

Sentiment Analysis And Topic Modeling of Stock Market Data

Akhil Bandhu Anand
Computer Science Department
Southern Utah University
Cedar City, United States
akhilanand@suuemail.net

Abstract—This paper covers an initial increment into stock market sentiment analytics. The purpose of this project is to explore the how good the models can get with predicting sentiment, the project will be an initial platform into modeling sentiment trading. If the models defined in the paper can get good enough in predicting sentiment, technical, fundamental, and algorithmic trading can be modelled with sentiment scores from this paper as a variable. Previous studies have shown that investor sentiment indicators can predict stock market change. This paper will be a mixture of previous work done in this field and then some more.

I. INTRODUCTION

Sentiment, that is one factor that may give financial modeling an edge on the market at least that is what the theory for this project is. To gain an understanding of financial markets we must first look at the Efficient Market Hypothesis (EMH). One fallacy of this theory is that all investors act rationally and as one can gather from intuition of the market, people don't always act rationally especially when finances are involved. Coming back to the theory, there are three factors in the EMH, the first one being weak form hypothesis. If one believes in this form of the hypothesis, the investor can never gain an edge on the market through technical analysis, that is, looking for patterns and trends in historical data to predict future stock market movements. The second leg of the theory is the semi-strong form of the hypothesis, investors cannot get an edge on the market through fundamental analysis, that is, by examining related economic and financial factors such as the balance sheet, strategic initiatives, micro-economic indicators, and consumer behavior, fundamental analysis attempts to identify a security's intrinsic value. And the final leg of the theory, the strong form of the hypothesis. This part of the theory highlights that investors cannot gain an edge on the market through insider trading, gaining information on a security before said information is made public. All this is possible because the theory states that a security's price reflects all the above highlighted information as soon as it is available to anyone in the market. Basically there is no way of gaining an edge on the market!

This paper will explore the analysis of sentiment, it won't reach the level of analysis needed for the bigger project in hand. This analysis is going to be a stepping stone for a polarity score to be used in the making of a financial trading

Efficient Market Hypothesis



model. The goal of this analysis is to be able to at least get a good hand at predicting the sentiment of the market using supervised and unsupervised learning methods.

The use of sentiment analysis (SA) is increasingly used to determine the feelings of social media users towards a subject. The most popular way of performing sentiment analysis is data mining. Our aim is to use Deep Learning and Supervised Learning technologies to assess investors' emotions. Our central idea is to adopt Deep Learning to determine investors' expectations about the price of stocks and the overall market based on their messages. The reason why we select Deep Learning methodology rather than data mining is that in data mining, identifying features and selecting the best of those features is the most challenging task to undertake especially in a Big Data.

In contrast to data mining, a Deep Learning model, learns features during the process of learning. Deep Learning algorithms lead to abstract representation, as a result, they can be invariant to the local change in the input data. In addition, Big Data problems including semantic indexing, data tagging, and fast information retrieval can be addressed better with the aid of Deep Learning. Deep Learning provides the opportunity to use a simpler model to accomplish complicated Artificial Intelligence tasks.

The remainder of this paper is organized as follows: "Related Work" section we look at previous work done in the field of sentiment analysis and the methods employed therein; "Methodology", in this section of the paper, we will look at the methodologies used in this paper. This paper has an emphasizes on the bag of words method and sentiment predictions but there is a section of topic modeling, especially

Latent Dirichlet Allocation, Word Clustering, Text mining, and Polar Subjectivity. In this section, we will also explore the data set and how it is prepared for the rest of the analysis. And the final part of the paper will discuss the results and conclusions from this project.

II. RELATED WORK

Following is some of the work that has been done in the area.

Pandey et al. (2017) suggested a new "metaheuristic scheme" based on the K-mean algorithm and the Cuckoo search method. The emotive information from the "Twitter dataset" was used to find the appropriate CH (Cluster heads). The suggested model's accuracy was tested on a variety of "Twitter datasets" and compared to other optimization approaches such as PSO ("Particle Swarm Optimization"), DE ("differential evolution"), CS ("cuckoo search"), and two n-gram schemes [1].

Wehrmann et al. (2017) presented a "language-agnostic translation-free" approach for "Twitter Sentiment analysis" using CNN. The polarity of tweets have been examined which has been written in various languages. It has been analyzed that the structure of deep neural network required less learner constraints. The experiment has been carried out on four different languages such as English, german, Portuguese and Spanish [2].

Sosa and P. M. (2017) proposed an integration of two NN (Neural networks) such as CNN (Convolution "Neural Network") -LSTM ('Long short term memory') and LSTM-CNN to perform sentiment analysis on Twitter data. For training 10,000 tweets have been tested that consists of positive and negative tweets equally. The average accuracy obtained by using LSTM-CNN is high about 75.2% [3].

To discover the feeling level of the News data, Shahare et al. (2017) offered techniques such as 'Naive Bayes' and 'Levenshtein algorithm.' The author uses different levels of processing to determine the emotion words. The author expertly extracts emotion from either a tiny or a large text. The author can also tell the difference between news events and text data [4].

This paper's experiments are focused on market sentiment. Market mood, according to [5], is the overall prevalent sentiment of investors in terms of anticipating price development in a market. This attitude is influenced by a variety of elements, including global events, history, economic news, seasonal conditions, and so on. Sentiment analysis, often known as opinion mining, is the practice of extracting a writer's attitude from source materials using natural language processing methods.

Multiple research initiatives have used supervised classification methods such as Support Vector Machines [7], Nave Bayes [6], or ensembles [7, 8] to perform sentiment analysis. The bag-of-words [52] paradigm is commonly used in machine learning approaches. A text is represented as a collection of its words in the bag-of-words model, irrespective of the sequence of those words in their sentences.

Loughran and McDonald's [10] work is one of the most well-known in this subject. They created a financial lexicon and manually created six-word lists using the US Securities and Exchange Commission portal from 1994 to 2008, including positive, negative, litigious, uncertainty, model strong and model weak.

III. METHODOLOGY

A. Dataset

For this paper, The first resource that would have provided a comprehensive dataset would have been StockTwits Inc. StockTwits is a financial social media platform that was founded in 2009. StockTwits users can get information on the stock market, such as current stock prices, price movement, stock exchange history, buying or selling recommendations, and so on. Furthermore, as a social network, it allows stock market traders to share their experience. Unfortunately, I was not able to gain access to their platform to gather a dataset. I was lucky enough to find a few datasets on Kaggle. Kaggle is a Google subsidiary, an online community of data scientist and machine learning practitioners, this community shares code and datasets about various different problems around the world.

In this experiment, I was able to find two different data sets. This first called "stock data.csv", this data set had two columns and 5,791 rows. The first column was a Sentiment score, "positive" and "negative", the other column was the actual text, this dataset was gathered from multiple twitter handles talking about market sentiment and economic news (web scraping). The next data set called "all data.csv", also had two columns one of Sentiment and one of the actual text. This data set had three categories "positive", "negative", and "neutral". Both these data set were combined to get a final dataset over 10,000 rows with three different categories of sentiment.

B. Text Mining

The first part of the analysis of this dataset is more of an exploratory data analysis, that is cleaning up the data and making some graphs to see how the data is structured. This next graph and word cloud is what the dataset looks like without any cleaning, without removing any kind of stop words or words like "eur" which will not be helpful in an american market.



Sentiment Word Counts

The figure consists of two horizontal bar charts side-by-side, both sharing a common y-axis of words. The left chart is titled 'negative' and shows counts for negative sentiment words. The right chart is titled 'positive' and shows counts for positive sentiment words. A legend on the right indicates that red bars represent negative sentiment and teal bars represent positive sentiment.

Word	Negative Count	Positive Count
loss	150	
share		250
fell	60	
production		130
lower	40	
corporation		130
demand	30	
technology		120
margin	30	
agreement		110
income	30	
increase		100
excluding	20	
building		90
government	20	
management		90
versus	20	
deal		80
vice	20	
real		70

cover is called Latent Dirichlet Allocation (LDA).

Latent Dirichlet Allocation (LDA) is a method that uses a corpus of text to represent a document as a random mix of hidden themes [11]. In natural language processing, the latent Dirichlet allocation (LDA) is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar. For example, if observations are words collected into documents, it posits that each document is a mixture of a small number of topics and that each word's presence is attributable to one of the document's topics. This next plot is from an LDA model that outputs and creates general topics through word probabilities.

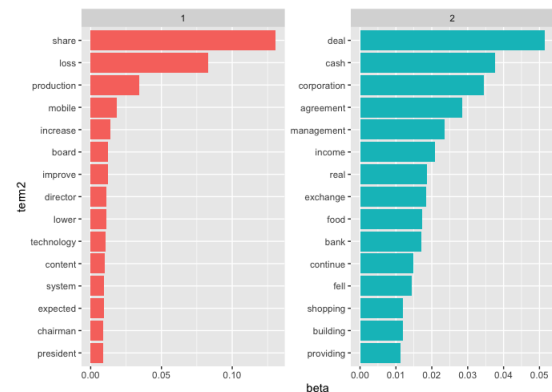


Figure 1 displays four horizontal bar charts, labeled 1, 2, 3, and 4, showing the distribution of beta values for various terms. The x-axis for all charts is 'beta'.

Cluster 1 (Red bars):

- deal
- cash
- production
- real
- food
- bank
- expected
- pay
- agreed
- mobile
- wireless
- cooperation
- demand
- organization
- system

Cluster 2 (Green bars):

- income
- technology
- continue
- fail
- president
- improve
- providing
- lower
- cutting
- content
- growth
- time
- delivery
- margin
- supply

Cluster 3 (Teal bars):

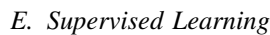
- loss
- corporation
- management
- exchange
- shopping
- building
- center
- provide
- resources
- credit
- chairman
- board
- network
- excluding
- public

Cluster 4 (Purple bars):

- share
- agreement
- increase
- result
- director
- planning
- mobile
- money
- level
- store
- vote
- music
- manufacturer
- completion
- versus

Word clustering is a technique for dividing a set of words into subsets of semantically comparable terms. It is rapidly being utilized in a variety of NLP applications, from word sense and structural disambiguation to information retrieval and filtering.

Hierarchical cluster analysis (also known as hierarchical clustering) is a type of cluster analysis in which the goal is to group objects or records that are "near" to one another. This next plot is a hierarchical cluster plot for the dataset, this shows us the word association for all the words.



1) *Data Examples:* Here are some examples of the text data, this will explain why further cleaning needs to be performed to get the data ready for some predictive modeling.

- 2) *Pre-Processing*: Some of the messages had a very weird symbols in them refer to the examples given above. Some preprocessing steps are performed to clean the text in the dataset, this was done after the above mentioned processing to get the data ready for modeling:

- All URLs were removed. Most URLs are short URLs and located at the end of the text.
- All punctuation was removed.
- All whitespace was removed from the text.
- The abbreviations were made into two separate words, like "hadn't" to "had not"
- Symbols were replaced from the text.
- All numbers were replaced with actual text. Example: "1" to "one", "5000" to "five thousand".

A word cloud of terms associated with the COVID-19 pandemic. The words are arranged in a circular pattern, with some words appearing larger than others. The words include: company, watch, helsinki, business, points, break, per, also, highest, back, services, global, finland, euromove, said, position, compared, volume, good, nifty, market, target, quarter, higher, goog, triangle, one, stocks, like, price, buy, total, value, ong, period, since, close, long, support, next, day, share, loss, bank, can, big, still, senex, line, weekly, cash, looking, two, percent, min, mobile, short, user, breakout, week, products, trade, lower, according, stop, operations, sales, earnings, group, well, net, first, see, new, coronavirus, now.

4) **ANOVA:** Analysis of variance (ANOVA) is a collection of statistical models and associated estimate processes (such as "variation" among and between groups) that are used to examine variations in means. Ronald Fisher, a statistician, invented ANOVA. ANOVA is based on the law of total variance, which divides observed variance in a variable into components attributed to various causes of variation.

5) *CART Modeling*: A Classification And Regression Tree (CART) is a predictive model that shows how the values of an outcome variable can be predicted based on the values of other variables. A CART output is a decision tree in which each fork represents a split in a predictor variable and each end node represents an outcome variable prediction.

6) *Random Forest*: Random forests, also known as random choice forests, are an ensemble learning method for classification, regression, and other tasks that works by building a large number of decision trees during training. For classification tasks, the random forest's output is the class chosen by the majority of trees. The mean or average prediction of the individual trees is returned for regression tasks.

7) *Naive Bayes*: In statistics, Naive Bayes classifiers are a group of straightforward "probabilistic classifiers" based on Bayes' theorem and strong (naive) independence assumptions between features (see Bayes classifier). They are one of the most basic Bayesian network models, but when combined with kernel density estimation, they can achieve great levels of accuracy.

8) *Support Vector Machines*: Support-vector machines (SVMs, also known as support-vector networks) are supervised learning models that examine data for classification and regression analysis in machine learning.

IV. RESULTS AND CONCLUSION

REFERENCES

- [1] A. Chandra Pandey, D. Singh Rajpoot and M. Saraswat, "Twitter Sentiment Analysis using Hybrid Cuckoo Search Method", *Information Process Management*, vol. 53, no. 4, pp. 764-779, 2017.
- [2] J. Wehrmann, W. Becker, H. E. L. Cagnini and R. C. Barros, "A Character-Based Convolutional Neural Network for Language-Agnostic Twitter Sentiment Analysis", *proceedings International Conference on Neural Networks*, pp. 2384-2391, 2017.
- [3] P. M. Sosa, "Twitter Sentiment Analysis using Combined LSTM-CNN Models", pp. 1-9, 2017, [online]
- [4] F. F. Shahare, "Sentiment Analysis for the News Data Based on the Social Media", *Proceedings 2017 IEEE International Conference on Intelligent Computer Control Systems*, pp. 1365-1370, 2017
- [5] Market Sentiment. <http://www.investopedia.com/>.
- [6] Saif H, He Y, Alani H. Semantic sentiment analysis of Twitter. *The semantic Web-ISWC 2012*. Berlin: Springer; 2012. p. 508–24.
- [7] Steinwart I, Christmann A. *Support vector machine*. Berlin: Springer; 2008.
- [8] Silva N, Hruschka E, Hruschka E. Tweet sentiment analysis with classifier ensembles. *Decis Support Syst*. 2014;66:170–9.
- [9] Fersini E, Messina E, Pozzi FA. Automatic construction of financial semantic orientation lexicon from large scale Chinese news corpus. *Decis Support Syst*. 2014;68:26–38.
- [10] Loughran T, McDonald B. When is a liability not a liability? Textual analysis, dictionaries. *J Finance*. 2011;66:35–65.
- [11] S. J. Putra, M. A. Aziz, and M. N. Gunawan, "Topic analysis of Indonesian comment text using the latent Dirichlet allocation," 2021 9th International Conference on Cyber and IT Service Management (CITSM), 2021.