# *LSTM MULTIVARIATE TIME SERIES Baseline Model*

## Chapter 4.1.1

```
1 from google.colab import drive
2 drive.mount('/content/gdrive')
```

    Mounted at /content/gdrive

```
 1 # Import librries
 2 import pandas as pd
 3 from pandas import DataFrame
 4 from pandas import concat
 5 from sklearn.preprocessing import LabelEncoder
 6 from sklearn.preprocessing import MinMaxScaler
 7 import glob #for maps
 8 import plotly.graph_objects as go
 9 import plotly.offline
10 import matplotlib.pyplot as plt
11 from datetime import datetime
12 import warnings
13 import itertools
14 import numpy as np
15 from statsmodels.tsa.stattools import adfuller
16 from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_erro
17 from math import sqrt
18 warnings.filterwarnings("ignore")
19 import plotly.express as px
20 import seaborn as sns
21 import pandas as pd
22 import plotly.graph_objects as go
23 import matplotlib.pyplot as plt
24 %matplotlib inline
25
26
27 from sklearn.model_selection import train_test_split
28 from tensorflow.keras.layers import Dense, Dropout, LSTM
29 import keras.models
30 import tensorflow
31 from keras.models import Sequential
32 from tensorflow.keras.layers import LSTM
33 from tensorflow.keras.layers import Dense
```
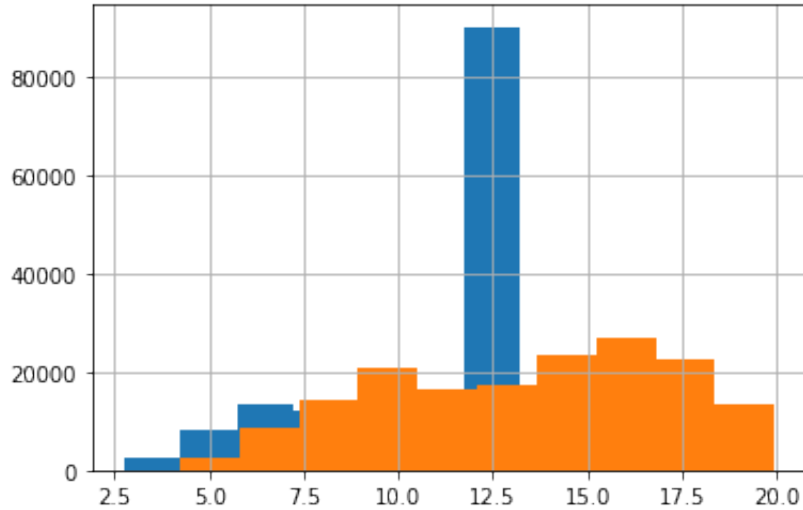
## ▾ Data Reading and Additional Processing

```
1 fin_data = pd.read_csv('/content/gdrive/MyDrive/Sent/final_data_May20_Sept22
```

```
1 fin_data.head(5)
```

| | latitude | longitute | Year | Month | Day | Hour | nitrogendioxide_tropospheric |
|---|---|---|---|---|---|---|---|
| **0** | 52.157350 | 13.716069 | 2020.0 | 5.0 | 1.0 | 11.0 | 1.45 |
| **1** | 52.192528 | 13.806399 | 2020.0 | 5.0 | 1.0 | 11.0 | -6.22 |
| **2** | 52.227226 | 13.895948 | 2020.0 | 5.0 | 1.0 | 11.0 | 5.75 |
| **3** | 52.167213 | 13.589217 | 2020.0 | 5.0 | 1.0 | 11.0 | 1.58 |
| **4** | 52.202910 | 13.680410 | 2020.0 | 5.0 | 1.0 | 11.0 | 1.68 |

```
1 print(final_data.shape)
2 print(fin_data['pm25'].hist())
3 print(fin_data['pm10'].hist())
```

```
(220179, 15)
AxesSubplot(0.125,0.125;0.775x0.755)
AxesSubplot(0.125,0.125;0.775x0.755)
```
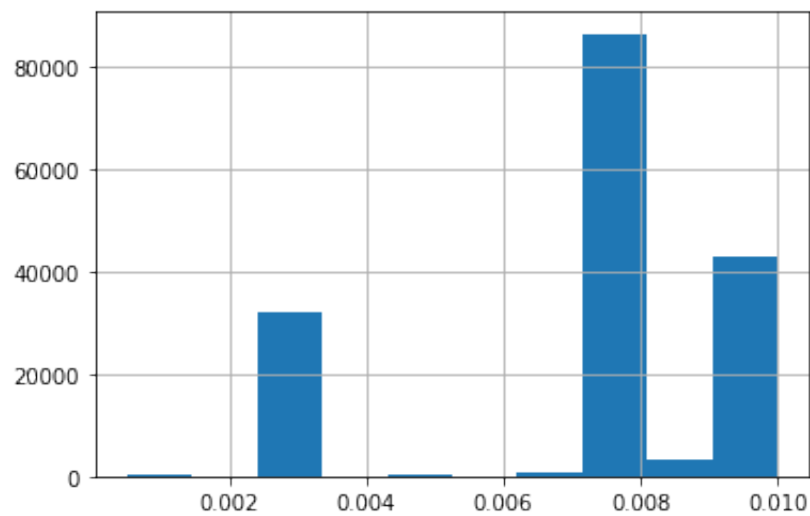


```
1 fin_data.shape
```

```
(166395, 17)
```

```
1 fin_data['qa_value'].unique() #quality value
```

```
array([0.0074, 0.01  , 0.0011, 0.0067, 0.009 , 0.0033, 0.0015, 0.0045,
       0.003 , 0.0014, 0.0005, 0.001 , 0.0073])
```

```
1 fin_data['qa_value'].hist()
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f0742f034f0>



```
1 fin_data = fin_data.loc[(fin_data['qa_value'] > 0.004)]
```

```
1 fin_data.shape
```

(133820, 17)

```
1 fin_data.sort_values(by='date', inplace=True)
2 fin_data.head()
```

| | latitude | longitute | Year | Month | Day | Hour | nitrogendioxide_tropospheri |
|---|---|---|---|---|---|---|---|
| 0 | 52.157350 | 13.716069 | 2020.0 | 5.0 | 1.0 | 11.0 | 1.4 |
| 95 | 52.548504 | 13.041852 | 2020.0 | 5.0 | 1.0 | 11.0 | 1.2 |
| 96 | 52.585503 | 13.135360 | 2020.0 | 5.0 | 1.0 | 11.0 | -4.1 |
| 97 | 52.621986 | 13.228033 | 2020.0 | 5.0 | 1.0 | 11.0 | 8.2 |
| 98 | 52.657963 | 13.319888 | 2020.0 | 5.0 | 1.0 | 11.0 | 1.6 |

```
1 data_LSTM_test1=fin_data.drop(columns=['date', 'utc', 'pm25','tm5_tropopause_
```

```
1 data_LSTM_test1.head(5)
```

| | nitrogendioxide_tropospheric_column | nitrogendioxide_tropospheric_colum |
|---|---|---|
| 0 | 1.459978e-05 | |
| 95 | 1.250596e-05 | |
| 96 | -4.180619e-07 | |
| 97 | 8.273078e-06 | |
| 98 | 1.682993e-05 | |

# ▾ Feature Engineering for LSTM

**Referenece:**

The model scrip was adopted from the following webpage:

https://machinelearningmastery.com/multivariate-time-series-forecasting-lstms-keras/

https://machinelearningmastery.com/time-series-prediction-lstm-recurrent-neural-networks-python-keras/

```
1
2 # convert series to supervised learning
3 def series_to_supervised(data, n_in=1, n_out=1, dropnan=True):
4     n_vars = 1 if type(data) is list else data.shape[1]
5     df = DataFrame(data)
6     cols, names = list(), list()
7     # input sequence (t-n, ... t-1)
8     for i in range(n_in, 0, -1):
9         cols.append(df.shift(i))
10        names += [('var%d(t-%d)' % (j+1, i)) for j in range(n_vars)]
11    # forecast sequence (t, t+1, ... t+n)
12    for i in range(0, n_out):
13        cols.append(df.shift(-i))
14        if i == 0:
15            names += [('var%d(t)' % (j+1)) for j in range(n_vars)]
16        else:
17            names += [('var%d(t+%d)' % (j+1, i)) for j in range(n_vars)]
18    # put it all together
19    agg = concat(cols, axis=1)
20    agg.columns = names
21    # drop rows with NaN values
22    if dropnan:
23        agg.dropna(inplace=True)
24    return agg
25
26 # load dataset
27 values = data_LSTM_test1.values
28 # integer encode direction
29 encoder = LabelEncoder()
30 #values[:,8] = encoder.fit_transform(values[:,8])
31 values
32 # ensure all data is float
33 values = values.astype('float32')
34 # normalize features
35 scaler = MinMaxScaler(feature_range=(0, 1))
36 scaled = scaler.fit_transform(values)
37 # frame as supervised learning
38 reframed = series_to_supervised(scaled, 1, 1)
39 reframed.drop(reframed.columns[[0,6, 8,9,10,11]], axis=1, inplace=True)
40 print(reframed.head())
```

```
      var2(t-1)  var3(t-1)  var4(t-1)  var5(t-1)  var6(t-1)   var2(t)
   1   0.018191   0.018763   0.243149   0.357339   0.129479  0.008212
   2   0.008212   0.008622   0.266456   0.371345   0.129479  0.008405
   3   0.008405   0.009252   0.232822   0.350280   0.129479  0.009880
   4   0.009880   0.010570   0.258888   0.361394   0.129479  0.012215
   5   0.012215   0.012459   0.232868   0.347342   0.129479  0.009161
```

```
1 #Train/Val/Test
2 values = reframed.values
3 x, x_test, y, y_test = train_test_split(values[:, :-1],values[:,-1],test_size
4 x_train, x_val, train_y, val_y = train_test_split(x,y,test_size = 0.2,train_s
```

```
1 # reshape input to be 3D [samples, timesteps, features]
2 train_X = x_train.reshape((x_train.shape[0], 1, x_train.shape[1]))
3 val_X = x_val.reshape((x_val.shape[0], 1, x_val.shape[1]))
4 test_X = x_test.reshape((x_test.shape[0], 1, x_test.shape[1]))
```

```
1 print(train_X.shape, train_y.shape, val_X.shape, val_y.shape, test_X.shape, y
```

```
(96349, 1, 5) (96349,) (24088, 1, 5) (24088,) (13382, 1, 5) (13382,)
```

## ▾ LSTM Baseline

```
 1 %%time
 2 model = Sequential()
 3 model.add(LSTM(500, input_shape=(train_X.shape[1], train_X.shape[2])))
 4 model.add(Dropout(0.2))
 5 model.add(Dense(1))
 6 model.compile(loss='mae', optimizer='adam') #mean absolute error
 7 history = model.fit(train_X, train_y, epochs=15, batch_size=64, validation_da
 8
 9 #plot history
10 plt.figure(figsize=(15, 10))
11
12 pyplot.plot(history.history['loss'], label='train')
13 pyplot.plot(history.history['val_loss'], label='val')
14 pyplot.legend()
15 pyplot.show()
```
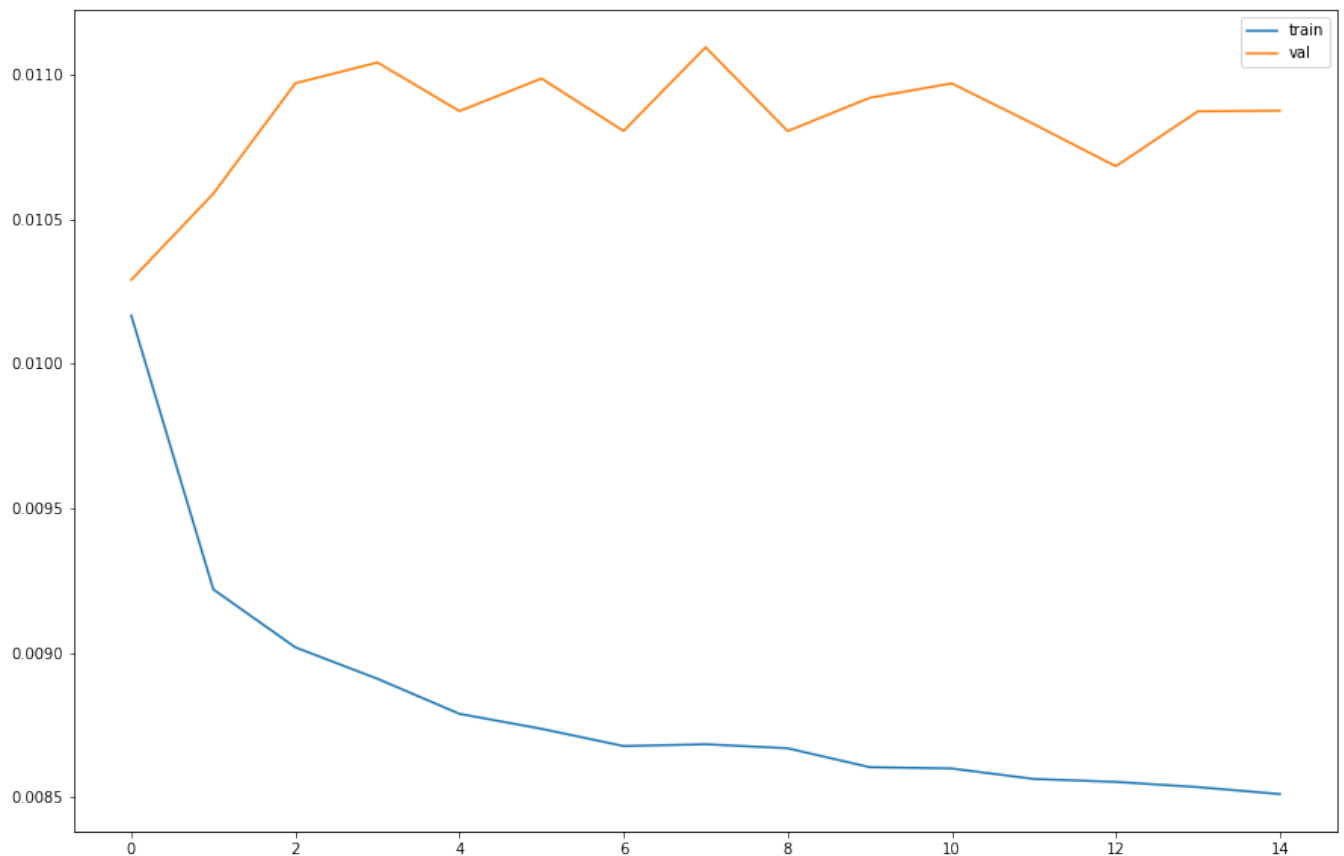
```
Epoch 1/15
1506/1506 - 39s - loss: 0.0102 - val_loss: 0.0103 - 39s/epoch - 26ms/step
Epoch 2/15
1506/1506 - 29s - loss: 0.0092 - val_loss: 0.0106 - 29s/epoch - 19ms/step
Epoch 3/15
1506/1506 - 29s - loss: 0.0090 - val_loss: 0.0110 - 29s/epoch - 19ms/step
Epoch 4/15
1506/1506 - 29s - loss: 0.0089 - val_loss: 0.0110 - 29s/epoch - 19ms/step
Epoch 5/15
1506/1506 - 29s - loss: 0.0088 - val_loss: 0.0109 - 29s/epoch - 19ms/step
Epoch 6/15
1506/1506 - 28s - loss: 0.0087 - val_loss: 0.0110 - 28s/epoch - 18ms/step
Epoch 7/15
1506/1506 - 28s - loss: 0.0087 - val_loss: 0.0108 - 28s/epoch - 18ms/step
```
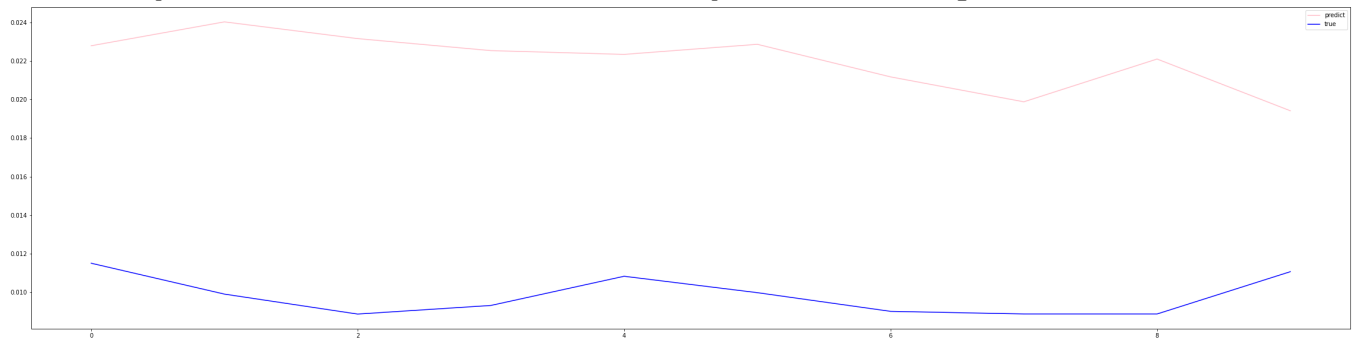
```
1506/1506 - 28s - loss: 0.0087 - val_loss: 0.0108 - 28s/epoch - 18ms/step
Epoch 8/15
1506/1506 - 28s - loss: 0.0087 - val_loss: 0.0111 - 28s/epoch - 19ms/step
Epoch 9/15
1506/1506 - 28s - loss: 0.0087 - val_loss: 0.0108 - 28s/epoch - 18ms/step
Epoch 10/15
1506/1506 - 28s - loss: 0.0086 - val_loss: 0.0109 - 28s/epoch - 18ms/step
Epoch 11/15
1506/1506 - 31s - loss: 0.0086 - val_loss: 0.0110 - 31s/epoch - 21ms/step
Epoch 12/15
1506/1506 - 28s - loss: 0.0086 - val_loss: 0.0108 - 28s/epoch - 19ms/step
Epoch 13/15
1506/1506 - 28s - loss: 0.0086 - val_loss: 0.0107 - 28s/epoch - 19ms/step
Epoch 14/15
1506/1506 - 28s - loss: 0.0085 - val_loss: 0.0109 - 28s/epoch - 19ms/step
Epoch 15/15
1506/1506 - 29s - loss: 0.0085 - val_loss: 0.0109 - 29s/epoch - 19ms/step
```



```
CPU times: user 11min 55s, sys: 27.3 s, total: 12min 22s
Wall time: 7min 18s
```
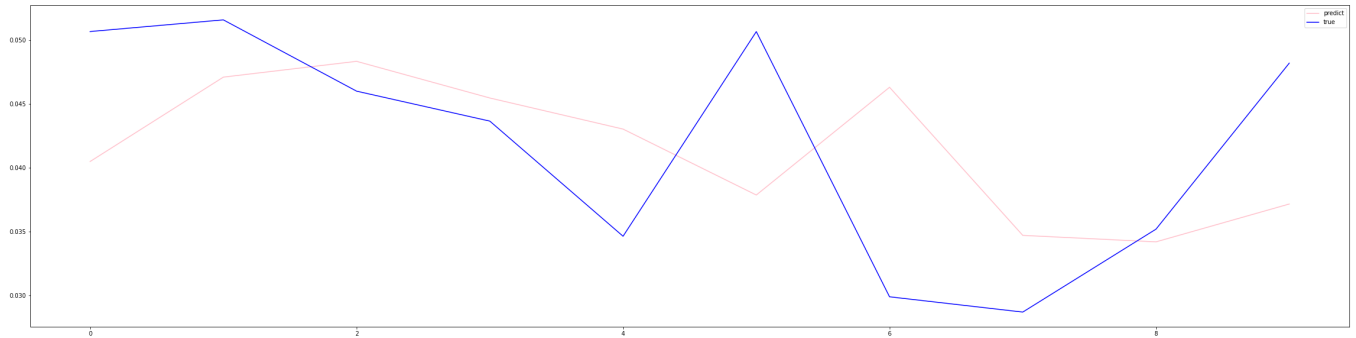
```
1 #Validation set (8 days forecast)
2
3 yhat = model.predict(val_X)
4 plt.figure(figsize=(40, 10))
5 pyplot.plot(yhat[-10:], label='predict', color = 'pink')
6 pyplot.plot(val_y[-10:], label='true', color = 'blue')
7 pyplot.legend()
8
9 pyplot.show()
```

753/753 [==============================] - 4s 5ms/step

```
1 #Train set (8 days forecast)
2
3 yhat = model.predict(train_X)
4 plt.figure(figsize=(40, 10))
5 pyplot.plot(yhat[-10:], label='predict', color = 'pink')
6 pyplot.plot(train_y[-10:], label='true', color = 'blue')
7 pyplot.legend()
8
9 pyplot.show()
```

```
3011/3011 [==============================] - 15s 5ms/step
```

```
1  #Test set (8 days forecast)
2
3  yhat = model.predict(test_X)
4  plt.figure(figsize=(40, 10))
5  pyplot.plot(yhat[-10:], label='predict', color = 'pink')
6  pyplot.plot(y_test[-10:], label='true', color = 'blue')
7  pyplot.legend()
8
9  pyplot.show()
```

419/419 [==============================] - 2s 6ms/step