

# Stock Price Prediction Using LSTM

STAT 6289 Statistical Deep Learning

Presenter: Xinyi Li

# Step1: Upload the Data

```
# Load the data and inspect them
data = pd.read_csv('TSLA.csv')
data.head()
```

|   | Date       | Open      | High      | Low       | Close     | Adj Close | Volume   |
|---|------------|-----------|-----------|-----------|-----------|-----------|----------|
| 0 | 2015-07-10 | 52.444000 | 52.599998 | 51.563999 | 51.830002 | 51.830002 | 13054500 |
| 1 | 2015-07-13 | 52.450001 | 52.509998 | 51.209999 | 52.431999 | 52.431999 | 14801500 |
| 2 | 2015-07-14 | 52.419998 | 53.198002 | 52.102001 | 53.130001 | 53.130001 | 9538000  |
| 3 | 2015-07-15 | 53.348000 | 53.498001 | 52.416000 | 52.627998 | 52.627998 | 10108000 |
| 4 | 2015-07-16 | 52.844002 | 53.439999 | 52.632000 | 53.335999 | 53.335999 | 8080000  |

Choosing to predict the **Closing Price** of each trading day

Other potential options: open price, or adjusted closing price

## Step2: Split the Data

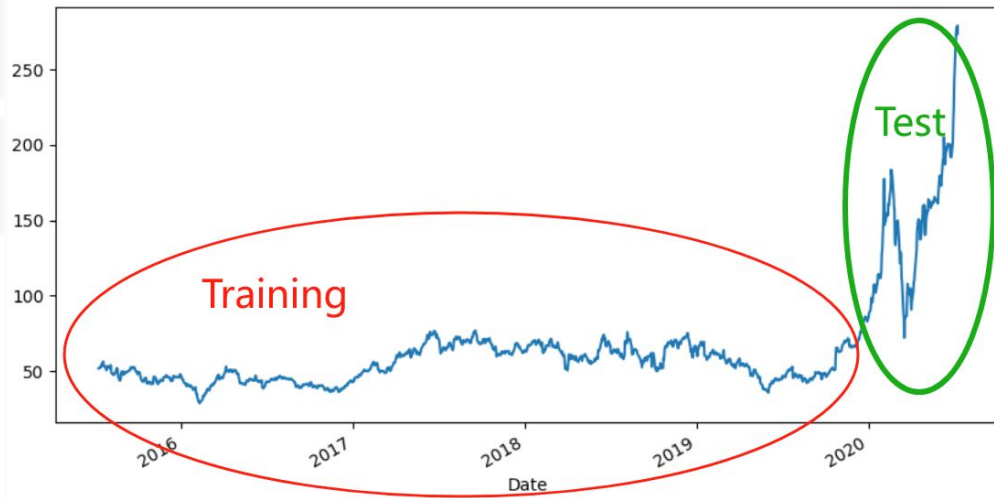
Why can't we randomly separate the data? Sequences matter; Unlike classification

```
#split the data into training and test sets
#set that training data is about 80%, test is about 20% (as normal setting)
training_size = int(len(data_normalized) * 0.8)
test_size= len(data_normalized)-training_size
```

```
train = data_normalized[0:training_size, :]
test = data_normalized[training_size:len(data_normalized), :]
```

```
print(training_size)
print(test_size)
train.shape, test.shape
```

```
1007
252
((1007, 1), (252, 1))
```



## Step3: Input the Time Lags

| Date      | Open   | High   | Low    | Close  | Adj Close | Volume   |
|-----------|--------|--------|--------|--------|-----------|----------|
| 7/10/2015 | 52.444 | 52.6   | 51.564 | 51.83  | 51.83     | 13054500 |
| 7/13/2015 | 52.45  | 52.51  | 51.21  | 52.432 | 52.432    | 14801500 |
| 7/14/2015 | 52.42  | 53.198 | 52.102 | 53.13  | 53.13     | 9538000  |
| 7/15/2015 | 53.348 | 53.498 | 52.416 | 52.628 | 52.628    | 10108000 |
| 7/16/2015 | 52.844 | 53.44  | 52.632 | 53.336 | 53.336    | 8080000  |
| 7/17/2015 | 54.5   | 55.108 | 53.65  | 54.932 | 54.932    | 25020500 |
| 7/20/2015 | 55     | 57.33  | 54.508 | 56.452 | 56.452    | 24892500 |
| 7/21/2015 | 54.01  | 54.7   | 53.31  | 53.354 | 53.354    | 30543500 |
| 7/22/2015 | 52.254 | 53.888 | 52.172 | 53.574 | 53.574    | 15525000 |
| 7/23/2015 | 53.93  | 53.98  | 53.054 | 53.44  | 53.44     | 11136000 |
| 7/24/2015 | 53.476 | 54.218 | 52.784 | 53.082 | 53.082    | 14182500 |

Not always constant;

Only trading days  
(weekdays) recorded

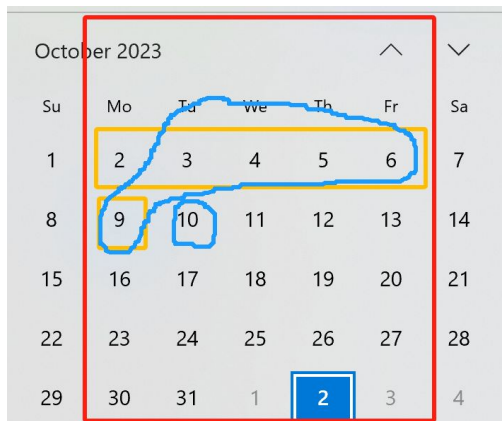
### Step3: Input the Time Lags (Time Lag=1)

| October 2023 |    |    |    |    |    |    |
|--------------|----|----|----|----|----|----|
| Su           | Mo | Tu | We | Th | Fr | Sa |
| 1            | 2  | 3  | 4  | 5  | 6  | 7  |
| 8            | 9  | 10 | 11 | 12 | 13 | 14 |
| 15           | 16 | 17 | 18 | 19 | 20 | 21 |
| 22           | 23 | 24 | 25 | 26 | 27 | 28 |
| 29           | 30 | 31 | 1  | 2  | 3  | 4  |

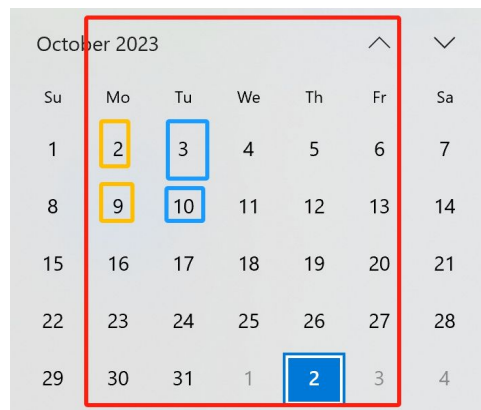
Use the price of yesterday to predict that of today; use the price of today to predict that of tomorrow  
(Predicting the future one day)

This is the simplest condition.

## Step3: Input the Time Lags (Time Lag=5)



“Updated...” file



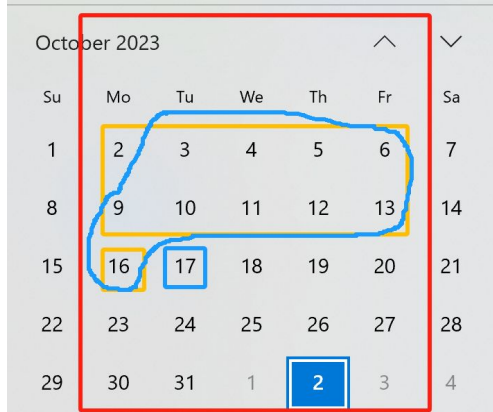
“Final...” file

When the time lag is not 1, there would be various understanding of time lag.

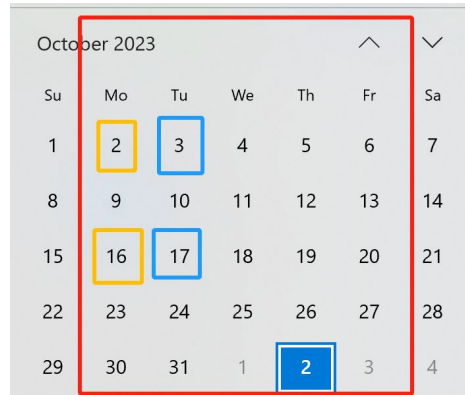
In the “Updated....ipynb” file, if time lag=5, that is to use the prices of the previous five days to predict that of today

In the “Final....ipynb” file, that is lead time =5, to use the price of the pervious fifth day to predict that of today.

## Step3: Input the Time Lags (Time Lag=10)



“Updated...” file



“Final...” file

Similarly, for time lag=10, there will be two understandings.

In the “Updated....ipynb” file, if time lag=10, that is to use the prices of the previous ten days to predict that of today

In the “Final....ipynb”file, that is lead time =10, to use the price of the pervious tenth day to predict that of today.

## Step3: Input the Time Lags

```
#build the input features with different time lags
def input_features(data, lag):
    X, Y = [], []
    for i in range(len(data) - lag):
        X.append(data[i:(i + lag), 0])
        Y.append(data[i + lag, 0])
    return np.array(X), np.array(Y)
```

“Updated...” file

```
#build the input features with different time lags
def input_features(data, lag):
    X, Y = [], []
    for i in range(len(data) - lag):
        X.append(data[i:(i+1), 0])
        Y.append(data[i + lag, 0])
    return np.array(X), np.array(Y)
```

“Final...” file



# Step4: Build the LSTM Model

Choosing the comparisons of **the number of hidden layers** and **drop out rate**

```
#create the LSTM model
def lstm(lstm_layer=4):
    model = Sequential()

    model.add(LSTM(units = 50, activation = 'relu', return_sequences=True,
                    input_shape = (X_train.shape[1], X_train.shape[2])))
    model.add(Dropout(0.2))

    for layer in range(lstm_layer-1):
        model.add(LSTM(50, return_sequences= False))
        model.add(Dropout(0.2))

    model.add(Dense(1))

    return model

for lstm_layer in [4,5,6]:
    model=lstm(lstm_layer)
```

```
#change another drop rate
def lstm2(lstm_layer=4):
    model = Sequential()

    model.add(LSTM(units = 50, activation = 'relu', return_sequences=True,
                    input_shape = (X_train.shape[1], X_train.shape[2])))
    model.add(Dropout(0.3))

    for layer in range(lstm_layer-1):
        model.add(LSTM(50, return_sequences= False))
        model.add(Dropout(0.3))

    model.add(Dense(1))

    return model

for lstm_layer in [4,5,6]:
    model=lstm2(lstm_layer)
```

```
models={'drop_0.2_4_h_layer': lstm(4), 'drop_0.2_5_h_layer': lstm(5), 'drop_0.2_6_h_layer': lstm(6),
        'drop_0.3_4_h_layer': lstm2(4), 'drop_0.3_5_h_layer': lstm2(5), 'drop_0.3_6_h_layer': lstm2(6)}
```

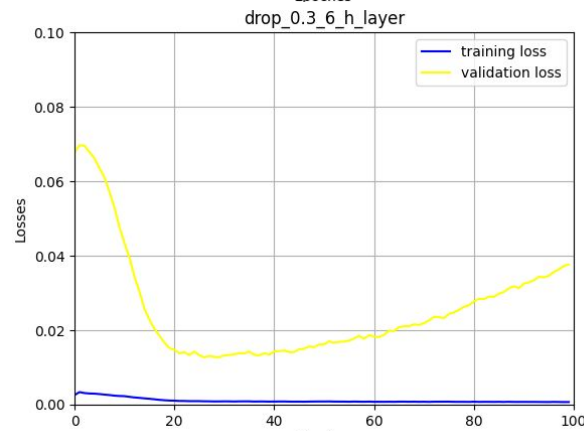
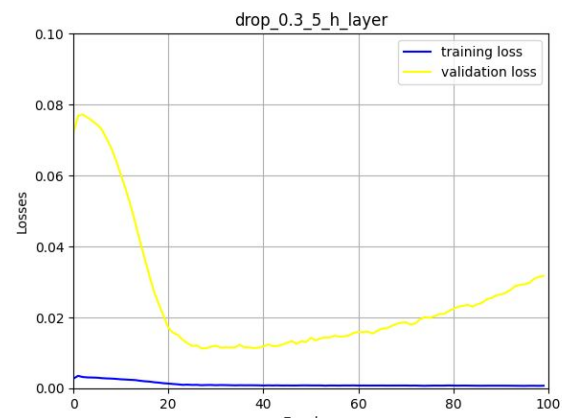
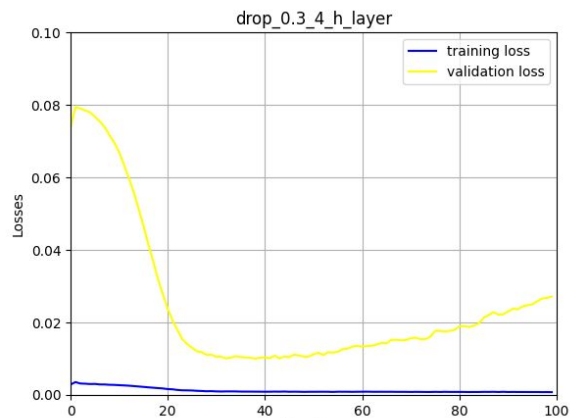
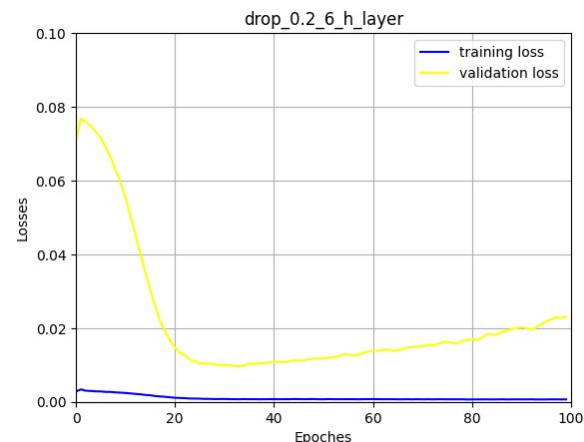
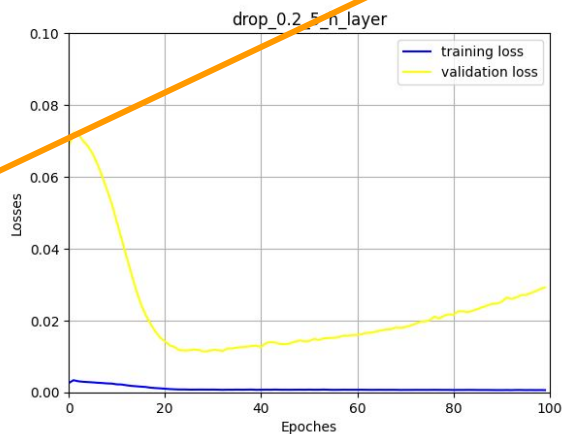
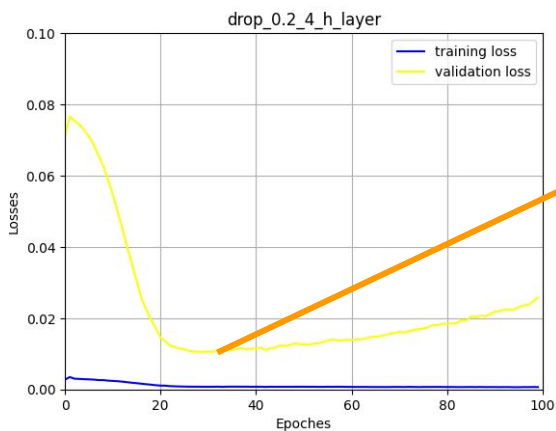
## Step4: Build the LSTM Model

Design optimizer and loss function; Set the number of epoches and batch sizes

```
for name, model in models.items():  
    model.compile(optimizer = 'adam', loss = 'mse' , metrics="mean_absolute_error")  
  
    history = model.fit(X_train, Y_train, epochs = 100, batch_size = 20,  
                        validation_data = (X_test, Y_test), verbose = 1, shuffle = False)
```

# Step4: Build the LSTM Model

Best model



## Step4: Build the LSTM Model

After finding the best model, go back to the model compile, and adjust the model to the best one

```
model=lstm(lstm_layer=6)
#lstm() is the model dropping out rate is 0.2, lstm2() is the rate of 0.3;
#also, lstm_layer=6, that's the number of hidden layers, 6
```

```
model.compile(optimizer = 'adam', loss = 'mse' , metrics="mean_absolute_error")
```

```
history = model.fit(X_train, Y_train, epochs = 30, batch_size = 20, validation_data = (X_test, Y_test), verbose = 1, shuffle = False) #select epoch=30 here
```

## Step5: Predict the stock

```
#prediction
Y_prediction = model.predict(X_test)
Y_prediction = scaler.inverse_transform(Y_prediction)
```

```
train_df = data.iloc[:training_size , :]
train_df.Date = pd.to_datetime(train_df.Date, format = '%Y/%m/%d')
test_df = data.iloc[training_size: , :]
test_df.Date = pd.to_datetime(test_df.Date, format = '%Y/%m/%d')
```

**test\_df.head()**



|      | Date       | Close     |
|------|------------|-----------|
| 1007 | 2019-07-11 | 47.720001 |
| 1008 | 2019-07-12 | 49.015999 |
| 1009 | 2019-07-15 | 50.700001 |
| 1010 | 2019-07-16 | 50.476002 |
| 1011 | 2019-07-17 | 50.972000 |

## Step5: Predict the stock

```
test_df['Prediction'] = np.nan # Initialize new column 'Prediction' with NaN values  
test_df.iloc[lag:, 2] = Y_prediction
```

```
test_df.head()
```

|      | Date       | Close     | Prediction |
|------|------------|-----------|------------|
| 1007 | 2019-07-11 | 47.720001 | NaN        |
| 1008 | 2019-07-12 | 49.015999 | 48.985081  |
| 1009 | 2019-07-15 | 50.700001 | 50.010769  |
| 1010 | 2019-07-16 | 50.476002 | 51.342285  |
| 1011 | 2019-07-17 | 50.972000 | 51.165260  |

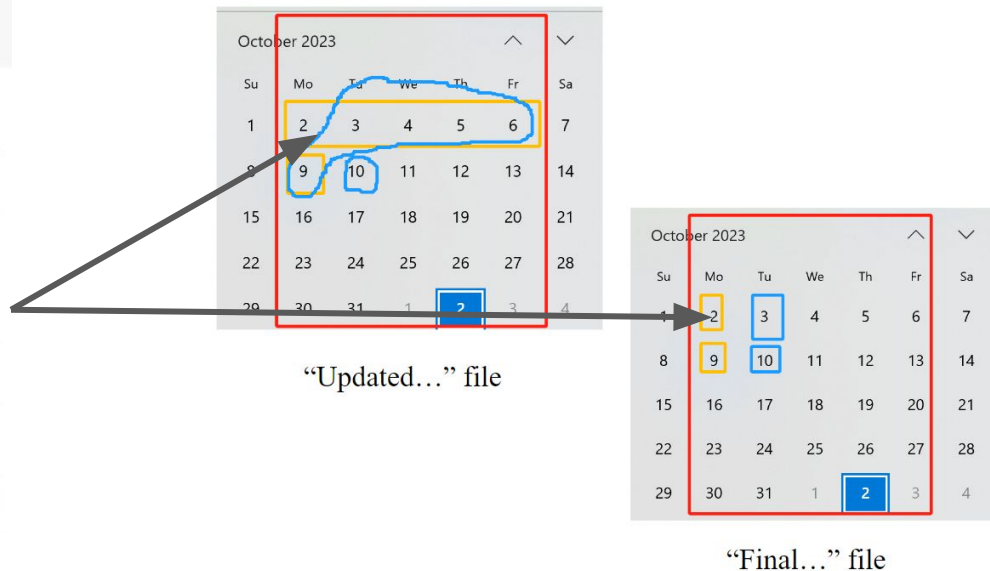
|              |    |    |    |    |    |    |
|--------------|----|----|----|----|----|----|
| October 2023 |    |    |    |    |    |    |
| Su           | Mo | Tu | We | Th | Fr | Sa |
| 1            | 2  | 3  | 4  | 5  | 6  | 7  |
| 8            | 9  | 10 | 11 | 12 | 13 | 14 |
| 15           | 16 | 17 | 18 | 19 | 20 | 21 |
| 22           | 23 | 24 | 25 | 26 | 27 | 28 |
| 29           | 30 | 31 | 1  | 2  | 3  | 4  |

## Step5: Predict the stock

```
test_df['Prediction'] = np.nan # Initialize new column 'Prediction' with NaN values  
  
test_df.iloc[lag:, 2] = Y_prediction
```

```
test_df.head(7)
```

|      | Date       | Close     | Prediction |
|------|------------|-----------|------------|
| 1007 | 2019-07-11 | 47.720001 | NaN        |
| 1008 | 2019-07-12 | 49.015999 | NaN        |
| 1009 | 2019-07-15 | 50.700001 | NaN        |
| 1010 | 2019-07-16 | 50.476002 | NaN        |
| 1011 | 2019-07-17 | 50.972000 | NaN        |
| 1012 | 2019-07-18 | 50.708000 | 50.954666  |
| 1013 | 2019-07-19 | 51.636002 | 51.234512  |

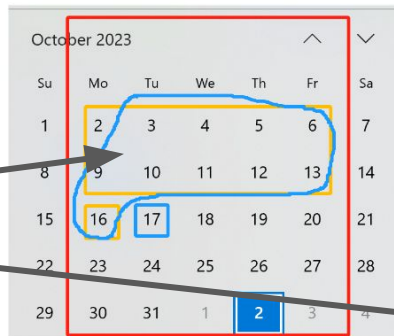




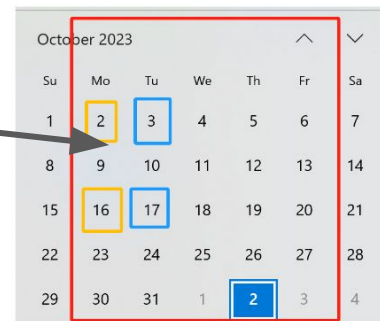
## Step5: Predict the stock

```
test_df['Prediction'] = np.nan # Initialize new column 'Prediction' with NaN values  
  
test_df.iloc[lag:, 2] = Y_prediction
```

|      | Date       | Close     | Prediction |
|------|------------|-----------|------------|
| 1007 | 2019-07-11 | 47.720001 | NaN        |
| 1008 | 2019-07-12 | 49.015999 | NaN        |
| 1009 | 2019-07-15 | 50.700001 | NaN        |
| 1010 | 2019-07-16 | 50.476002 | NaN        |
| 1011 | 2019-07-17 | 50.972000 | NaN        |
| 1012 | 2019-07-18 | 50.708000 | NaN        |
| 1013 | 2019-07-19 | 51.636002 | NaN        |
| 1014 | 2019-07-22 | 51.136002 | NaN        |
| 1015 | 2019-07-23 | 52.034000 | NaN        |
| 1016 | 2019-07-24 | 52.976002 | NaN        |
| 1017 | 2019-07-25 | 45.764000 | 51.409676  |
| 1018 | 2019-07-26 | 45.608002 | 52.347935  |
| 1019 | 2019-07-29 | 47.153999 | 53.574364  |
| 1020 | 2019-07-30 | 48.452000 | 53.410755  |
| 1021 | 2019-07-31 | 48.321999 | 53.773228  |



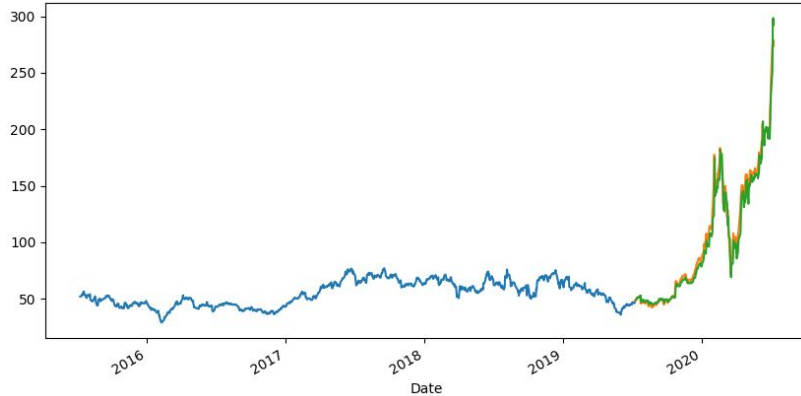
“Updated...” file



“Final...” file

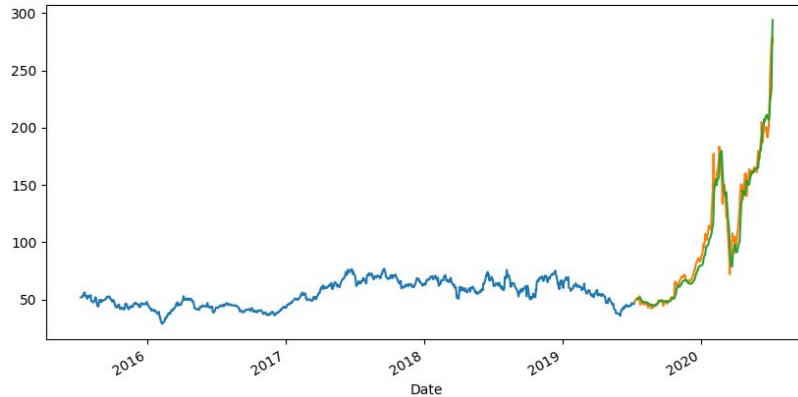


## Step5: Predict the stock (Results in Updated...file)

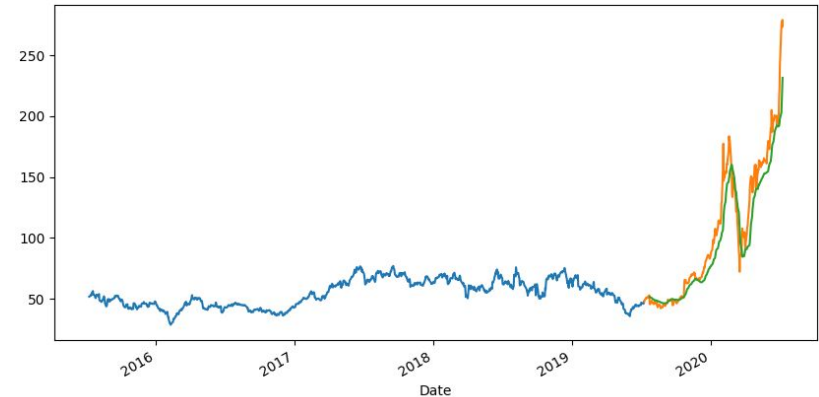


Future One Day Prediction (Time lag=1)

**Orange line is the test set.**  
**Green line is the prediction.**

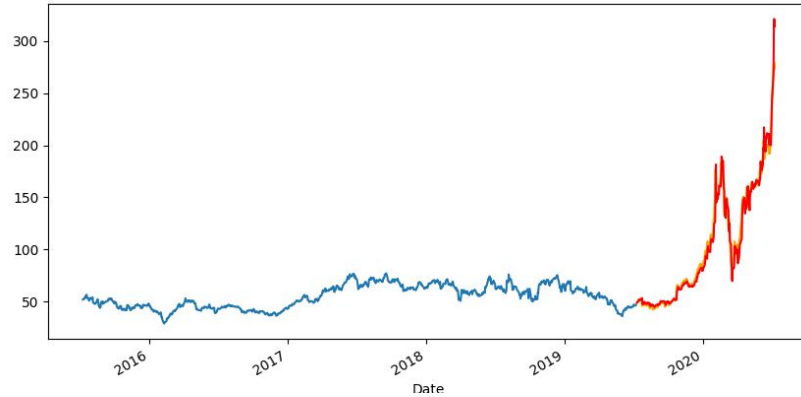


Future One Week Prediction (Time lag=5)



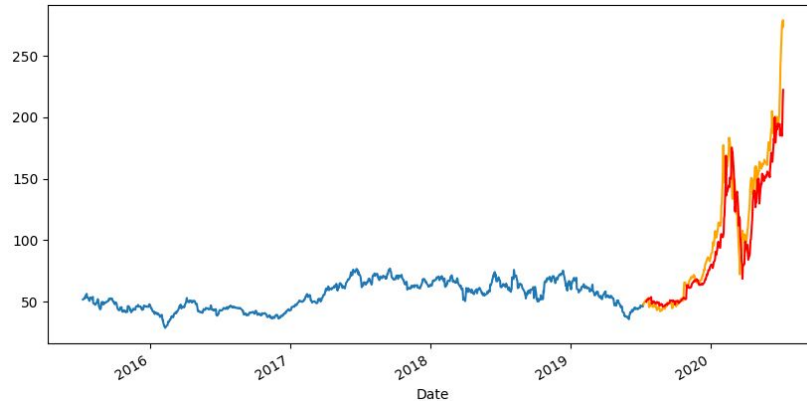
Future Two Weeks Prediction (Time lag=10)

## Step5: Predict the stock (Results in Final...file)

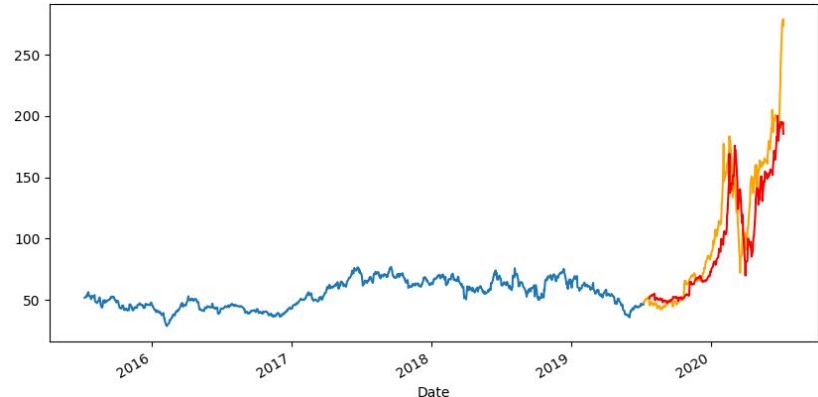


Future One Day Prediction (Time lag=1)

**Orange line** is the test set.  
**Red line** is the prediction.



Future One Week Prediction (Time lag=5)



Future Two Weeks Prediction (Time lag=10)