

# Trading Stock with Deep Reinforcement Learning

Yixiao Song, Zhongjie Li, Stephen Guan Feng Yang

*Stanford, California, United States*

---

---

## 1. Task definition

The stock market has largely been a mystery to many investors. Few have been able to consistently and reliably predict its movements. We hope to solve this challenge by building an AI agent that can accurately predict stock trends. We hope to apply several variants of the technique of Reinforcement Learning, including deep Q learning, and double Q learning, to build a stock trading agent. We will compare the total rewards obtained by picking actions according to each of the generated policies.

## 2. Literature Review

Stock market prediction have always been one of the most interesting topic for artificial intelligence. As it is impacted by many different factors, it's very hard to predict, no matter if it's "fundamental analysis" or "technical analysis".

Malkiel [1] concluded that no investment system can consistently yield average returns exceeding the average returns of a market as a whole after testing these systems on 40 years' data. Chen et al built a LSTM model to predict stock price in Chinese stock market and it effectively improved the return. [2] As social media become more and more popular, the twitter data can be a good source of information for stock market prediction.

In 2011, Bollen et al implemented a Self-Organizing Fuzzy Neural Network based model to predict stock market trend based on twitter data.[3] Neuneier proposed to use Q learning to help allocate assets for stock trading. [4] Lee et al proposed a daily stock trading policy based on Q learning, which showed promising results in Korean stock market. [5]

### 3. Infrastructure

#### 3.1. Model

We will be running deep Q-learning with additional feature sets, which is a model-free approach. We'll define the following to model the market:

State: (last ten days' stock price, last ten days' index price, cash amount, stock amount)

Action: Buy or Sell

$Q_{opt} : \max Q(s, a; w)$

$\gamma$  : interest rate

Reward: increase of asset(based on current cash and stock holding )

Cost: Stock transaction fee, time delay (not all transaction can be fulfilled immediately)

Output: final reward value (Cash + # Stock\*price per stock)

Feature set:

LSTM time series forecasting

Sentiment analysis from social media data

Index (SP 500, Dow, Nasdaq etc)

#### 3.2. Dataset

The dataset we used for our results is the daily(Open and Close) stock price history of SAP from 1/1/2016 to 10/1/2017. We are training our data from 1/1/2016 to 1/1/2017, and then run prediction from 1/1/2017 to 10/1/2017.

We also have about 13 million tweets from twitter tracked by the company's name that we used to do sentiment analysis, which we hope to incorporate into features of some of the Q learning variants that we do.

For our current model, we introduced a \$ 2 cost with each action taken. We defined the total reward to be the total cash at the end of the test period, starting initially with \$1000.

Stock price and indexes data from Panda data reader:

<https://pandas-datareader.readthedocs.io/en/latest/>

Twitter data from live tweets adaptor, tracking single company's name over 2 years period around 12 M+ data:

Data structure: id,created(time stamp),text,language,

60 username,replyUser,retweeteduser  
61



Figure 1: Some sample twitters with positive and negative sentiment for SAP's stock

62 3.3. *Evaluation metric*

63 Total asset at the end of test period compared to baseline and oracle.

64 4. **Approach**

65 4.1. *Baseline*

66 For the baseline of our project, we are doing a greedy algorithm as follows:  
67 if we have shares in hand, sell when price rises, and if not, we buy when price  
68 falls. (compare today's price with yesterday's). With an initial investment  
69 of 1000 into the market, between 2017-1-1 and 2017-10-1, our total cash in  
70 the end was 1046.09.

71 4.2. *Oracle*

72 For the Oracle of our project, we will do a linear search over the period  
73 2017-1-1 to current date, which will find the optimal buying and selling date  
74 between the period. We will compare the data output after learning against  
75 this data to compute the loss. 41,698.21 USD is the total reward after running  
76 the linear search.

### 77 4.3. Stock price prediction

78 To predict the future stock price, we implemented two methods, including  
79 linear regression and deep neural network.

#### 80 4.3.1. Problem setup

81 We used SAP stock prices between 2015 and 2017 as training dataset and  
82 the prices beyond 2017 as eval dataset. The input is the last 10 days' stock  
83 prices and the prediction is next 10 days' stock prices. The cost is calculated  
84 as following:

$$85 \text{cost} = tf.reduce\_sum(tf.pow(prediction - true\_prices, 2)) / total\_input\_qty$$

#### 86 4.3.2. Linear regression model

87 The most straightforward method to predict stock prices is linear regres-  
88 sion.

$$89 Prediction = W * X + b$$

90 With 1500 iterations of training, the training loss is decreased from 49766 to  
91 83. And the eval loss is 53. Following figure showed the prediction (in blue)  
92 vs the ground truth (in red).

#### 93 4.3.3. Fully connected neural network model

94 We also constructed a 4 layer neural network to predict future stock  
95 prices. The network is structured with 4 hidden layers with size of 1024, 512,  
96 256, 128, respectively. And the final output size is 10.

97  
98 With 1500 iterations of training, the training loss is decreased from 102  
99 to 46. The eval loss is 38.

#### 101 4.3.4. LSTM neural network model

102 Considering stock prices are sequential data, the LSTM model should  
103 work better for that. The model setup is shown as following.

104  
105 Following figure showed the comparison between prediction (in blue) and  
106 ground truth (in red). We can see that the predicted data overall follows the  
107 trend of the ground truth data.

108

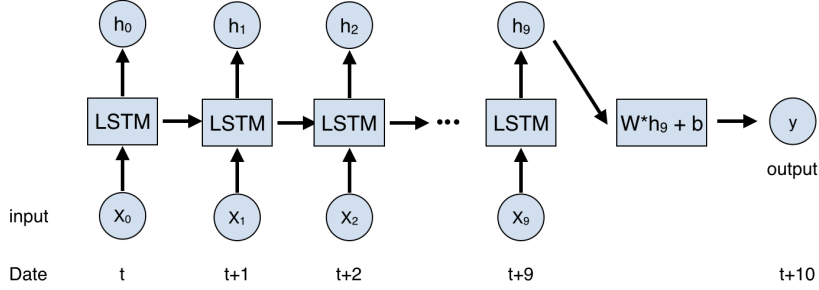


Figure 2: LSTM model structure

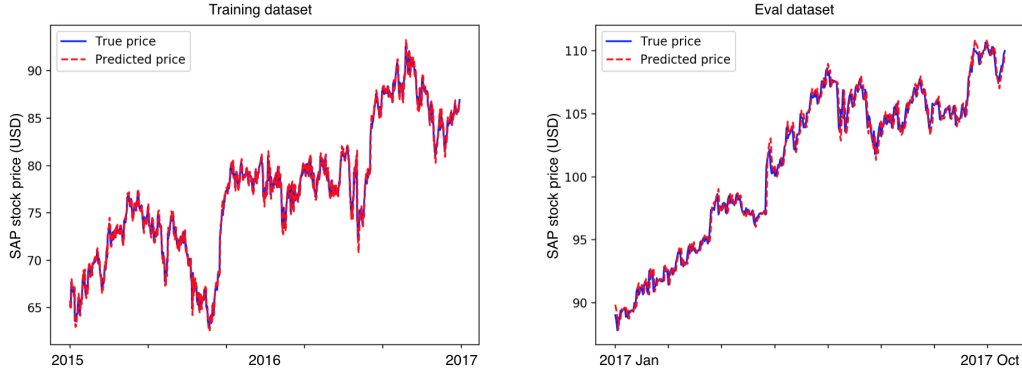


Figure 3: Comparison between ground truth and prediction, left: training data, right: eval data

#### 4.4. Sentiment classification

Social media data has produced a crescent interest in the task of sentiment analysis. In this task, we will need process 13M social media data from Twitter, which tracking by single company name. However, there are few challenges we found out during the implementation.

1. From the dataset, there are over 60 different languages: en: 4718581, tr: 4034451, und: 1053715, pt: 536049 etc. In order to process all the languages, we will need pre-trained word vectors or labeled source training set.

2. Tweet from twitter is limited to 140 characters. The limitations bring a challenge to the regular bag-of-words model by providing very few contextual data. (n-grams only have 60% accuracy )

Based on the scope of the project and timeline, leverage NLP APIs to achieve the task is more advisable.

We are using google NLP API(<https://cloud.google.com/natural-language/>), Sentiment Analysis to process tweet of supported languages (English, Spanish, Japanese, Chinese (Simplified and Traditional), French, German, Italian, Korean and Portuguese.)

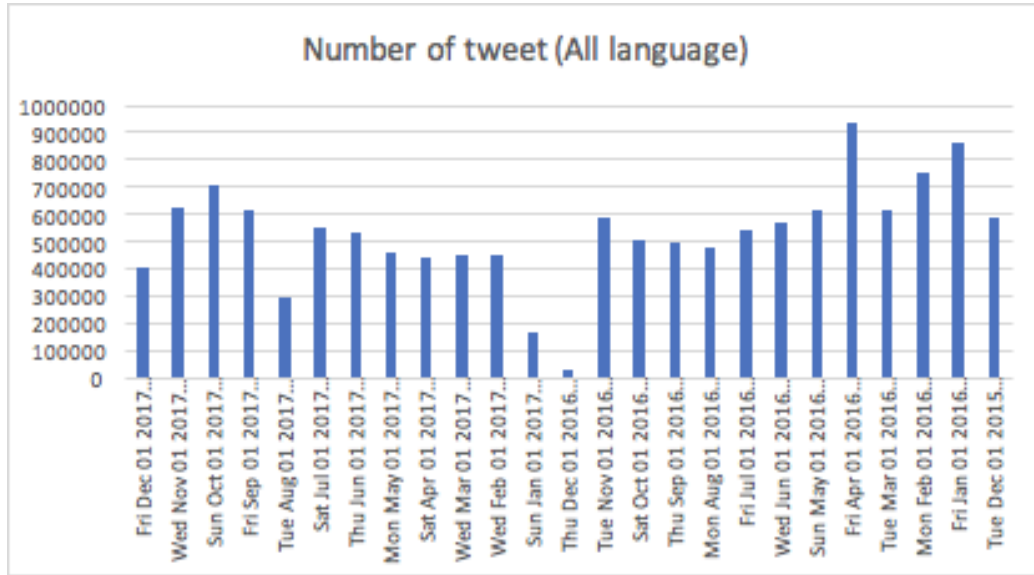


Figure 4: Number of tweets on 'SAP' by date

For this task, we build a scheduler agent to load data from social media table and process through sentiment API and store the result in the sentiment column.

Sentiment score for each tweet range from -1.0(Strong Negative) to 1.0(Strong Positive) and we are using sum of the tweets sentiment as well as average of the sentiment for our feature:

Result chart



Figure 5: The average sentiment of 'SAP'

#### 4.5. *Q-learning*

We've defined the model we are going to use for the agent above. We are interested in finding the maximum reward for a stock, so for each state we just want to find the maximum expected utility,  $V_\pi(s)$ , and it's corresponding action. This would in turn give us a policy  $\pi$  to predict future actions.

Due to the huge state space of the stock market(infinite), we need to use function approximation to estimate  $Q_{opt}$  for unseen (state, action) pair. Instead of having a simple linear  $\phi(x) \cdot w$ , we decided that our  $Q_{opt}$  will be approximated using a neural network that is 3 layers (hence deep Q learning), with the lower most layer outputting actions { "buy", "sell", "hold" }. The equation is as follows:

$$Q(s, a) \leftarrow r + \gamma Q(s', \operatorname{argmax}_a(Q(s', a)))$$

The problem with the above approach is that Q learning by itself is noisy because it picks the highest estimated  $Q(s', a')$  to update the current estimated  $Q(s, a)$ , which when next updated will further amplify the affect of the estimated  $Q(s', a')$ , whether or not it is actually the  $Q_{opt}$ . Therefore, we introduce a second Q as described in Hado van Hasselt's Double Q-Learning[6], as follows:

$$Q_1(s, a) \leftarrow r + \gamma Q_2(s', \operatorname{argmax}_a(Q_1(s', a)))$$

161 and  
162  $Q_2(s, a) \leftarrow r + \gamma Q_1(s', \operatorname{argmax}_a(Q_2(s', a)))$   
163 which van Hasselt proved to solve the problem described.

164  
165 To address the exploration vs the exploitation problem, we don't fully  
166 rely on the estimated Q values, and use a  $\epsilon$ -greedy approach to make sure  
167 we see as many states-action pairs as possible. (Which we label in our graph  
168 RAM(randomized Agent memory). We also added in one of our models opti-  
169 mizations PER(Prioritized Action Replay) to our Q learning Algorithm. The  
170 main idea behind Prioritized Action Replay is that we would pick cached re-  
171 plays that are worse fits of our model than better fits, so our model will learn  
172 the most from those replays.

173 Result we have comparing all the implementations:  
174

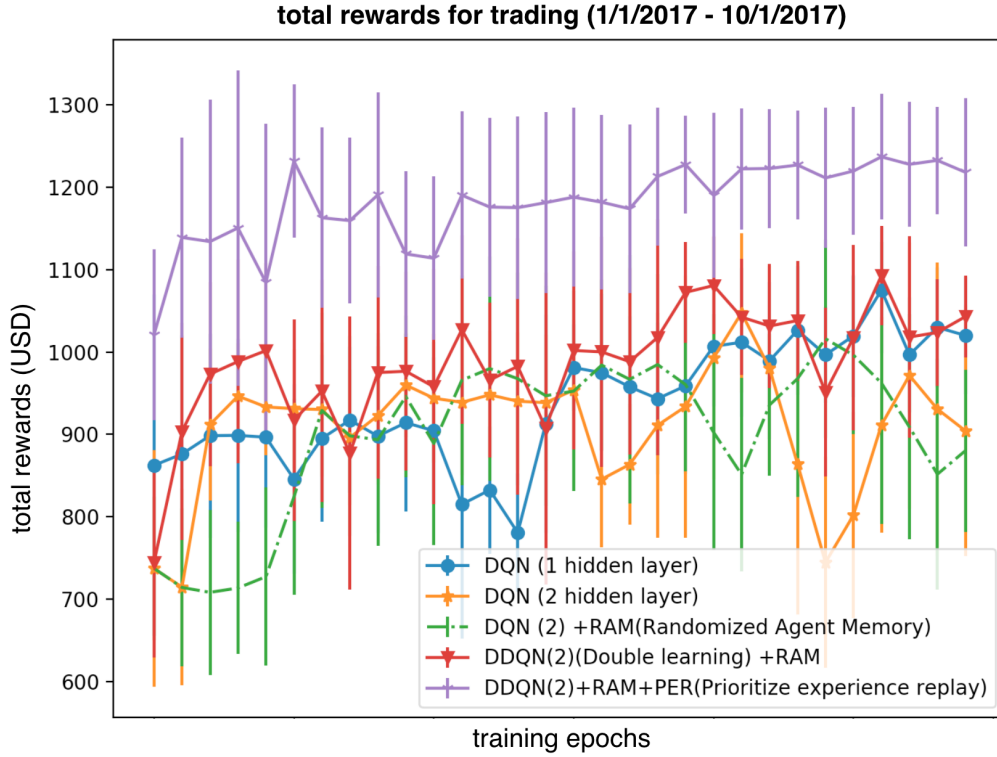


Figure 6: The comparison of several different Q learning algorithms



## 175 5. Error Analysis

176 1. For all the Q algorithms we implement during learning phase  
 177

| DQN(1) | DQN(2) | DQN(2)_RAM | DDQN(2)_RAM | DDQN(2)_ADV2 |
|--------|--------|------------|-------------|--------------|
| 100%   | 97%    | 96%        | 105%        | 126%         |

Figure 7: The average profit for different Q learning

178 we found the DQN with single layer have more robust performance in av-  
 179 erage reward, and DQN with two layer become less stable. This could because  
 180 of the unappropriated network size we had for DQN with two layer(64+128).  
 181

182 2. Mean square error loss function could lead to network weights change  
 183 substantially when the predict value have large bias from estimated value.  
 184 We found Huber loss is more robust and less sensitive to outliers. Which will  
 185 help the performance become more stable.

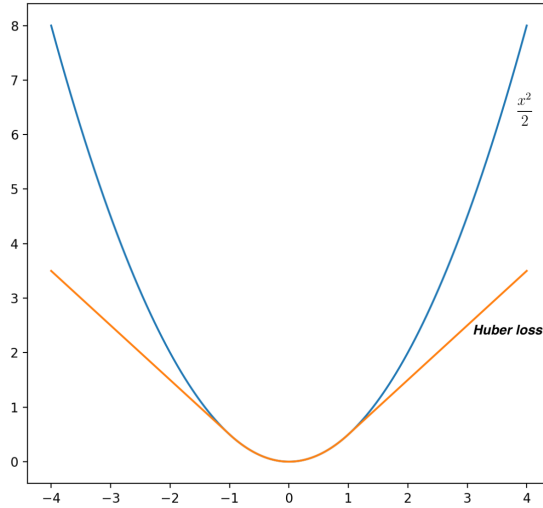


Figure 8: Huber Loss vs MSE

186 3. We tried to use longer time to learn our Q policy (from 2000-2017),  
 187 but the performance is very unstable event in the learning phase (range from  
 188 10,000+ to less than 500 base on 1000 initial investment and 2\$ transaction  
 189 cost). The possible reason could be the memory limitation (We can not store  
 190 all the possible states for very large data set).

191

192 4. Social media data keyword tracking to the unrelated tweet. Because  
 193 we are using single key word (SAP) to pull data from twitter. And this key-  
 194 word could be use differently in other language or irrelevant words. This  
 195 could have bias on the sentiment result.

196

197

## 198 6. Reference

- 199 [1] B. G. Malkiel, A random walk down Wall Street: including a life-cycle  
 200 guide to personal investing, WW Norton & Company, 1999.
- 201 [2] K. Chen, Y. Zhou, F. Dai, A lstm-based method for stock returns pre-  
 202 diction: A case study of china stock market, in: Big Data (Big Data),  
 203 2015 IEEE International Conference on, IEEE, pp. 2823–2824.
- 204 [3] J. Bollen, H. Mao, X. Zeng, Twitter mood predicts the stock market,  
 205 Journal of computational science 2 (2011) 1–8.
- 206 [4] R. Neuneier, Enhancing q-learning for optimal asset allocation, in: Ad-  
 207 vances in neural information processing systems, pp. 936–942.
- 208 [5] J. W. Lee, J. Park, O. Jangmin, J. Lee, E. Hong, A multiagent approach  
 209 to  $q$ -learning for daily stock trading, IEEE Transactions on Systems,  
 210 Man, and Cybernetics-Part A: Systems and Humans 37 (2007) 864–877.
- 211 [6] H. V. Hasselt, Double q-learning, in: J. D. Lafferty, C. K. I. Williams,  
 212 J. Shawe-Taylor, R. S. Zemel, A. Culotta (Eds.), Advances in Neural  
 213 Information Processing Systems 23, Curran Associates, Inc., 2010, pp.  
 214 2613–2621.