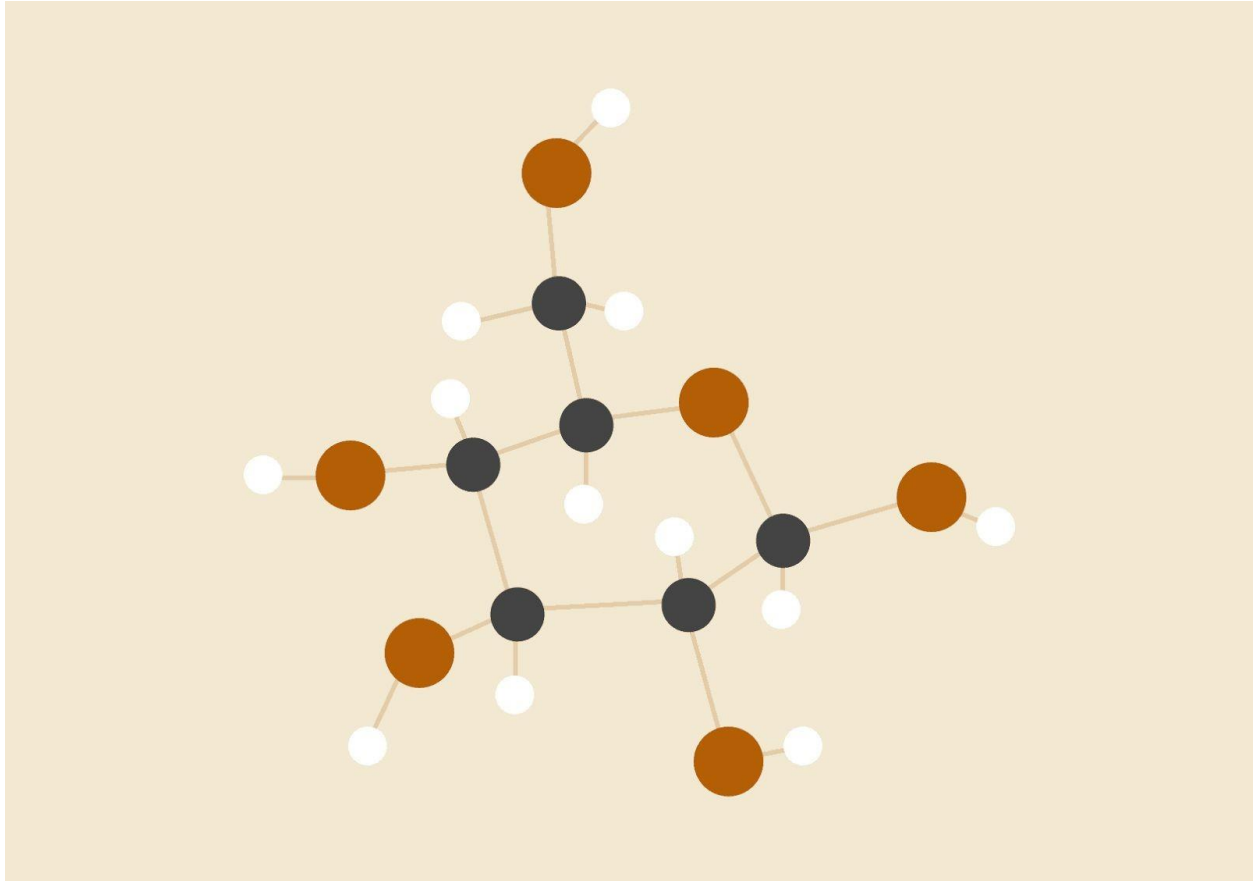


STOCK PRICE PREDICTION

Using Long-Short Term Memory Networks



Venkatesh Venkatraman

15.06.2019

Machine Learning Nanodegree

Table of Content

Table of Content	1
DEFINITION	2
Project Overview	2
Problem Statement	3
Metrics	3
ANALYSIS	4
Data Exploration	4
Exploratory Visualization	5
Algorithms and Techniques	6
Benchmark Model	7
METHODOLOGY	8
Data Preprocessing	8
Implementation	10
Refinement	11
RESULT	14
Model Evaluation and Validation	14
Justification	17
CONCLUSION	18
Free-Form Visualization	18
Reflection	19
Improvement	20

DEFINITION

Project Overview

Investment firms, hedge funds and even individuals have been using financial models to better understand market behavior and make profitable investments and trades. A wealth of information is available in the form of historical stock prices and company performance data, suitable for machine learning algorithms to process.

Can we actually predict stock prices with machine learning? Investors make educated guesses by analyzing data. They'll read the news, study the company history, industry trends and other lots of data points that go into making a prediction. The prevailing theories is that stock prices are totally random and unpredictable but that raises the question why top firms like Morgan Stanley and Citigroup hire quantitative analysts to build predictive models. We have this idea of a trading floor being filled with adrenaline infused men with loose ties running around yelling something into a phone but these days they're more likely to see rows of machine learning experts quietly sitting in front of computer screens. In fact about 70% of all orders on Wall Street are now placed by software, we're now living in the age of the algorithm.

This project seeks to utilize Deep Learning models, Long-Short Term Memory (LSTM) Neural Network algorithm, to predict stock prices. For data with timeframes recurrent neural networks (RNNs) come in handy but recent researches have shown that LSTM, networks are the most popular and useful variants of RNNs.

I will use Keras to build a LSTM to predict stock prices using historical closing price and trading volume and visualize both the predicted price values over time and the optimal parameters for the model.

Problem Statement

The challenge of this project is to accurately predict the future closing value of a given stock across a given period of time in the future. For this project I will use a **Long Short**

Term Memory networks¹ – usually just called “LSTMs” to predict the closing price of the Bank Nifty using a dataset of past prices

GOALS

1. Explore stock prices.
2. Implement basic model using linear regression.
3. Implement LSTM using keras library.
4. Compare the results and submit the report.

Metrics

For this project measure of performance will be using the Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) calculated as the difference between predicted and actual values of the target stock at adjusted close price and the delta between the performance of the benchmark model (Linear Regression) and our primary model (Deep Learning).

ANALYSIS

Data Exploration

The data used in this project is of the **Bank Nifty** from **January 1, 2008 to May 31, 2019**, this is a series of data points indexed in time order or a time series. My goal was to predict the closing price for any given date after training. The data was downloaded from the **NSE(National Stock Exchange)**

The prediction has to be made for Closing price of the data.

The dataset is of following form :

	Date	Open	High	Low	Close	Volume	Turnover
0	01-Jan-08	6136.75	6165.35	6109.85	6144.35	81305366	4410.97
1	02-Jan-08	6144.70	6197.00	6060.85	6179.40	110353305	7184.35
2	03-Jan-08	6184.25	6230.15	6126.40	6178.55	156427760	9475.37
3	04-Jan-08	6179.10	6300.05	6179.10	6274.30	147642063	8779.23
4	07-Jan-08	6271.00	6289.80	6193.35	6279.10	139778977	8873.37

Table: The whole data can be found out in 'v_bank_nifty.csv' in the project root folder⁶

Note: I did not observe any abnormality in datasets, i.e, no feature is empty and does not contains any incorrect value as negative values.

The mean, standard deviation, maximum and minimum of the data was found to be following:

Open	---	mean : 6963.446213729645	Std: 2290.955577376435	Max: 12052.65	Min: 2526.2
High	---	mean : 7004.592392073593	Std: 2287.0091030936446	Max: 12103.05	Min: 2585.3
Low	---	mean : 6911.586394196764	Std: 2289.5635729660344	Max: 12005.85	Min: 2252.75
Close	---	mean : 6958.5147912243565	Std: 2287.960306242476	Max: 12088.55	Min: 2524.2
Volume	---	mean : 195436167.0279547	Std: 84129125.37035412	Max: 741153209	Min: 2768292
Turn.0	---	mean : 8464.916206652517	Std: 3873.694247427281	Max: 35131.19	Min: 113.99

We can infer from this dataset that date, high and low values are not important features of the data. As it does not matter at what was the highest prices of the stock for a particular day or what was the lowest trading prices. What matters is the opening price of the stock and closing prices of the stock. If at the end of the day we have higher closing prices than the opening prices that we have some profit otherwise we saw losses. Also volume of share is important as a rising market should see rising volume, i.e, increasing price and decreasing volume show lack of interest, and this is a warning of a potential reversal. A price drop (or rise) on large volume is a stronger signal that something in the stock has fundamentally changed.

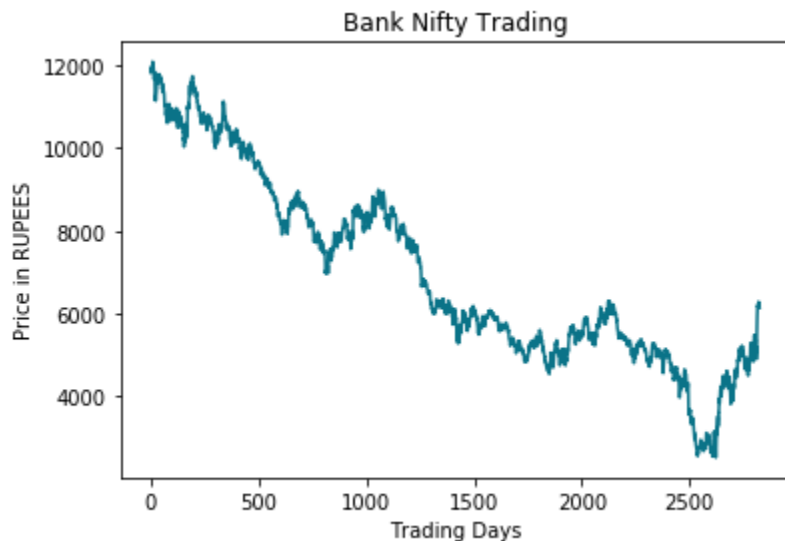
Therefore i have removed Date, High and low features from data set at preprocessing step. The mean, standard deviation, maximum and minimum of the preprocessed data was found to be following:

Open	---	mean : 0.4657817144612799	Std: 0.2404836615293666	Max: 0.9999999999999998	Min: 0.0
Close	---	mean : 0.4636294982120443	Std: 0.2392175428798064	Max: 1.0000000000000002	Min: 0.0
Volume	---	mean : 0.2609314878894725	Std: 0.11393667914041924	Max: 1.0	Min: 0.0

Exploratory Visualization

To visualize the data i have used **matplotlib**⁷ library. I have plotted Closing stock price of the data with the no of items(no of days) available.

Following is the snapshot of the plotted data :



X-axis: Represents Tradings Days
Y-axis: Represents Closing Price In Rupee

Through this data we can see a price range of bank nifty. The major fall in the prices between 12000-4000 might be because of the Global Financial Crisis of 2008-2009.

Algorithms and Techniques

The goal of this project was to study time-series data and explore as many options as possible to accurately predict the Stock Price. Through my research i came to know about **Recurrent Neural Nets (RNN)**⁸ which are used specifically for sequence and pattern learning. As they are networks with loops in them, allowing information to persist and thus ability to memorise the data accurately. But Recurrent Neural Nets have vanishing Gradient descent problem which does not allow it to learn from past data as was expected. The remedy of this problem was solved in **Long-Short Term Memory Networks**⁹, usually referred as LSTMs. These are a special kind of RNN, capable of learning long-term dependencies.

In addition to adjusting the architecture of the Neural Network, the following full set of parameters can be tuned to optimize the prediction model:

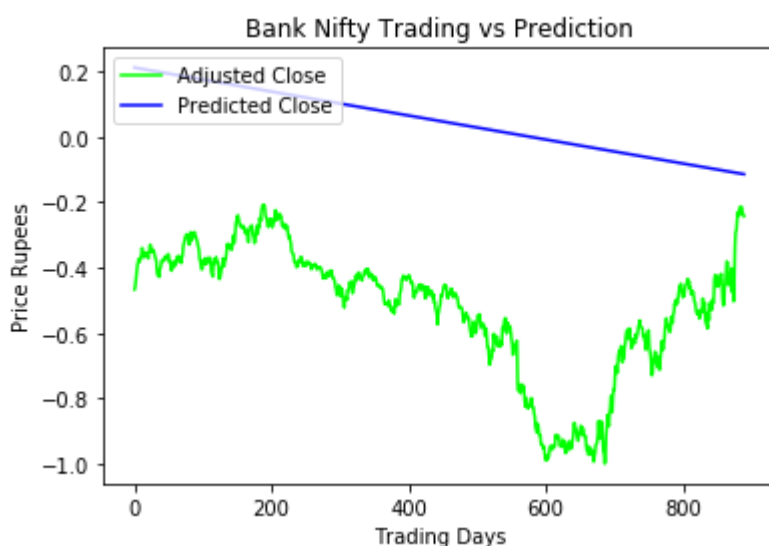
- Input Parameters
 - Preprocessing and Normalization (see Data Preprocessing Section)
- Neural Network Architecture
 - Number of Layers (how many layers of nodes in the model; used 3)
 - Number of Nodes (how many nodes per layer; tested 1,3,8, 16, 32, 64, 100,128)
- Training Parameters
 - Training / Test Split (how much of dataset to train versus test model on; kept constant at 82.95% and 17.05% for benchmarks and lstm model)
 - Validation Sets (kept constant at 0.05% of training sets)
 - Batch Size (how many time steps to include during a single training step; kept at 1 for basic lstm model and at 512 for improved lstm model)
 - Optimizer Function (which function to optimize by minimizing error; used “Adam” throughout)
 - Epochs (how many times to run through the training process; kept at 1 for base

model and at 5 for improved LSTM)

Benchmark Model

For this project i have used a Linear Regression model as its primary benchmark. As one of my goals is to understand the relative performance and implementation differences of machine learning versus deep learning models.

Following is the predicted results that i got from my benchmark model :



X-axis: Represents Tradings Days
Y-axis: Represents Closing Price Rupee
Green line: Close price
Blue Line: Predicted Close price

Train Score: 0.8627 MSE (0.9288 RMSE)
Test Score: 0.35408284 MSE (0.59504860 RMSE)

METHODOLOGY

Data Preprocessing

Acquiring and preprocessing the data for this project occurs in following sequence, much of which has been modularized into the **preprocess.py** file for importing and use across all notebooks:

- Downloaded the data from the National Stock Exchange and saved it in **v_bank_nifty.csv**.
- Remove unimportant features(date, high and low) :

	Item	Open	Close	Volume
0	0	11910.10	11823.30	390294634
1	1	11873.90	11914.05	447901993
2	2	11962.45	11906.20	283718253
3	3	11959.85	11965.60	332181958
4	4	11934.90	11922.70	303329445

	Item	Open	Close	Volume
2821	2821	6271.00	6279.10	139778977
2822	2822	6179.10	6274.30	147642063
2823	2823	6184.25	6178.55	156427760
2824	2824	6144.70	6179.40	110353305
2825	2825	6136.75	6144.35	81305366

- Normalised the data using **MinMaxScaler** helper function from Scikit-Learn.

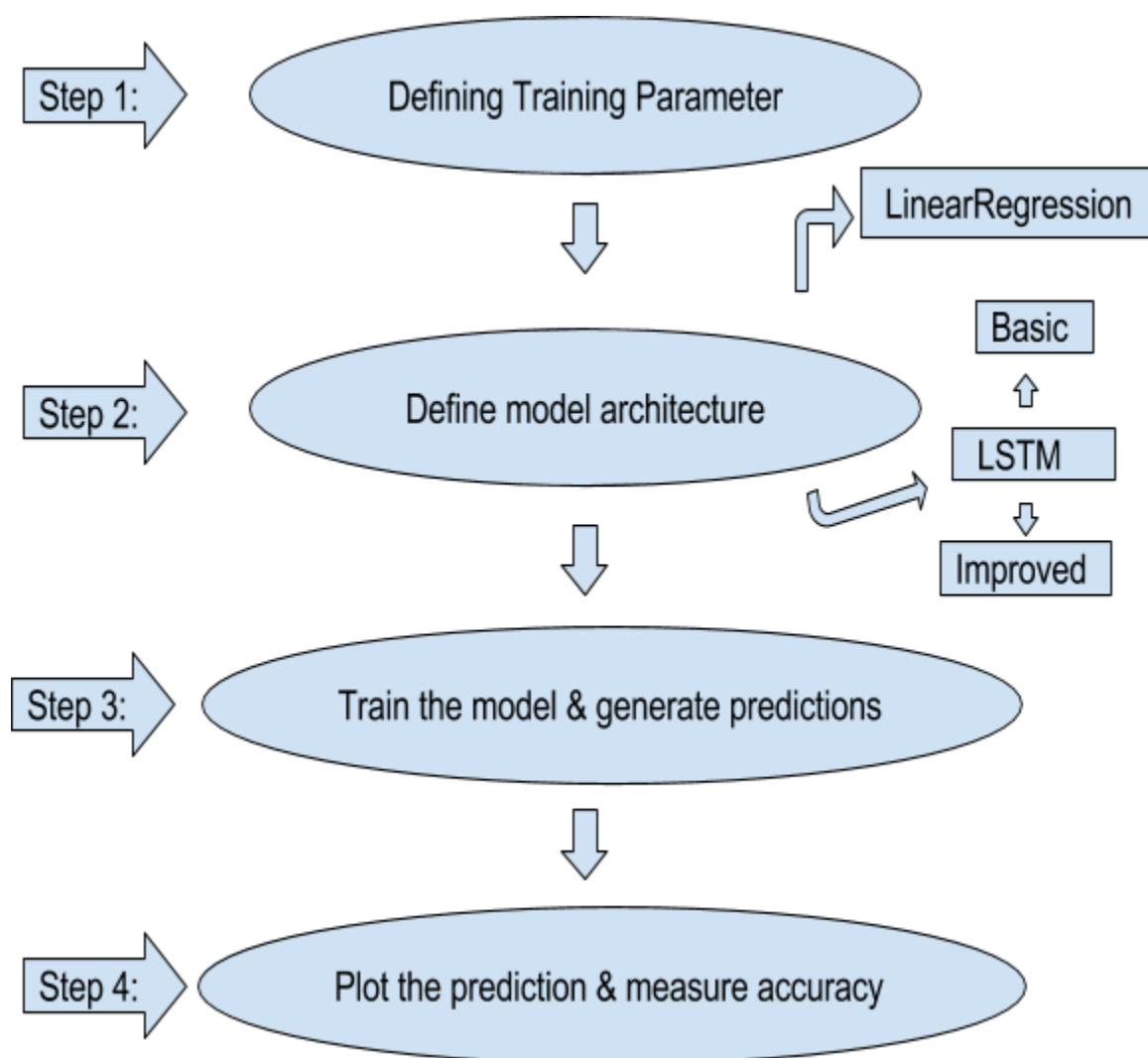
	Item	Open	Close	Volume
0	0	0.985036	0.972267	0.524830
1	1	0.981236	0.981755	0.602848
2	2	0.990532	0.980934	0.380493
3	3	0.990259	0.987145	0.446127
4	4	0.987640	0.982660	0.407052

- Stored the normalised data v_bank_nifty_preprocessed.csv file for future reusability.
- Splitted the dataset into the training (68.53%) and test (31.47%) datasets for linear regression model. The split was of following shape :
 - x_train (1936, 1)
 - y_train (1936, 1)
 - x_test (890, 1)
 - y_test (890, 1)

- Splitted the dataset into the training (82.95%) and test (17.05%) datasets for LSTM model. The Split was of following shape:
 - x_train (2320, 50, 3)
 - y_train (2320,)
 - x_test (396, 50, 3)
 - y_test (396,)

Implementation

Once the data has been downloaded and preprocessed, the implementation process occurs consistently through all three models as follow:



I have thoroughly specified all the steps to build, train and test model and its predictions in the notebook itself.

Some code implementation insight:

Benchmark model :

Step 1 : Split into train and test model :

```
X_train, X_test, y_train, y_test, label_range= sd.train_test_split_linear_regression(stocks)
```

Here I am calling a function defined in '**stock_data.py**' which splits the data for linear regression model. The function is as follows :

```
def train_test_split_linear_regression(stocks):
    """
    Split the data set into training and testing feature for Linear Regression Model
    :param stocks: whole data set containing ['Open','Close','Volume'] features
    :return: X_train : training sets of feature
    :return: X_test : test sets of feature
    :return: y_train: training sets of label
    :return: y_test: test sets of label
    :return: label_range: scaled range of label used in predicting price,
    """
    # Create numpy arrays for features and targets
    feature = []
    label = []

    # Convert dataframe columns to numpy arrays for scikit learn
    for index, row in stocks.iterrows():
        # print([np.array(row['Item'])])
        feature.append([row['Item']])
        label.append([row['Close']])

    # Regularize the feature and target arrays and store min/max of input data for rescaling later
    feature_bounds = [min(feature), max(feature)]
    feature_bounds = [feature_bounds[0][0], feature_bounds[1][0]]
    label_bounds = [min(label), max(label)]
    label_bounds = [label_bounds[0][0], label_bounds[1][0]]

    feature_scaled, feature_range = scale_range(np.array(feature), input_range=feature_bounds, target_range=[-1.0, 1.0])
    label_scaled, label_range = scale_range(np.array(label), input_range=label_bounds, target_range=[-1.0, 1.0])

    # Define Test/Train Split 80/20
    split = .315
    split = int(math.floor(len(stocks['Item']) * split))

    # Set up training and test sets
    X_train = feature_scaled[:-split]
    X_test = feature_scaled[-split:]

    y_train = label_scaled[:-split]
    y_test = label_scaled[-split:]

    return X_train, X_test, y_train, y_test, label_range
```

Step 2: In this step model is built using **scikit-learn linear_model**¹⁰ library.

```
model = LinearRegressionModel.build_model(X_train,y_train)
```

Here I am calling a function defined in '**LinearRegressionModel.py**' which builds the model for the project. The screenshot of the function is as follows:

```
def build_model(X, y):
    """
    build a linear regression model using sklearn.linear_model
    :param X: Feature dataset
    :param y: label dataset
    :return: a linear regression model
    """
    linear_mod = linear_model.LinearRegression() # defining the linear regression model
    X = np.reshape(X, (X.shape[0], 1))
    y = np.reshape(y, (y.shape[0], 1))
    linear_mod.fit(X, y) # fitting the data points in the model

    return linear_mod
```

Step 3: Now it's time to predict the prices for given test datasets.

```
predictions = LinearRegressionModel.predict_prices(model,X_test, label_range)
```

The screenshot of the function is as follows, it is defined in '**LinearRegressionModel.py**':

```
def predict_prices(model, x, label_range):
    """
    Predict the label for given test sets
    :param model: a linear regression model
    :param x: testing features
    :param label_range: normalised range of label data
    :return: predicted labels for given features
    """
    x = np.reshape(x, (x.shape[0], 1))
    predicted_price = model.predict(x)
    predictions_rescaled, re_range = sd.scale_range(predicted_price, input_range=[-1.0, 1.0], target_range=label_range)

    return predictions_rescaled.flatten()
```

Step 4: Finally calculate the test score and plot the results of benchmark model.

```
: trainScore = mean_squared_error(X_train, y_train)
print('Train Score: %.4f MSE (%.4f RMSE)' % (trainScore, math.sqrt(trainScore)))

testScore = mean_squared_error(predictions, y_test)
print('Test Score: %.8f MSE (%.8f RMSE)' % (testScore, math.sqrt(testScore)))

Train Score: 0.8627 MSE (0.9288 RMSE)
Test Score: 0.35408284 MSE (0.59504860 RMSE)
```

Improved LSTM model :

Step 1 : Split into train and test model :

Note : The same set of training and testing data is used for improved LSTM as is used with basic LSTM.

Step 2 : Build an improved LSTM model :

```
# Set up hyperparameters
batch_size = 100
epochs = 5

# build improved lstm model
model = lstm.build_improved_model( X_train.shape[-1], output_dim = unroll_length, return_sequences=True)
```

Here I am calling a function defined in '**lstm.py**' which builds the improved lstm model for the project. The screenshot of the function is as follows:

NOTE: The function uses **keras Longshort term memory**¹¹ library to implement LSTM model.

I have increased the batch_size to 512 from 1

Epochs from 1 to 5 for my improved LSTM model

Also in the function i have add increased the no of nodes in hidden layer to 128 from 100 and have added a drop out of 0.2 to all the layers.

```
def build_improved_model(input_dim, output_dim, return_sequences):
    """
    Builds an improved Long Short term memory model using keras.layers.recurrent.lstm
    :param input_dim: input dimension of model
    :param output_dim: output dimension of model
    :param return_sequences: return sequence for the model
    :return: a 3 layered LSTM model
    """
    model = Sequential()
    model.add(LSTM(
        input_shape=(None, input_dim),
        units=output_dim,
        return_sequences=return_sequences))

    model.add(Dropout(0.2))

    model.add(LSTM(
        128,
        return_sequences=False))

    model.add(Dropout(0.2))

    model.add(Dense(
        units=1))
    model.add(Activation('linear'))

    return model
```

Step 3: We now need to train our model.

```
model.fit(X_train,
        y_train,
        batch_size=batch_size,
        epochs=epochs,
        verbose=2,
        validation_split=0.05
    )
```

```
Train on 2204 samples, validate on 116 samples
Epoch 1/5
  - 9s - loss: 0.0422 - val_loss: 0.0112
Epoch 2/5
  - 6s - loss: 0.0054 - val_loss: 0.0023
Epoch 3/5
  - 6s - loss: 0.0023 - val_loss: 3.0196e-04
Epoch 4/5
  - 6s - loss: 0.0022 - val_loss: 3.1358e-04
Epoch 5/5
  - 6s - loss: 0.0019 - val_loss: 2.9985e-04
<keras.callbacks.History at 0x3707438>
```


I have used here a built in library function to train the model.

Step 4: Now it's time to predict the prices for given test datasets.

```
# Generate predictions
predictions = model.predict(X_test, batch_size=batch_size)
```

I have used a built-in function to predict the outcomes of the model.

Step 5: Finally calculate the test score and plot the results of improved LSTM model.

```
trainScore = model.evaluate(X_train, y_train, verbose=0)
print('Train Score: %.8f MSE (%.8f RMSE)' % (trainScore, math.sqrt(trainScore)))

testScore = model.evaluate(X_test, y_test, verbose=0)
print('Test Score: %.8f MSE (%.8f RMSE)' % (testScore, math.sqrt(testScore)))

Train Score: 0.00048111 MSE (0.02193420 RMSE)
Test Score: 0.00111130 MSE (0.03333616 RMSE)
```

Refinement

For this project i have worked on fine tuning parameters of LSTM to get better predictions. I did the improvement by testing and analysing each parameter and then selecting the final value for each of them.

To improve LSTM i have done following:

- Increased the number of hidden node from 100 to 128.
- Added Dropout of 0.2 at each layer of LSTM
- Increased batch size from 1 to 512
- Increased epochs from 1 to 5
- Added verbose = 2
- Made prediction with the batch size

Thus improved my mean squared error, for testing sets, from **0.35408284 MSE** to **0.00111130 MSE**.

The predicted plot difference can be seen as follows:

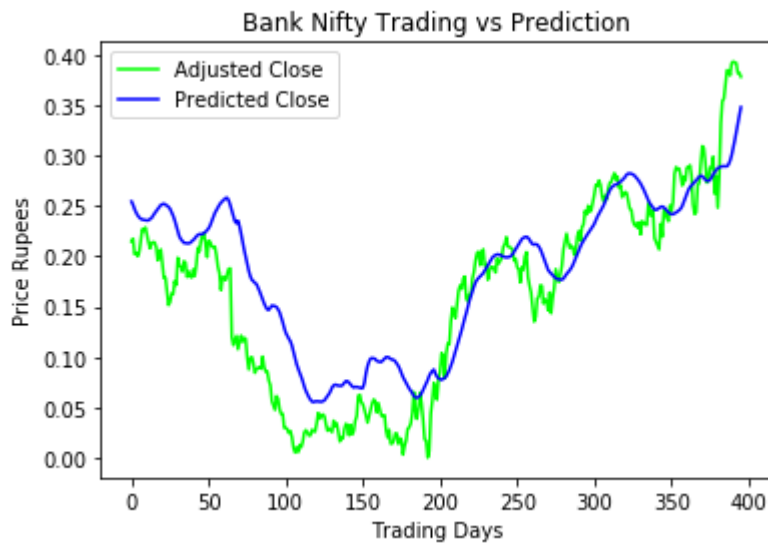


Fig : Plot For Adjusted Close and Predicted Close Prices for basic LSTM model

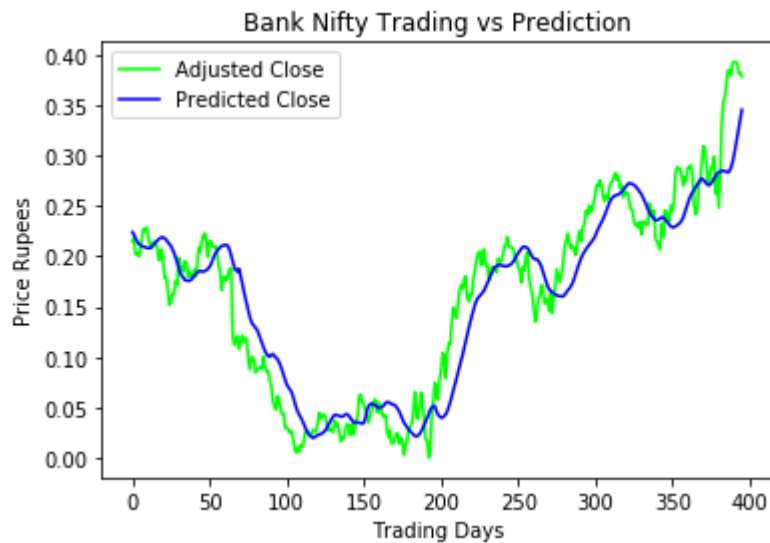


Fig : Plot For Adjusted Close and Predicted Close Prices for improved LSTM model

RESULT

Model Evaluation and Validation

With each model i have refined and fined tune my predictions and have reduced mean squared error significantly.

- For my first model using linear regression model:
 - Train Score: 0.8627 MSE (0.9288 RMSE)
 - Test Score: 0.35408284 MSE (0.59504860 RMSE)

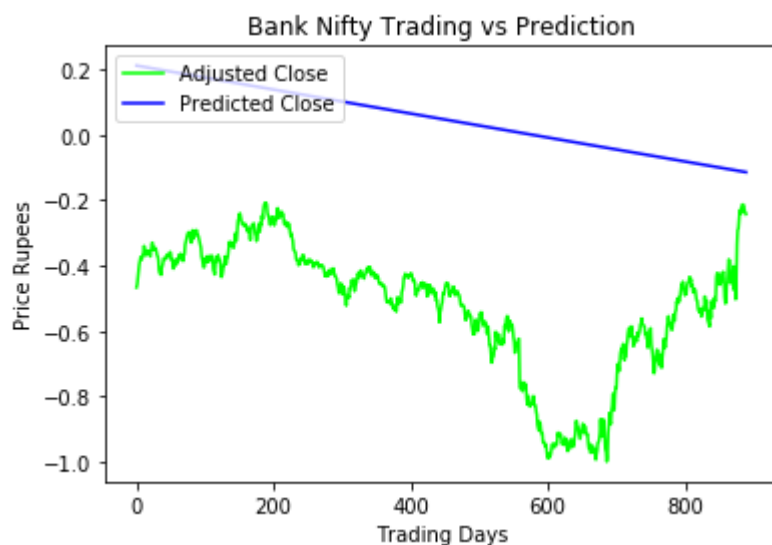


Fig: Plot of Linear Regression Model

- For my second model using basic Long-Short Term memory model:
 - Train Score: 0.00057929 MSE (0.02406843 RMSE)
 - Test Score: 0.00234837 MSE (0.04845994 RMSE)

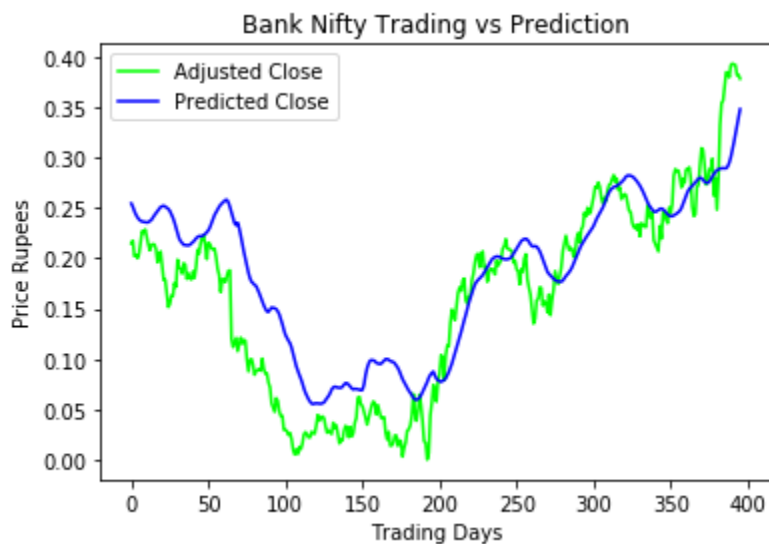


Fig: Plot of basic Long-Short Term Memory model

- For my third and final model, using improved Long-Short Term memory model:
 - Train Score: 0.00048111 MSE (0.02193420 RMSE)
 - Test Score: 0.00111130 MSE (0.03333616 RMSE)

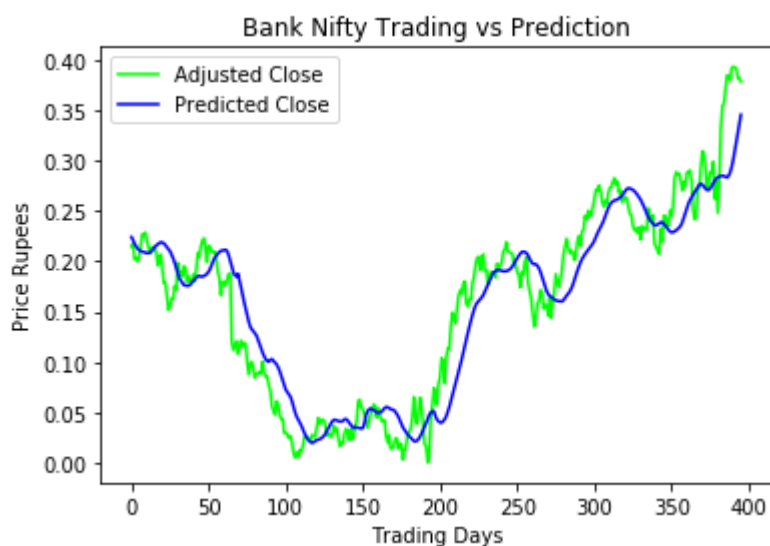


Fig: Plot of Improved Long-Short Term Memory Model

Robustness Check :

For checking the robustness of my final model i used an unseen data, i.e, data of Alphabet Inc. from May 1, 2019 to May 31, 2019. On predicting the values of unseen data i

got a decent result for the data. The results are as follows:

Test Score: 0.3308 MSE (0.5752 RMSE)

Justification

Comparing the benchmark model - Linear Regression to the final improved LSTM model, the Mean Squared Error improvement ranges. This significant decrease in error rate clearly shows that my final model has surpassed the basic and benchmark model.

Also the Average Delta Price between actual and predicted Closing Price values was:

Delta Price: 0.001111 - RMSE * Close Range

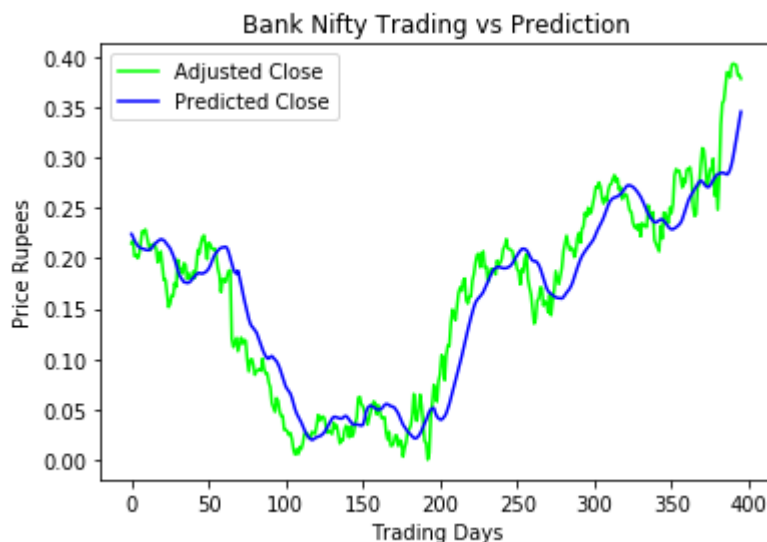
Which is less than one cent :)

CONCLUSION

Free-Form Visualization

I have already discussed all the important features of the datasets and their visualisation in one of the above sections. But to conclude my report i would choose my final model visualization, which is improved version of LSTM by fine tuning parameters. As i was very impressed on seeing how close i have gotten to the actual data, with a mean square error of just 0.0009. It was an 'Aha!' moment for me as i had to poke around a lot (really ALOT !! :P). But it was fun working on this project.

Fig: Plot of Improved Long-Short Term Memory Model



Reflection

To recap, the process undertaken in this project:

- Set Up Infrastructure
 - iPython Notebook
 - Incorporate required Libraries (Keras, Tensor flow, Pandas, Matplotlib, Sklearn, Numpy)
 - Git project organization
- Prepare Dataset
 - Incorporate data of Alphabet Inc company
 - Process the requested data into Pandas Dataframe
 - Develop function for normalizing data
 - Dataset used with a 80/20 split on training and test data across all models
- Develop Benchmark Model
 - Set up basic Linear Regression model with Scikit-Learn
 - Calibrate parameters
- Develop Basic LSTM Model
 - Set up basic LSTM model with Keras utilizing parameters from Benchmark Model
- Improve LSTM Model
 - Develop, document, and compare results using additional labels for the LSMT model
- Document and Visualize Results
- Plot Actual, Benchmark Predicted Values, and LSTM Predicted Values per time series
- Analyze and describe results for report.

I started this project with the hope to learn a completely new algorithm, i.e, Long-Short Term Memory and also to explore a real time series data sets. The final model really exceeded my expectation and have worked remarkably well. I am greatly satisfied with these results.

The major problem i faced during the implementation of project was exploring the data. It was toughest task. To convert data from raw format to preprocess data and then to split them into training and test data. All of these steps require a great deal of patience and very precise approach. Also i had to work around a lot to successfully use the data for 2 models, i.e, Linear Regression and Long-Short Term Memory, as both of them have

different inputs sizes. I read many research papers to get this final model right and i think it was all worth it :)

Improvement

Before starting my journey as Machine Learning Nanodegree Graduate i had no prior experience in python. In the beginning of this course to do everything with python, i had to google it. But now i have not only made 7 projects in python, i have explored many libraries along the ways and can use them very comfortably. This is all because of highly interactive videos and forum provided by Udacity. I am really happy and satisfied taking up this course.

And as there is scope of improvement in each individual so is the case with this project. This project though predicts closing prices with very minimum Mean Squared Error, still there are many things that are lagging in this project. Two of most important things are :

- There is no user interaction or interface provided in this project. A UI can be provided where user can check the value for future dates.
- The stocks used for this project are only of Bank nifty, we can surely add more stocks in the list so as to make this project more comprehensive.

I would definitely like to add these improvement to this project in future.

Thank You