

# **Twitter posts as indicator for future price of Bitcoin**

by

**Temirlan Ulugbek uulu**

Bachelor Thesis in Computer Science

Prof. Michael Sedlmair  
Bachelor Thesis Supervisor

Date of Submission: May 10, 2018

With my signature, I certify that this thesis has been written by me using only the indicates resources and materials. Where I have presented data and results, the data and results are complete, genuine, and have been obtained by me unless otherwise acknowledged; where my results derive from computer programs, these computer programs have been written by me unless otherwise acknowledged. I further confirm that this thesis has not been submitted, either in part or as a whole, for any other academic degree at this or another institution.

Signature

Place, Date

## **Abstract**

Consider this a separate document, although it is submitted together with the rest. The abstract aims at another audience than the rest of the proposal. It is directed at the final decision maker or generalist, who typically is not an expert at all in your field, but more a manager kind of person. Thus, don't go into any technical description in the abstract, but use it to motivate the work and to highlight the importance of your project.

(target size: 15-20 lines)

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Goal and Motivation	1
1.2	Understanding Bitcoin	1
1.3	Difference between crypto-market and other financial markets	2
1.4	Twitter as source of emotional state	2
<b>2</b>	<b>Related Work</b>	<b>3</b>
2.1	Unrelated prediction algorithms	3
2.2	Closely related work	3
2.2.1	Predicting Bitcoin price fluctuations with Twitter sentiment analysis (2017) [7]	3
2.2.2	Algorithmic Trading of Cryptocurrency Based on Twitter Sentiment Analysis [8]	3
<b>3</b>	<b>Design and Implementation</b>	<b>4</b>
3.1	Data collection	4
3.2	Sentiment analysis	5
3.3	Feature vectors	5
3.3.1	Price vectors	5
3.3.2	Tweet vectors	6
3.4	Performing linear regression	7
3.5	Result evaluation	7
<b>4</b>	<b>Evaluation</b>	<b>7</b>
4.1	Results	7
4.2	Discussion	7
<b>5</b>	<b>Conclusions</b>	<b>7</b>
<b>6</b>	<b>Future Work</b>	<b>7</b>

# 1 Introduction

## 1.1 Goal and Motivation

Since early ages of financial markets, people have been trying to predict the future price of traded assets. With knowledge of future price or at least the direction of the price movement, one can gain profits or even more - manipulate the market. Knowing the price for the coming year, month, day or hour gives you by far more opportunities than simple buy low and sell high. This way or the other, knowing the future price is always good. This paper describes another attempt to predict future price of an asset. However, in this paper we don't try to predict the future price of ordinary financial market with usual stocks or physical commodities. We deal with cryptocurrency, in particular Bitcoin. To understand the core idea of prediction, one has to understand what is cryptocurrency and how it is traded.

## 1.2 Understanding Bitcoin

Initially, Bitcoin was invented as a currency for peer-to-peer electronic cash system, where users don't have to trust anybody. One can think of bitcoin as electronic money. All the transactions and creations of this coin are stored in the blockchain data-structure. The idea of blockchain was proposed nearly a decade ago in 2009 by someone who named himself as Satoshi Nakamoto [1]. There are many good sources out there that try to explain the technology behind this currency, but that knowledge is not necessary for the understanding of this paper. Here are few things that we need to understand about bitcoin.

Satoshi Nakamoto wrote in his paper: We have proposed a system for electronic transactions without relying on trust [1]. And by that he meant that one doesn't have to trust anyone to manage his bitcoins. No bank, government or other third party in between. As long as you keep your private keys safe, no one will be able to steal your bitcoins. That is achieved by a huge computational power needed to maintain the bitcoin blockchain. Huge amount of computers all over the world contribute to this network, i.e. they do computations to verify transactions and if one successfully finds the solution, he/she gets some amount of bitcoin for his work. At the moment of writing this, the total computational power is around 30 000 000 tera hashes per second [2]. If someone wants to do illegal transactions or steal someones bitcoins, he has to have the majority of the CPU power, which is nearly impossible to achieve. The difficulty of verifying transaction brings us to one of the major problems of this blockchain system. Transactions are costly and slow. Average transactions take couple of hours and may be faster if you are willing to pay more [3]. And there is no guarantee that your transaction will be confirmed at all. Theoretically it may be stuck in the pool of transactions forever. This makes it a bit inconvenient for bitcoin to be the substitution for the current payment ways. However, the great technology and the idea of trustless payment system made bitcoin a great asset for investments. Bitcoin became a digital gold and is currently traded in tens, if not hundreds, of exchanges all over the world. Anyone who has some money and is willing to invest can do so. Over the past years, there was tremendous growth of bitcoin price [4]. This growth leads to increase of people investing, which in its turn leads back to price growth.

With time, thousands alternatives of bitcoin were created [5]. Mostly they have the same

idea of blockchain and decentralization, but have their own features and advantages. Those alternative coins are shortly called alt-coins. Even though bitcoin is not the only cryptocurrency out there, we will analyze only bitcoin price in this paper, since it was the first one of a kind and is still leading among all altcoins according to the market capitalization [5].

### **1.3 Difference between crypto-market and other financial markets**

From a first glance, cryptomarket can seem nothing else but another type of financial market. But they are completely different. The major difference is that there is no real physical commodity or some company stock in cryptomarket. Bits in the agreed blockchain is all what is traded. The price is purely on a level of supply and demand. The highest price that the people are willing to pay and the lowest price for which others are willing to sell. No real world factors affect the price, except the ones that affect the mood of traders. Some people might believe that the price of bitcoin is controlled by people or group of people with a lot of money, also called "whales". But mostly whales are asleep, or in other words not actively trading in the markets. No doubt that they can change the direction of the market completely if they want to. But mostly, the price depends on the smaller active traders, at least that is what we assume here. We make a claim that bitcoin price depends on the mood of people all over the world. Thus, if this claim is true and we have the mood of all bitcoin traders in numbers, we could predict the future price.

### **1.4 Twitter as source of emotional state**

Now, to see if the claim in previous section is true, we need a way to extract the emotional state of people in digital format in order to be able to perform computations and predict the future price. And here we involve the Twitter. We make another claim and assume that Twitter has the emotional state of people trading the Bitcoin. We extract all the Twitter posts related to bitcoin and perform a linear regression on them together with history price to find the best predicting function. After getting the predicting function, we can apply it on the current state of peoples mood all over the world and get the future prices for bitcoin.

In order to check the above claims, we extract all the bitcoin related tweets and do sentiment analysis on them. VADER [6] was used in order to extract the polarity of Twitter posts, to be exact, it determined how negative or positive they are. Those sentiment scores and the historical bitcoin price data then have been fed into linear regression algorithm and the optimal predicting function have been obtained.

However Twitter doesn't allow to extract tweets older than on week and we don't have the full flexibility here. Even though if the algorithm was fed only with one month of data, the results were pretty impressive. The weekly predictions using the Twitter had accuracy between 66-100%. Provided results should strongly support the claims made earlier and encourage people to investigate more in the direction of considering the Twitter as one of the significant indicators for the price predictions of cryptocurrencies.

## **2 Related Work**

### **2.1 Unrelated prediction algorithms**

One has to understand the difference between previous works that have the same goal - predicting the future market state. Many good algorithms were developed in predicting the future price in other financial markets. Despite the fact that some of them achieved very good results, those works have very little to do with the described in this paper. That is mostly because the markets are completely different as stated in section 1.3. Thus, we only should look at the works related to the cryptocurrency market. Even if the field is relatively new, a lot of good attempts were done in developing predicting algorithms. Different algorithms use different input data. In most of the cases, in cryptocurrency markets, people try to analyze the historical price of the coin. We will not be looking at those kind of works and will focus only the ones that are very similar to what we tried to do here.

### **2.2 Closely related work**

#### **2.2.1 Predicting Bitcoin price fluctuations with Twitter sentiment analysis [7]**

As the paper states itself, the work is ... done by a naive method of solely attributing rise or fall based on the severity of aggregated Twitter sentiment change. The work tries a very naive approach and sees if the increase of aggregated Twitter sentiment leads to some direct changes in price. As a result of the work, they get fluctuations of the price in the future. [7] Knowing the fluctuations of the price is nearly useless to the traders, unless the accuracy is 100%. Imagine the case if the output of the algorithm for the next three hours is [decrease, increase, increase]. This might be the case for [-10%, +1%, +1%] as well. Thus trader wont be able to get the most profit, unlike the case if he knew the exact magnitude of the price movement.

#### **2.2.2 Algorithmic Trading of Cryptocurrency Based on Twitter Sentiment Analysis [8]**

This work does some more advanced work in terms of data processing. It performs multiple ML techniques on the Twitter and Bitcoin data that has been collected in around 20 days [8]. First of all, when comparing to the big history that Bitcoin already has, 20 days of data is a very small dataset to train your system. Moreover, the data has been collected from November 15th to December 3rd in 2015. As Bitcoin usually close to end of the year, the price was very unstable during that period. Somewhat steady down-trend of the price abruptly turned into fast growth. This might be not the very best period to collect the data, but we can only do assumptions here since their work hasnt been tested on other data than the one collected during that period. The main goal described in that paper was to again get the binary result or, to be more precise, to know if the price will increase or decrease, which, as we discussed above, is not very useful. Two approaches were made to achieve the result: feeding words from tweets separately into the learning algorithms and feeding sentiment scores of the tweets into the algorithms. We dont cover the first approach, since it is not closely related to the our work. However, the second

approach they tried come very close to what has been done here. They took the whole tweet, did sentiment analysis on it and used those labels as input for Naive Bayes. The main flaw here was that they used text-processing.com as their sentiment analyzer, which would calculate the positivity, negativity, and neutrality scores [8]. As has been stated in the web-site itself, the sentiment analyzer is composed of 2 classifiers trained on movie reviews. If your text is not similar to movie reviews, then its less likely to make a correct guess [9].

Unlike the papers described above, we try to predict the exact price in the future instead of predicting just the binary increase/decrease of the price. This information will be much more useful for traders. Moreover, we use linear regression to get the optimal predicting function instead of naive approach. As an input to our linear regression, we use polarity scores of tweets. We use VADER to do sentiment analysis on Twitter posts and will later see how VADER outperforms other algorithms. Data for one month has been collected, which is still relatively small dataset, but both price downtrend and uptrend were observed during the period of data collection.

### 3 Design and Implementation

Very general idea about the design of this work is this: we take all the bitcoin related posts from Twitter, do sentiment analysis on them using Vader, feed this sentiment scores into linear regression, get the optimal predicting function, and finally apply it to the current state to get the future prices.

#### 3.1 Data collection

Data collection is hard when you are dealing with Twitter, because one cant extract old tweets using the official Twitter API. The Twitter Search API searches against a sampling of recent Tweets published in the past 7 days [10]. One has to pay to get more privileges. Using some third party APIs, like GetOldTweets-python [11], wouldnt do the job as well, because they use the Standard Search API in that back. ... its important to know that the standard search API is focused on relevance and not completeness. [10]. Thus, using third party APIs cant get you complete data as has been seen during the process of data collection. Using GetOldTweets [11], we tried to get Bitcoin related tweets since the beginning of 2018, but couldnt get a satisfactory results. It returned only around 1 million tweets, which cant be a complete dataset by any means, given the fact that we were getting more than 100 000 bitcoin related tweets per day when getting the live streaming. Thus, obtaining the complete dataset of old Twitter posts was not possible. Instead, tweets were collected using live stream. Tweepy [12] was used to interact with Twitter API. It is python library, which gets the job done. Every tweet, which contained btc/bitcoin in it, was streamed and stored. In order to be able to extract tweets 24/7, the web-server was rented from Aruba Cloud [13] with minimal requirements: 1GB RAM, 20GB space. MySQL database was set-up on the same server to store the data. The data collection has started on March 30th and around 5 million Twitter posts were stored in approximately 1 month.



## 3.2 Sentiment analysis

What we wanted from tweets is to know how positive or negative they are. Assuming that we have very huge stream of tweets, we need this analysis to be very quick and preferably as accurate as possible. VADER (Valence Aware Dictionary for sEntiment Reasoning) [6] is the perfect tool we could use for this. It is very computationally cheap and accurate at the same time [6]. Moreover, in the latest update, they made a major improvement in performance: Reconstructing for much improved speed/performance, reducing the time complexity from something like  $O(N^4)$  to  $O(N)$  [14]. With hundreds of tweets arriving to our server every day, performance was very important when making decision in which tool to use. ... a corpus that takes a fraction of a second to analyze with VADER can take hours when using more complex models like SVM (if training is required) or tens of minutes if the model has been previously trained. [6] Other than performance advantages, it has transparency. VADER is an open-source tool that can easily be inspected and extended [6] [14]. This is important if we would like to introduce it a cryptocurrency related lexicon. VADER also doesn't require any training dataset and is completely self-contained [6]. Given all these conveniences, this tool still outperforms other well-known tools and techniques. When looking at social media context VADER outperforms seven well-established sentiment analysis lexicons (LIWC, GI, ANEW, SWN, SCN, WSD, and Hu-Liu04 opinion lexicon) and four machine learning algorithms (Naive Bayes, Maximum Entropy, SVM-Classification, and SVM-Regression) [6]. Thus, by considering all mentioned factors, VADER was chosen to do the sentiment analysis for this work.

## 3.3 Feature vectors

In order to be able to use linear regression, we need to think of what should the input and output vectors should contain to meet our problem well. Linear regression was performed on different input data or to be more precise, three types of input: historical price, tweets vector, and combination of these two. Output of our predicting algorithm will be the price vector. Predictions in this work are done for different frame-widths and intervals. This terminology was chosen for this paper with the following definitions. Frame-width is the length of the period of time for which the prediction is being made. For example if the frame-width is 1 hour, we are trying to predict the price over the next coming hour. Interval is the distance between predictions. For example if the frame-width is 1 hour and the interval is 1 minute, we are trying to predict the price for each minute of the coming hour, or in other words, we do 60 predictions for the coming hour. The terms frame-width and interval will be now used without further explanations.

### 3.3.1 Price vectors

Price vector are the vectors representing price of Bitcoin over certain period in the timeline or in our case the whole frame-width. In all cases of linear regression done in this work, the output vector is the price vector for the future. In some cases, price vectors of the past are used as input vectors as well. The structure of the price vector is straightforward. If we have frame-width of 1 day and interval of 1 hour, our price vector would be of dimension 24 or in other words, it would contain hourly price information for the whole 24 hours or 1 complete day. Here, we make a small change. Instead of storing simply the prices, we

store the price change in percentage from the beginning of the frame. Thus, the first value of the price vector would always be zero since that is the point of time which we take as our baseline. The next value of the vector would be the price change after one interval. If we continue our example, and the price changed from 100 to 110 in one interval, first two values of the price vector would be 0 and 10. This way we ensure that we are looking at the trend instead of looking at the concrete price changes and being bound to some time periods. The price vector from 2010 can now be similar to 2018 and our prediction algorithm would recognize it better. Imagine if we had prices instead of price changes. Then, the price vector from 2010 would contain couple of dollars, if not cents, unlike the vector from 2018 which would contain thousands of dollars. Using our version of price vectors, we reduce the size of vector by 1, since the first value will always be zero. That place can be reserved for future use. The very straight forward way to use it would be to store the initial price as the first value of the vector.

### 3.3.2 Tweet vectors

When talking about hundreds of thousands of tweets, extracting the feature vector from it is not that obvious. Keeping in mind the idea/claim, that the price depends on the mood of people all around the world, all we need is a feature vector that, simply said, can tell how happy or sad is Twitter Bitcoin community. The raw dataset we have after processing tweets with VADER is list of values (or sentiment scores of tweets) ranging between -1 and 1, where -1 is completely negative and 1 is the opposite. Now we want to divide them into  $M$  partitions and calculate the percentage of tweets coming to each partition. For example if  $M$  would be 2, we would just have two partitions:  $[-1,0)$ ,  $[0,1]$ . Or in other words, we just calculate what part of tweets is negative and what part is positive. Thus, we our tweet vector would have dimension of 2. Here is what was described shortly:

$$M_{partitions} \Rightarrow VTweet \in \mathbb{R}^M \quad (1)$$

$$VTweet_i \in [0, 1], \text{ where } i = 1, \dots, M \quad (2)$$

$$\sum_1^M VTweet_i = 1 \quad (3)$$

However, it might be tricky to choose the optimal  $M$  Its not obvious what resolution would be enough to preserve enough data. Moreover, the value of optimal  $M$  may depend on the frame-width and interval. For example if we are trying to predict the price for the coming year, just the small  $M$  may be sufficient to find out the general mood for long term predictions. Thus, the testing data was divided into parts and cross validation has been performed for each frame-width and interval to find the optimal  $M$ . The error was computed for different values for  $M$  starting with 2 and the error was calculated on each iteration.  $M$  was being increased until the testing error stopped decreasing or the decrease was below certain threshold. The following optimal  $M$ s were discovered by doing cross-validation:

// Include the chart here

### **3.4 Performing linear regression**

Linear regression is not the most advanced machine learning technique and most people might say that other techniques, like neural networks, are much better at approximating functions or doing price predictions. However, linear regression is the perfect starting point. If the correlation between the mood of Twitter community and Bitcoin price is defined by some non-linear function, linear regression would most probably fail at making predictions. But due to its simplicity and computational cheapness, linear regression was chosen as a starting point. If mood can be mapped to price by some linear affine mapping we would already get good predictions and there would be no need to do computationally expensive and advanced techniques. Even if predictions won't be as good as we wish they were, linear regression could make the further direction of research clearer and getting even somewhat good (e.g. better than random) could motivate people to investigate more into this field.

As mentioned above, we have three types of inputs to feed into the linear regression: price vectors, tweet vectors, and combined price-tweet vectors. In all three cases, we perform standard linear regression procedure. We divide our data instances into training and testing (validation) sets. We learn the optimal function or optimal matrix from the set of training instances and we evaluate the error of this optimal function by applying it on the validation set. By doing this on different datasets, or to be more precise on only prices, only tweets, and price-tweets combined, we could learn how tweets could be useful when predicting future price using linear regression.

### **3.5 Result evaluation**

Bla-bla.

## **4 Evaluation**

### **4.1 Results**

### **4.2 Discussion**

## **5 Conclusions**

Summarize the main aspects and results of the research project. Provide an answer to the research questions stated earlier.

(target size: 1/2 page)

## **6 Future Work**

## References

- [1] Satoshi Nakamoto. Bitcoin: A Peer-to-Peer Electronic Cash System.
- [2] Blockchain.info. <https://blockchain.info/charts/hash-rate>, May 2018.
- [3] Blockchain.info. <https://blockchain.info/charts/avg-confirmation-time>, May 2018.
- [4] Coinmarketcap. <https://coinmarketcap.com/currencies/bitcoin/>, May 2018.
- [5] Coinmarketcap. <https://coinmarketcap.com/>, May 2018.
- [6] Eric Gilbert C.J.Hutto. Vader: A parsimonious rule-based model for sentiment analysis of social media text, 2014.
- [7] Jacob Loenneke Evita Stenqvist. Predicting bitcoin price fluctuation with twitter sentiment analysis, June 2017.
- [8] Stephanie Rosales Stuart Colianni and Michael Signorotti. Algorithmic trading of cryptocurrency based on twitter sentiment analysis.
- [9] Text-Processing.com FAQ. <http://text-processing.com/docs/faq.html>, May 2018.
- [10] Official Twitter API. <https://developer.twitter.com/en/docs/tweets/search/overview/standard>, May 2018.
- [11] Jefferson-Henrique. <https://github.com/jefferson-henrique/getoldtweets-python>, April 2018.
- [12] Tweepy. <http://www.tweepy.org/>, May 2018.
- [13] Aruba Cloud. <https://www.arubacloud.com/>, May 2018.
- [14] Vader Github repository. <https://github.com/cjhutto/vadersentiment>, 2018.