

Air Pollution Forecasting using LSTM

Xiaoyang Chen, Julia Jin

George Washington University

DATS 6203: Machine Learning II

Professor: Amir Jafari

December 12, 2022

Abstract

In this report, we discuss models that use LSTMs to predict air quality. In section 1, we talk about the background of our topic and introduce the LSTM model. In section 2, we discuss our dataset description and how we preprocess the data. Also in section 2, we visualized the data, looked at the graph of different characteristic data changing according to time and the relationship diagram between each characteristic data.

In section 3, we discussed the multivariate model and univariate model established in the project. In section 4, we discuss the data results from multivariate model and univariate model and compare and summarize in section 5.

1 Introduction

This project aims to predict air pollution using the Air Quality dataset from Kaggle. This is a dataset that reports on the weather and the level of pollution each hour for five years at the US embassy in Beijing, China. The data includes date and time, our target which is the pollution called PM2.5 concentration, and the weather information as features, including dew point, temperature, pressure, wind direction, wind speed and the cumulative number of hours of snow and rain. In the subsection 1.1, we provide background information about air pollution. In the subsection 1.2, the report outline is presented.

1.1 Background

Air pollution is one of the most concerns for urban areas. There are countries around the world have built a variety of sensing devices for monitoring PM2.5 concentrations. There were also many studies have been constructed to predict and forecast various air pollution [1].

Due to its proven track record of success with time-series data, a Long Short-Term Memory (LSTM) Recurrent Neural Network (RNN) model was chosen to perform the task of air quality forecasting. The data provides hourly average concentrations of various air pollutants around the U.S. embassy in Beijing over five years. We chose the Long Short-Term Memory (LSTM) Recurrent Neural Network (RNN) model for the air quality prediction task due to its successful track record with time-series data. This project compares the difference between the multivariate input model and the single variable input

model in the LSTM model by establishing a multivariate input model and a single variable input model.

1.2 LSTM

According to the Wikipedia's definition, long short-term memory (LSTM) is an artificial neural network used in the fields of artificial intelligence and deep learning [2]. LSTM units are a building unit for layers of a (RNN). A RNN composed of LSTM units is often called an LSTM network.

The difference between LSTM and traditional RNN neural networks is that each neuron in LSTM is a memory cell. The LSTM links the previous data information to the current neurons.

Each neuron contains three gates as shown in Figure 1: input gate, forget gate, and output gate. Using the internal gate, the LSTM can solve the problem of long-term dependence of the data [1].

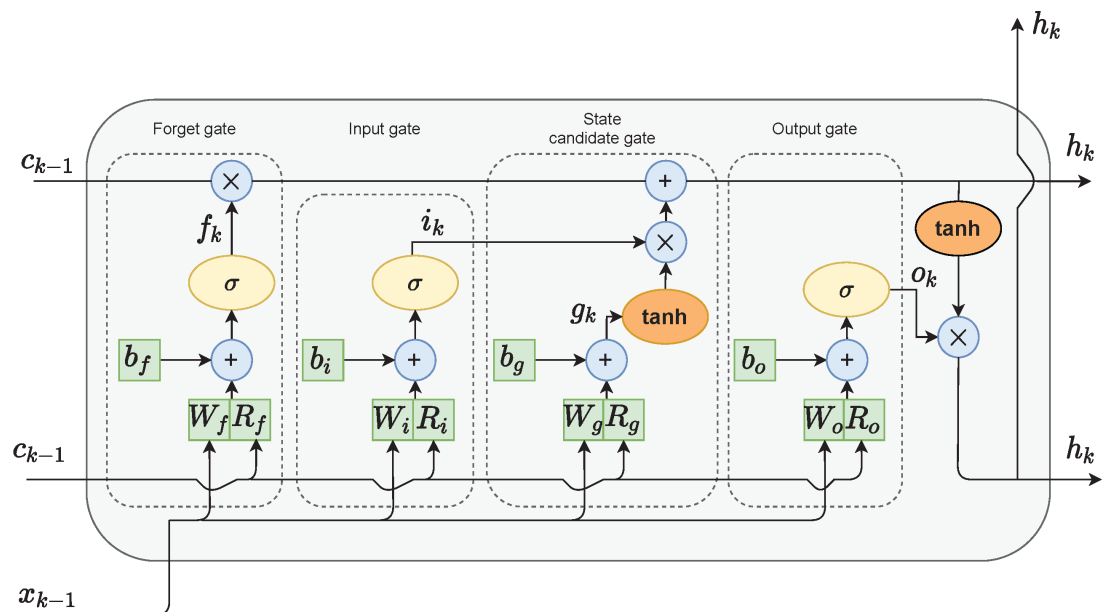


Figure 1. LSTM Architecture

1) Forget gate

One of the main properties of the LSTM is to memorize and recognize the information coming inside the network and to discard the information which is not required to the network to learn the data and predictions. This gate is responsible for this feature of the LSTM.

2) Input gate

Input gate helps in deciding the importance of the information by updating the cell state. where the forget gate helps in the elimination of the information from the network input gate decides the measure of the importance of the information and helps the forget function in elimination of the not important information and other layers to learn the information which is important for making predictions.

3) Output gate

It is the last gate of the circuit that helps in deciding the next hidden state of the network in which information goes through the sigmoid function. Updated cell from the cell state goes to the tanh function then it gets multiplied by the sigmoid function of the output state. Which helps the hidden state to carry the information [3].

In a long short-term memory (LSTM) network, the activation function used for the gates (input, forget, and output gates) is the sigmoid function, while the activation function used for the cell state and output is the hyperbolic tangent (tanh) function.

The sigmoid function is used for the gates because it has a range of 0 to 1, which makes it well-suited for binary classification tasks. It also has a smooth derivative, which makes it easy to compute the gradient during training.

The tanh function is used for the cell state and output because it has a range

of -1 to 1, which makes it well-suited for representing numerical values. It also has a non-linear shape, which allows the LSTM network to capture complex patterns in the data.

In summary, the choice of activation functions in an LSTM network is a trade-off between range, smoothness, and non-linearity, and the sigmoid and tanh functions are well-suited for the tasks performed by the gates and cell state/output, respectively.

2 Data Mining

2.1 Dataset Description

The dataset we used in this project is the air quality dataset which was reported on the weather and the level of pollution each hour for five years at the US embassy in Beijing, China.

Target

pm2.5: PM2.5 concentration ($\mu\text{g}/\text{m}^3$)

Features

dew: Dew point

temperature: Temperature around the embassy (F)

pressure: Air pressure

wind_dir: Combined wind direction

wind_speed: Cumulated wind speed (m/s)

snow: Cumulated hours of snow

rain: Cumulated hours of rain

2.2 Dataset Preprocessing

- Visualization of features over time

Figure 2 is the information about different features in dataset.

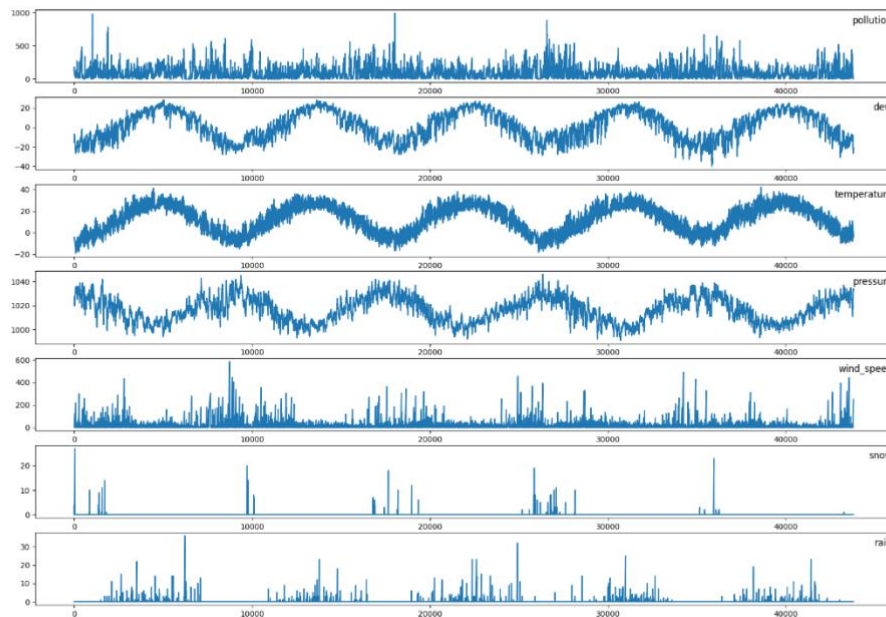


Figure 2: Columns Information

In Figure 2, there are 8 columns in our dataset, but only 7 are shown. “wind_dir” is a string. We will explain later how we convert it to a numerical value.

All features display high seasonality. “dew”, “temperature” and “pressure” also display high cyclicity.

- Correlation Matrix

As shown in Figure 3, we compute the Pearson correlation coefficient of the features with the built-in function `corr()` in pandas and display the correlation matrix with matplotlib.

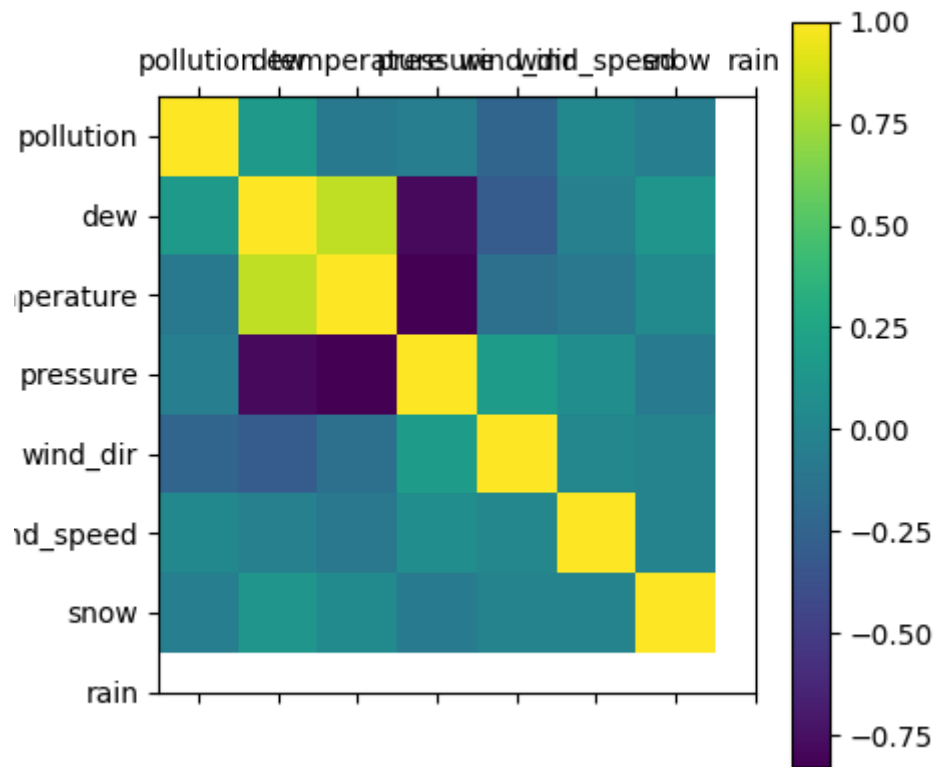


Figure 3: Correlation Matrix

- Data Preprocessing

In addition to doing data visualization, we also perform some preprocessing on the dataset.

As mentioned above, we convert “wind_dir”, a string, to a numerical value. We use the function `LabelEncoder` from `sklearn.preprocessing` to change its type to float [4].

We used `MinMaxScaler` to normalize the data. The benefit of normalizing the data is that it can increase the convergence speed of the iterative solution and improve the accuracy of the iterative solution.

Our dataset is rather small for LSTM. Before machine learning can be used, our time series forecasting problem must be reframed as a supervised

learning problem, which is, from one sequence to pairs of input and output sequences. Then we use the function `series_to_supervised` as shown in Figure 4. The function is defined with default parameters so that it will construct a data frame with $t - 1$ as X and t as y [5].

```
def series_to_supervised(data, n_in=1, n_out=1, dropnan=True):
    n_vars = 1 if type(data) is list else data.shape[1]

    df = pd.DataFrame(data)
    cols, names = list(), list()

    # input sequence (t-n, ... t-1)
    for i in range(n_in, 0, -1):
        cols.append(df.shift(i))
        names += [('var%d(t-%d)' % (j + 1, i)) for j in range(n_vars)]

    # forecast sequence (t, t+1, ... t+n)
    for i in range(0, n_out):
        cols.append(df.shift(-i))
        if i == 0:
            names += [('var%d(t)' % (j + 1)) for j in range(n_vars)]
        else:
            names += [('var%d(t+%d)' % (j + 1, i)) for j in range(n_vars)]

    # put it all together
    agg = pd.concat(cols, axis=1)
    agg.columns = names
    # drop rows with NaN values
    if dropnan:
        agg.dropna(inplace=True)

    return agg
```

Figure 4: Convert Time Series to Supervised Problem

3 Model Description

3.1 Multivariate Model

To build this model, we divided our dataset into two parts. The dataset from 2010 to 2013 is used as the training part, and the dataset from 2014 is used

as the validation part.

For our model, we use batch normalization which is used to improve the performance and stability of neural networks. This normalization helps to reduce the internal covariate shift, which is the change in the distribution of the inputs to a layer caused by the change in the parameters of the previous layer. It can improve the convergence rate and overall performance of the network. Our mini batch size is 128.

As shown in Figure 5, we also use dropout layer. It is used to prevent overfitting in our networks. It works by randomly dropping out, or setting to 0, a certain number of outputs from a layer during training. It reduces the number of connections between the layers, which in turn reduces the complexity of the model and prevents overfitting.

In tensorflow.keras library, there is a function called “LSTM” which can be easily applied [6]. The activation function is already set as “tanh” and the recurrent activation is set as “sigmoid” as we talked in previous section.

```
Model: "sequential"
```

Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 256)	271360
dense (Dense)	(None, 64)	16448
dropout (Dropout)	(None, 64)	0
batch_normalization (Batch Normalization)	(None, 64)	256
dense_1 (Dense)	(None, 1)	65

```

Total params: 288,129
Trainable params: 288,001
Non-trainable params: 128

```

Figure 5: Multivariate Model Summary

3.2 Univariate Model

The univariate model differs from the multivariate model in that multivariate model reads in all the features. In the data preparation for the univariate model, we only read the data of the pollution feature and preprocess it.

Again, we divided dataset into two parts. The dataset from 2010 to 2013 is used as the training part, and the dataset from 2014 is used as the validation set.

As shown in the Figure 6, the parameters in LSTM layer are not the same as the ones in the multivariate model.

```

Model: "sequential"
-----
Layer (type)                Output Shape              Param #
-----
lstm (LSTM)                  (None, 256)              264192
dense (Dense)                (None, 64)               16448
dropout (Dropout)           (None, 64)               0
batch_normalization (BatchN (None, 64)              256
ormalization)
dense_1 (Dense)              (None, 1)                65
-----
Total params: 280,961
Trainable params: 280,833
Non-trainable params: 128
-----

```

Figure 6: Univariate Model Summary

4 Numerical Results

4.1 Multivariate Model

Here, we start with the results of the multivariate model. The training loss vs. the validation loss plot for the multivariate model is shown in Figure 7.

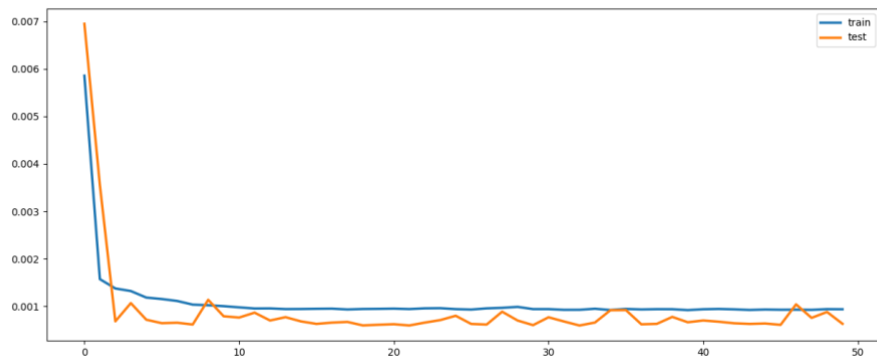


Figure 7: Training loss vs validation loss for Multivariate Model

We have our first increase in validation loss at epoch 2. Then at epoch 8, we have our second increase. After that, we can see some minimal increase in loss, but the trend seems stable overall. From epoch 24 we can see some more noticeable ups and downs, but overall, the validation loss oscillates within a very small range.

In Figure 8, we compare the multivariate model output with the actual test data.

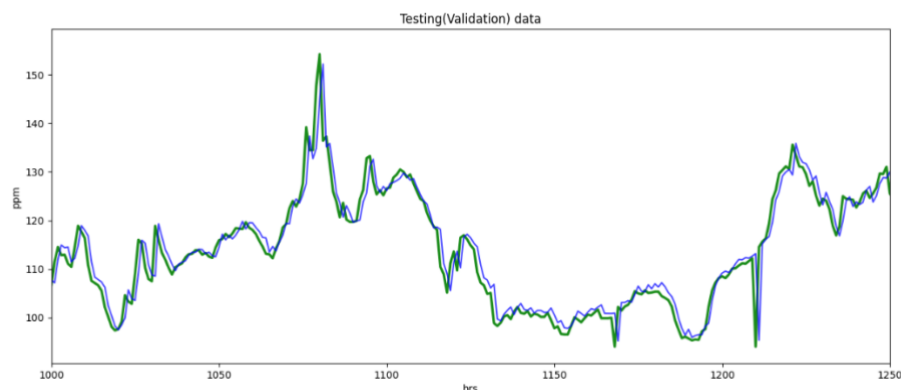


Figure 8: Test Result for Multivariate Model

The multivariate model predicts the air quality closely. Most of the time, the predicted data is not much different from the actual data. The most noticeable discrepancy happens around the interval of 1120 to 1180.

We choose Root Mean Squared Error (RMSE) as our evaluation metric. For multivariate model, the RMSE is 2.32079.

4.2 Univariate Model

The training loss vs. the validation loss plot for the univariate model is shown in Figure 9.

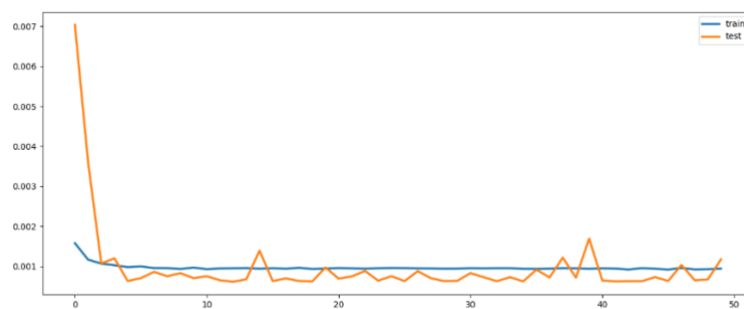


Figure 9: Training loss vs test loss for Univariate Model

In univariate model, the first increase in error is also at epoch 2. Then at epoch 6, we have our second increase. At epoch 19, we have a sudden huge increase. After that we notice ups and downs in the next 20 epochs, until another sudden huge increase at epoch 39. The validation loss does not converge and continues to fluctuate. It varies in a much bigger range compared to the multivariate model.

In Figure 10, we compare the univariate model output and the actual test data.

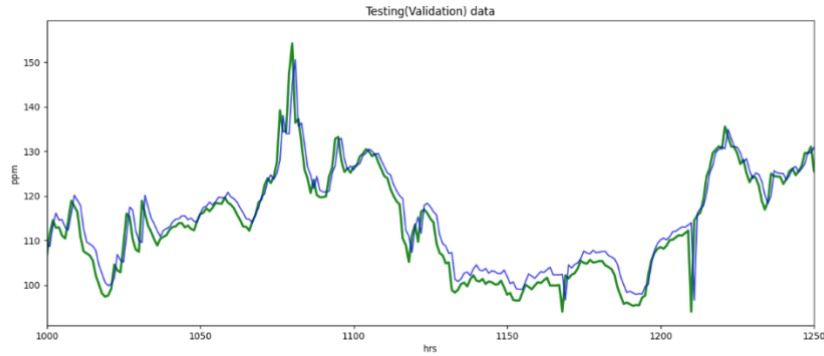


Figure 10: Test result for Univariate Model

It shows that the univariate model predicts the air quality well. Although the trend of the predicted data is roughly consistent with the trend of the actual data, we have more noticeable discrepancies, around interval 1000 to 1025, interval 1040 to 1065 and interval 1120 to 1180.

The RMSE value is 3.15849, higher than the multivariate model.

5 Conclusion

In this project, we implement LSTM network to predict the air quality situation. We construct a multivariate model and a univariate model and compared the performance of the two models. With a closer match to the actual dataset and lower RMSE, the multivariate model is the best model we select for this problem. It can capture more information about the data and make more accurate predictions.

The biggest problem of both models is that the validation loss is not steadily decreasing as we expected. Increase in loss is most likely due to overfitting, although we have taken actions to prevent that from happening as we mentioned in the model description. Having said that, our dataset is relatively small for an LSTM problem. If the variance is large enough, overfitting can still

happen.

Another point worth mentioning is that in the interval 1120 to 1180 in the test set, both models' outputs deviate significantly from the actual data. We will have to further investigate in this range and see if there are other factors that lead to this difference.

Future work could be targeted at using a deeper and wider LSTM network, using a learning rate schedule, and using a different optimization algorithm. Moreover, we also can use GRU to compare the performance between GRU and LSTM.

Reference

- [1] Tsai, Zeng, Y.-R., & Chang, Y.-S. (2018). Air Pollution Forecasting Using RNN with LSTM. 2018 16TH IEEE INT CONF ON DEPENDABLE, AUTONOM AND SECURE COMP, 16TH IEEE INT CONF ON PERVAS INTELLIGENCE AND COMP, 4TH IEEE INT CONF ON BIG DATA INTELLIGENCE AND COMP, 3RD IEEE CYBER SCI AND TECHNOL CONGRESS (DASC/PICOM/DATACOM/CYBERSCITECH), 1074–1079. <https://doi.org/10.1109/DASC/PiCom/DataCom/CyberSciTec.2018.00178>
- [2] Wikipedia, "Long short-term memory," 2022. [Online] Available at: https://en.wikipedia.org/wiki/Long_short-term_memory, Accessed: 2022.
- [3] "A Complete Guide to LSTM Architecture and its Use in Text Classification," 2022. [Online] Available at: <https://analyticsindiamag.com/a-complete-guide-to-lstm-architecture-and-its-use-in-text-classification/>, Accessed: 2022.
- [4] "sklearn.preprocessing.LabelEncoder." [Online] Available at: <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.LabelEncoder.html>, Accessed: 2022.
- [5] "How to Convert a Time Series to a Supervised Learning Problem in Python," 2022. [Online] Available at: <https://machinelearningmastery.com/convert-time-series-supervised-learning-problem-python/>, Accessed: 2022.
- [6] "tf.keras.layers.LSTM" [Online] Available at: https://www.tensorflow.org/api_docs/python/tf/keras/layers/LSTM, Accessed: 2022.
- [7] Air Pollution Forecasting. Kaggle. <https://www.kaggle.com/datasets/rupakroy/lstm-datasets-multivariate-univariate>

Appendix

Code for Multivariate Model

```
import warnings
warnings.filterwarnings('ignore')

import tensorflow as tf
import os
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from tensorflow.keras.models import Sequential, load_model
from tensorflow.keras.layers import *
from sklearn.preprocessing import LabelEncoder, MinMaxScaler
from sklearn.metrics import mean_squared_error as mse
from tensorflow.keras.callbacks import ModelCheckpoint
from tensorflow.keras.losses import MeanSquaredError
from tensorflow.keras.metrics import RootMeanSquaredError
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.layers import LSTM, Dense, Dropout,
BatchNormalization
from tensorflow.keras.models import Sequential

print('Imports Complete')
# ----- read data -----
--
df = pd.read_csv('LSTM-Multivariate_pollution.csv')
# print(df.head())
# print(df.shape)
# print(df.info())
# print(df.describe())
# ----- data pre-processing -----
-----
col_names = ['pollution', 'dew', 'temperature', 'pressure',
'wind_dir', 'wind_speed', 'snow', 'rain']
df.drop_duplicates(inplace=True)
df.dropna(inplace=True)
# print(df.shape)

df.index = pd.to_datetime(df['date'], format='%Y.%m.%d %H:%M:%S')
# print(df.head())
df.drop('date', axis=1, inplace=True)
df.columns = col_names
features = df.values
# print(df.head())

columns = [0, 1, 2, 3, 5, 6, 7]
plt.figure(figsize=(20,14))
for i, c in enumerate(columns, 1):
    plt.subplot(len(columns), 1, i)
    plt.plot(features[:, c])
    plt.title(df.columns[c], y=0.75, loc="right")
# plt.show()

plt.matshow(df.corr())
plt.xticks(range(len(col_names)), col_names)
plt.yticks(range(len(col_names)), col_names)
```

```

plt.colorbar()
# plt.show()

# print(df["wind_dir"].unique())
wind_dir_encoder = LabelEncoder()
df["wind_dir"] = wind_dir_encoder.fit_transform(df["wind_dir"])
df["wind_dir"] = df["wind_dir"].astype(float)
# print(df.head())

values = df.values
target = df['pollution']
# plt.plot(target)
# plt.show()

# How to Convert a Time Series to a Supervised Learning Problem in Python
# https://machinelearningmastery.com/convert-time-series-supervised-learning-problem-python/
# The function is defined with default parameters so that if you call it with just your data, it will construct a DataFrame with t-1 as X and t as y
def series_to_supervised(data, n_in=1, n_out=1, dropnan=True):
    n_vars = 1 if type(data) is list else data.shape[1]

    df = pd.DataFrame(data)
    cols, names = list(), list()

    # input sequence (t-n, ... t-1)
    for i in range(n_in, 0, -1):
        cols.append(df.shift(i))
        names += [('var%d(t-%d)' % (j + 1, i)) for j in range(n_vars)]

    # forecast sequence (t, t+1, ... t+n)
    for i in range(0, n_out):
        cols.append(df.shift(-i))
        if i == 0:
            names += [('var%d(t)' % (j + 1)) for j in range(n_vars)]
        else:
            names += [('var%d(t+%d)' % (j + 1, i)) for j in range(n_vars)]

    # put it all together
    agg = pd.concat(cols, axis=1)
    agg.columns = names
    # drop rows with NaN values
    if dropnan:
        agg.dropna(inplace=True)

    return agg

scaler = MinMaxScaler(feature_range=(0, 1))
scaled_dataset = scaler.fit_transform(values)
# print(scaled_dataset.shape[1])
reframed = series_to_supervised(scaled_dataset, 1, 1)
# print(reframed.head())
# print(reframed.shape)

reframed.drop(reframed.columns[[9, 10, 11, 12, 13, 14, 15]], axis=1, inplace=True)
# print(reframed.head())

```

```

# print(reframed.shape)

values = reframed.values
# First 4 years data
n_train_hours = 365 * 24 * 4

train = values[:n_train_hours, :]
test = values[n_train_hours:, :]

# split into input and outputs
train_X, train_y = train[:, :-1], train[:, -1]
test_X, test_y = test[:, :-1], test[:, -1]

# reshape input to be 3D :- (no.of samples, no.of timesteps, no.of
features)
train_X = train_X.reshape((train_X.shape[0], 1, train_X.shape[1]))
test_X = test_X.reshape((test_X.shape[0], 1, test_X.shape[1]))
# print(train_X.shape, train_y.shape, test_X.shape, test_y.shape)

#----- Design Model -----
---
model = Sequential()
model.add(LSTM(256, input_shape=(train_X.shape[1], train_X.shape[2])))
model.add(Dense(64))
model.add(Dropout(0.25))
model.add(BatchNormalization())
model.add(Dense(1))

model.summary()

model.compile(loss='mse', optimizer='adam')

history = model.fit(train_X, train_y, epochs=50, batch_size=128,
validation_data=(test_X, test_y))

plt.figure(figsize=(15,6))
plt.plot(history.history['loss'], label='train', linewidth = 2.5)
plt.plot(history.history['val_loss'], label='test', linewidth = 2.5)
plt.legend()
plt.show()

# ----- Pridict -----
--
prediction = model.predict(test_X)
prediction = prediction.ravel()

ture_test = test[:, 8]

poll = np.array(df["pollution"])

meanop = poll.mean()
stdop = poll.std()

ture_test = ture_test * stdop + meanop
prediction = prediction * stdop + meanop

plt.figure(figsize=(15, 6))
plt.xlim([1000, 1250])
plt.ylabel("ppm")
plt.xlabel("hrs")

```

```
plt.plot(ture_test, c="g", alpha=0.90, linewidth=2.5)
plt.plot(prediction, c="b", alpha=0.75)
plt.title("Testing(Validation) data")
plt.show()

rmse = np.sqrt(mse(ture_test, prediction))
print("Test(Validation) RMSE =", rmse)
```

Code for Univariate Model

```
import warnings
warnings.filterwarnings('ignore')

import tensorflow as tf
import os
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from tensorflow.keras.models import Sequential, load_model
from tensorflow.keras.layers import *
from sklearn.preprocessing import LabelEncoder, MinMaxScaler
from sklearn.metrics import mean_squared_error as mse
from tensorflow.keras.callbacks import ModelCheckpoint
from tensorflow.keras.losses import MeanSquaredError
from tensorflow.keras.metrics import RootMeanSquaredError
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.layers import LSTM, Dense, Dropout,
BatchNormalization
from tensorflow.keras.models import Sequential

print('Imports Complete')
# ----- read data -----
--
df = pd.read_csv('LSTM-Multivariate_pollution.csv')
# print(df.head())
# print(df.shape)
# print(df.info())
# print(df.describe())
# ----- data pre-processing -----
-----
col_names = ['pollution', 'dew', 'temperature', 'pressure',
'wind_dir', 'wind_speed', 'snow', 'rain']
df.drop_duplicates(inplace=True)
df.dropna(inplace=True)
# print(df.shape)

df.index = pd.to_datetime(df['date'], format='%Y.%m.%d %H:%M:%S')
# print(df.head())
df.drop('date', axis=1, inplace=True)
df.columns = col_names
df.drop(['dew', 'temperature', 'pressure', 'wind_dir', 'wind_speed',
'snow', 'rain'], axis=1, inplace=True)
values = df.values
features = df.values
# print(df.head())

plt.figure(figsize=(20, 14))
plt.plot(df['pollution'])
plt.title('pollution', y=0.75, loc="right")
plt.show()

col_names = df.columns.tolist()
print(col_names)

# How to Convert a Time Series to a Supervised Learning Problem in
Python
# https://machinelearningmastery.com/convert-time-series-supervised-learning-problem-python/
```

```

# The function is defined with default parameters so that if you call
it with just your data, it will construct a DataFrame with t-1
# as X and t as y
def series_to_supervised(data, n_in=1, n_out=1, dropnan=True):
    n_vars = 1 if type(data) is list else data.shape[1]

    df = pd.DataFrame(data)
    cols, names = list(), list()

    # input sequence (t-n, ... t-1)
    for i in range(n_in, 0, -1):
        cols.append(df.shift(i))
        names += [('var%d(t-%d)' % (j + 1, i)) for j in range(n_vars)]

    # forecast sequence (t, t+1, ... t+n)
    for i in range(0, n_out):
        cols.append(df.shift(-i))
        if i == 0:
            names += [('var%d(t)' % (j + 1)) for j in range(n_vars)]
        else:
            names += [('var%d(t+%d)' % (j + 1, i)) for j in
range(n_vars)]

    # put it all together
    agg = pd.concat(cols, axis=1)
    agg.columns = names
    # drop rows with NaN values
    if dropnan:
        agg.dropna(inplace=True)

    return agg

scaler = MinMaxScaler(feature_range=(0, 1))
scaled_dataset = scaler.fit_transform(values)
# print(scaled_dataset.shape[1])
reframed = series_to_supervised(scaled_dataset, 1, 1)
print(reframed.head())
print(reframed.shape)

# print(reframed.head())
# print(reframed.shape)

values = reframed.values
# First 4 years data
n_train_hours = 365 * 24 * 4

train = values[:n_train_hours, :]
test = values[n_train_hours:, :]
# print(train)
# print(test)

# split into input and outputs
train_X, train_y = train[:, :-1], train[:, -1]
test_X, test_y = test[:, :-1], test[:, -1]

# reshape input to be 3D :- (no.of samples, no.of timesteps, no.of
features)
train_X = train_X.reshape((train_X.shape[0], 1, train_X.shape[1]))
test_X = test_X.reshape((test_X.shape[0], 1, test_X.shape[1]))
# print(train_X.shape, train y.shape, test X.shape, test y.shape)

```

```

#----- Design Model -----
---
model = Sequential()
model.add(LSTM(256, input_shape=(train_X.shape[1], train_X.shape[2])))
model.add(Dense(64))
model.add(Dropout(0.25))
model.add(BatchNormalization())
model.add(Dense(1))

model.summary()

model.compile(loss='mse', optimizer='adam')

history = model.fit(train_X, train_y, epochs=50, batch_size=128,
validation_data=(test_X, test_y))

plt.figure(figsize=(15,6))
plt.plot(history.history['loss'], label='train', linewidth = 2.5)
plt.plot(history.history['val_loss'], label='test', linewidth = 2.5)
plt.legend()
plt.show()

# ----- Pridict -----
--
prediction = model.predict(test_X)
prediction = prediction.ravel()

ture_test = test[:, 1]

poll = np.array(df["pollution"])

meanop = poll.mean()
stdop = poll.std()

ture_test = ture_test * stdop + meanop
prediction = prediction * stdop + meanop

plt.figure(figsize=(15, 6))
plt.xlim([1000, 1250])
plt.ylabel("ppm")
plt.xlabel("hrs")
plt.plot(ture_test, c="g", alpha=0.90, linewidth=2.5)
plt.plot(prediction, c="b", alpha=0.75)
plt.title("Testing(Validation) data")
plt.show()

rmse = np.sqrt(mse(ture_test, prediction))
print("Test(Validation) RMSE =", rmse)

```