
Hybrid Attention Networks for predicting Chinese stock trend by news data

Anonymous Author(s)

Affiliation

Address

email

1 (Implementation Track, group of 4)

Abstract

2 Making a prediction of the trend of the stock market is a very important part for
3 investors to maximize profit. However, predicting the trend of the stock market
4 is a very difficult due to the sensitivity of the stock market. The stock market
5 is very sensitive to the real-world, there are many factors that can influence the
6 stock market, for example, the sudden break of the pandemic, the bankruptcy
7 of a relative company, etc. Internet, as a huge source of information, can be an
8 exhaustive description of what is happening in the world. However, it is important
9 to distinguish the high quality and real knowledge in the lake of information. In
10 order to address these challenges, we implemented and improved methods from
11 Peking University and Microsoft research team which used news data to predict the
12 stock trends and improve the annualized return. They purposed a framework with a
13 Hybrid Attention Networks (HAN) and a self-paced learning mechanism to address
14 three main problems which are *sequential content dependency*, *diverse influence*,
15 *and effective and efficient learning* when dealing with online stock news. In this
16 paper, we are going to replace the bi-directional GRU-based RNN method with
17 BERT and train with both stock news data and previous stock price data, which
18 can ensure a better prediction performance.

19 1 Introduction or Problem statement

20 Predicting the future trend of the stock price is always a goal that all investors doing in order to make
21 the profit maximized. However, this is a very challenging task due to the high sensitivity of the market
22 which made it very fluctuated. There are conventional ways that trying to use information from price,
23 volume, and other fields to predict the future price, but these methods cannot reveal the ground truth
24 beneath the fluctuation of the stock price, which made them limited. So we come up with an idea
25 that, if we can gather all information about certain companies, we can make a strong prediction about
26 their stock prices. The Internet, as a huge pool of information, can be considered to hold all recent
27 information in the world. In other words, if we want some public accessible information, it must be
28 able to be found somewhere on the internet. This job requires people to search and extract useful and
29 reliable information from the internet. Fortunately, as the Natural Language Processing techniques
30 and Web crawling techniques improved, which made our goal become realizable. The research team
31 stated that there are several ML tools based on analyzing the Wall Street journal. So they decided to
32 develop a tool called Hybrid Attention Networks (HAN) to predict stock price by relevant news.[1]
33 Our project is focusing on implement HAN in the Chinese stock market and evaluates its performance.
34 Since BERT is one of the state-of-art nlp techniques being widely used today, we also want to replace
35 the bi-directional GRU structure in the HAN framework with BERT to enhance the model prediction
36 performance.

37 2 Method

38 The core idea is to imitate the learning process of human when facing chaotic online news which can
39 be concluded into three characteristics:

- 40 • Sequential Context Dependency: It refers to the fact that a single news is more informative
41 within a broader context than isolated. Therefore, when we take a sequence of related news
42 as a unified content, we can make better prediction.
- 43 • Diverse Influence: Different online news have different influence. It is common that one
44 critical news can affect the stock price for weeks, whereas a trivial one may have zero effect.
45 So when we predict the subsequent stock trend, we need to comprehensively identify the
46 estimated impact of each online news.
- 47 • Effective and efficient Learning: It is not surprising that the content of news contains
48 vague information on stock trend which leads to the fact that we can not make any reliable
49 new-oriented prediction. In order to address this problem, people usually tend to first gain
50 knowledge from the more common situations before turning to more complicated situations.
51 And if there is limited number of stock news in a period, we may discard those data due to
52 the efficiency requirement.

53 Ziniu Hu et al. proposed a Hybrid Attention Networks (HAN) [1] to capture first two characteristics,
54 Sequential Context Dependency and Diverse Influence. The Hybrid Attention Networks (HAN)
55 includes temporal-level attention-based recurrent neural networks (RNN) to identify more influential
56 time periods of the sequence and news-level attention-based recurrent neural networks (RNN) to iden-
57 tify more influential news at the same time point. With structure of recurrent neural networks (RNN),
58 we are able to consider a sequence of related recent news as a unified context and the attention mech-
59 anism enables the process of identifying diverse influence of different news and different time periods.
60

61 To perform a more effective and efficient learning, they implemented a Self-Paced Learning (SPL)
62 algorithm. Kuma et al. designed Self-Paced Learning (SPL) which is based on Curriculum learning
63 but it can optimize the original objective and curriculum design at the same time. It skips certain
64 data points that are considered to be yet too hard and gradually increase the complexity of training
65 samples thus the model performs better predictions.

66 Specifically, the framework of Hybrid Attention Networks (HAN) includes:

- 67 1. a pre-trained unsupervised Word2Vec as the word embedding layer;
- 68 2. an attention layer to catch the unevenly distributed importance level of a sequence of stock
69 news in a day;
- 70 3. a bi-directional GRU network to incorporate information of a sequence of stock news to get
71 a latent vector for one day;
- 72 4. Another attention layer to adjust the importance of each day's latent vector;
- 73 5. A standard Multi-layer Perceptron (MLP) to output the prediction results;

74 3 Related work

75 In the traditional finance theory, Brownian motion and Efficient Market Hypothesis (EMH) [2] are
76 the cornerstones of the stock market prediction and strategies. However, many empirical results
77 indicate that EMH is not applicable to the stock market. To capture the characteristics of time-series
78 historic market data, neural networks have been considered to be a very effective learning algorithm
79 for decoding those nonlinear time series data [3]. Among different widely used deep learning neural
80 networks architectures, Recurrent Neural Networks (RNN) are best fitted for time series modeling
81 problems in financial market. In 1996, Roman et al. [4] used back propagation and RNN models to
82 predict stock index among different stock market. Recently, Long-Short Term Memory (LSTM)
83 is considered as one of the most popular approaches for stock market prediction since it has the
84 capability of holding past information [5].
85

86 Ziniu Hu et al. stated that although there has been a recent surge of interest on development of Deep
87 Neural Network (DNN) architectures based on news or public data to predict stock trend movements
88 [6], most of them are devoted to reliability and quality of information sources. In addition, Ronaghi
89 et al. [7] designed a Noisy Deep Stock Movement Prediction Fusion framework (ND-SMPF) for
90 stock price movement prediction to address this problem.

91 **4 Experiments**

92 **4.1 Implementation**

93 To collect new data for our project, we use an amazing stock data crawling tool called Tushare. For
94 model implementation part, in general, we are going to implement the Hybrid Attention Networks
95 (HAN) framework mainly with PyTorch module. And regarding the self-paced learning (SPL)
96 algorithm, we might need to write our own code to implement the parameter optimization steps in
97 Python.

98
99 More specifically, for the deep learning network construction part, we need to define all the
100 layer listed above and try running it on Google Colab which offers us a free GPU. For the self-paced
101 learning part(SPL), we need to implment the iterative algorithm for parameters optimization such as
102 back-propagation on our own. So we have to be careful about deriving the math formula correctly
103 before coding.

104 **4.2 Dataset**

105 The data we are going to use for this project can be mainly divided into two parts: one is the relevant
106 stock market news information, the other is the previous stock market prices. We retrieve our data via
107 the crawling tool called Tushare, basically, it includes:

- 108 1. Historical stock open price, close price, all interesting factors for each stock
- 109 2. Corresponding news for each stock with published date
- 110 3. all the quick news from 2019-01-01 till now, still need to assign to corresponding stocks
111 with names or ids
- 112 4. Overall news about finance or stock market with published date
- 113 5. all the general news from 2019-01-01 till now

114 This process is challenging because first we need to extract the publication timestamp, title, and
115 content for each news. Then we need to correlate each of the collected news to a specific stock if the
116 news mentioned the name of the stock in the title or content. We then filter out the news without any
117 correlation to stocks. After such a process, for each stock, we then aggregate all the news in a certain
118 date to construct the daily news corpus.

119 **4.3 Evaluation**

120 To specify the label of the tri-classification problem (UP, DOWN, STAY), we set up two particular
121 thresholds to bin the rise percents (rise percent is defined as the increasing rate between next day's
122 and today's open price), We define the thresholds so that the three categories are approximately
123 even. We also split the dataset into a training set (66.7%), and a test set (33.3%). Then we further
124 randomly sample a validation set from the training set with 10% size of it, in order to optimize the
125 hyper-parameters and choose the best epoch.

126
127 In the experimental setting of a tri-label classification problem, as the three label is approx-
128 imately evenly split, we choose accuracy, which means the proportion of true results among the total
129 number of testing samples, as the evaluation metric.

5 Plan of Project

5.1 Anticipated division of work over team

For the model construction part: we are going to split our work mainly into 4 parts:

1. data crawling
2. HAN framework implementation(replace bi-directional GRU with BERT)
3. SPL iteration math formula derivation
4. SPL algorithm implementation in Python

For the actual experiment part: we are going to try 4 other models or deep learning frameworks to compare the performance with the improved HAN+SPL model described above. Each member would take care of one model or framework.

5.2 Experience or expertise that we might have to complete the project successfully

Two of our members have taken the NLP course which could help implement the BERT model here. And all of us have prior experience coding with PyTorch module. Besides, one of our members majoring financial engineering and risk management degree also has lots of experience working with stock market data.

5.3 Expected challenges and difficulties

The first challenge is to retrieve the data from the internet by using crawling tool. We have to learn how to use the new tool Tushare in a short time. And then we need to process the data by relevant each stock market news to a certain date, which is also a challenging work. For the model construction part, we need to not only implement the HAN framework but also replace the bi-directional GRU with BERT. Last but not least, we need to write our own code to implement the SPL algorithm.

5.4 milestones and contingency plans to get positive results

We are going to hold a weekly recurrent meeting to mainly focus on each part of work described above.

1. data crawling and pre-processing
2. HAN framework implementation(replace bi-directional GRU with BERT)
3. SPL iteration math formula derivation
4. SPL algorithm implementation in Python
5. experiment with new data and compare to other models' performance

References

- [1] Hu, Ziniu, et al. "Listening to chaotic whispers: A deep learning framework for news-oriented stock trend prediction." Proceedings of the eleventh ACM international conference on web search and data mining. 2018.
- [2] Malkiel, Burton G. "The efficient market hypothesis and its critics." Journal of economic perspectives 17.1 (2003): 59-82.
- [3] Batres-Estrada, Bilberto. "Deep learning for multivariate financial time series." (2015).
- [4] Roman, Jovina, and Akhtar Jameel. "Backpropagation and recurrent neural networks in financial analysis of multiple stock market returns." Proceedings of HICSS-29: 29th Hawaii International Conference on System Sciences. Vol. 2. IEEE, 1996.
- [5] Jia, Hengjian. "Investigation into the effectiveness of long short term memory networks for stock price prediction." arXiv preprint arXiv:1603.07893 (2016).

- 171 [6] Ding, Xiao, et al. "Deep learning for event-driven stock prediction." Twenty-fourth international
172 joint conference on artificial intelligence. 2015.
- 173 [7] Ronaghi, Farnoush, et al. "ND-SMPF: A Noisy Deep Neural Network Fusion Framework for
174 Stock Price Movement Prediction." 2020 IEEE 23rd International Conference on Information Fusion
175 (FUSION). IEEE, 2020.