

Sentiment Analysis: Forecasting Bitcoin Price Changes with Twitter Data

Masterarbeit

zur Erlangung des akademischen Grades „Master of Science (M. Sc.)“ im Studiengang Wirtschaftswissenschaft der Wirtschaftswissenschaftlichen Fakultät der Leibniz Universität Hannover

vorgelegt von

Manuel Hilmer



Prüfer: Prof. Dr. Michael H. Breitner

Hannover, den 26. September 2018

Contents

List of Figures	V
List of Tables	VI
List of Abbreviations	VII
1. Introduction & Motivation	1
2. Background Theory	4
2.1. Bitcoin	4
2.2. Twitter Sentiment Analysis	5
2.3. Machine Learning	8
2.4. Feature Extraction	17
2.5. Classification Evaluation Metrics	21
2.6. Related Work	22
3. Experimental Methodology	24
3.1. General Idea	24
3.2. Experimental Tool and Platform	25
3.3. Train-/Test-Dataset	26
3.4. Preprocessing Methods	27
3.5. Classifiers	29
4. Econometric Methodology	36
4.1. General Idea	36
4.2. Description of the Datasets	37
4.3. Twitter-based Prediction of Financial Bitcoin Indicators	39
5. Experimental Results and Discussion	42
5.1. Classifier Comparison	42
5.2. Sentiment of Tweets as a Predictor	51
6. Conclusions, Limitations and Further Research	55
References	57
A. Figures	63
B. Tables	65

Chapter 1.

Introduction & Motivation

Considering virtual currencies, the most prominent "cryptocurrency" Bitcoin comes in mind. A virtual currency is an alternative currency which is electronic and thus has no physical form. The economic and social impact of virtual currencies in the recent past continued to grow very quickly. Bitcoin achieves great success and is the leading cryptocurrency in the world. It graduated from unfamiliarity to mainstream mostly because of the insane increase in value it saw in recent years. Big companies and brands started accepting Bitcoin, banks integrated Bitcoin accounts and exchanges gained public interest. From then on cryptocurrencies are used as investment, trading object and real-world payment method. The Bitcoin ecosystem is peer-to-peer and trustless. Bitcoins are not issued by any bank or other financial institution. Cryptography is used and the system depends on software algorithms. Transactions are verified decentralized. This structure leads to the independence of other financial markets, gold or fiat currency. "The value of Bitcoin is derived from the value that people assign to it" (Brito and Castillo, 2013). Based upon this the bitcoin price should have a connection with news articles and social media.

Just like the impact of the Bitcoin, grows the frequency of related news articles and posts on microblogging websites. Publishing good or bad information regarding a particular company or product can be done by using a microblog such as Twitter. Users create messages and thus interact with others. These messages are known as "tweets" and are the chosen tool to express thoughts or feelings about several different topics. Opinion mining or its alternate term "Sentiment Analysis" finds critical use in areas where organizations or individuals gain benefits from knowing the general sentiment related to a particular object, i.e. a product, person or a movie. Sentiment analysis has many areas of application and is often used to study the effect of sentiment. For stock market returns there exists evidence that the sentiment of people is capable of predicting returns (Da et al., 2014). Bollen et al. (2011a) uses twitter sentiment analysis to identify quantifiable relationships between overall public sentiment and social, economic or other events in the media. Consumers and merchants conduct sentiment analysis on tweets to gather knowledge of products or to perform market analysis. Advanced machine learning algorithms and higher computing

power allows sentiment analysis on large datasets provided by social media websites. These offer a relatively new source of information to gather sentiment in real time and on a large scale. The amount of data being published on twitter or other social media sites far exceeds the manual classifiable amount. Thus the grown impact of social media leads to a growth in research using sentiment analysis (Liu, 2012).

While there are still many other causes of Bitcoin price changes than sentiment, it is reasonable to explore whether social media data can inform predictions on the price changes. For example, Kristoufek (2013) analyzes the dynamic connection between the Bitcoin price and the overall attention of the Bitcoin measured by search queries on Google Trends and on Wikipedia. He finds a strong causal relationship between the prices and searched terms. Apart from research of correlations between sentiment and other indicators, it is important to get accurate classification models. This is realized with various complex machine learning algorithms and allows to predict the sentiment of a document or a sentence with high accuracy.

This thesis aims to implement sentiment analysis on tweets using various state-of-the-art machine learning algorithms. Because accuracy of diverse classification methods differs greatly, current methods with respect to their quality need to be compared. This grants a deeper insight into strength and weaknesses of investigated models. The attempt is to classify the polarity of the tweet in either positive or negative. I conduct a classifier comparison using 1.6 million pre-labeled tweets for training and testing. The data comes with usernames, links, hashtags and other twitter specific formatting which are in need of processing and conversion into a standard form. Useful features; which represent the tweet in a certain way have to be extracted from the text. The best classification model is selected to conduct sentiment analysis on tweets regarding "bitcoin". In addition to that I investigate the dynamics affecting different market indicators of Bitcoin by focusing on Twitter sentiment as an explanatory variable. I explore these quantifiable relationships by using vector autoregressive (VAR) models.

Recapitulating, this thesis is concerned with the following research questions:

RQ 1. *Which machine learning models are best suited for Twitter Sentiment analysis?*

RQ 2. *Is there a correlation between Twitter sentiment and Bitcoin market indicators which can be used to forecast Bitcoin price changes?*

In order to answer the research questions, the thesis is organized as follows: After this introduction and motivation follows a theoretical chapter which illustrates the Bitcoin system and gives insight in theoretical foundations of sentiment analysis. In this context

machine learning algorithms and evaluation metrics used in this work are explained in detail. To conclude this chapter, I present related research in the field of twitter sentiment analysis and Bitcoin price correlation. Chapter 3 describes the approach of conducted classifier comparison in order to get sentiment data for "Bitcoin"-tweets. Subsequently, Chapter 4 specifies the procedure of correlation analysis between Bitcoin market indicators and sentiment indicators. Section 5 provides results and discussion of the previous conducted experiments. Section 6 concludes. The appendix provides additional figures, tables and codings which have a reference in the text.

Chapter 6.

Conclusions, Limitations and Further Research

In this thesis I investigated the relationship between twitter sentiment and Bitcoin price. On that account an empirical analysis was conducted assimilating different economic Bitcoin and Twitter-related variables. At first I examined machine learning algorithms for twitter sentiment analysis. I conducted a classifier comparison to study the effect of different classifiers on the accuracy of sentiment analysis models. Systematically, six machine learning algorithms (Logistic Regression, Multinomial Naive Bayes, Support Vector Machine, Feed-forward neural networks, Recurrent neural networks, Convolutional neural networks) were used to train different models for text classification and thus to measure the sentiment of twitter posts regarding Bitcoin. The models were trained and evaluated on a dataset containing 1.6 million pre-labeled tweets. I tested different feature extraction methods with the algorithms to make the comparison even more manifold. The comparison reveals, that no machine learning algorithm performs badly. The accuracy ranges from 80.04 % to 83.13 % validation accuracy. For algorithms using sparse feature vectors, tf-idf with combined unigram, bigram and trigram features result in the highest accuracy. For the classifier using dense feature vectors a pre-trained GloVe word vector with updated weights during the training process performed best. I determined that recursive neural networks with LSTM layer performs most accurate and which answered RQ 1.

The RNN model was then used to classify two different sets of tweets concerning Bitcoin. The predicted sentiment for every tweet was used to calculate a sentiment score. In order to perform the empirical analysis, which answers RQ 2., I had to handle non-stationary data and co-integration. Pairwise correlation of all relevant variables was performed, indicating, that both datasets result in different correlations. The following estimated OLS regression models confirm prior correlation results. The Twitter sentiment ratio for tweets from big news sites and selected twitter-users with special insights in the Bitcoin system has a positively influence on Bitcoin price in the short-run (daily frequency). For hourly twitter and Bitcoin data, the total number of daily tweets have a negative impact on the Bitcoin price. The findings suggest that tweets can be used to make assumptions

of the future Bitcoin prices. The following VECM analysis results in the statement, that the sentiment score helps to 1-day-predict the price of Bitcoin (dataset 1). Days with an increase of positive tweets lead days with an rise of the Bitcoin price. VECM for hourly data shows that a greater number of tweets per hour lead hours with an increase of the Bitcoin price. If the dataset is limited to tweets from users with the most followers this relationship becomes insignificant. However, the results attest the forecasting capabilities of Twitter sentiment. Twitter can be used as information source for Bitcoin market participants, as it offers significant details on Bitcoin price changes.

There are several ways the above sentiment and empirical analysis could be extended. Even though Twitter Sentiment Analysis works effective and achieves accurate results, there are still limitations. Classifiers hardly can perceive linguistic details like irony, humor or sarcasm. Since Tweets often show exactly such details, false classifications are not infrequently. Already in pre-labeled training data, like the data I used to train the classifier, this is not taken account of. Even if the machine learning algorithms were able to learn irony or sarcasm, they cannot because of mislabeled training data. As a consequence of this, algorithms need better labeled training data, which are hard to obtain, since complex machine learning algorithms require huge amounts of training data. Another problem of the training data I used is, that it is not domain specific at all. Training data with reference to Bitcoin might yield a more effective classifier, when trained on Bitcoin related tweets. The tweets collected via a webscraper for dataset 2 contain a lot of spam. A effective way to remove spam-tweets from the data could improve robustness of the empirical findings. This is also the case for the length of the considered time period.

Both Twitter and Bitcoin offer large amounts of data for further research. Results of this thesis suggest several extensions to develop more precise sentiment classifier models. In the context of Bitcoin it would be interesting to see, if semi-supervised learning (Dorado and Ratté, 2016) can surpass the problem of non-domain specific training data. Another approach that could lead to more accurate classifiers is the use of complex deep neural networks (dos Santos and Gatti, 2014). Incorporating not only Twitter but also other sources like news pages or Reddit could provide further insights in Bitcoin price correlation to sentiment.