

An Application Predicting Stock Price Trends using Neural Network

Yi Man
yi.man@colorado.edu

Han Ngo
han.ngo@colorado.edu

1. Problem Motivation

Stock price trend prediction is an active research area for data scientists. Therefore, in recent years, significant efforts have been put into developing models that can predict for future trend of a specific stock or overall market. Also, some of the researchers showed that there is a strong relationship between news article about a company and its stock prices fluctuations. We create Stockalyzer to help us project the tendency in stock price changes with regards to past datasets and public reactions. Using deep learning model and text mining techniques to forecast the interactivity between change in price and released news will help us foresee stock trend movement.

The overall purpose of our work will help clarify such following questions:

1. Given a set of symbols (company indexes), how would the model predict stock price change in a specified time frame?
2. How efficient does the neural network work towards different sets of data? Does it cause too much over-fitting or under-fitting?
3. What would be a good decision in different situations, e.g. if we decided to previously buy the stock and there is an increase in the stock price?
4. Once decision is made, based on historical data source, how does the program generate analysis for us to see if the action yields any profits?

2. Literature Survey

Financial stock analysis is a well-studied area that has seen much prior work. There are approaches that have been utilized to build efficient predictive models from Deep Learning technologies such as Neural Network [1], Support Vector Machines [2], Q-Learning, etc. to identify the correlation between past datasets and current public opinions' impact on stock movements. A much more detailed and well-developed version has been developed by Facebook team using additive models to predict certain sets of stocks in current years.

Moreover, collaborative open-source work has also been performed by Kaggle community members and competition candidates interested in exploring this domain. While many of these projects provide bases analysis of stock trends as a function of time, Facebook's Prophet model seems to be the most popular library to be used currently for its animating visualization and high validation accuracy.

3. Proposed Work

Data Collection and Integration

We have collected constituent financials and download the full S&P 500 data from Yahoo! Finance-GSPC to our dataset directory. The Python script load.py will

fetch the prices of individual stocks in S&P 500 for 4000 days maximum, each saved to our dataset directory.

Collecting and analyzing tweets

We will fetch the filtered news from Reuters via the Twitter streaming API. The sentiment analysis takes the input from package then looks up the sentiment using a sentiment analysis API and stores the scores as a new feature in our tables.

Training data preprocessing

We will use content in one sliding windows to make prediction for the next and ensuring there is no overlap between two consecutive windows. Layer inputs and outputs will form an unrolled version of recurrent neural network for training on Tensorflow. After that, for model validation and accuracy evaluation, we will use the latest 10% of data as the testing data.

Data normalization

We will normalize the close prices in each sliding window. The task becomes predicting the relative change rates instead of the absolute values. In a normalized sliding window at time t , all the values are divided by the previous price to help prevent out-of-scale issues.

Model construction

The model is expected to learn the price sequences of different stocks in time. Due to the different underlying patterns, we will use embedding. We will create a Tensorflow graph using input variables from datasets that we have and the available optimizing algorithm built in TF library for gradient descent optimization.

For our focus in data with time frames, we found that recurrent neural network come accordingly useful. We will construct the model first without no hyper-parameters being initiated. We will mostly run some trials and errors with different set of variables and depend on Keras library for creating LSTM (Long Short Term memory) layers combined with a number of fully connected ones.

Training/Testing

We will initiate different hyperparameters, start a running tf session, and saved our train model at the end which will be later used to evaluate on the last 10% of our datasets. There are a lot of approaches that machine learning enthusiasts pursue to forecast stock trends as the aforementioned in literature survey section. For the sake of learning the whole process of data mining, we decided to try Neural Network approach. Moreover, reading researches about the relevance between public opinions and stock movement motivates us to discover add on a new feature to the training data sets.

4. Data Source

Index file:

<https://raw.githubusercontent.com/datasets/s-and-p-500-companies/master/data/constituents-financials.csv/>

Data file for 500 companies:

<https://finance.yahoo.com/quote/%5EGSPC/history?p=%5EGSPC/>

5. Evaluation Methods

We will use mean square error as the loss metric and the RMSPropOptimizer algorithm

which is built in Tensorflow for gradient descent optimization

6. Tools

For communication and sharing purposes, our group will work on the project through Github. The code will be mostly based on Python. We plan to utilize Pandas and Numpy libraries to aid data cleaning process. In addition, we plan to use WEKA for text analysis, Tableau and Matplotlib for basic graphing visualizations. Most interactions will be based on Ubuntu CLI and data training process will also make use of our GPU instance on AWS.

7. Milestones

Date	Responsibilities	Updated status
March 15th	Data collection, preprocessing + cleaning	Done
March 30th	Data integration + normalization	Done. Yi still took care of the tweets analyzing part, Hannie moved on with the raw data from SP500 in the meantime.
April 7th	Model construction + set up TF environment Work on progress report	Installed Anaconda environment to use Keras/ Tensorflow. Constructed a simple model but fitting take quite a bit long time.

April 14th	Training + testing on our dataset	Almost. Better improve the network structure to get more accordance between prediction and ground truth. Probably need to set up GPU to optimize time
April 21th	Finalize visualizations, write up our findings	To-do Visualize on difference SP500 companies. Try with different number of epochs
April 28th	Finish group presentation slides and report	To-do

9. Results so far

We wrote a script to download the needed data up-to-date dated back to 4000 days ago for each of the 500 S&P companies. They are all stored in folder 'data'. We also completed the code to load the training data from csv files into the appropriately shaped numpy array and split them in batches of training/testing data with proportion of 9/1.

After we have the data ready, we noticed that all the data scale in each column differ a lot. So before throwing batches of data into the Recurrent Neural Network, we used Standard Scaler in sklearn to help us complete the normalization preprocessing step.

Now that we had the normalized dataset, we wanted the LSTM to learn the sin wave from a set window size of data that we would feed it and then hopefully we could ask the LSTM to predict the next N-steps in the series and it will keep spitting out the sin wave. After some researches, we noted that one of the important things was also to normalize the sliding window. Because running the adjusted returns of a stock index through a network would make the optimization process not converge to any sort of optimums for such large numbers. So to combat this, we have researched ways and decided to take each n-sized window of training/testing data and normalize each one to reflect percentage changes from the start of that window (so the data at point $i=0$ will always be 0). We used the following equations to normalise and subsequently de-normalise at the end of the prediction process to get a real world number out of the prediction.

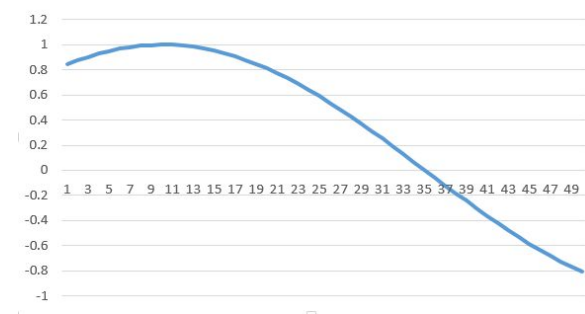
Normalization formula:

$$n_i = \left(\frac{p_i}{p_0} \right) - 1$$

Denormalization formula:

$$p_i = p_0(n_i + 1)$$

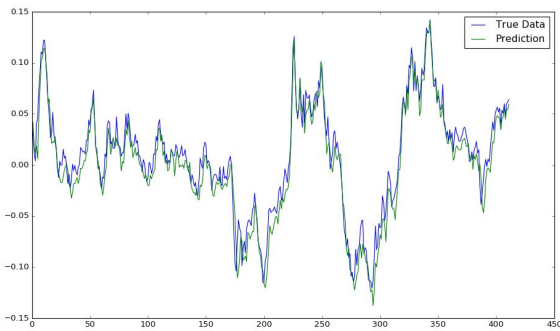
About the model, the way Keras LSTM layers work is by taking in a numpy array of 3 dimensions (N, W, F) where N is the number of training sequences, W is the sequence length and F is the number of features of each sequence. We chose to go with a sequence length (read window size) of 50 which allows for the network so get glimpses of the shape of the sin wave at each sequence and hence hopefully teach itself to build up a pattern of the sequences based on the prior window received. The sequences themselves are sliding windows and hence shift by 1 each time, causing a constant overlap with the prior windows.



Example of 'High' data sequence from A.csv file

For the model construction, we have had some time to learn Tensorflow and Keras. So we started with using a simple network structure where we initially have one input layer (consisting a sequence of size 50 data

values) which feeds into an LSTM layer with 50 neurons, that in turn feeds into another LSTM layer with 100 neurons which then feed into a fully connected layer with a linear activation function which will be used to give the prediction of the next time step.



Epochs = 20, window size = 50

After running our simple model, the last summary we achieved is as follows:

Simple model result summary

<i>Compilation time</i>	0.02s
<i>Training MSE score</i>	~0.02
<i>Testing MSE score</i>	~0.08

REFERENCES

- [1] Monik V.2017. RNN-LSTM in stock prediction. (2017).
<http://monik.in/a-noobs-guide-to-implementing-rnn-lstm-using-tensorflow/>
- [2] Ruchi Desai. 2015. Stock Market Prediction Using Data Mining. (2015).
<https://www.ijedr.org/papers/IJEDR1402237.pdf>
- [3] Mark Dunne. Stock Market Prediction,
- [4] Sarvajanik Sudeep. 2017. A simple deep learning model for stock price prediction using tensorflow. Medium(2017).