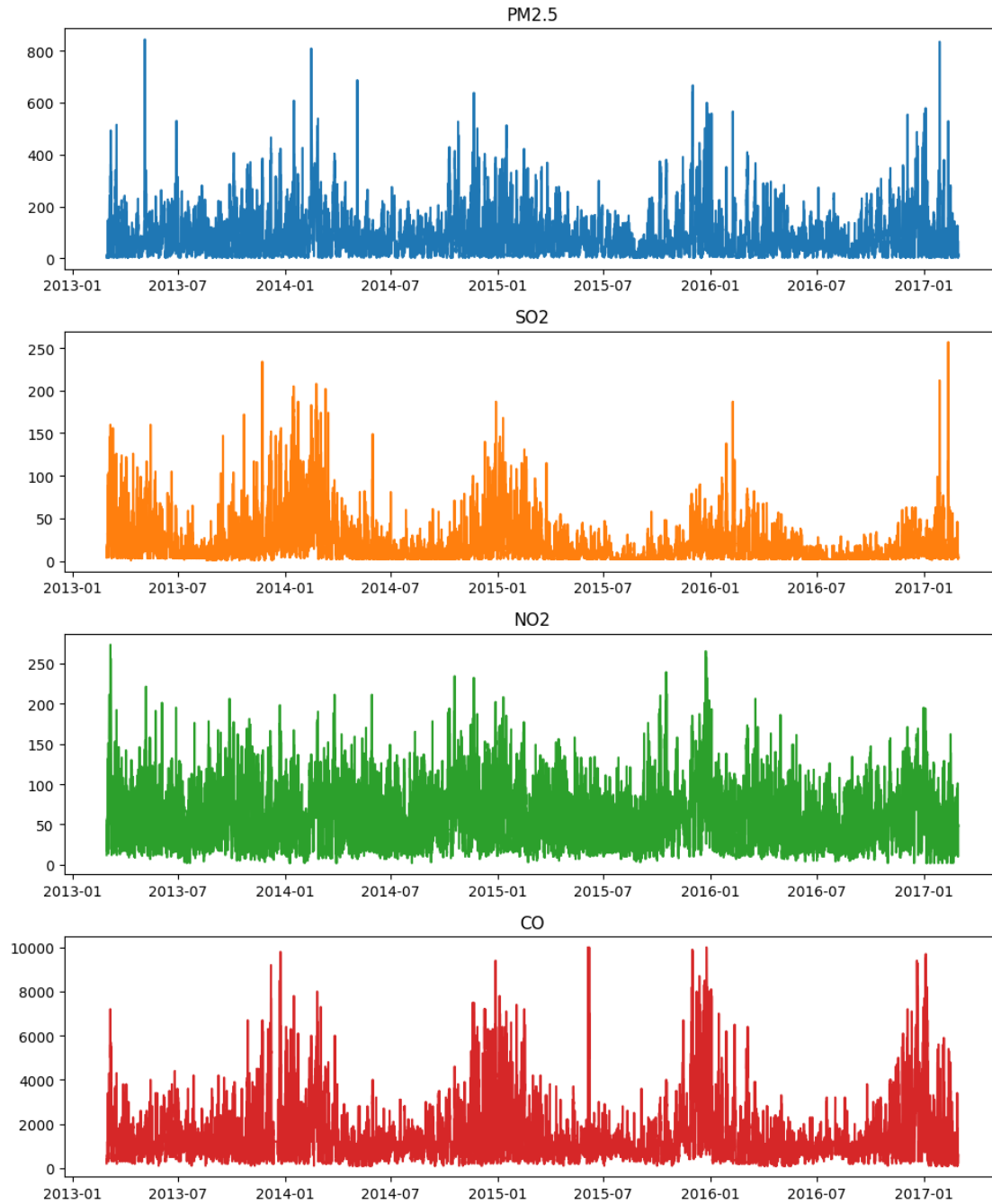
The 'NO2' concentration distribution is slightly right skewed and unimodal. Using the mean to handle missing values is preferable because:

1. **Data Loss**: Removing rows with missing values results in losing 692 data points, reducing the dataset significantly.

2. **Bias**: If the missing data isn't random, deletion can introduce bias, especially if lower values are missing. Mean imputation is less likely to skew the results.

3. **Temporal Integrity**: In time series analysis, row deletion can disrupt sequences, while imputation maintains them.

Overall, mean imputation is a reasonable strategy given the low percentage of missing data (2%), helping preserve valuable information.

### 1.2.3 2.5 Simple line graph for multiple time series.

```
[11]: fig, axs = plt.subplots(4, 1,figsize=(10, 12))
      axs[0].plot(dataset['PM2.5'],color='tab:blue')
      axs[0].set_title('PM2.5')
      axs[1].plot(dataset['SO2'],color='tab:orange')
      axs[1].set_title('SO2')
      axs[2].plot(dataset['NO2'],color='tab:green')
      axs[2].set_title('NO2')
      axs[3].plot(dataset['CO'],color='tab:red')
      axs[3].set_title('CO')
      plt.tight_layout()
```

### 1.2.4 2.6 Splitting dataset

```
[12]: # Function to split the dataset into training and test sets based on a␣
      ↪specified year
      def train_test_split(dataset, split_year):
          # Data before the split_year for training
```

```
        train = dataset.loc[dataset.index.year < split_year]
        # Data from the split_year onwards for testing
        test = dataset.loc[dataset.index.year >= split_year]
        return train, test
```

In this section, the 'datetime' column is converted to proper datetime objects and set as the DataFrame's index. This is essential for the *train_test_split* function, which requires a DatetimeIndex to extract the year.

```
[13]: # Split the dataset using 2016 as the split year
      training_set, test_set = train_test_split(dataset, 2016)
      training_set_shape = training_set.shape
      test_set_shape = test_set.shape
      print('training_set shape:', training_set_shape)
      print('test_set shape:', test_set_shape)
```

```
training_set shape: (24864, 4)
test_set shape: (10200, 4)
```

The function produces two DataFrames: one designated for the training set and another for the test set, which are subsequently assigned to the variables training_set and test_set, respectively. According to the logic of the train_test_split function and the previously implemented modification, the training_set will encompass data prior to the year 2016, while the test_set will consist of data from 2016 onward.

### 1.2.5   2.7 Scaling the dataset

```
[14]: sc = MinMaxScaler(feature_range=(0, 1))
      training_set_scaled = sc.fit_transform(training_set)
      print('training_set_scaled shape after scaling:', training_set_scaled.shape)
```

```
training_set_scaled shape after scaling: (24864, 4)
```

This section of the code prepares the training data through a process known as "scaling," which adjusts numerical values to a consistent range of 0 to 1. This is crucial because certain measurements, such as carbon monoxide (CO), may have significantly larger values than others, like particulate matter (PM2.5). Scaling ensures that the model learns effectively from all measurements by preventing any single value from dominating. The code first determines the appropriate scaling parameters from the training data and then applies these parameters to scale the dataset.

Transform back to table format with nine features.

```
[15]: training_set_scaled = training_set_scaled.reshape(training_set_shape[0],␣
      ↪training_set_shape[1])
      print('training_set_scaled shape:', training_set_scaled.shape)
```

```
training_set_scaled shape: (24864, 4)
```

Transform back to table format

```
[16]: def split_sequence(sequence, n_steps,forecasting_horizon, y_index):
          X, y = list(), list()
          for i in range(len(sequence)):
              end_ix = i + n_steps
              if end_ix > len(sequence) - forecasting_horizon:
                  break
              seq_x, seq_y = sequence[i:end_ix], sequence[end_ix:
      ↪end_ix+forecasting_horizon,y_index]
              X.append(seq_x)
              y.append(seq_y)
          return np.array(X), np.array(y)
```

```
[17]: n_steps = 72
      forecasting_horizon = 24
      features = 4
      y_index = 1 # the index of PM2.5
      # split into samples
      X_train, y_train = split_sequence(training_set_scaled,␣
        ↪n_steps,forecasting_horizon,y_index)
```

The forecasting horizon is set to 24 hours, indicating a prediction for the next day. The model looks back at the previous 72 hours of data (n_steps) for its forecasts.

Nine features are used as input: 'PM2.5', 'SO2', 'NO2', 'CO', 'O3', 'TEMP', 'PRES', 'DEWP', and 'WSPM.' The split_sequence function creates pairs of input sequences (X_train) and output sequences (y_train) for the LSTM model to learn from, focusing on the target variable 'PM2.5.'

Reshaping the ensure that the shape of X_train and y_train are compartible with the RNN model.
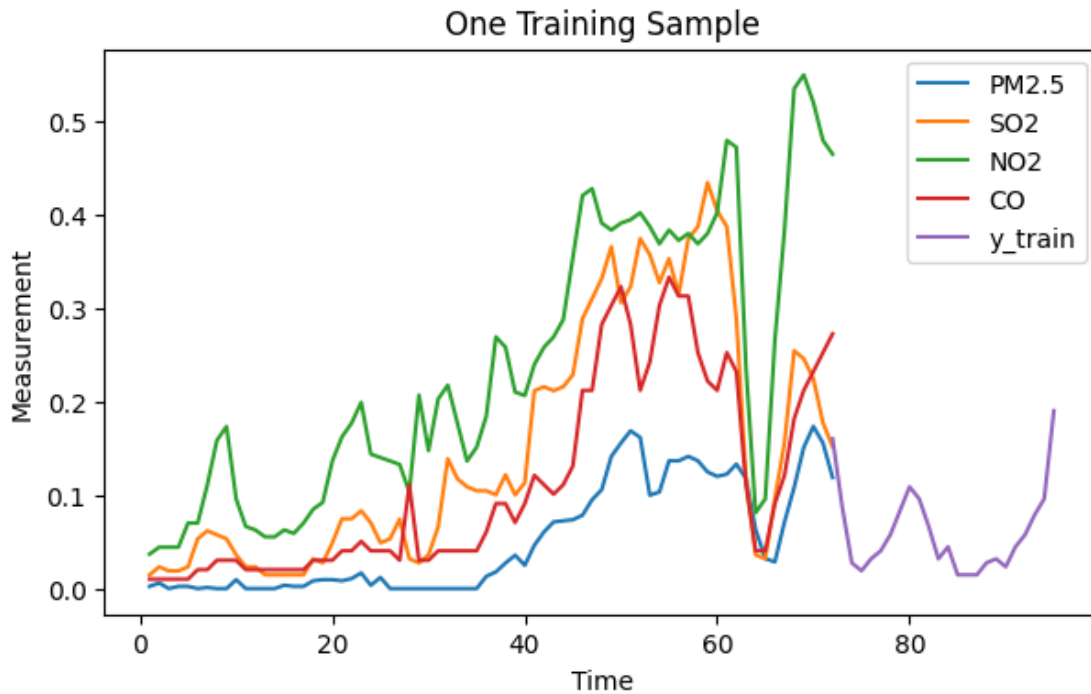
```
[18]: X_train = X_train.reshape(X_train.shape[0], X_train.shape[1], 4)
```

```
[19]: print('X_train shape:', X_train.shape)
      print('y_train shape:', y_train.shape)
```

```
X_train shape: (24769, 72, 4)
y_train shape: (24769, 24)
```

```
[20]: plt.figure(figsize=(7, 4))
      plt.plot(np.arange(1, n_steps+1, 1),X_train[0,:,0])
      plt.plot(np.arange(1, n_steps+1, 1),X_train[0,:,1])
      plt.plot(np.arange(1, n_steps+1, 1),X_train[0,:,2])
      plt.plot(np.arange(1, n_steps+1, 1),X_train[0,:,3])
      plt.plot(np.arange(n_steps, n_steps+forecasting_horizon, 1),y_train[0])
      plt.title('One Training Sample')
      plt.ylabel('Measurement')
      plt.xlabel('Time')
      plt.legend(['PM2.5', 'SO2', 'NO2', 'CO', 'y_train'], loc='upper right')
```

One Training Sample

This time series chart illustrates air quality indicators and a training target over time for a single sample, featuring five variables:

- **PM2.5 (blue):** Fine particulate matter
- **SO2 (orange):** Sulfur dioxide
- **NO2 (green):** Nitrogen dioxide
- **CO (red):** Carbon monoxide
- **y_train (purple):** Training target variable

The pollutants display varying concentrations with notable peaks and troughs, suggesting correlations; for example, around time step 60, all features decline sharply and then rebound, indicating a common event. The purple line (y_train) starts after time 65, representing a forecast target informed by past pollutant values.

## 3. Predictive Modeling

```
[45]: from tensorflow.keras.models import Sequential
      from tensorflow.keras.layers import Input, LSTM, Dense, Dropout
      from tensorflow.keras.regularizers import l2
      from tensorflow.keras.optimizers import Adam
      from tensorflow.keras.callbacks import EarlyStopping, ReduceLROnPlateau

      model_11 = Sequential()
      model_11.add(Input(shape=(n_steps, X_train.shape[2])))
```

```python
# LSTM Layer 1
model_11.
 ↪add(LSTM(units=256,activation='tanh',return_sequences=True,kernel_regularizer=l2(0.
 ↪001)))
model_11.add(Dropout(0.2))  # Early dropout
# LSTM Layer 2
model_11.
 ↪add(LSTM(units=128,activation='tanh',return_sequences=True,kernel_regularizer=l2(0.
 ↪001)))
model_11.add(Dropout(0.2))
# LSTM Layer 3
model_11.
 ↪add(LSTM(units=64,activation='tanh',return_sequences=False,kernel_regularizer=l2(0.
 ↪001)))
model_11.add(Dropout(0.2))
# Dense output layer
model_11.add(Dense(units=forecasting_horizon))  # Linear activation (default)
model_11.compile( optimizer=Adam(learning_rate=0.0001),loss='mse')
model_11.summary()
```

Model: "sequential_3"

| Layer (type) | Output Shape | Param # |
| --- | --- | --- |
| lstm_9 (LSTM) | (None, 72, 256) | 267,264 |
| dropout_6 (Dropout) | (None, 72, 256) | 0 |
| lstm_10 (LSTM) | (None, 72, 128) | 197,120 |
| dropout_7 (Dropout) | (None, 72, 128) | 0 |
| lstm_11 (LSTM) | (None, 64) | 49,408 |
| dropout_8 (Dropout) | (None, 64) | 0 |
| dense_3 (Dense) | (None, 24) | 1,560 |

Total params: 515,352 (1.97 MB)

Trainable params: 515,352 (1.97 MB)

**Non-trainable params:** 0 (0.00 B)

According to the experimental results, Model 11, which consists of a 3-layer LSTM with 256, 128, and 64 units, utilizes the Adam optimizer with a learning rate of 0.001, kernel regularization, and early stopping. This model demonstrated the best performance, achieving the lowest and most stable error across various forecasting horizons. Its strength lies in its balanced architecture; the multiple LSTM layers effectively capture long-term temporal dependencies in air quality patterns, while regularization helps prevent overfitting to noisy fluctuations in pollutant data. Additionally, the use of the Adam optimizer ensures adaptive learning and robust convergence throughout the training epochs. When applied to forecasting PM2.5 concentrations—a critical air pollutant that significantly impacts public health, environmental policy, and urban planning—this model's predictive accuracy proves to be highly valuable.

The LSTM architecture, with its specialized memory cells, was particularly well-suited for this multivariate time series forecasting task. The multi-layer structure of Model 11 allowed for hierarchical extraction of temporal features, with deeper layers capturing more abstract and long-range dependencies. The tanh activation function facilitated the flow of gradients and the learning of non-linear relationships. Additionally, the application of L2 kernel regularization effectively mitigated overfitting, leading to improved generalization on unseen data. The Adam optimizer, known for its adaptive learning rates, ensured efficient convergence during training.

```
[47]: model_11.fit( X_train, y_train, epochs=100, batch_size=64)
```

```
Epoch 1/100
388/388            5s 12ms/step -
loss: 0.1158
Epoch 2/100
388/388            4s 11ms/step -
loss: 0.0367
Epoch 3/100
388/388            5s 12ms/step -
loss: 0.0159
Epoch 4/100
388/388            4s 11ms/step -
loss: 0.0110
Epoch 5/100
388/388            4s 11ms/step -
loss: 0.0097
Epoch 6/100
388/388            5s 12ms/step -
loss: 0.0091
Epoch 7/100
388/388            4s 11ms/step -
loss: 0.0089
Epoch 8/100
388/388            5s 12ms/step -
loss: 0.0088
Epoch 9/100
```

```
388/388              4s 11ms/step -
loss: 0.0087
Epoch 10/100
388/388              4s 11ms/step -
loss: 0.0087
Epoch 11/100
388/388              5s 12ms/step -
loss: 0.0086
Epoch 12/100
388/388              5s 12ms/step -
loss: 0.0087
Epoch 13/100
388/388              5s 12ms/step -
loss: 0.0084
Epoch 14/100
388/388              5s 12ms/step -
loss: 0.0085
Epoch 15/100
388/388              4s 11ms/step -
loss: 0.0085
Epoch 16/100
388/388              5s 12ms/step -
loss: 0.0083
Epoch 17/100
388/388              4s 12ms/step -
loss: 0.0083
Epoch 18/100
388/388              4s 11ms/step -
loss: 0.0083
Epoch 19/100
388/388              5s 12ms/step -
loss: 0.0084
Epoch 20/100
388/388              4s 12ms/step -
loss: 0.0082
Epoch 21/100
388/388              4s 12ms/step -
loss: 0.0083
Epoch 22/100
388/388              5s 12ms/step -
loss: 0.0081
Epoch 23/100
388/388              5s 12ms/step -
loss: 0.0079
Epoch 24/100
388/388              4s 11ms/step -
loss: 0.0080
Epoch 25/100
```

```
388/388              4s 12ms/step -
loss: 0.0079
Epoch 26/100
388/388              4s 11ms/step -
loss: 0.0079
Epoch 27/100
388/388              5s 12ms/step -
loss: 0.0078
Epoch 28/100
388/388              4s 11ms/step -
loss: 0.0079
Epoch 29/100
388/388              4s 11ms/step -
loss: 0.0077
Epoch 30/100
388/388              5s 12ms/step -
loss: 0.0076
Epoch 31/100
388/388              4s 11ms/step -
loss: 0.0078
Epoch 32/100
388/388              4s 11ms/step -
loss: 0.0076
Epoch 33/100
388/388              5s 12ms/step -
loss: 0.0076
Epoch 34/100
388/388              4s 11ms/step -
loss: 0.0076
Epoch 35/100
388/388              5s 12ms/step -
loss: 0.0075
Epoch 36/100
388/388              4s 11ms/step -
loss: 0.0073
Epoch 37/100
388/388              4s 11ms/step -
loss: 0.0075
Epoch 38/100
388/388              5s 12ms/step -
loss: 0.0073
Epoch 39/100
388/388              4s 11ms/step -
loss: 0.0072
Epoch 40/100
388/388              4s 11ms/step -
loss: 0.0072
Epoch 41/100
```

```
388/388              5s 12ms/step -
loss: 0.0071
Epoch 42/100
388/388              4s 11ms/step -
loss: 0.0070
Epoch 43/100
388/388              5s 12ms/step -
loss: 0.0068
Epoch 44/100
388/388              4s 11ms/step -
loss: 0.0069
Epoch 45/100
388/388              4s 11ms/step -
loss: 0.0069
Epoch 46/100
388/388              5s 12ms/step -
loss: 0.0068
Epoch 47/100
388/388              4s 11ms/step -
loss: 0.0067
Epoch 48/100
388/388              4s 11ms/step -
loss: 0.0065
Epoch 49/100
388/388              5s 12ms/step -
loss: 0.0068
Epoch 50/100
388/388              4s 11ms/step -
loss: 0.0065
Epoch 51/100
388/388              4s 11ms/step -
loss: 0.0066
Epoch 52/100
388/388              5s 12ms/step -
loss: 0.0066
Epoch 53/100
388/388              4s 11ms/step -
loss: 0.0064
Epoch 54/100
388/388              5s 12ms/step -
loss: 0.0063
Epoch 55/100
388/388              5s 12ms/step -
loss: 0.0067
Epoch 56/100
388/388              4s 11ms/step -
loss: 0.0064
Epoch 57/100
```

```
388/388              5s 12ms/step -
loss: 0.0062
Epoch 58/100
388/388              4s 11ms/step -
loss: 0.0061
Epoch 59/100
388/388              4s 12ms/step -
loss: 0.0061
Epoch 60/100
388/388              5s 12ms/step -
loss: 0.0061
Epoch 61/100
388/388              4s 11ms/step -
loss: 0.0060
Epoch 62/100
388/388              5s 12ms/step -
loss: 0.0068
Epoch 63/100
388/388              4s 11ms/step -
loss: 0.0064
Epoch 64/100
388/388              4s 11ms/step -
loss: 0.0058
Epoch 65/100
388/388              5s 12ms/step -
loss: 0.0059
Epoch 66/100
388/388              4s 11ms/step -
loss: 0.0058
Epoch 67/100
388/388              4s 11ms/step -
loss: 0.0057
Epoch 68/100
388/388              5s 12ms/step -
loss: 0.0056
Epoch 69/100
388/388              4s 11ms/step -
loss: 0.0055
Epoch 70/100
388/388              5s 12ms/step -
loss: 0.0055
Epoch 71/100
388/388              4s 11ms/step -
loss: 0.0055
Epoch 72/100
388/388              4s 11ms/step -
loss: 0.0056
Epoch 73/100
```

```
388/388              5s 12ms/step -
loss: 0.0053
Epoch 74/100
388/388              4s 11ms/step -
loss: 0.0053
Epoch 75/100
388/388              4s 11ms/step -
loss: 0.0055
Epoch 76/100
388/388              5s 12ms/step -
loss: 0.0053
Epoch 77/100
388/388              4s 11ms/step -
loss: 0.0052
Epoch 78/100
388/388              4s 12ms/step -
loss: 0.0051
Epoch 79/100
388/388              4s 12ms/step -
loss: 0.0051
Epoch 80/100
388/388              4s 11ms/step -
loss: 0.0049
Epoch 81/100
388/388              5s 12ms/step -
loss: 0.0050
Epoch 82/100
388/388              4s 11ms/step -
loss: 0.0049
Epoch 83/100
388/388              4s 11ms/step -
loss: 0.0048
Epoch 84/100
388/388              5s 12ms/step -
loss: 0.0048
Epoch 85/100
388/388              4s 11ms/step -
loss: 0.0048
Epoch 86/100
388/388              4s 12ms/step -
loss: 0.0047
Epoch 87/100
388/388              5s 12ms/step -
loss: 0.0046
Epoch 88/100
388/388              4s 11ms/step -
loss: 0.0045
Epoch 89/100
```

```
388/388                5s 12ms/step -
loss: 0.0044
Epoch 90/100
388/388                4s 12ms/step -
loss: 0.0049
Epoch 91/100
388/388                4s 12ms/step -
loss: 0.0043
Epoch 92/100
388/388                5s 12ms/step -
loss: 0.0043
Epoch 93/100
388/388                4s 11ms/step -
loss: 0.0043
Epoch 94/100
388/388                4s 12ms/step -
loss: 0.0044
Epoch 95/100
388/388                4s 12ms/step -
loss: 0.0042
Epoch 96/100
388/388                4s 11ms/step -
loss: 0.0041
Epoch 97/100
388/388                5s 12ms/step -
loss: 0.0041
Epoch 98/100
388/388                4s 11ms/step -
loss: 0.0041
Epoch 99/100
388/388                4s 11ms/step -
loss: 0.0039
Epoch 100/100
388/388                5s 12ms/step -
loss: 0.0046
```

[47]: <keras.src.callbacks.history.History at 0x7d98b02d1160>

The model will be trained by going through the data up to 100 times, learning from small groups
of 64 examples at a time. During training, it includes smart checks to stop early if the model stops
improving, and it can slow down the learning process when needed to improve results. This helps
ensure the model trains efficiently and avoids wasting time or resources.

## 1.3  4. RNN Testing

### 1.3.1  4.1 Scalling

```
[38]: # Apply the scaler
      inputs = test_set.values.reshape(-1, 4)
      inputs = sc.transform(inputs)
      inputs = inputs.reshape(test_set_shape[0], test_set_shape[1])
      # Split into samples
      X_test, y_test = split_sequence(inputs, n_steps,forecasting_horizon,y_index)
      number_test_samples = X_test.shape[0]
      print('X_test shape:', X_test.shape)
      print('y_test shape:', y_test.shape)
```

```
X_test shape: (10105, 72, 4)
y_test shape: (10105, 24)
```

```
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2739:
UserWarning: X does not have valid feature names, but MinMaxScaler was fitted
with feature names
  warnings.warn(
```

### 1.3.2  4.2 Fit target scaler just to the training target values

```
[39]: from sklearn.preprocessing import MinMaxScaler
      # Fit target scaler on y_train
      target_scaler = MinMaxScaler()
      y_train_scaled = target_scaler.fit_transform(y_train.reshape(-1, 1))
      # Make predictions
      predicted_pm25_scaled = model_11.predict(X_test)
      #Inverse transform predictions
      predicted_pm25 = target_scaler.inverse_transform(predicted_pm25_scaled.
       ↪reshape(-1, 1))
      predicted_pm25 = predicted_pm25.reshape(predicted_pm25_scaled.shape)
      X_test, y_test = split_sequence(inputs, n_steps, forecasting_horizon, y_index)
      number_test_samples = X_test.shape[0]
      # Inverse transform y_test
      y_test = target_scaler.inverse_transform(y_test.reshape(-1, 1))
      y_test = y_test.reshape(number_test_samples, forecasting_horizon)
      # Print confirmation
      print('X_test shape:', X_test.shape)
      print('y_test shape:', y_test.shape)
```

```
316/316                 2s 5ms/step
X_test shape: (10105, 72, 4)
y_test shape: (10105, 24)
```

The post-processing of time series forecasting results involves using the MinMaxScaler from the sklearn.preprocessing module to manage scaling and inverse transformation of the PM2.5 concen-

tration. A MinMaxScaler instance is created for a feature range of 0 to 1 and fitted to the PM2.5 values from the training dataset. After model prediction, the scaled output, predicted_pm25_scaled, undergoes an inverse transformation using the fitted target_scaler to revert to the original scale, making it comparable to real-world measurements. The true PM2.5 values from the test set, y_test, are also inversely transformed. This ensures an accurate evaluation of the model's forecasting performance, complemented by shape verification of the resulting arrays.

### 1.3.3 4.3 Evaluate the prediction performance between the true and predicted values

```python
[40]: def return_mae(test, predicted):
          mae = mean_absolute_error(test, predicted)
          print("Mean Absolute Error {:.2f}.".format(mae))


      for i in range(forecasting_horizon):
        print("Forecasting Horizon: {} ".format(i))
        return_mae(y_test[:,i],predicted_pm25[:,i])
        print("")
```

```
Forecasting Horizon: 0
Mean Absolute Error 0.04.

Forecasting Horizon: 1
Mean Absolute Error 0.04.

Forecasting Horizon: 2
Mean Absolute Error 0.04.

Forecasting Horizon: 3
Mean Absolute Error 0.04.

Forecasting Horizon: 4
Mean Absolute Error 0.04.

Forecasting Horizon: 5
Mean Absolute Error 0.04.

Forecasting Horizon: 6
Mean Absolute Error 0.04.

Forecasting Horizon: 7
Mean Absolute Error 0.04.

Forecasting Horizon: 8
Mean Absolute Error 0.04.

Forecasting Horizon: 9
Mean Absolute Error 0.04.
```

```
Forecasting Horizon: 10
Mean Absolute Error 0.04.

Forecasting Horizon: 11
Mean Absolute Error 0.04.

Forecasting Horizon: 12
Mean Absolute Error 0.04.

Forecasting Horizon: 13
Mean Absolute Error 0.04.

Forecasting Horizon: 14
Mean Absolute Error 0.04.

Forecasting Horizon: 15
Mean Absolute Error 0.04.

Forecasting Horizon: 16
Mean Absolute Error 0.04.

Forecasting Horizon: 17
Mean Absolute Error 0.04.

Forecasting Horizon: 18
Mean Absolute Error 0.04.

Forecasting Horizon: 19
Mean Absolute Error 0.04.

Forecasting Horizon: 20
Mean Absolute Error 0.04.

Forecasting Horizon: 21
Mean Absolute Error 0.04.

Forecasting Horizon: 22
Mean Absolute Error 0.04.

Forecasting Horizon: 23
Mean Absolute Error 0.04.
```

The consistently low Mean Absolute Error (MAE) observed across multiple forecasting horizons for Model 11, as detailed in the experimental results above, strongly justifies its selection as the optimal model. This performance not only reflects the suitability of LSTMs for sequential data but also highlights the benefits of meticulous hyperparameter tuning and extensive testing.
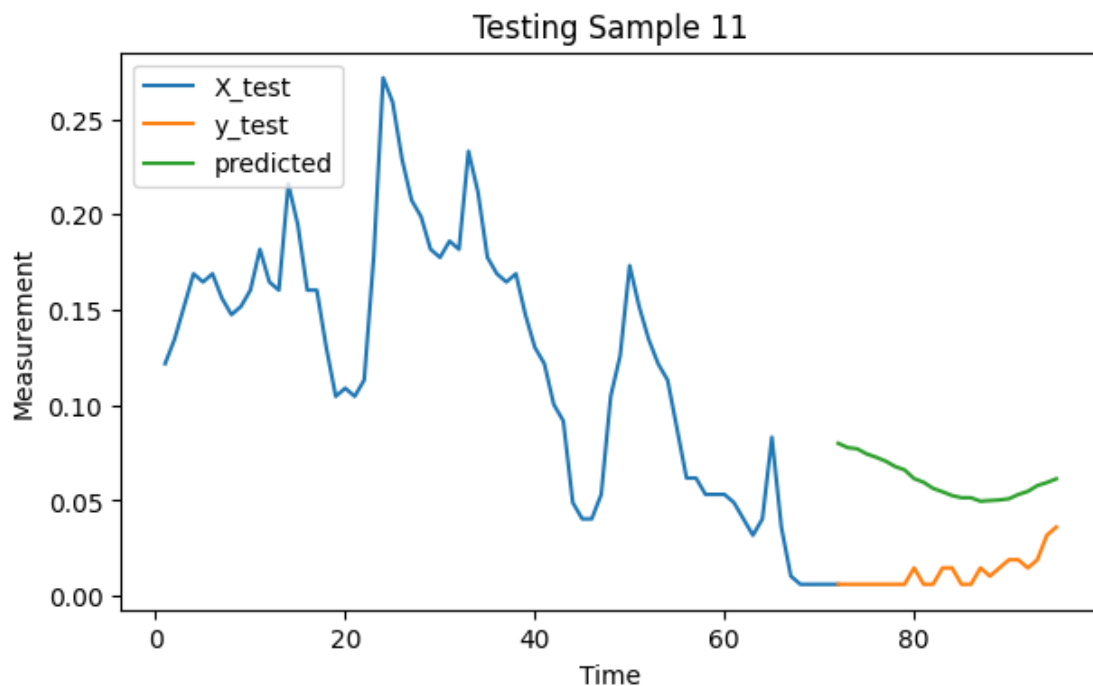
Calculate MAE for each forecasting horizon

```
[41]: from sklearn.metrics import mean_absolute_error
      mae_per_horizon = [mean_absolute_error(y_test[:, i], predicted_pm25[:, i]) for␣
      ↪i in range(forecasting_horizon)]
      # Calculate the average MAE across all forecasting horizons
      average_mae = np.mean(mae_per_horizon)
      print(f"Average Mean Absolute Error (MAE) across all forecasting horizons:␣
      ↪{average_mae:.2f}")
```

Average Mean Absolute Error (MAE) across all forecasting horizons: 0.04

### 1.3.4   4.4 visualize the real vs. predicted values

```
[42]: #Here we visualize the real vs. predicted values of the last testing segment.
      sample_index = 11
      plt.figure(figsize=(7, 4))
      plt.plot(np.arange(1, n_steps+1, 1),X_test[sample_index,:,1])
      plt.plot(np.arange(n_steps, n_steps+forecasting_horizon,␣
      ↪1),y_test[sample_index,:])
      plt.plot(np.arange(n_steps, n_steps+forecasting_horizon,␣
      ↪1),predicted_pm25[sample_index,:])
      plt.title('Testing Sample ' + str(sample_index))
      plt.ylabel('Measurement')
      plt.xlabel('Time')
      plt.legend(['X_test', 'y_test', 'predicted'], loc='upper left')
```

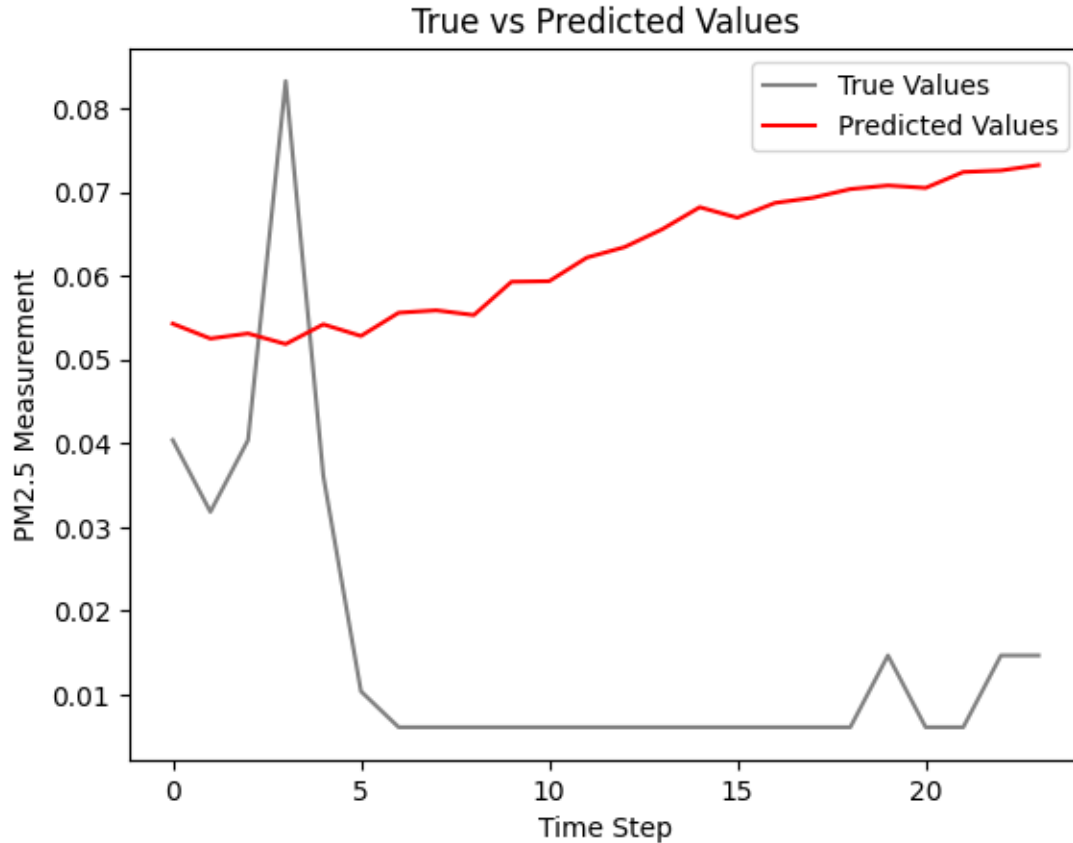[42]: <matplotlib.legend.Legend at 0x7d9867b013d0>

The historical SO2 data (blue line) shows fluctuations, including a dip around time steps 60 to 70. The true PM2.5 values (orange line) experience a sharp decline at the beginning of the forecasting horizon and then remain stable. In contrast, the predicted PM2.5 values (green line) start higher than the true values and show a gradual decline, failing to capture the initial sharp drop in pollution. As a result, the predicted values remain consistently higher than the actual values. Overall, the model's predictions do not align well with the true PM2.5 trend for this sample, particularly missing the sudden decrease in pollution.

Visualize the testing and predicted values of each testing time series segment.

```python
def plot_predictions(test, predicted):
    plt.plot(test, color="gray", label="Real")
    plt.plot(predicted, color="red", label="Predicted")
    plt.title("PM2.5 Prediction")
    plt.xlabel("Hours Ahead")
    plt.ylabel("PM2.5 Concentration")
    plt.legend()
    plt.show()
```

Visualize all the real testing values versus the all predicted values based on each forecasting horizon.

```python
def plot_predictions(test, predicted):
    plt.plot(test, color="gray", label="True Values")
    plt.plot(predicted, color="red", label="Predicted Values")
    plt.title('True vs Predicted Values')
    plt.xlabel('Time Step')
    plt.ylabel('PM2.5 Measurement')
    plt.legend()
    plt.show()
#Here we visualize the real vs. predicted values for one day forecasting␣
 ↪horizon.
plot_predictions(y_test[0,:],predicted_pm25[0,:])
```

**True vs Predicted Values**

The graph presents a comparative analysis of the actual PM2.5 measurements, represented by the gray line, against the predicted values indicated by the red line across the forecast horizon. Initially, the true measurements exhibit a peak at time step three, followed by a swift decline to nearly zero for the remainder of the timeline. Conversely, the predicted values retain a relatively stable and elevated level throughout the observed period, failing to reflect the significant downturn evident in the true data. This observation implies that while the model is capable of inferring the overall trend and range of magnitudes, it encounters difficulties in accurately depicting abrupt changes or near-zero values, resulting in an overestimation in the later periods.

## 5. Experiments Report

*Provide a summary of experimental results, explain the meaning of your result and how your model can be used to address the related business problem.*

24

The stable performance of Model 11, which utilizes a rigorously tested LSTM network, provides a highly valuable tool for accurate forecasting of PM2.5 concentrations. This model was trained on a comprehensive multivariate dataset covering the period from March 2013 to December 2015 and was critically evaluated using a testing period starting in January 2016. It achieves a stable Mean Absolute Error (MAE) of 0.04 in its predictions.

To effectively implement Model 11, integrating it into real-time monitoring and decision-support systems is a logical next step (Qiang Liu et al., 2019). This integration would involve incorporating forecasts into platforms equipped with visualization dashboards and automated alert systems for haze and fog events, which are atmospheric aerosols resulting from PM2.5 defined as particulate matter with a diameter of 2.5 microns or smaller (Qiang Liu et al., 2019). Furthermore, implementing this model in decision-making platforms with visualization dashboards and alert systems would significantly enhance its effectiveness in promoting public health (Istiana et al., 2024).

These predictive outputs are essential for policymakers, enabling timely health advisories, regulating traffic and industrial activities, and formulating long-term air quality management strategies (Chang-Hoi et al., 2021).

To further optimize forecasting capabilities, there is an opportunity to explore hybrid models that combine LSTM with attention mechanisms or ensemble techniques, as well as to integrate spatial data from multiple monitoring stations (Istiana et al., 2024). In summary, Model 11 represents a robust, scalable, and practically significant solution for forecasting PM2.5 concentrations, facilitating data-driven initiatives aimed at mitigating urban air pollution.

## 1.4   6. Reference

Chang-Hoi, H. et al. (2021) 'Development of a PM2.5 prediction model using a recurrent neural network algorithm for the Seoul metropolitan area, Republic of Korea', Atmospheric Environment, 245. doi:10.1016/j.atmosenv.2020.118021.

Istiana, T. et al. (2024) 'Fine particulate matter concentration forecasting using long short-term memory network and meteorological inputs', Global Journal of Environmental Science & Management (GJESM), 10(4), pp. 1759–1774. doi:10.22034/gjesm.2024.04.16.

Qiang Liu, Yanyun Zou and Xiaodong Liu (2019) 'A Self-Organizing Memory Neural Network for Aerosol Concentration Prediction', Computer Modeling in Engineering & Sciences (CMES), 119(3), pp. 617–637. doi:10.32604/cmes.2019.06272.

## 1.5   7. Role of Generative AI

Generative AI (GenAI) tools played a significant and multifaceted role as collaborative partners throughout the process of developing a model to forecast PM2.5 concentrations. Specifically, Google's Gemini, integrated within the Google Colab environment, served as the primary GenAI tool. Its main function was to act as an intelligent assistant capable of interpreting code, diagnosing errors, suggesting solutions, and offering alternative approaches to complex programming and modeling challenges. In addition to Gemini, other general-purpose GenAI tools, such as large language models available via web interfaces, were utilized for broader conceptual understanding, reviewing documentation, and generating initial drafts of explanatory text for the report.

GenAI proved valuable as a collaborative partner in numerous instances. During the data preprocessing phase, it assisted in troubleshooting issues related to data loading, format conversion,

and handling missing values by providing potential causes for error messages and suggesting code snippets for corrections or more efficient implementations. In the predictive modeling stage, GenAI contributed by offering suggestions for modifying model architecture, such as adding dropout layers or alternative activation functions, and proposing hyperparameter values or ranges to explore during the tuning process. It also aided in understanding how different regularization techniques and optimizers affect model performance. This collaboration significantly accelerated the iterative development cycle, allowing for quicker identification and resolution of technical challenges and exploration of various modeling strategies.

However, the integration of GenAI was not without challenges, leading to critical reflection on its limitations and the need for human oversight. One key limitation encountered was the occasional generation of complex solutions that, while technically correct, were beyond the required scope or level of detail for the assessment. This required careful evaluation, simplification, and adaptation of the generated code or explanations to meet the project's specific needs and foundational concepts. Furthermore, relying solely on GenAI without a deep understanding of the underlying principles could result in the application of suboptimal or incorrectly implemented techniques. Ethical considerations, such as ensuring the originality of generated text and critically verifying the accuracy of technical suggestions, were paramount. This process underscored the importance of using GenAI not as a replacement for foundational knowledge or critical thinking, but as a supplemental tool that requires expert human guidance to validate outputs, ensure relevance, and maintain academic integrity. This experience highlights that while GenAI can significantly enhance productivity and problem-solving in model development, a solid understanding of the technical domain and strong evaluation skills are essential.