

Using Sentiment Analysis to Predict Cryptocurrency Prices

Mao Guan

mg3844@columbia.edu

David Lee

jl4397@columbia.edu

Carlo Provinciali

cp2984@columbia.edu

Weiqi Tong

wt2247@columbia.edu

Abstract

We investigate the relationship between community sentiment and the price of cryptocurrency. Our goal is to prove whether features that capture the sentiment of online discussion boards can help predict Bitcoin price movements. We collected daily discussions threads on a Reddit community where users share their expectations of Bitcoin markets. We extracted sentiment features from these corpora using domain specific lexicons (VADER and SocialSent) and trained an LSTM model to predict hourly price variations. Our results indicate that the model prediction power benefits from sentiment features when predicting large price shifts. In addition, we discovered that our model is most effective at predicting short term price movements (10 hours in the future).

1 Introduction

The cryptocurrency world has recently been at the center of attention after the price of one Bitcoin reached \$19,000, which in turn attracted a new wave of investors eager to make profit out of the fast growing market. Despite the growth in popularity, the technology behind many cryptocurrencies is often not well documented and generally not well understood by the public aside from a small fraction of experts. For this reason, the vast majority of investors rely on online communities to gather information and base their investment decisions off the opinions shared by other users. Therefore, it is widely believed that general excitement and confidence from the community toward a cryptocurrency

project can lead to more investments and increase the price while fear, doubt, and uncertainty can cause the price to sink. We propose using sentiment as a way to predict price fluctuations caused by shift of community opinions toward a particular cryptocurrency and analyze its impact on its price.

This paper documents our efforts of measuring sentiment across a number of subreddits that discuss specific cryptocurrency markets, and using the extracted sentiment metrics to predict the price's movements of the respective cryptocurrency. In the initial phase, we identified daily threads on Reddit where users exchanged their opinion and sentiment for the future price of a specific cryptocurrency. After extracting this content, we used the VADER algorithm to calculate the sentiment score using both the VADER lexicon and the SocialSent lexicon pre-trained on r/btc. We then input these features along with other information such as user flair and average comment popularity to an LSTM model and used it to predict the cryptocurrency's future price movements.

2 Related Work

Several prior papers have analyzed the relationship between sentiment analysis and the price movement of various financial assets. Bollen et. al (2010) investigated a causative relation between public mood from a collection of tweets and showed that calmness and happiness mood are causative of the Dow Jones Industrial Average by three to four days. The findings from Su et. Al (Hong et al., 2017) demonstrate that dissemination speed and

visibility of a post can also further enrich stock market predictions from stock market data.

More recently, given public interest in cryptocurrency, researchers have begun to investigate the relationship between social media sentiment and the price movement of various cryptocurrencies (Mai et al., 2015; Bin et al., 2016; Stenqvist and Lnn, 2017). Because of the promising results from these papers, we have decided to try and take a novel approach to predicting prices that combines forum data tagged with sentiment values with an LSTM model.

2.1 Why Reddit?

Prior research has shown that forum data is useful in predicting cryptocurrency price fluctuations than other forms of social media that are popularly used such as Twitter. Mai et. al (Mai et al., 2015) used vector error correction models (VECM) to capture relationships between time series data and Granger causality to show that when aggregated at the interday level, the sentiments on forum messages are more telling indicators of future Bitcoin returns than are Twitter messages. In addition, Kim et. al (2015) applied a model that uses sentiment analysis on gaming forum data to predict market price fluctuations of gaming currencies traded in cash to Bitcoin and found that sentiment data from Bitcoin forums are highly indicative of the price movement in the next one to seven days.

Other papers also mentions the reason why online communities might be more helpful for predicting price movements of cryptocurrencies compared to social platforms such as Twitter. Online communities such as Reddit serve as forums where people can share opinions regarding topics of common interest (Bernstein et al., 2011; Yong and Kim, 2011). Such communities reflect the responses of many users to certain cryptocurrencies on a daily basis, while social media contains mainly expression of current feeling that are shorter and more fragmentary. More specifically, Phillips and Gorse (2017) argue that discussion boards within Reddit do a better job at representing communities that share a common interest. In particular, the user activity in specific subreddits (i.e. Reddit communi-

ties) can be a good proxy of the excitement of the community for a particular cryptocurrency, which is then reflected in its price movement. In fact, cryptocurrencies are largely traded online, where many users rely on information on the Web community to make decisions about trading on them (Maurer et al., 2013). For this reason, Kim et. al (2016)[6] propose a method to predict price movements and number of transactions of cryptocurrencies based on user comments on online cryptocurrency communities such as Bitcoin Talk. Daily topics and relevant comments as well as following replies are analyzed to determine how the opinions of community users are associated with price and number of transactions of cryptocurrencies.

2.2 Why Domain-Specific Lexicons?

Texts on online communities involve frequent use of neologisms, slang, and emoticons that transcend grammatical usage, therefore, this paper applies the VADER method that C.J. Hutto and Eric Gilbert introduced (Gilbert, 2014) to parse such expressions and analyze social media texts by drawing on a rule-based model. VADER (for Valence Aware Dictionary for Sentiment Reasoning) is a gold-standard sentiment lexicon that is especially attuned to microblog-like contexts. VADER retains (and even improves on) the benefits of traditional lexicons like LIWC. VADER is bigger, yet just as simply inspected, understood, quickly applied (without a need for extensive learning/training) and easily extended.

Based on the learning data at the time of higher accuracy, the types of comments that most significantly influenced fluctuations in the price and the number of transactions of each cryptocurrency were identified. Opinions affecting price fluctuations varied across cryptocurrencies. Positive user comments significantly affected price fluctuations of Bitcoin, whereas those of the other two currencies, Ethereum and Ripple, were significantly influenced by negative user comments and replies. Moreover, user opinions proved useful to predict the fluctuations in 6 to 7 days.

Although this method produced good results, it did not include domain specific lexicon, which according to Hamilton et al. (Hamilton et al., 2016)

can improve sentiment analysis tasks. For this reason, we plan to leverage SocialSent, the package they developed that uses a random walk model and a small collection of seed words to generate domain specific sentiment embeddings. This gives us an opportunity to improve the performance of our model without requiring huge corporas for training, which makes it easily adaptable to new domains.

2.3 Why LSTM?

Traditionally, the Time Series Methods like the ARMA, VECM and etc. (Mai et al., 2015) are popular for trading models. However, such models rely heavily on the specification and assumptions of the data which, in reality could be wrong. The LSTM model, on the contrary, outperforms the classic time series models in that it is more flexible on the assumptions and more intelligent in its forgetting and remembering the past information.

Recently, there have been many efforts focused on exploring how variations in sentiment are related to cryptocurrencies prices changes (Bin et al., 2016). Many of these projects apply the traditional econometrics methods to forecast future prices; However, we believe there is an opportunity to apply advanced machine learning algorithms such as the Long Short-term Memory Model (LSTM). LSTM is famous for its advantage in solving the problems of gradient vanishing in the traditional RNNs. LSTM models have successfully been applied in similar domains such as algorithm trading in equity, futures and etc.

3 Dataset

3.1 Overview

To build our corpus, we scraped 337,190 comments from Daily Discussion posts on a subreddit on www.reddit.com called BitcoinMarkets¹. This subreddit features Daily Discussion posts which are auto-generated, daily post where users are free to have open-ended discussions about Bitcoin trading, including discussions related to the day's events, technical analyses, trading strategies, and quick questions. Although there are a number of Bitcoin-related subreddits on Reddit, we chose to scrape comments from the Daily Discussion posts of the

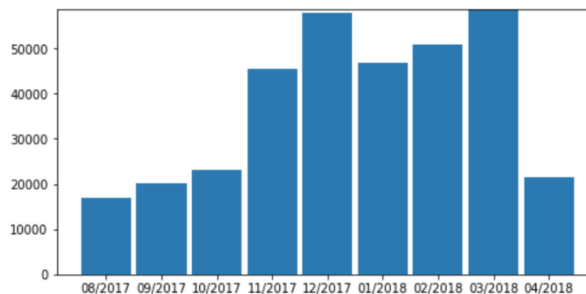


Figure 1: Distribution of comments by month.

BitcoinMarkets subreddit because it provided a forum for a general discussion of Bitcoin and has higher quality content relative to the other Bitcoin communities. However, one of the most important reasons why we chose this subreddit is that users have the option to add a textual "flair" so that their trading position on Bitcoin is advertised to the community. Users can choose from these four flairs, sorted from most optimistic to least optimistic about the prospects of Bitcoin: "Long-term Holder", "Bullish", "Bearish", and "Bitcoin Skeptic." This provides us with additional information that can be used to infer a specific user's sentiment toward the market.

3.2 Dataset Analysis

The comment ID, parent comment ID, number of upvotes and downvotes, body, time of creation, day of daily discussion, and author's flair were scraped from Daily Discussion posts from August 2017 to April 2018. Figure 1 shows the the number of comments by month. There is a linear correlation between the number of comments and the volatility of Bitcoin that month. Figure 2 shows the distribution of flairs among the 21.5% of comments whose users showed their opinion. The disproportionate number of bullish positions ("Long-term Holder" and "Bullish") suggests the presence of a skewed dataset. One explanation for this might be that people who are skeptical about Bitcoin are less likely to participate to discussions around Bitcoin investments.

Further analysis of the comments where the author's opinion was written was done using Scattertext², a tool to find distinguishing unigrams in

¹www.reddit.com/r/BitcoinMarkets

²www.github.com/JasonKessler/scattertext

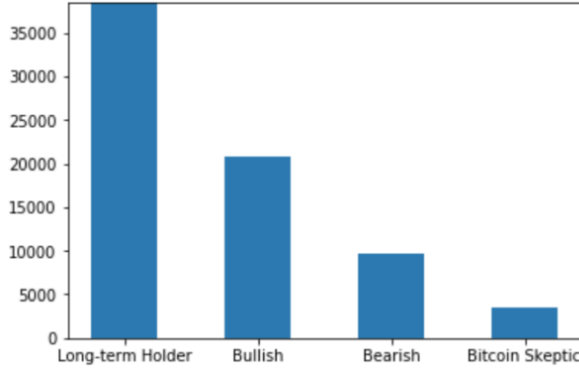


Figure 2: Distribution of flairs among flaired authors.

Bull	Bear	Whole Corpus
09	12h	bitcoin
01	chikou	btc
2019	kumo	bch
08	3d	gdax
02	despair	bullish
performance	ill	crypto
04	dcb	bitfinex
present	6h	alts
05	4h	bearish
06	ish	coinbase

Table 1: Top 10 characteristic terms for each class of user flairs and the whole corpus.

small-to-medium-sized corpora. The library uses a scaled F-score to determine which words are the most characteristic of a class. In this dataset, we split the data into two classes, bulls and bears. Table 1 shows the top characteristic unigrams for both classes.

The "bull" terms (which represent users with confidence in the market) are comprised mostly of numerical unigrams. One hypothesis as to why that might be is because "bulls" tend mostly to talk about the long-term future. The integer unigrams might correspond to months and years. In contrast, the numerical unigrams on the "bearish" or skeptical side measure days and hours. In addition, the "bear" terms also contain unigrams such as *chikou* and *kumo* that refer to technical analysis methods. This suggests that technical analysis³ techniques might be common among "bears" or people predicting future Bitcoin price decreases. Figure 3 shows the full result of the plotting.

³investopedia.com/terms/t/technicalanalysis.asp

4 Evaluation

As mentioned in the introduction, the main goal of the project is to prove that sentiment features can help predict Bitcoin's future price movements. For this reason, most of the features that power our model are related to the sentiment on daily discussion threads. Although, we decided to include additional features that were proven useful for this task by previous literature such as previous price information and dissemination (Su et. Al [12])

4.1 Framework

Since the goal of our project is to determine whether there exists a positive effect of the sentiment data to the prediction of bitcoin prices and how different groups of features affect the results. We designed our experiment framework as follows. We evaluated our data following a two step approach. As figure 4 suggests, we determined the contribution of each feature group to the accuracy of the prediction. Finally, we trained our final model with all set of features and tested the best time interval for future predictions. All the analysis is based on the parameters and metrics including, f-score, precision, recall and auc.

4.2 Features

4.2.1 Sentiment Features

For the sentiment features, we leveraged the VADER sentiment analysis tool⁴ to generate sentiment polarity score for each comment in the daily discussion thread. This tool gives a positive, negative, and compound sentiment for a given corpus based off an input sentiment lexicon. For the latter, we used both the VADER sentiment lexicon itself and the pre-trained Socialsent lexicon⁵ for r/btc⁶. Using the sentiment scores for each comment, we constructed an hourly average sentiment score for each category. In addition, in order to capture the degree of influence that each comment might have had on the community, we built hourly average sentiments scores by weighing each comment by the number of votes and number of children comments respectively. The intuition is that comments with

⁴<https://github.com/cjhutto/vaderSentiment>

⁵<https://nlp.stanford.edu/projects/socialsent/>

⁶[reddit for general Bitcoin-related discussion](https://www.reddit.com/r/bitcoin/)

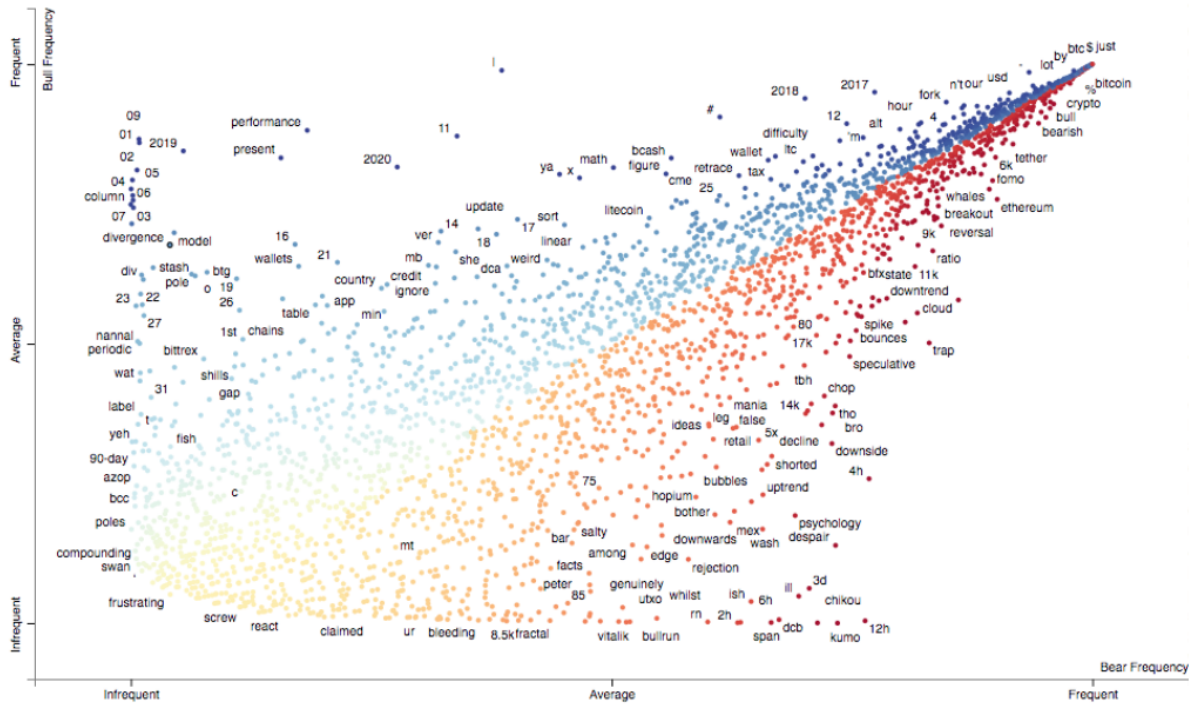


Figure 3: Unigrams plotted against two axes corresponding to frequency for two classes. The further up the word is, the more often it is used by bulls; the further right the word is, the more often it is used by bears.

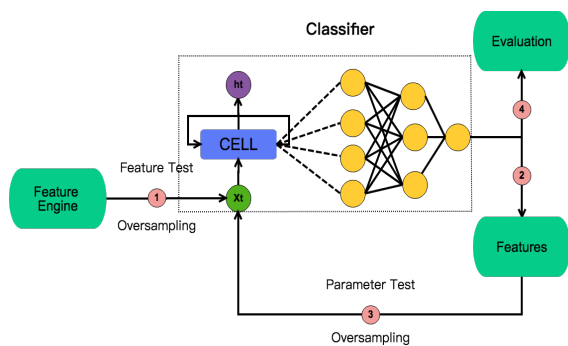
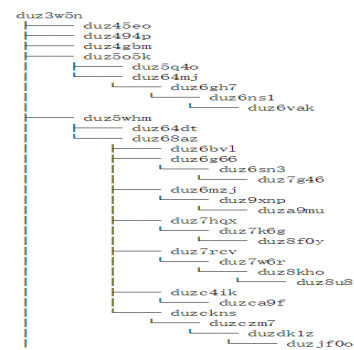
**Figure 4: Framework**

Figure 5: Tree-like structure in reddit

many "upvotes" or that prompted a longer discussion (measured by the number of "children" comments) convey sentiments that are more likely to be shared across the community, hence the larger weight assigned to these comments when computing the hourly sentiment. The comments on a Reddit post form a tree-like structure, as shown in the figure 5, where each comment id refers to a comment posted and child of a note refers to a comment reply.

Figure 6 shows the trend of VADER negative score weighted by different methods. The simple

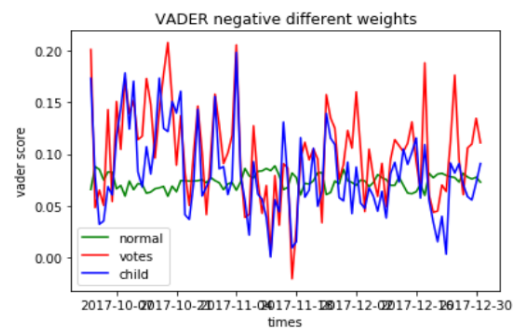


Figure 6: Trend of negative vader scores

Features	Number(33)	Notes
VADER_avg	4	pos,neg,neutral, compound
SocialSent_avg	4	pos,neg,neutral, compound
VADER_vote_avg	4	weighted by num of votes
SocialSent_vote_avg	4	
VADER_child_avg	4	weighted by num of child
SocialSent_child_avg	4	
Proportion of Flairs	5	5 types
Num of word	1	
Num of comment	1	
Num of replies	1	
Past Price	1	past price

Table 2: List of Features.

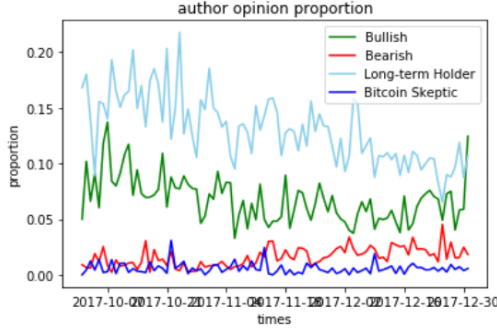


Figure 7: flair type distribution

average score is relatively stable and has a small variance, while two other methods are much more volatile.

In total there are 24 sentiment features.

4.2.2 Other Features

We also included basic statistics by hour such as average number of words in each comment, total number of comments, and average number of comment replies by hour. We include the proportion that each user flair appeared in the discussion. The figure 7 shows the trend of distribution of four flair types. All the hourly averages for these group of features are simple averages and are not weighted by the "influence" of each comment. In total there are 8 statistical features.

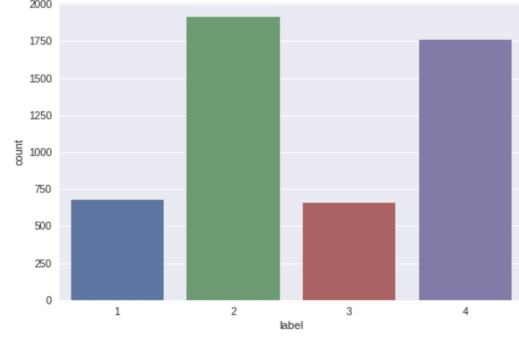


Figure 8: Label Distribution

4.3 Labels

We extracted the historical Bitcoin prices using the GDAX API⁷. We then grouped the prices into four categories: Large Rise (1), Minor Rise (2), Large Drop (3), Minor Drop (4). We defined a "minor" change as one within one standard deviation of all hourly price changes in the previous 30 hours, while "large" changes are outside of one standard deviation. The reason for this distinction is that we believe minor changes might be due to market corrections rather than macro-trends and overall changes in community sentiment. The figure below shows the proportion of each category in the whole dataset. As the figure show, most of the change of price is minor change and it is about the same ratio between significant and minor for both rise and drop.

4.4 Model and Metrics

Due to the imbalance of labels in our dataset, we adopted SMOTE (Chawla et al., 2002) as an oversampling method. The oversampling process allowed us to use a balanced dataset for the classification task.

The models we use is a Long Short Term Memory Model (LSTM) (Hochreiter and Schmidhuber, 1997) connected by two dense layers. LSTM is famous for its ability to capture the historic information for both long period and short period. We believe the LSTM was the best option for determine the effect on classification of a time series such as Bitcoin price changes. Besides, keeping the model structure simple helps us to neglect other model parameters effects on the final results.

The evaluation metrics is determined based on

⁷<https://docs.gdax.com/#get-historic-rates>

Table 3: Feature Category

Category	Description
Base	Sentiment(VADER,Socialsent)
Base	+ Statistical Features
Votes	Votes Features(VADER, Socialsent)
Children	Child Features (VADER ,Socialsent)
Price	Price Features (Last close price change)

Table 4: Feature Test Results for: window 10, 50 epochs

Features	B	B,P	B,V
auc	0.500	0.547	0.509
precision_micro	0.289	0.277	0.300
precision_macro	0.277	0.211	0.333
recall_macro	0.285	0.275	0.298
recall_micro	0.289	0.277	0.301
f1_score_macro	0.230	0.228	0.252
f1_score_micro	0.289	0.277	0.300
Features	B,C	B,P,V	All
auc	0.530	0.537	0.496
precision_micro	0.308	0.286	0.304
precision_macro	0.345	0.295	0.308
recall_macro	0.310	0.287	0.303
recall_micro	0.308	0.286	0.304
f1_score_macro	0.273	0.283	0.294
f1_score_micro	0.308	0.286	0.304

the multi-classification task. Since we have 4 class, we use the categorical crossentropy as the cost function in Keras. And for final evaluation, we use the metrics: auc, precision-micro,precision-macro,recall-micro,recall-macro,f1score-micro, and f1score-macro.

4.5 Results

For notation, we classify the sentiment features into Base, Children, Votes. Table .3

We choose the look-back window for the LSTM to be 10 and trained it for 50 epochs and we tested six groups: Base, Base + Price, Base + Votes, Base + Children, Base + Price + Votes, Base + Price + Votes + Children. From the result table, we observe that all of these combinations yield positive results in terms of accuracy compared to a random classifier (0.25 accuracy). This suggests that all the feature groups are helpful when predicting the classification. As is noted in table 4(B,P,V,C represents Base, Price, Votes, Children Features), Base +

Table 5: Price Feature 500 epochs

Metric	Price 500 epochs
auc	0.521
precision_micro	0.420
precision_macro	0.399
recall_macro	0.419
recall_micro	0.420
f1_score_macro	0.399
f1_score_micro	0.420

Table 6: Confusion Matrix:Price Feature 500 epochs

Actual/Predicted	1	2	3	4
1	338	90	101	45
2	181	103	158	137
3	88	67	356	72
4	141	98	150	165

Price is worse than Base alone, though better than benchmark. There are multiple reasons why this is the case; the price might be a confounding factor and harm the prediction our our model - which is counter-intuitive - or 50 epochs might not be enough for the model to learn the dynamics of daily price changes. Since we believe the second option is the more likely, we decided to run 500 epochs on the price feature alone.

After examining the classification results on our validation set (Table 5 ,6),we find that the price feature alone is helpful for prediction and that the model can more accurately predict large price changes than small variations. This finding aligns with our intuition that minor change could be due to market corrections that generally do not reflect into the community's sentiment. Furthermore, the fact that the price feature itself seems to have a similar performance reinforces the idea that smaller variations in prices might be more random than larger shifts. On the other hand, the model that included the feature groups Base + Children + Votes (Table 7, 8) did not perform as well as the model that was also trained on price information.

For the second stage, we focused on determining the best time interval for prediction by testing our final models for different time intervals. In this scenario, we only evaluated the predictions of the base model that was only trained on the sentiment related features.(figure 9) We compared results for the fol-

Table 7: Base + votes + children 500 epochs

Metric	Price 500 epochs
auc	0.519
precision_micro	0.345
precision_macro	0.343
recall_macro	0.344
recall_micro	0.345
f1_score_macro	0.339
f1_score_micro	0.345

Table 8: Confusion Matrix:Base + votes + children 500 epochs

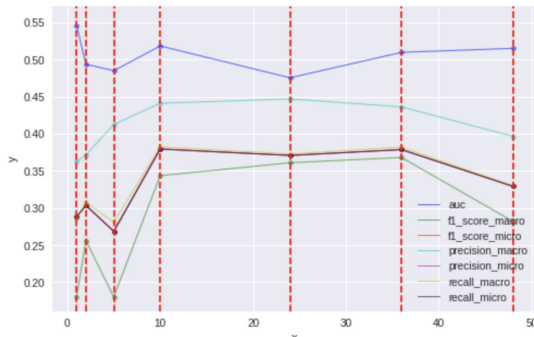
Actual/Predicted	1	2	3	4
1	212	120	173	69
2	124	172	173	110
3	127	106	276	74
4	135	150	138	131

lowing time hour intervals: 1,2,5,10,24,36,48. As the result shows, the model performs best when predicting price changes 10 hours ahead. We concluded that a look-back window size of 10 hours is the best for prediction.

5 Conclusion and Further Research

5.1 Conclusion

Our analysis shows that sentiment features extracted from r/bitcoinmarkets improve the prediction of future short term price movements of Bitcoins, especially in the case large price variations. On the other hand, the model trained using sentiment features did not perform well when classifying small price variation, which we hypothesize might in part due to the weaker role of community sentiment in causing these fluctuations. When training a LSTM model

**Figure 9:** Parameter:windowsize Test Result

using sentiment features, the best performance is obtained when predicting price changes 10 hours in advance. Other features such as prior price movements also proved very useful at predicting future prices.

5.2 Further Improvement

We acknowledge there are still a number of opportunities for improvements and further research which we could not explore due to time constraints.

5.2.1 Valuation Methods

We report our results on a validation set which was collected along with our training set. In order to prove the model is effective at predicting future price movements, we should track its performance on Bitcoin price that that we have not yet observed. Moreover, we could use the model's predictions as an input to an automated investment strategy and evaluate the quality of the predictions based on our portfolio's gains or losses - similarly to what was done by Phillips and Gorse (Phillips and Gorse, 2017).

5.2.2 Sentiment Embedding

Word embedding methods such as word2vec and glove have been widely used in academia and industry, and there has been research on using embedding features as inputs to neural networks, which can be used to determine the sentiment of a sentence. However, word2vec embedding only encodes the context but not the sentiment information. (Tang D, 2014) proposes a method for combining sentiments with word2vec, which results in word embedding in sentiment domain. Their approach includes not only the loss function from context words, but also sentiment polarity annotated for the sentence (twitter in their case). So they use annotated twitter dataset to train a word embedding that encodes sentiment information. Such embedding features can be regarded as high-dimensional sentiment scores. We can just average the sentence embedding vectors (for instance 100 dimensions) as input features, or utilize some dimension reduction methods. We can test the effectiveness of these features compared to lexicon-based features such as VADER and SocialSent.

5.2.3 Tree-structural or user-based Features

Reddit comments have a tree-like structure. The complicate structure actually encodes all the children comments replied to a first-level comment, and

there could be many levels in a parent comment. We tried to parse such a tree but it takes large running time. We can compare the influence of sentiments from large trees that involve more discussion in the community than small trees or even one-node trees. We can also extract user information from the original data set, and try to find who are the influential users in this subreddit, and whether their comments have many children which indicates more discussion.

5.2.4 Generalization for Other Instruments

One limitation in our research is that we only focus on cryptocurrency. For other financial instruments such as equity and commodity, there exist some communities where investors discuss and exchange their ideas on their opinions. We could apply a similar approach to conduct sentiment analysis on other markets and test the capability of our model. Although, since those markets depend largely on other factors such as news and global markets, sentiments might only be partially indicative of future prices.

6 Team Work Distribution

David: I was mostly responsible for everything related to finding and scraping the initial dataset. I collected all of the data using Reddit's API for /r/BitcoinMarkets as well as /r/ethtrader, which we did not end up using. I found the scattertext module to generate the visualizations and I helped with some of the prior research into why we should use Reddit. For the paper, I was responsible for writing the Dataset section, generating all accompanying graphs and tables, and helped out with creating other tables in the paper as well.

Mao: I was responsible for data analysis, including framework design, data labeling, model construction, and result analysis. I constructed neural-network model and use it to test different sentiment feature groups, different parameters and did comparison among them to testify the assumption concerned. I also did some research about sentiment embedding with Weiqi. For the paper, I am responsible for the Evaluation section, including Framework, Labels, Methods and Metrics, Results and all the related figures and tables in these parts.

Weiqi: I was responsible mainly for feature engineering and preparing for the dataset as input into the model. I also researched on methodology such as VADER. Meanwhile, I used the API mentioned to obtain price information of the Bitcoin. I contributed mainly to the code of extracting features and writings of this part in the final paper. I also discussed some advanced method such as sentiment embedding with Mao and I wrote the further research parts in the final paper.

Carlo: I was responsible for creating the sentiment features from the data provided by David. I researched the VADER algorithm and the SocialSent lexicon and worked closely with Weiqi on the feature engineering. I was responsible for the Abstract, Introduction, Related Work, Conclusion, and general review of all sections of the paper.

References

- Michael S. Bernstein, Andrs Monroy-Hernandez, Drew Harry, Paul Andr, Katrina Panovich, and Gregory G. Vargas. 2011. 4chan and /b/: An analysis of anonymity and ephemerality in a large online community. In *International Conference on Weblogs and Social Media, Barcelona, Catalonia, Spain, July*.
- Kim Young Bin, Kim Jun Gi, Kim Wook, Im Jae Ho, Kim Tae Hyeong, Kang Shin Jin, and Kim Chang Hun. 2016. Predicting fluctuations in cryptocurrency transactions based on user comments and replies:. *Plos One*, 11(8):e0161197.
- Johan Bollen, Huina Mao, and Xiao-Jun Zeng. 2010. Twitter mood predicts the stock market. *CoRR*, abs/1010.3003.
- Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer. 2002. Smote: synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16(1):321–357.
- CJ Hutto Eric Gilbert. 2014. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Eighth International Conference on Weblogs and Social Media (ICWSM-14)*. Available at (20/04/16) <http://comp. social. gatech. edu/papers/icwsml4. vader. hutto. pdf>.
- William L Hamilton, Kevin Clark, Jure Leskovec, and Dan Jurafsky. 2016. Inducing domain-specific sentiment lexicons from unlabeled corpora. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing. Conference on Empirical Methods in Natural Language Processing*, volume 2016, page 595. NIH Public Access.

- Sepp Hochreiter and Jrgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- Kee Sul Hong, Alan R. Dennis, and Lingyao Yuan. 2017. Trading on twitter: Using social media sentiment to predict stock returns. *Decision Sciences*, 48(3):454–488.
- Young Bin Kim, Hyeok Lee Sang, Shin Jin Kang, Myung Jin Choi, Jung Lee, and Hun Kim Chang. 2015. Virtual world currency value fluctuation prediction system based on user sentiment analysis. *Plos One*, 10(8):e0132944.
- Young Bin Kim, Jun Gi Kim, Wook Kim, Jae Ho Im, Tae Hyeong Kim, Shin Jin Kang, and Chang Hun Kim. 2016. Predicting fluctuations in cryptocurrency transactions based on user comments and replies. *PloS one*, 11(8):e0161197.
- Feng Mai, Qing Bai, Zhe Shan, Xin Wang, and Roger H. L. Chiang. 2015. The impacts of social media on bitcoin performance. *Social Science Electronic Publishing*.
- Bill Maurer, Taylor C Nelms, and Lana Swartz. 2013. when perhaps the real problem is money itself!: the practical materiality of bitcoin. *Social Semiotics*, 23(2):261–277.
- Ross C Phillips and Denise Gorse. 2017. Predicting cryptocurrency price bubbles using social media data and epidemic modelling. In *Computational Intelligence (SSCI), 2017 IEEE Symposium Series on*, pages 1–7. IEEE.
- Evita Stenqvist and Jacob Lnn. 2017. Predicting bitcoin price fluctuation with twitter sentiment analysis.
- Yang N Tang D, Wei F. 2014. Learning sentiment-specific word embedding for twitter sentiment classification. *Annual Meeting of the Association for Computational Linguistics*, 1(1):1555–1565.
- Sauk Hau Yong and Young Gul Kim. 2011. *Why would online gamers share their innovation-conducive knowledge in the online game user community? Integrating individual motivations and social capital perspectives*. Elsevier Science Publishers B. V.