

Domain Background:

Records of prices for traded commodities go back thousands of years. Merchants along popular silk routes would keep records of traded goods to try and predict price trends, so that they could benefit from them.

In finance the field of quantitative analysis is about 25 years old and even now it's not fully accepted, understood or widely used. Quantitative analysis is the study of how certain variables correlate with stock price behavior.

One of the first attempts at this was made in 1970 by two British statisticians named Box and Jenkins using mainframe computers. The only historical data they had access to were prices and volume. They called their model Arima and at the time it was extremely slow and expensive to run but by the 80s things started to get interesting. Spreadsheets were invented so that the firms could model company's financial performance and automated data collection became a reality. And with improvements in computing power, model could analyze data much faster, it was a renaissance on Wall Street, people were excited about the new possibilities.

Box and Jenkins work:

https://en.wikipedia.org/wiki/Box%E2%80%93Jenkins_method

Investor makes educated guess by analysing data, they read the news, study the company history, company trends and there are lot more points that go into stock price prediction. The prevailing theory is that stock prices are totally random and unpredictable.

As Burton Malkiel said,

"A blindfolded monkey throwing darts at a newspaper's financial pages could select a portfolio that would do just as well as one carefully selected by experts"

But that rises a question that why do top companies like Morgan Stanley and Citigroup hire a quantitative analyst to build predictive models? We have this idea of a trading floor being filled with adrenaline-infused man with loose ties running around yelling something into a phone. But these days we are more likely to see rows of machine learning experts quietly sitting in front of computers. Infact about 70% of all order on Wall Street are now placed by a software.

We are now living in the age of algorithms.

In the past few years we have seen lots of academic papers published using neural nets to predict stock prices with varying degrees of success. But until recently the ability to build these models has been restricted to academics who spend their days writing very complex code. Now

with libraries like Tensorflow, keras etc anyone can build powerful models trained on massive data sets.

Research papers on stock prediction using nn:

https://people.eecs.berkeley.edu/~akar/IITK_website/EE671/report_stock.pdf

<https://nseindia.com/content/research/FinalPaper206.pdf>

Motivation behind the project:

This project is about building a deep learning model to predict the stock prices. For this project I'll be using Keras framework with Tensorflow backend to predict the stock price of Google.

Problem statement :

To build a deep learning model using LST network to predict the stock price of Google using daily closing price of the S&P from 2004-08-19 to 2013-06-25 for training data. This is a series of data points indexed in time series. Our goal would be to predict the closing price for any given date after training. This is basically a regression problem and the reason for using LSTM along with Recurrent neural network is that it should be able to remember a sequence rather than just the preceding the value.

Mean squared error is used as a cost function or classification error and root mean squared error is used as a error metric. Now let's see why we have to use MSE instead of absolute mean squared error .

Mean Absolute error(MAE) and Root mean squared error are the most common metrics used to measure the accuracy for a continuous variable. Let's define these metrics

Mean Absolute Error: MAE measures the average magnitude of the errors in a set of predictions, without considering their direction. It's the average over the test sample of the absolute differences between prediction and actual observation where all individual differences have equal weights.

$$\text{MAE} = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

If the absolute value is not taken (the signs of the errors are not removed), the average error becomes the Mean Bias Error(MBE) and is usually intended to measure average model

bias. MBE can convey useful information, but should be interpreted cautiously because positive and negative errors will cancel out.

Root mean squared error (RMSE): RMSE is a quadratic scoring rule that also measures the average magnitude of the error. It's the square root of the average differences between prediction and actual observation

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

Comparison

Similarities: Both MAE and RMSE express average model prediction error in units of the variable of interest. Both metrics can range from 0 to ∞ and are indifferent to the direction of errors. They are negatively-oriented scores, which means lower values are better.

Differences: Taking the square root of the average squared errors has some interesting implications for RMSE. Since the errors are squared before they are averaged, the RMSE gives a relatively high weight to large errors. This means the RMSE should be more useful when large errors are particularly undesirable. The three tables below show examples where MAE is steady and RMSE increases as the variance associated with the frequency distribution of error magnitudes also increases.

CASE 1: Evenly distributed errors				CASE 2: Small variance in errors				CASE 3: Large error outlier			
ID	Error	Error	Error^2	ID	Error	Error	Error^2	ID	Error	Error	Error^2
1	2	2	4	1	1	1	1	1	0	0	0
2	2	2	4	2	1	1	1	2	0	0	0
3	2	2	4	3	1	1	1	3	0	0	0
4	2	2	4	4	1	1	1	4	0	0	0
5	2	2	4	5	1	1	1	5	0	0	0
6	2	2	4	6	3	3	9	6	0	0	0
7	2	2	4	7	3	3	9	7	0	0	0
8	2	2	4	8	3	3	9	8	0	0	0
9	2	2	4	9	3	3	9	9	0	0	0
10	2	2	4	10	3	3	9	10	20	20	400
MAE		RMSE		MAE		RMSE		MAE		RMSE	
2.000		2.000		2.000		2.236		2.000		6.325	

The last sentence is a little bit of a mouthful but I think is often incorrectly interpreted and important to highlight. RMSE does not necessarily increase with the variance of the errors. RMSE increases with the variance of the frequency distribution of error magnitudes.

To demonstrate, consider Case 4 and Case 5 in the tables below. Case 4 has an equal number of test errors of 0 and 5 and Case 5 has an equal number of test errors of 3 and 4. The variance of the errors is greater in Case 4 but the RMSE is the same for Case 4 and Case 5.

CASE 4: Errors = 0 or 5

ID	Error	Error	Error^2
1	5	5	25
2	0	0	0
3	5	5	25
4	0	0	0
5	5	5	25
6	0	0	0
7	5	5	25
8	0	0	0
9	5	5	25
10	0	0	0

var	MAE	RMSE
6.944	2.500	3.536

CASE 5: Errors = 3 or 4

ID	Error	Error	Error^2
1	3	3	9
2	4	4	16
3	3	3	9
4	4	4	16
5	3	3	9
6	4	4	16
7	3	3	9
8	4	4	16
9	3	3	9
10	4	4	16

var	MAE	RMSE
0.278	3.500	3.536

There may be cases where the variance of the frequency distribution of error magnitudes (still a mouthful) is of interest but in most cases (that I can think of) the variance of the errors is of more interest.

Another implication of the RMSE formula that is not often discussed has to do with sample size. Using MAE, we can put a lower and upper bound on RMSE.

1. $[MAE] \leq [RMSE]$. The RMSE result will always be larger or equal to the MAE. If all of the errors have the same magnitude, then $RMSE=MAE$.

2. $[RMSE] \leq [MAE * \sqrt{n}]$, where n is the number of test samples. The difference between RMSE and MAE is greatest when all of the prediction error comes from a single test sample. The squared error then equals to $[MAE^2 * n]$ for that single test sample and 0 for all other samples. Taking the square root, RMSE then equals to $[MAE * \sqrt{n}]$.

Focusing on the upper bound, this means that RMSE has a tendency to be increasingly larger than MAE as the test sample size increases. This can be problematic when comparing RMSE results calculated on different sized test samples, which is frequently the case in real world modeling.

Conclusion

RMSE has the benefit of penalizing large errors more so can be more appropriate in some cases, for example, if being off by 10 is more than twice as bad as being off by 5. But if being off by 10 is just twice as bad as being off by 5, then MAE is more appropriate.

From an interpretation standpoint, MAE is clearly the winner. RMSE does not describe average error alone and has other implications that are more difficult to tease out and understand.

On the other hand, one distinct advantage of RMSE over MAE is that RMSE avoids the use of taking the absolute value, which is undesirable in many mathematical calculations

And the main reason behind using RMSE over MSE is that RMSE returns values in their standard units.

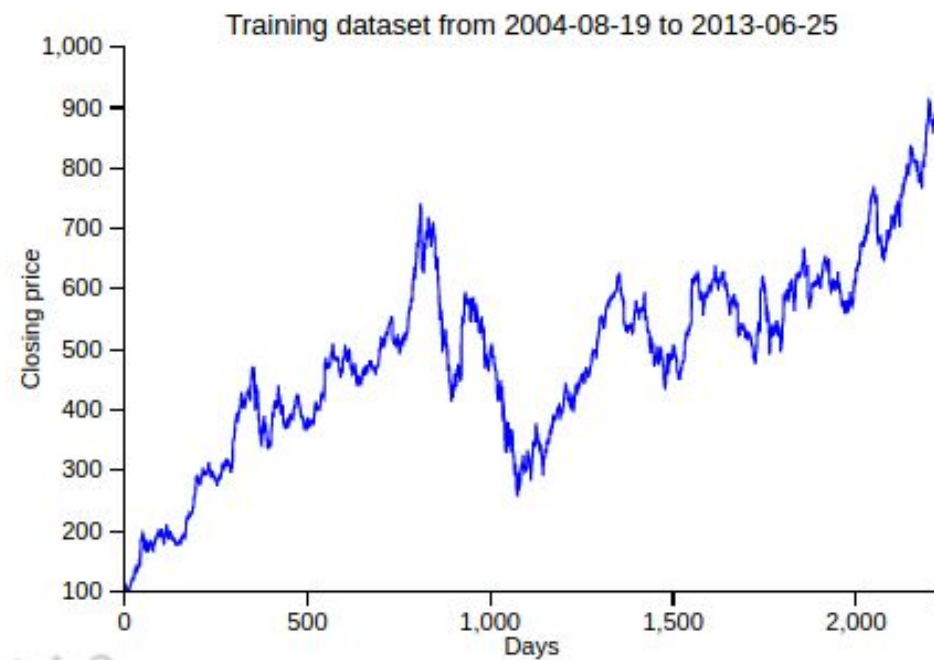
Data sets and inputs:

The data set consists of daily closing prices of Google from 2004-08-19 to 2013-06-25 for training purpose and Closing stock price from 2013-06-26 to 2016-01-05 is used as a testing dataset for testing the model's performance .Closing stock price from 2016-01-06 to 2017-04-11 is used as an another test case containing most recent google stock prices,used to ensure how the model can perform on unseen data. Here the only feature that we are passing to the learning model is the closing price of the Google stock price on various days. These features can extracted from other google finance page.

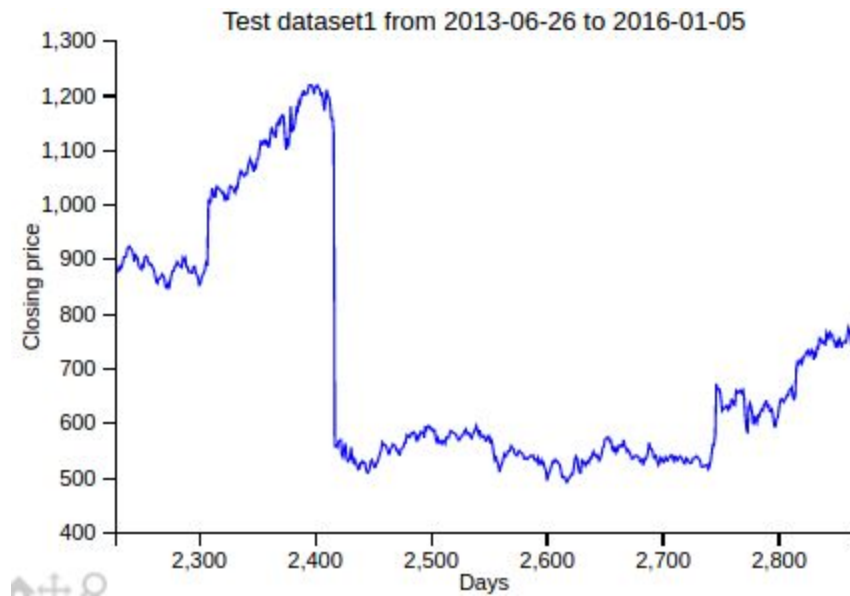
Google Finance:

<https://www.google.com/finance?q=google&ei=GJDnWJC4McaFuwTVgZrgCg>

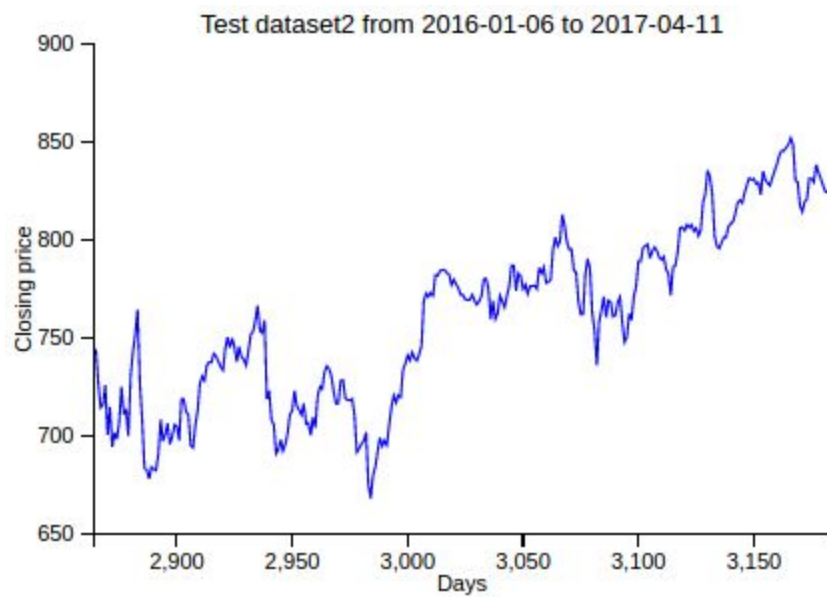
Training dataset:



Test dataset1:



Test dataset2:



[Dataset](https://github.com/satishjasthi/Udacity-MLNDP-Submissions/blob/master/CapstoneProject/Google2004-2017April11.csv)

[:https://github.com/satishjasthi/Udacity-MLNDP-Submissions/blob/master/CapstoneProject/Google2004-2017April11.csv](https://github.com/satishjasthi/Udacity-MLNDP-Submissions/blob/master/CapstoneProject/Google2004-2017April11.csv)

Let's take a look at the first 20 Closing stock prices of Google

	Close
0	100.340176
1	108.310183
2	109.400185
3	104.870176
4	106.000184
5	107.910182
6	106.150181
7	102.010172
8	102.370175
9	100.250171
10	101.510173
11	100.010169
12	101.580175
13	102.300173
14	102.310174
15	105.330177
16	107.500188
17	111.490192
18	112.000192
19	113.970198

What are the maximum, minimum closing prices and other statistics of dataset?

Maximum price: 1220.172036 units

Minimum price: 100.010169 units

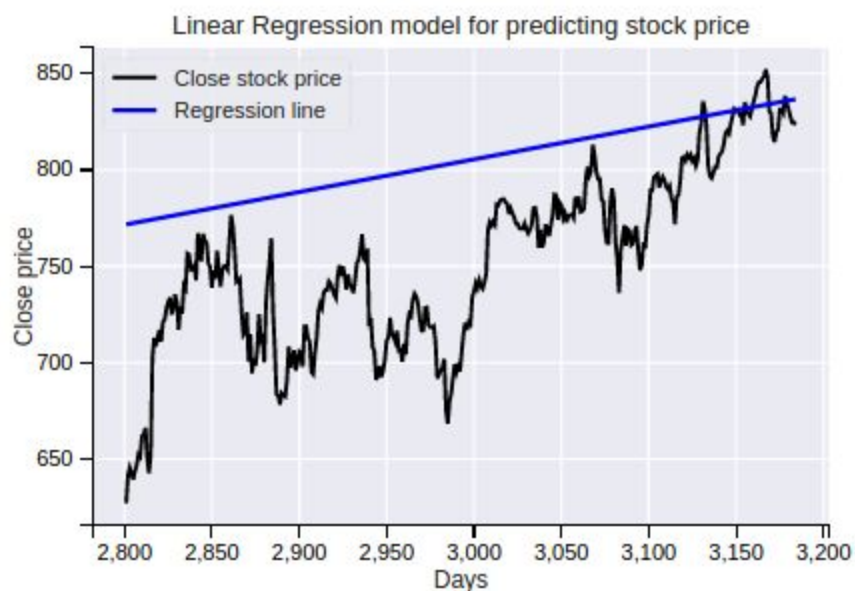
Mean price: 560.402579 units

Standard deviation : 198.763233 units

Abnormalities within the data: As we can see from the Closing price of Google stock figure that there few abnormalities within the data. For example consider the lowest price point on April 26 2006.

Benchmark model:

As a benchmark model I'll be using Linear Regression model against the deep learning model. In both cases error function or cost function and the error metric both is root mean square error itself.



Linear regression model for this dataset can be found here:

<https://github.com/satishjasthi/Udacity-MLNDP-Submissions/blob/master/CapstoneProject/Linear.ipynb>

Workflow of project:

1. Data preprocessing: Initially the input data is normalised, rather than directly using the raw values normalised values helps in faster convergence of the algorithm. When our model predicts the values we'll then denormalise the data to get a real world number out of it.

2. After we model our data and estimate the skill of our model on the training dataset, we need to get an idea of the skill of the model on new unseen data. For a normal classification or regression problem, we would do this using cross validation.

However with time series data, the sequence of values is important. So we use a function called `create_dataset` which will create a dataset where X is the stock price at any given time (t) and Y is the stock price at the next time step (t + 1).

For example:

If dataset = 112,118,132,129,121,135

Then

X:112,118,132,129,121

Y:118,132,129,121,135

3. Now we'll rearrange the data so that it matches with the format of the input of a LSTM network, which is [samples, time step, features]. Once the data is in the right format we can train the LSTM network with the training data.

Configuration of LSTM network:

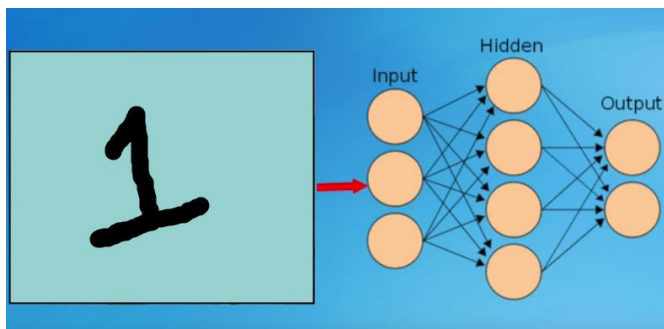
Number of LSTM units = 5

Number of epochs = 10

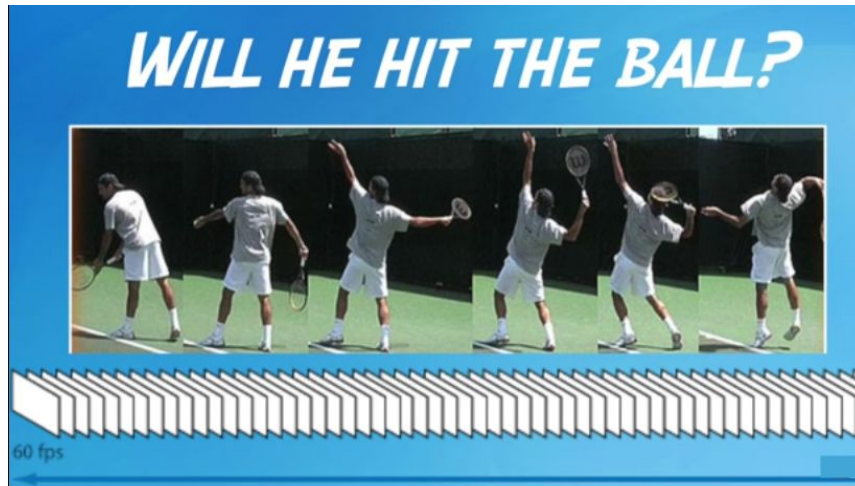
Batch_size = 1

Let's see how LSTM can learn to forecast time-series data :

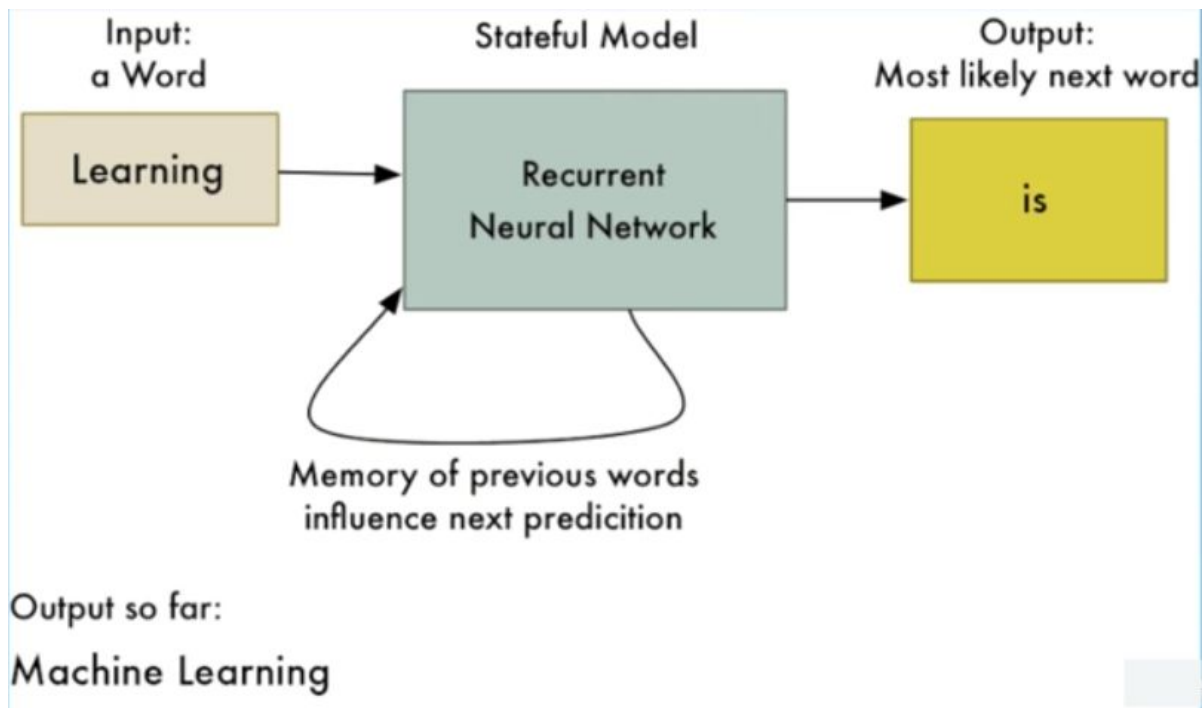
In case of sequences, memory matters because it uses conditional memory. However feedforward neural networks don't have this property. They accept a fixed size vector as an input, like an image.



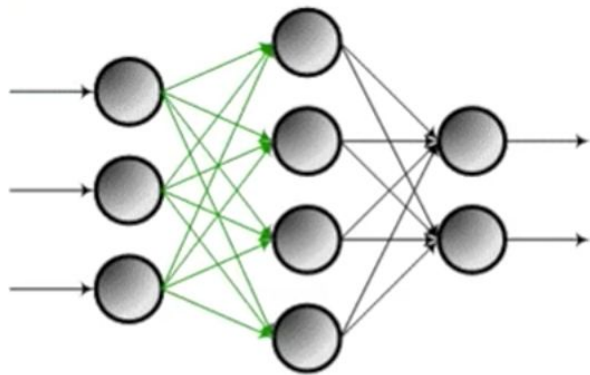
So we couldn't use it to predict the next frame in a movie because that would require a sequence of image vectors as inputs. And not just one since the probability of a certain event happening would depend on what happened every frame before it.



So we need a way to allow information to persist and that's why we can use a recurrent neural net. Recurrent neural nets can accept sequences of vectors as inputs.

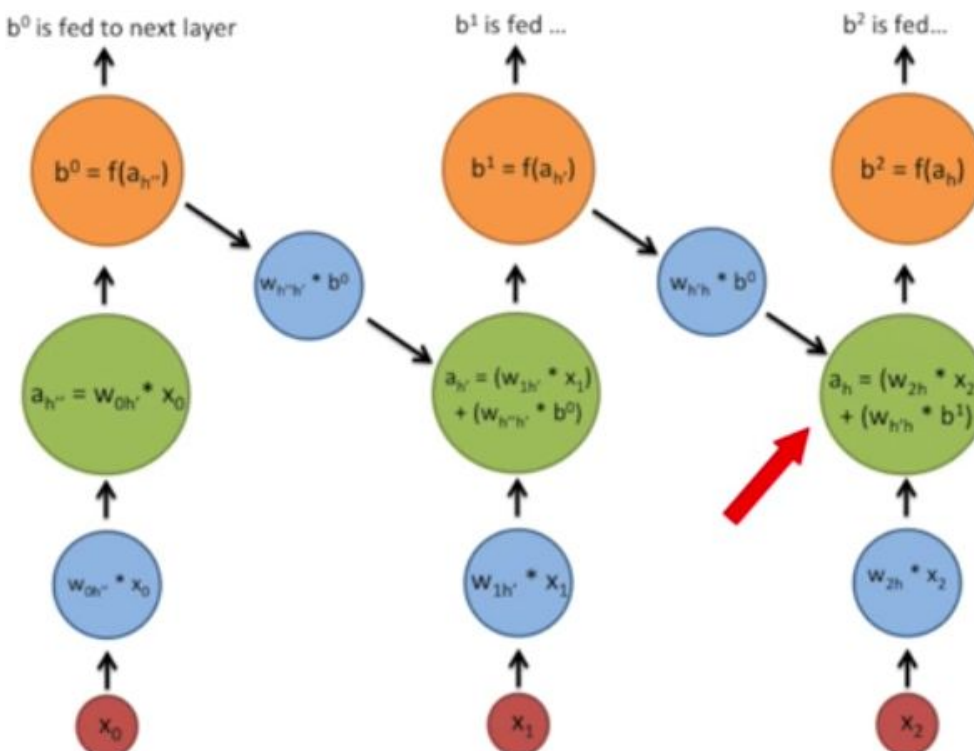


So as we discussed before, in feedforward neural nets, the hidden layer's weights are based only on the input data.



input -> hidden -> output

But in a recurrent net, the hidden layer is a combination of the input data at the current time step and the hidden layer at a previous time step. The hidden layer is constantly changing as it gets more inputs and the only way to reach these hidden states is with correct sequence inputs. This is how memory is incorporated in and we can model this process mathematically as



$$h_t = \phi (W * X_t + U * h(t-1))$$

Where h_t = hidden state at a given time step

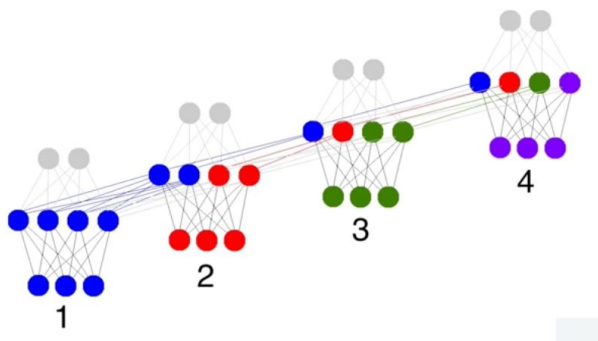
X_t = input at the given time step

W = weight matrix similar to one used in feedforward neural nets

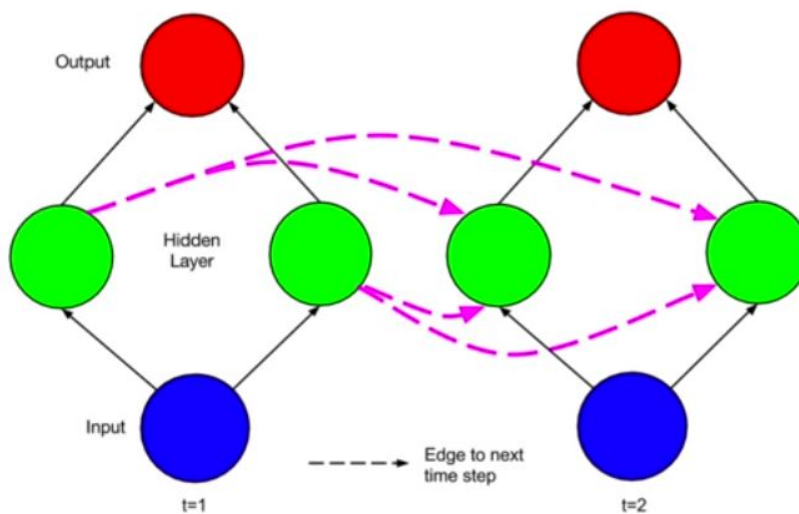
h_{t-1} = hidden state of previous time step

U = transition matrix

And because this feedback loop is occurring at every time step in the series each hidden state has traces of not only the previous hidden state but also of all of those that preceded it. That's why we call it Recurrent Neural net.



In a way, we can think of it as copies of the same network each passing a message to the next.

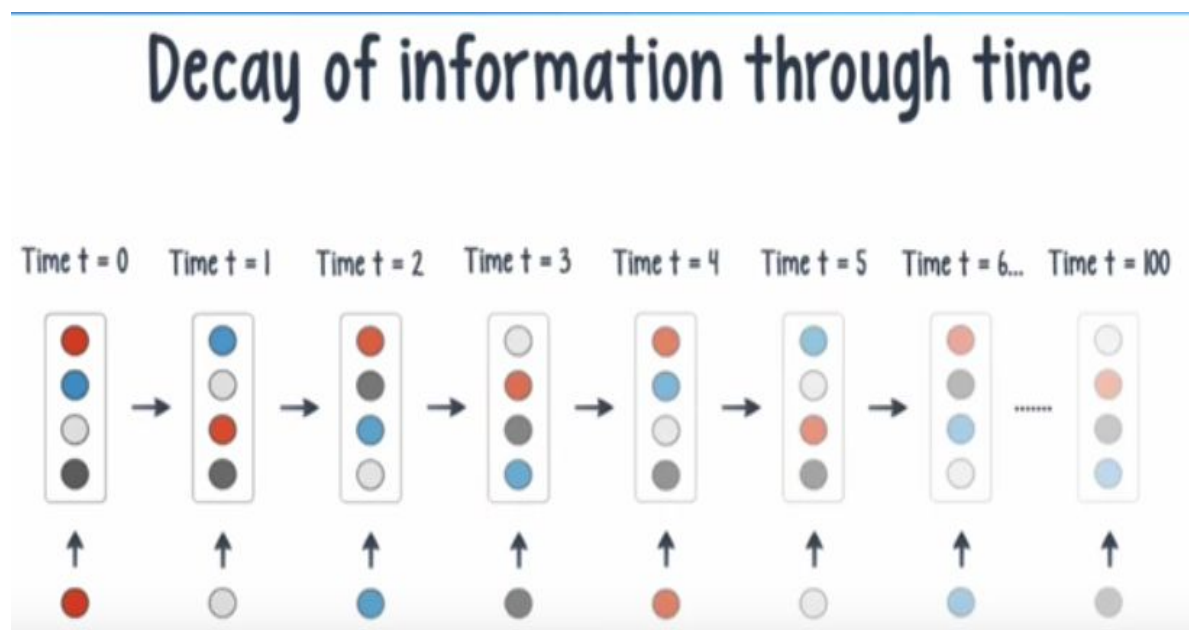


So Recurrent neural nets can connect the previous data with the present task. But there is problem in using RNN, Consider an example

“ I hope Senpai will notice me. Like walking through a barren street in a crumbling ghost town, isolation can feel melancholy and hopeless. Yet, all it takes is an ordinary flower bud amidst the desolation to show life really can exist anywhere. This is similar to Stephen’s journey in The Samurai’s Garden. This novel is about an ailing Chinese boy named Stephen who goes to Japanese village during a time of war between Japan and China to recover from his disease. She is my friend. He is my Senpai.””

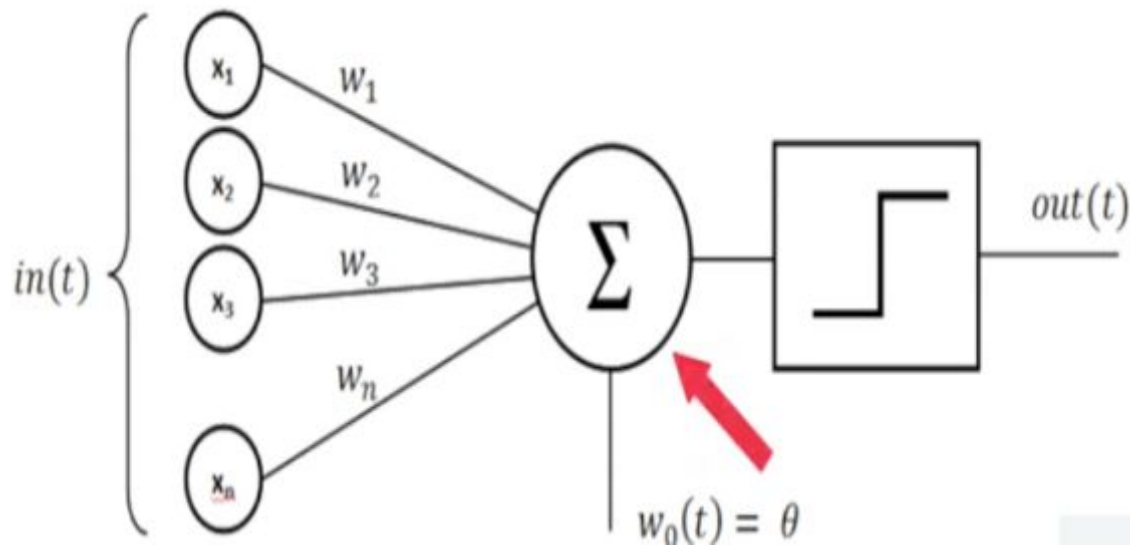
Let’s say we want to train a model predict the last word Senpai given all previous words. We need the context from the very beginning of the sequence to know that this word is probably Senpai.

In a regular neural net memories become more subtle as they fade into the past since the error signal from the later time steps doesn’t make it far enough back in time to influence the network at early time steps during the back propagation. This is what we call vanishing gradient problem.

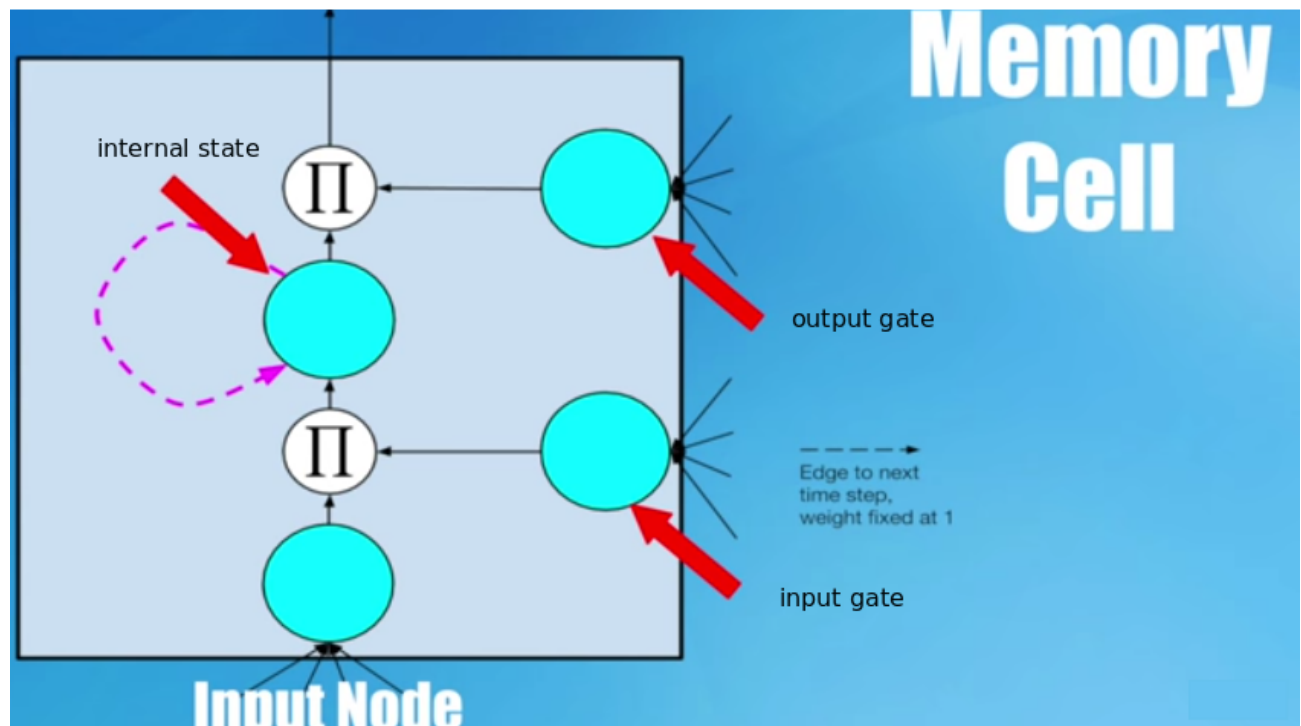


Vanishing gradient problem by Yoshua Bengio:
<http://www-dsi.ing.unifi.it/~paolo/ps/tnn-94-gradient.pdf>

A popular solution to this is a modification to recurrent neural nets called Long Short Term Memory (LSTM). Normally neurons are units that apply an activation function, like a sigmoid to a linear combination of their inputs.



In an LSTM recurrent net, we instead replace these neurons with memory cells. Each cell has an input gate, an output gate and an internal state that feeds into itself across time step with a constant weight of 1. This eliminates the vanishing gradient problem since any gradient that flows into this self recurring unit during back propagation is preserved indefinitely, So this helps RNN to remember long term dependencies.



Ok now that we are clear why we add LSTM in code, now we can proceed further.

4. Hyper parameter tuning: Hyperparameters of a LSTM network are:

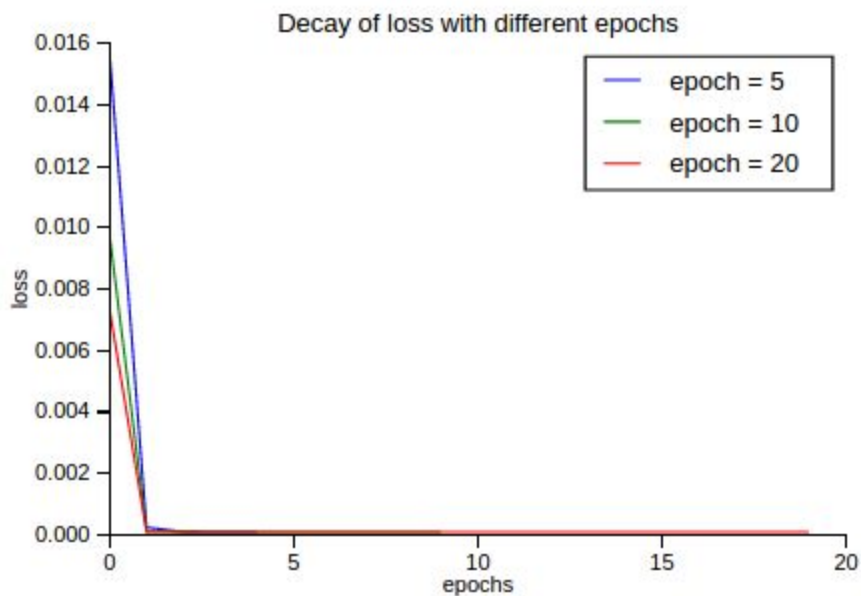
1. Epochs
2. Number of LSTM units
3. Batch_size

Epochs : one epoch equals one forward pass and one backward pass of *all* the training dataset through the network.

Number of LSTM units: is the number of LSTM neurons within a single hidden layer .

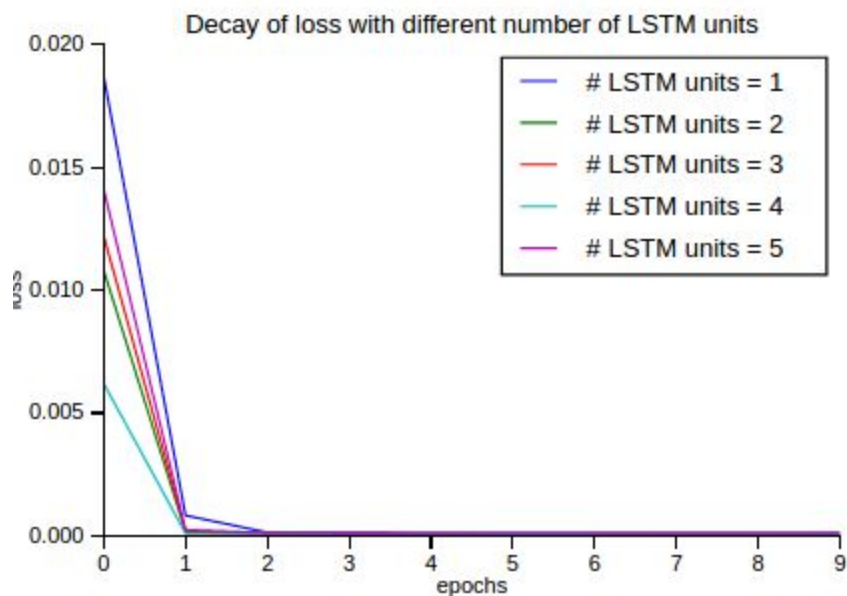
Batch_size: is the number of training examples in one forward/backward pass. The higher the batch size, the more memory space you'll need.

Now let's see the affect on epochs on model's performance:



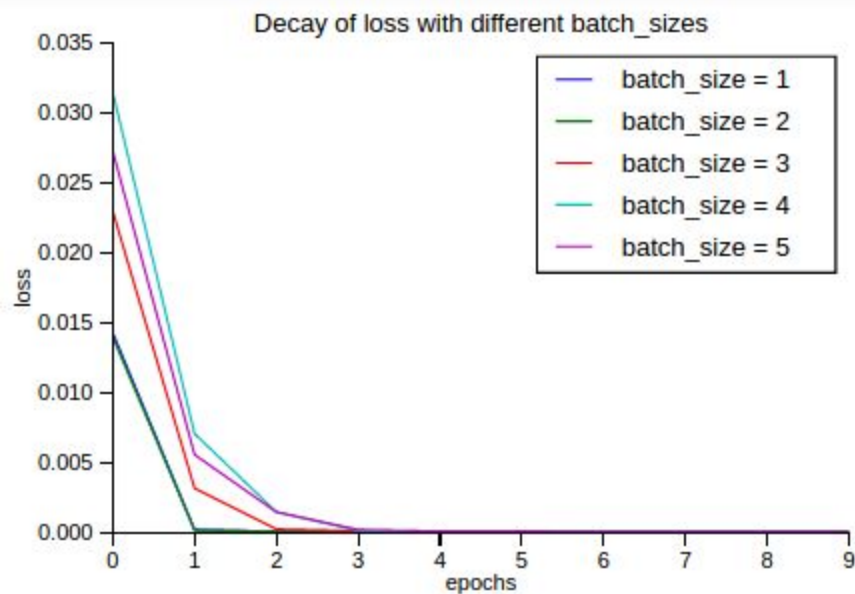
As we can see from the graph that loss function becomes almost equal to zero after epochs= 1, Since loss value doesn't change with the increase in epochs, we can choose any epochs > 3, I'll be choosing epochs = 10.

Now let's see how number of LSTM units affects the model's performance



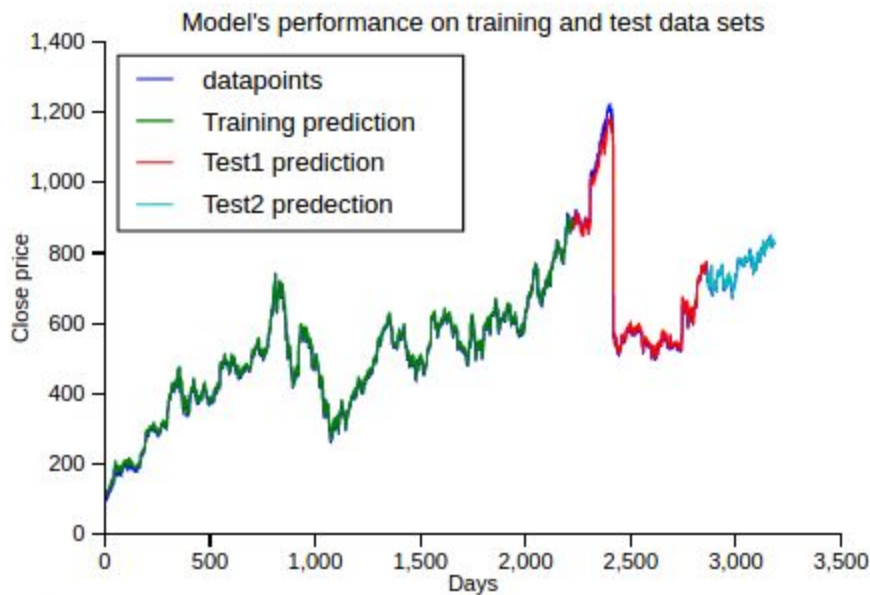
As we can see from the graph that all LSTM unit's losses are almost equal to zero after 2 epochs. So I'll be choosing 5 LSTM units.

Now let's see how batch_size affects the model's performance

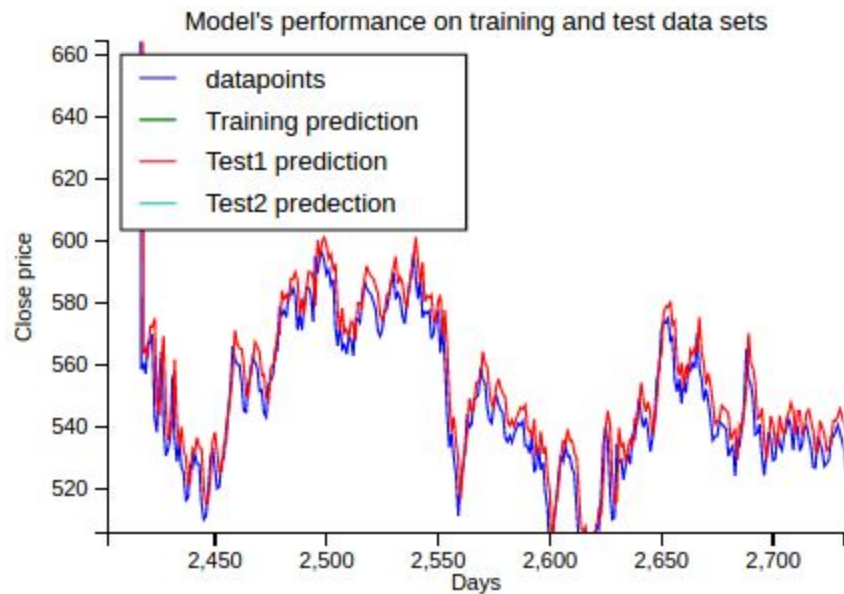


As we can see from the graph that as batch_size increases the model is taking more epochs to converge loss value to zero. So I'd be using a batch_size = 1, since its loss value becomes almost zero within 2 epochs.

5. Once the model is trained, its performance will be tested on the test data1 dataset and again on test data2 dataset to ensure the model's performance is good even on unseen recent stock price data.

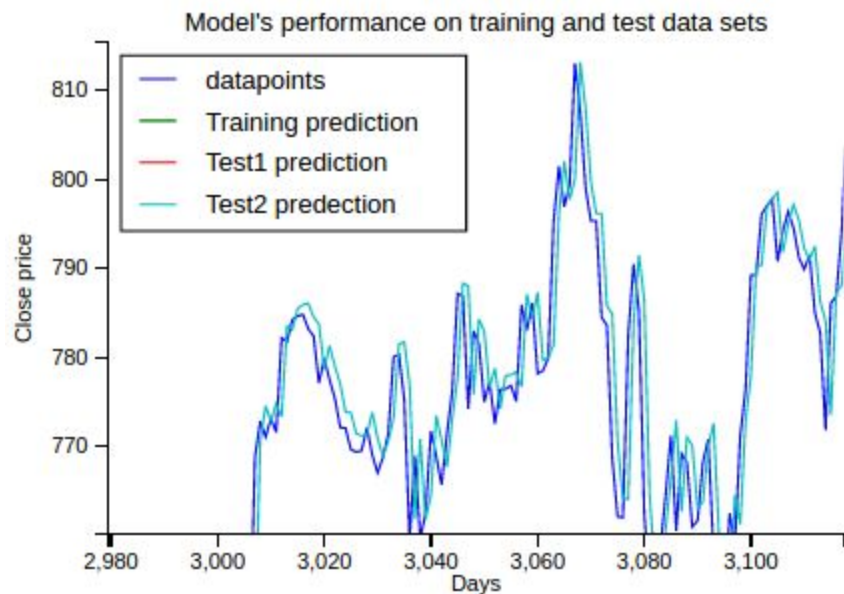


Closer look at the model's prediction on test data1:



As we can see from the graph above, the model's prediction is matching approximately with the actual data.

Closer look at the model's prediction on test data2:



Even though the model has new seen the test data2 set it is able to predict the stock prices which are very near to actual prices

Evaluation metrics:

To track and reduce the loss value I'm using mean squared error as a cost function and Gradient descent as an optimizer. Here root means squared error is used as an error metric.

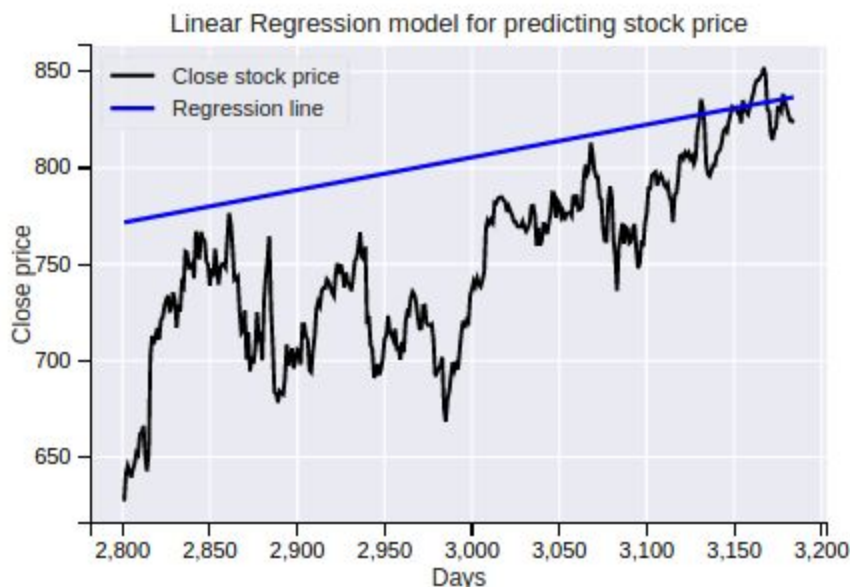
If we take a look at the RMSE value for training dataset it is 10.82 units and for the first test dataset ie test data1 RMSE value is 27.5 units. The reason for test data1 score to be higher than any other score is that, within the range of samples that test1 contains there is sudden increase in the stock price in the year 2013, I think the model has approximated the increase in stock price and the difference between this approximated value and actual value might have contributed to increase in test data1 score.

And the RMSE value for test data2 is 8.84 units, this value is lesser because the model actually performing really well since it's predictions are pretty closer to the actual prices.

Comparison with Benchmark model:

As a benchmark model I have used Linear Regression model with root mean squared error as a cost function as well as a error metric.

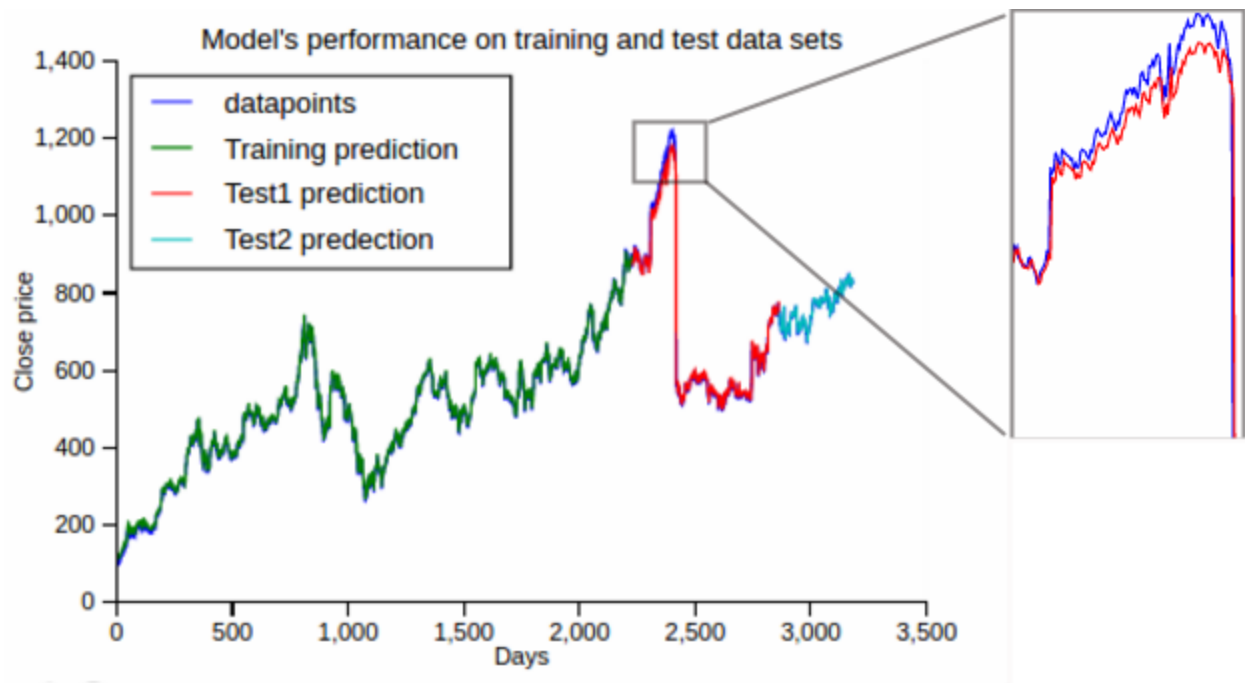
Performance of Linear Regression model in predicting the stock prices is shown below:



As we can see the performance of Linear regression model is nowhere near the performance of the LSTM network. This difference can be inferred from the RMSE value of Linear Regression model, which is 60.89 units. Even if we consider the highest RMSE values out of two test sets, ie test data1 and test data2 in case of LSTM networks, RMSE of Linear regression is more than twice the RMSE of LSTM network which clearly indicates the poor performance of Linear regression on this dataset.

Conclusion :

LSTM networks can be used to forecast on time series data like the stock price data that I have used with in this project. Whenever we analyse a stock data the key point is to maintain a long term sequence dependence, ie to connect between previous data and the present data. LSTM networks are really good at this since they are not affected by vanishing gradient problem.



As we can see from the graph the model performs well even when there is a sudden abnormality within the data. Even though the prediction is approximate, the model can capture the sudden trend (like rise or drop) within the stock data.

Further improvements:

- One can experiment with variation in LSTM network architecture to improve the model's overall performance. Various versions of LSTM for regression are:
 - LSTM for regression using the window method
 - LSTM for regression with Time steps
 - LSTM with Memory between batches
 - Stacked LSTM with memory between batches
- One can extend the project to an API which can use the LSTM model to provide appropriate recommendations to stockholders.

References:

Frame works:

*Tensor flow: <https://www.tensorflow.org/>

*Scikit learn: <http://scikit-learn.org/stable/>

*Keras: <https://keras.io/>

Research Papers:

*Learning long-term dependencies with gradient descent is difficult by Y Bengio:

<http://www.dsi.ing.unifi.it/~paolo/ps/tnn-94-gradient.pdf>

Blogs:

*Understanding LSTM networks by Colah:

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

*Deep learning documentation: <http://deeplearning.net/tutorial/lstm.html>