

Udacity Machine Learning Nanodegree 2020

Capstone Project

Stock Prediction Using Deep Learning

Nitu Nivedita

June 2020

Definition

Project Overview

As per my Capstone Project Proposal, this project aims at using Deep Learning and improved Recurrent Neural Networks in order to predict the stock market. Accurate stock market prediction is of great interest to investors; however, stock markets are driven by volatile factors among myriad other factors. For achieving maximum prediction accuracy, this Stock Prediction model tries to use deep learning and time series modelling methodologies that can look at the history of a sequence of data and correctly predict what the future elements are going to be.

Machine learning is being increasingly used to help make trading decisions and financial institutions are using artificial intelligence extensively. Although there is an abundance of stock data for machine learning models to train on, a high noise to signal ratio and the multitude of factors that affect stock prices are among the several reasons that predicting the market extremely difficult. However, these models don't need high levels of accuracy because even 60% accuracy can deliver solid returns. One really good method for predicting stock prices is using a long short-term memory neural network (LSTM) for times series forecasting, which I will be using in this project.

My motivation in solving this problem is both a) to understand the dynamics of the stock market and b) to use the power of Machine Learning to discover buying opportunities. My goal is to develop a structured LSTM network for time series predictive modeling and use the LSTM model to predict stock price and stock price movement for taking Buy/Sell/Hold decisions.

This project uses the concepts of Recurrent Neural Networks published in the research paper [Comparative Study of Stock Trend Prediction using Time Delay, Recurrent and Probabilistic Neural Networks](#) [1] D. C. Wunsch et al., "Comparative Study of Stock Trend Prediction using Time Delay, Recurrent and Probabilistic Neural Networks," IEEE Transactions on Neural Networks, Institute of Electrical and Electronics Engineers (IEEE) as well as the more recent [2] [An ensemble of LSTM neural networks for high-frequency stock market classification](#) long-short-term memory (LSTM) neural networks for stock prediction.

Problem Statement

The main purpose of the prediction is to reduce uncertainty associated with investment decision making. Without a doubt, forecasting stock prices is very difficult because of the market volatility that needs accurate forecasting models. The stock market prices are highly fluctuating and are considered to be very dynamic because of the underlying nature of the financial domain and in part because of the mix of parameters (Previous day's closing price, P/E ratio etc.) and the unknown factors (like Election Results etc.).

Metrics

Metrics used to measure the performance of this model are described below:

- 1) For evaluation of the model, RMSE will be used. Root mean squared error (RMSE) can provide what is the average deviation of the prediction from the true value, and it can be compared with the mean of the true value to see whether the deviation is large or small. We can plot the predicted and actual data. Valid and predicted prices visualized can be used to evaluate if the model passed.
- 2) To analyze our model's effectiveness, we can validate our model against untested data. One metric we can create is a prediction interval (similar to a confidence interval) which will tell us with some high degree of confidence in the true range of our prediction. We can use the root mean squared error for getting to a prediction interval
- 3) I will use R-square and RMSE to check the model performance, on both benchmark model and LSTM model
- 4) I will use scikit-learn LinearRegression and use Keras's LSTM function to build the neural networks. The networks will have two or three layers

Analysis

Data Exploration

As part of the exploration, I will perform **5 steps as described in detail below**

(1) Data Exploration: Description of Dataset

Attribute values are in floating point except to date and volume.

Date: — Trading date of the stock.

Open: — This price of stock's opening price which means the very beginning price of particular trading day, but which is not be the same price of previous's day ending price.

High: — This is the highest price of the stock on a particular trading day.

Low: — This is the lowest stock price during trade day.

Close: — This is the closing price of the stock during trade-in particular day.

Adj Close: — This is the ending or closing price of the stock which was changed to contain any corporations' actions and distribution that is occurred during trade time of the day.

Volume: — This is the number of stocks traded on a particular day.

(2) Data Exploration: Data Preparation

The first stage we need to import all necessary libraries. The gathered data set was read using the panda library in python and display the records for understanding the data set for pre-processing. At this point, I was able to identify the behavior and characteristics of the data set.

Data Preparation

Using last 10 years stock prices for data exploration and analysis

```
In [125]: # The tech stocks we'll use for analysis
tech_list = ['AAPL', 'GOOG', 'MSFT', 'AMZN']

# Set up End and Start times for data pull.
end = datetime.now()
start = datetime(end.year - 10, end.month, end.day)

#For loop for grabbing yahoo finance data and setting as a dataframe
for stock in tech_list:
    # Set DataFrame as the Stock Ticker
    globals()[stock] = DataReader(stock, 'yahoo', start, end)
```

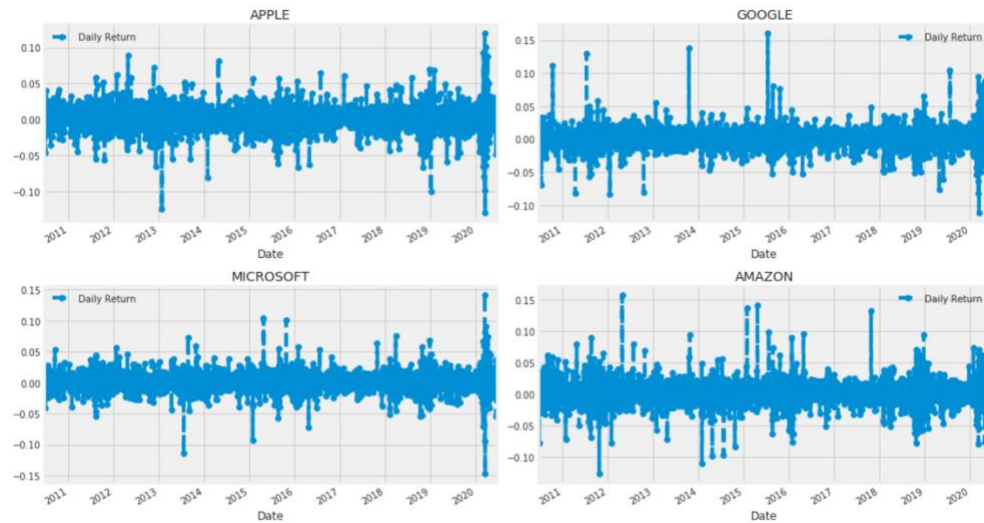
	High	Low	Open	Close	Volume	Adj Close	company_name
Date							
2020-05-08	2387.239990	2357.000000	2372.139893	2379.610107	3211200.0	2379.610107	AMAZON
2020-05-11	2419.669922	2372.110107	2374.699951	2409.000000	3253700.0	2409.000000	AMAZON
2020-05-12	2419.000000	2355.000000	2411.850098	2356.949951	3074900.0	2356.949951	AMAZON
2020-05-13	2407.699951	2337.800049	2366.800049	2367.919922	4782900.0	2367.919922	AMAZON
2020-05-14	2391.370117	2353.209961	2361.010010	2388.850098	3648100.0	2388.850098	AMAZON
2020-05-15	2411.000000	2356.370117	2368.520020	2409.780029	4235000.0	2409.780029	AMAZON
2020-05-18	2433.000000	2384.010010	2404.350098	2426.260010	4366600.0	2426.260010	AMAZON
2020-05-19	2485.000000	2428.969971	2429.830078	2449.330078	4320500.0	2449.330078	AMAZON
2020-05-20	2500.000000	2428.969971	2429.830078	2449.330078	4320500.0	2449.330078	AMAZON

```
AAPL.describe()
```

	High	Low	Open	Close	Volume	Adj Close
count	2518.000000	2518.000000	2518.000000	2518.000000	2.518000e+03	2518.000000
mean	123.854524	121.507019	122.659366	122.731401	6.899855e+07	116.095715
std	67.178624	65.702543	66.339431	66.522625	5.332056e+07	68.393280
min	34.658573	33.651428	34.005714	34.275715	1.136200e+07	29.672398
25%	74.250357	72.619999	73.388569	73.573570	3.014805e+07	64.786510
50%	108.545002	106.355000	107.460003	107.480003	5.043535e+07	99.784672
75%	165.299999	161.542496	163.722500	163.574997	9.397745e+07	157.845711
max	356.559998	351.089996	355.149994	352.839996	4.702495e+08	352.839996

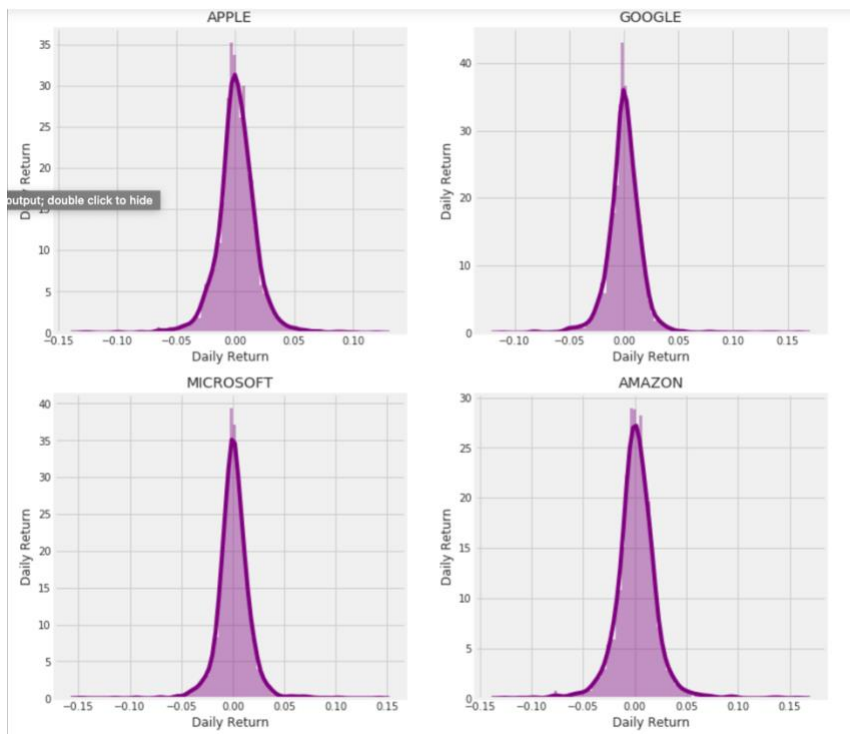
(3) Data Exploration and Analysis: Determination of risk factors of stocks

Here, I will try to analyze the risk factor of stocks and for that I will take a deeper look at the daily changes of the stock. I will use pandas to retrieve the daily returns for the Tech stocks



(4) Data Exploration and Analysis: Determination of Daily Return of stocks

I plotted the daily return of stocks based on my above analysis of risk factors



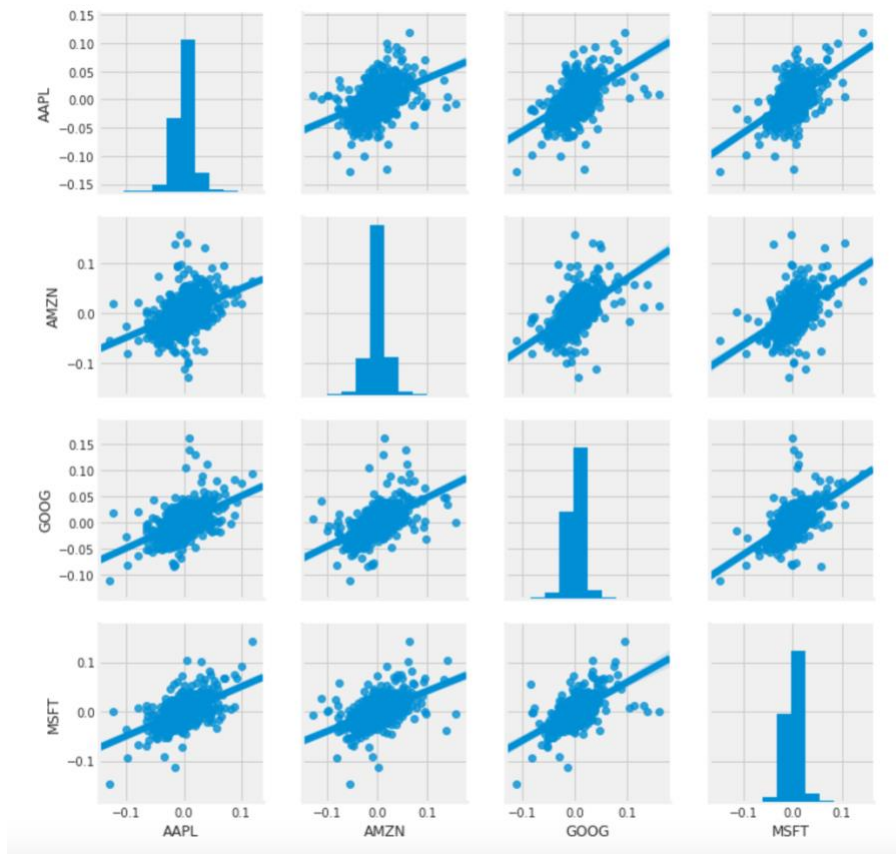
(5) [Data Exploration and Visualization: Determination of Correlation among stocks](#)

I want to analyze the returns of all the stocks in my list above so I will build a Dataframe with all the ['Close'] columns for each of the stocks.

```
closing_df = DataReader(tech_list, 'yahoo', start, end)['Adj Close']  
  
# Let's check  
closing_df.head()
```

Symbols	AAPL	AMZN	GOOG	MSFT
Date				
2010-06-21	33.412212	122.550003	243.367798	20.500095
2010-06-22	33.867310	122.309998	242.217117	20.357901
2010-06-23	33.511150	121.449997	240.124954	19.994505
2010-06-24	33.267521	118.330002	236.662933	19.749605
2010-06-25	32.983074	121.000000	235.457443	19.378315

One important parameter in data exploration is that I want to plot all the relationships on daily returns between all the stocks. This shows an interesting correlation between Tech stocks



Exploratory Visualization

Below figure illustrates a visualization of the feature distribution. Through this graph, I am getting a better understanding of features and coefficient of the determination score. A heat map is created by over-index features. Here, I was able to see the co-relation in between the given features of the data set such as Open, High, Low, Close and Adj Close. I also plotted a heat map for analysis

Exploratory Visualization

```
# Let's see a historical view of the closing price
```

```
plt.figure(figsize=(12, 8))
plt.subplots_adjust(top=1.25, bottom=1.2)

for i, company in enumerate(company_list, 1):
    plt.subplot(2, 2, i)
    company['Adj Close'].plot()
    plt.ylabel('Adj Close')
    plt.xlabel(None)
    plt.title(f"{tech_list[i - 1]}")
```



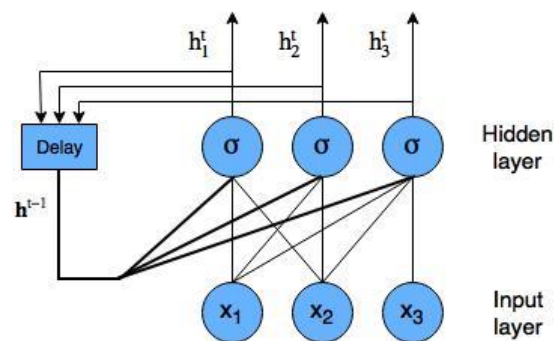
Algorithms and Techniques

1. Reading and analyzing data. (Pandas)
2. Normalizing the data (Sci-kit Learn)
3. Converting data to time-series and supervised learning problem.
4. Creating model (Keras)
5. Fine tuning the model
6. Training, predicting and visualizing the result.

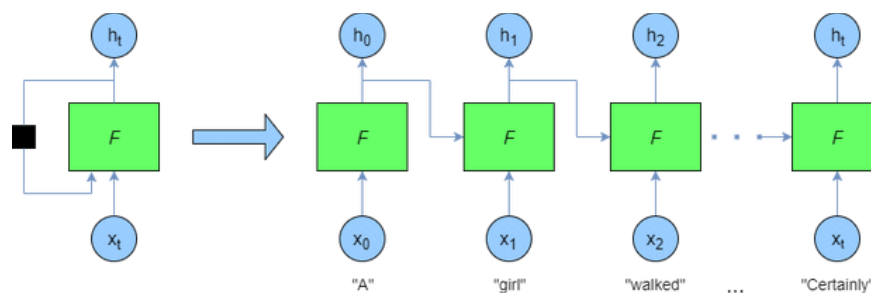
Analyzing Data: Determining Moving Average of Stocks

A moving average for a list of numbers is like arithmetic average but instead of calculating the average of all the numbers, we calculate the average of the first 'n' numbers (n is referred as window size or time period) and then move (or slide) the window by 1 index, thus excluding the first element and including the n+1 element and calculate their average. This process continues.

To explain the algorithm and technique used to solve this problem, I will take the help of a diagram. A recurrent neural network is a neural network that attempts to model time or sequence dependent behavior. This is performed by feeding back the output of a neural network layer at time t to the input of the same network layer at time $t + 1$.

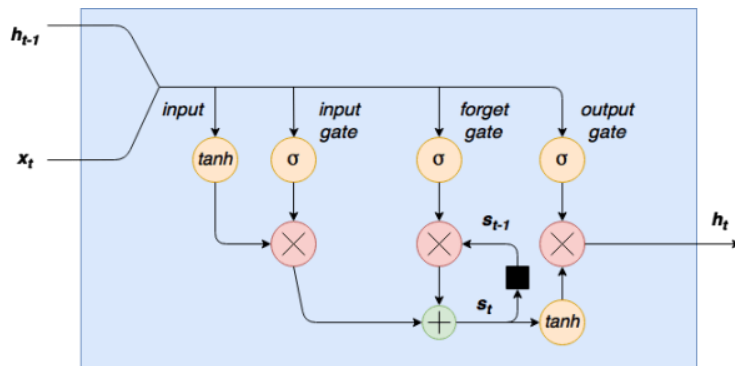


Recurrent neural networks are “unrolled” programmatically during training and prediction, so we get something like the following:



An LSTM network is a recurrent neural network that has LSTM cell blocks in place of our standard neural network layers. These cells have various components called the input gate, the forget gate and the output gate. Here is a graphical representation of the LSTM cell. The new neuron has 3 gates, each with a different goal. The gates are:

- *Input Gate*
- *Output Gate*
- *Forget Gate*



Data Cleansing

Scrape the data set from Yahoo Finance and cleanse the acquired data. b) We have prices of different stocks, they are in different scales and hence for the neural network to converge, we need to scale. c) Normalize the data for the LSTM networks which will be done using the function MinMaxScaler of the sklearn library. Create the scaled training data set d) Split the dataset into the training (80%) and test (20%) datasets

Model Design

The First step is to import the libraries needed - sklearn library, Keras, seaborn and matplotlib for visualization b) Build the LSTM model, compile and train LSTM has a similar control flow as a recurrent neural network. It processes data passing on information as it propagates forward. The differences are the operations within the LSTM's cells. These operations are used to allow the LSTM to keep or forget information.

Algorithm

- This model will be implemented using Keras module.
- Create an object for the model with Sequential function.
- Add layers to the model
- In this network first a LSTM layer will be added which takes the 3-dimensional array as input and has dimension 50 i.e the number of neurons.
- Then we compile the function using a loss function parameter, the optimizer and the metrics we want to use to check the model's efficiency.
- Call function fit to train the model
- Use the model to predict closing price of any given stock (ex. AMZN, AAPL, GOOG etc)

Benchmark

The benchmark model for this project is linear regression. My main goal to compare the accuracies of the various machine learning models. This benchmark will use exact the same input as our LSTM network model and provide a benchmark performance for the LSTM.

Benchmark Model

```
[191]: from sklearn.preprocessing import MinMaxScaler
import csv
scaler = MinMaxScaler()
data=dataset.iloc[:,1:7].values
print(data[1])
X = data[:,1:2]
y = dataset['Close']

y_new = []
for i in y:
    y_new.append([i])
y_new = np.array(y_new)

X = np.array(X)

X=scaler.fit_transform(X)
y_new=scaler.fit_transform(y_new)

[5.03200e+01 2.29428e+07 5.43400e+01 5.00600e+01 5.39500e+01]
```

Next, divide the test data “goog.csv” containing 10 years of Google closing prices into training and test data to be fed into the Benchmark model.

Linear Regression Model

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y_new, test_size = 0.3, shuffle = False, random_state = 0)
print(X_train[1])
print(y_train[1])
print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)

[0.27788561]
[0.00279645]
(2783, 1)
(1194, 1)
(2783, 1)
(1194, 1)
```

Fit the Model

```
from sklearn.linear_model import LinearRegression
model= LinearRegression()
model = model.fit(X_train,y_train)
predicted_price = model.predict(X_test)
predicted_train = model.predict(X_train)
print(model.coef_)
print(model.intercept_)

[[-0.50343115]]
[0.22706923]
```

Evaluate the Benchmark Model Performance

Evaluate Benchmark Model Performance

```
In [196]: from sklearn.metrics import r2_score, mean_squared_error
r2_score(y_test,predicted_price),np.sqrt(mean_squared_error(y_test,predicted_price))

Out[196]: (-8.446253522156992, 0.4534881927598406)
```

Methodology

Data Preprocessing

First I imported the basic libraries and got the data of 4 Tech stocks as below from [Yahoo Finance for last 10 years](#)

```
# The tech stocks we'll use for analysis
tech_list = ['AAPL', 'GOOG', 'MSFT', 'AMZN']
|
# Set up End and Start times for data pull.
end = datetime.now()
start = datetime(end.year - 10, end.month, end.day)

#For loop for grabbing yahoo finance data and setting as a dataframe
for stock in tech_list:
    # Set DataFrame as the Stock Ticker
    globals()[stock] = DataReader(stock, 'yahoo', start, end)
```

Implementation

(1) [Explore the data](#)

In this step I am exploring the datasets first – last 10 years of price history of stocks below

```
company_list = [AAPL, GOOG, MSFT, AMZN]
company_name = ["APPLE", "GOOGLE", "MICROSOFT", "AMAZON"]

for company, com_name in zip(company_list, company_name):
    company["company_name"] = com_name

df = pd.concat(company_list, axis=0)
df.tail(30)
```

	High	Low	Open	Close	Volume	Adj Close	company_name
Date							
2020-05-08	2387.239990	2357.000000	2372.139893	2379.610107	3211200.0	2379.610107	AMAZON
2020-05-11	2419.669922	2372.110107	2374.699951	2409.000000	3253700.0	2409.000000	AMAZON
2020-05-12	2419.000000	2355.000000	2411.850098	2356.949951	3074900.0	2356.949951	AMAZON
2020-05-13	2407.699951	2337.800049	2366.800049	2367.919922	4782900.0	2367.919922	AMAZON
2020-05-14	2391.370117	2353.209961	2361.010010	2388.850098	3648100.0	2388.850098	AMAZON
2020-05-15	2411.000000	2356.370117	2368.520020	2409.780029	4235000.0	2409.780029	AMAZON
2020-05-18	2433.000000	2384.010010	2404.350098	2426.260010	4366600.0	2426.260010	AMAZON
2020-05-19	2485.000000	2429.060071	2429.830078	2440.330078	4320500.0	2440.330078	AMAZON

(2) [Building the LSTM](#)

In order to build the LSTM, we need to import a couple of modules from [Keras](#):

1. Sequential for initializing the neural network
2. Dense for adding a densely connected neural network layer

3. LSTM for adding the Long Short-Term Memory layer
4. Dropout for adding dropout layers that prevent overfitting

I have added the LSTM layer and later added a few Dropout layers to prevent overfitting. The LSTM layer arguments are:

1. *50 units* which is the dimensionality of the output space
2. *return_sequences=True* which determines whether to return the last output in the output sequence, or the full sequence
3. *input_shape* as the shape of our training set.

When defining the Dropout layers, we specify 0.2, meaning that 20% of the layers will be dropped. Thereafter, we add the Dense layer that specifies the output of 1 unit. After this, we compile our model using the popular adam optimizer and set the loss as the *mean_squared_error*. This will compute the mean of the squared errors.

Our data is now ready and we can build the RNN LSTM model and fit the data. Importing the required keras libraries to build the LSTM network

LSTM Model

```
In [156]: from keras.models import Sequential
          from keras.layers import Dense, LSTM

          #Build the LSTM model
          model = Sequential()
          model.add(LSTM(50, return_sequences=True, input_shape= (x_train.shape[1], 1)))
          model.add(LSTM(50, return_sequences=False))
          model.add(Dense(25))
          model.add(Dense(1))
```

I implemented an efficient stacked RNN. *return_sequences* is to return the last output in the output sequence. *input_shape* will be of the 3D format of test sample size, time steps, no. of input features.

(3) [Normalization](#)

Normalizing the input data using *MinMaxScaler* so that all the input features are on the scale from 0 to 1

```
In [152]: #Scale the data
          from sklearn.preprocessing import MinMaxScaler

          scaler = MinMaxScaler(feature_range=(0,1))
          scaled_data = scaler.fit_transform(dataset)

          scaled_data
```

Feature scaling and normalizing data are the best way to reduce the error rate and improve the accuracy of the model. There are various types of data in a given data set. Here, I am putting all selected features on the same scale. Therefore, none of the features are dominating others. Below code, a snapshot is used to scale the data.

After the pre-processing stage, I have created the implementation of LSTM algorithm. I have created a training data set that contains the past 8 years of close price of stock. Split the training data set into two parts Xtrain and Ytrain. Here, Xtrain dataset is an independent dataset and Ytrain data set is a dependent data set. I have converted the Xtrain and Ytrain data set into NumPy array because it is needed for training the LSTM model.

(4) [Convert to numpy arrays and reshape the data](#)

```
In [154]: # Convert the x_train and y_train to numpy arrays
x_train, y_train = np.array(x_train), np.array(y_train)

In [155]: #Reshape the data
x_train = np.reshape(x_train, (x_train.shape[0], x_train.shape[1], 1))
x_train.shape

Out[155]: (1644, 60, 1)
```

(5) [Compiling the Keras LSTM model](#)

Next I compiled the model. In this step, the type of loss that Keras should use to train the model needs to be specified. In this case, we are using 'mean_squared_error'. Next, in this example, the optimizer that will be used is the Adam optimizer – an effective “all round” optimizer with adaptive stepping. I also used Keras *callback* – callbacks are certain functions which Keras can optionally call, usually after the end of a training epoch. The callback that is used in this example is a model checkpoint callback – this callback saves the model after each epoch, which can be handy for when we are running long-term training.

(6) [Train and Fit the model](#)

```
In [157]: # Compile the model
model.compile(optimizer='adam', loss='mean_squared_error')

In [158]: #Train the model
model.fit(x_train, y_train, batch_size=1, epochs=1)

Epoch 1/1
1644/1644 [=====] - 361s 220ms/step - loss: 6.1706e-04

Out[158]: <keras.callbacks.History at 0x7f550003cf98>
```

Refinement

I have followed the below methodology to increase the accuracy of the model in the Long-Short Term Memory (LSTM) algorithm.

- Allocating more time-series data for training and testing (80% and 20%). It is a very important strategy to get higher accuracy results.
- Feature Scaling and Normalizing data to represent the independent features into a fixed range. Here, I put all the features on the same scale.
- Identifying the most important features from the given data set before training. It would be increasing the accuracy of the model.

Results

Model Evaluation and Validation

In this project, feature columns used are ['adjclose', 'volume', 'open', 'high', 'low']. As we see the model is produced predictive results very well with high accuracy and low error!



RMSE is quite low. Model has high accuracy.

RMSE – code and output

```
In [162]: # Get the models predicted price values
          predictions = model.predict(x_test)
          predictions = scaler.inverse_transform(predictions)

In [163]: # Get the root mean squared error (RMSE)
          rmse = np.sqrt(np.mean(((predictions - y_test) ** 2)))
          rmse

Out[163]: 9.51579715537576
```

Justification

I have compared the final results to the benchmark result with the help of statistical analysis. The LSTM solution is significant and has adequately solved the problem. As we see above, LSTMs are very powerful in sequence prediction problems because they're able to store past information.

This is important in my case because the previous price of a stock is crucial in predicting its future price. LSTM is a very great choice to handle with time-series data rather than traditional Recurrent Neural Network (RNN). In RNN, there is a gradient vanishing/exploding problem, and the problem comes from updating the weights by only multiplications. To solve the problem, LSTM considers another way to updating the weights not only by multiplications but also by additions.

In my work and project above, I have created a high-accuracy LSTM Stock Predictor model and evaluated it to be of high accuracy by testing over multiple stocks

References :

- <https://en.wikipedia.org/wiki/NASDAQ>
- https://en.wikipedia.org/wiki/New_York_Stock_Exchange
- <http://www.investopedia.com/terms>
- https://en.wikipedia.org/wiki/Root-mean-square_deviation
- https://en.wikipedia.org/wiki/Coefficient_of_determination
- http://scikit-learn.org/stable/modules/model_evaluation.html
- <https://www.kaggle.com/wiki/RootMeanSquaredError>