



# Sentiment-Based Prediction of Alternative Cryptocurrency Price Fluctuations Using A Gradient Boosting Tree Model

Tianyu "Ray" Li, Anup Chamrajnagar, Xander Fong, Nick Rizik  
Math 106/87: Professor Feng Fu



## ABSTRACT

In this paper, we analyze Twitter signals as a medium for user sentiment to predict the price fluctuations of a small-cap alternative cryptocurrency called ZClassic. We extracted tweets on an hourly basis for a period of 3.5 weeks, classifying each tweet as positive, neutral, or negative. We then compiled these tweets into an hourly sentiment index, creating an unweighted and weighted index, with the latter giving larger weight to retweets. These two indices, alongside the raw summations of positive, negative, and neutral sentiment were juxtaposed to ~400 data points of hourly pricing data to train an Extreme Gradient Boosting Regression Tree Model. Price predictions produced from this model were compared to historical price data, with the resulting predictions having a 0.81 correlation with the testing data. Our model's predictive data yielded statistical significance at the  $p < 0.0001$  level. Our model is the first academic proof of concept that social media platforms such as Twitter can serve as powerful social signals for predicting price movements in the highly speculative alternative cryptocurrency, or "alt-coin", market.

## METHODS

We began by researching different alternative cryptocurrencies to decide which would be best suited within the confines of our analysis. We decided to choose ZClassic (ZCL), a private, decentralized, fast, open-source community driven virtual currency, as the primary target of our academic focus. The technological nature of the ZClassic cryptocurrency lends itself to a high level of predictability via tweet analysis, as it is set to "hard fork" into Bitcoin Private on February 28th, 2018. A hardfork is a major change to blockchain protocol which makes previously invalid blocks or transactions valid [19]. Typically, in the time leading up to the "hard fork," a coin's price reflects investor outlook on the success of the resulting fork.

To collect the tweets, we decided to base our program in RStudio, using open-sourced tweet package [20], which accesses Twitter's REST API. We were able to use the tweet package to retrieve, from each of the last seven days, searching from midnight backwards, tweets that had the terms "ZClassic," "ZCL," and "BTCP" while simultaneously eliminating repeat tweets. This collection process was repeated 3 times over the course of three and a half weeks to provide sufficient data for our analysis. In the end, we garnered a final data set of 130,000 unique tweets.

We then created an algorithm to classify each tweet as positive, negative, or neutral sentiment using natural language processing. The dictionary, primarily sourced from the Python package "TextBlob", that assigns impactful words and phrases a polarity value (e.g., "top" and "not great" have values of 0.5 and -0.4, respectively), which we view as sentiment

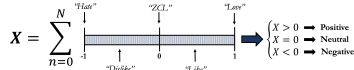


Figure 1: Summation of polarity values to Sentiment Conversion ( $N = \text{total words}$ )

For model selection, we employed 10-fold cross validation on 589 data points to choose an optimal model framework among linear regression, logistic regression, polynomial regression, exponential regression, tree model, and support vector machine regression. A tree model called the Extreme Gradient Boosting Regression (also known as XGBoost [21]), exhibited the smallest loss, or inaccuracy, and was thus chosen to train the model on our data. The XGBoost model, as well as other tree-based models, is particularly suited for applications on our data for the following reasons:

1. Tree models are not sensitive to the arithmetic range of the data and features. Thus, we do not need to normalize the data and possibly prevent loss due to normalization.
2. Tree models are by far the most scalable machine learning model due to their construction processes—simply adding more children nodes to the pre-existing tree nodes will update the tree and allow our strategy to continue to accurately predict price as our collection of price and tweet data increases into the future. It also makes the model adaptable for currencies with larger daily tweet volumes.
3. On the abstract level, the tree model is a rule-based learning method which, unlike a traditional regression learning method, has more potential to unveil insightful relationships between features.

For completeness, we sketch the key ideas behind XGBoost as follows. Let us define:

$$\hat{y}_i = \phi(x_i) = \sum_{k=1}^K f_k(x_i), \quad f_k \in \mathcal{F},$$

- $\phi(x)$  is our predicting function
- $\mathcal{F}$  is the prediction from our model for the  $i$ -th observation
- Each  $f_k$  representing a tree in our regression tree forest  $\mathcal{F}$ .

Our goal is to minimize the objective function  $L$ , defined below:

$$L(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k),$$

$$\Omega(f) = \gamma T + \frac{1}{2} \|w\|^2$$

Manual Sentiment Decision	Algorithm Sentiment Prediction		
	Positive	Neutral	Negative
Positive	79%	21%	0%
Neutral	34%	51%	15%
Negative	25%	0%	75%

Figure 2: Validation analysis of algorithm sentiment prediction by manual inspection.

**Algorithm 1:** Greedy algorithm for split finding [21] used in our price prediction model.  
Input:  $L$ , matrix of feature nodes  
Input:  $f$ , feature node  
Output:  $g$ , best split  
1.  $g \leftarrow \text{null}$   
2.  $g \leftarrow \text{find\_best\_split}(f)$   
3.  $g \leftarrow \text{find\_best\_split}(f)$   
4.  $g \leftarrow \text{find\_best\_split}(f)$   
5.  $g \leftarrow \text{find\_best\_split}(f)$   
6.  $g \leftarrow \text{find\_best\_split}(f)$   
7.  $g \leftarrow \text{find\_best\_split}(f)$   
8.  $g \leftarrow \text{find\_best\_split}(f)$   
9.  $g \leftarrow \text{find\_best\_split}(f)$   
10.  $g \leftarrow \text{find\_best\_split}(f)$   
11.  $g \leftarrow \text{find\_best\_split}(f)$   
12.  $g \leftarrow \text{find\_best\_split}(f)$   
13.  $g \leftarrow \text{find\_best\_split}(f)$   
14.  $g \leftarrow \text{find\_best\_split}(f)$   
15.  $g \leftarrow \text{find\_best\_split}(f)$   
16.  $g \leftarrow \text{find\_best\_split}(f)$   
17.  $g \leftarrow \text{find\_best\_split}(f)$   
18.  $g \leftarrow \text{find\_best\_split}(f)$   
19.  $g \leftarrow \text{find\_best\_split}(f)$   
20.  $g \leftarrow \text{find\_best\_split}(f)$   
21.  $g \leftarrow \text{find\_best\_split}(f)$   
22.  $g \leftarrow \text{find\_best\_split}(f)$   
23.  $g \leftarrow \text{find\_best\_split}(f)$   
24.  $g \leftarrow \text{find\_best\_split}(f)$   
25.  $g \leftarrow \text{find\_best\_split}(f)$   
26.  $g \leftarrow \text{find\_best\_split}(f)$   
27.  $g \leftarrow \text{find\_best\_split}(f)$   
28.  $g \leftarrow \text{find\_best\_split}(f)$   
29.  $g \leftarrow \text{find\_best\_split}(f)$   
30.  $g \leftarrow \text{find\_best\_split}(f)$   
31.  $g \leftarrow \text{find\_best\_split}(f)$   
32.  $g \leftarrow \text{find\_best\_split}(f)$   
33.  $g \leftarrow \text{find\_best\_split}(f)$   
34.  $g \leftarrow \text{find\_best\_split}(f)$   
35.  $g \leftarrow \text{find\_best\_split}(f)$   
36.  $g \leftarrow \text{find\_best\_split}(f)$   
37.  $g \leftarrow \text{find\_best\_split}(f)$   
38.  $g \leftarrow \text{find\_best\_split}(f)$   
39.  $g \leftarrow \text{find\_best\_split}(f)$   
40.  $g \leftarrow \text{find\_best\_split}(f)$   
41.  $g \leftarrow \text{find\_best\_split}(f)$   
42.  $g \leftarrow \text{find\_best\_split}(f)$   
43.  $g \leftarrow \text{find\_best\_split}(f)$   
44.  $g \leftarrow \text{find\_best\_split}(f)$   
45.  $g \leftarrow \text{find\_best\_split}(f)$   
46.  $g \leftarrow \text{find\_best\_split}(f)$   
47.  $g \leftarrow \text{find\_best\_split}(f)$   
48.  $g \leftarrow \text{find\_best\_split}(f)$   
49.  $g \leftarrow \text{find\_best\_split}(f)$   
50.  $g \leftarrow \text{find\_best\_split}(f)$   
51.  $g \leftarrow \text{find\_best\_split}(f)$   
52.  $g \leftarrow \text{find\_best\_split}(f)$   
53.  $g \leftarrow \text{find\_best\_split}(f)$   
54.  $g \leftarrow \text{find\_best\_split}(f)$   
55.  $g \leftarrow \text{find\_best\_split}(f)$   
56.  $g \leftarrow \text{find\_best\_split}(f)$   
57.  $g \leftarrow \text{find\_best\_split}(f)$   
58.  $g \leftarrow \text{find\_best\_split}(f)$   
59.  $g \leftarrow \text{find\_best\_split}(f)$   
60.  $g \leftarrow \text{find\_best\_split}(f)$   
61.  $g \leftarrow \text{find\_best\_split}(f)$   
62.  $g \leftarrow \text{find\_best\_split}(f)$   
63.  $g \leftarrow \text{find\_best\_split}(f)$   
64.  $g \leftarrow \text{find\_best\_split}(f)$   
65.  $g \leftarrow \text{find\_best\_split}(f)$   
66.  $g \leftarrow \text{find\_best\_split}(f)$   
67.  $g \leftarrow \text{find\_best\_split}(f)$   
68.  $g \leftarrow \text{find\_best\_split}(f)$   
69.  $g \leftarrow \text{find\_best\_split}(f)$   
70.  $g \leftarrow \text{find\_best\_split}(f)$   
71.  $g \leftarrow \text{find\_best\_split}(f)$   
72.  $g \leftarrow \text{find\_best\_split}(f)$   
73.  $g \leftarrow \text{find\_best\_split}(f)$   
74.  $g \leftarrow \text{find\_best\_split}(f)$   
75.  $g \leftarrow \text{find\_best\_split}(f)$   
76.  $g \leftarrow \text{find\_best\_split}(f)$   
77.  $g \leftarrow \text{find\_best\_split}(f)$   
78.  $g \leftarrow \text{find\_best\_split}(f)$   
79.  $g \leftarrow \text{find\_best\_split}(f)$   
80.  $g \leftarrow \text{find\_best\_split}(f)$   
81.  $g \leftarrow \text{find\_best\_split}(f)$   
82.  $g \leftarrow \text{find\_best\_split}(f)$   
83.  $g \leftarrow \text{find\_best\_split}(f)$   
84.  $g \leftarrow \text{find\_best\_split}(f)$   
85.  $g \leftarrow \text{find\_best\_split}(f)$   
86.  $g \leftarrow \text{find\_best\_split}(f)$   
87.  $g \leftarrow \text{find\_best\_split}(f)$   
88.  $g \leftarrow \text{find\_best\_split}(f)$   
89.  $g \leftarrow \text{find\_best\_split}(f)$   
90.  $g \leftarrow \text{find\_best\_split}(f)$   
91.  $g \leftarrow \text{find\_best\_split}(f)$   
92.  $g \leftarrow \text{find\_best\_split}(f)$   
93.  $g \leftarrow \text{find\_best\_split}(f)$   
94.  $g \leftarrow \text{find\_best\_split}(f)$   
95.  $g \leftarrow \text{find\_best\_split}(f)$   
96.  $g \leftarrow \text{find\_best\_split}(f)$   
97.  $g \leftarrow \text{find\_best\_split}(f)$   
98.  $g \leftarrow \text{find\_best\_split}(f)$   
99.  $g \leftarrow \text{find\_best\_split}(f)$   
100.  $g \leftarrow \text{find\_best\_split}(f)$

To minimize the above objective function, we employed a greedy Algorithm (1) to create our regression tree forest  $\mathcal{F}$  as originally implemented in Ref. [21]. One-third of the 589 data points is separated as the testing data, and the remainder is used as the training set as we built our Extreme Gradient Boosting Regression model. The model also tests different lead-lag on the range of [0, 1, 2, 3, 4, 5 hours] since we do not know how quickly the public would react to the market update or the social media sentiment. Based on the testing result, we decided that there is a 3-hour lag effect between social media information and price effects.

## RESULTS

Sentiment	Tweets
Positive	"Bitcoin Private Sbtcp is coming! better hop in \$zcl to grab some before the snapshot on the 28th! 15 days is a times for a crypto to be going up." "#ZCL going to give Grandma her retirement back!" "If you haven't got \$ZCL yet, now is the time to grab some."
Neutral	"Do you plan on selling your \$zcl before the fork or holding on to it for the \$btcp?" "don't miss out registering on Binance, before they close registration again." "RT @cryptobid: When everyone realizes you get FREE #BitcoinPrivate SBTCP if you buy \$ZCL #CryptoCurrency #FreeCoins ..."
Negative	"RT@Crypto.Bitlord: Bitcoin private is just another scam coin knock off like bitcoin cash trying to ride the satoshi's original vision." "2017 end - most thought 2018 would b the year crypto would explode. I said it b4 will say it again - gonna b a bad year." "What a joke this \$zcl game is. Good news, drops like a rock." "#cryptocurrency ..."

Figure 3: Examples of Tweets with positive, neutral and negative sentiment classifications in our dataset.

Having set the sentiment classification algorithms in place, we decided to train our model using six different features: Pure Positive Sentiment, Pure Negative Sentiment, Neutral Sentiment, An Weighted Sentiment Index, and Hourly Trading Volume. These six features proved to be varied enough to train the model effectively on a variety of different trading points and resulted in the best and most accurate overall correlation with the testing data (as shown in Figure 5). The detailed co-plots of the different features versus the price curve over the study period is shown in Figure 4.

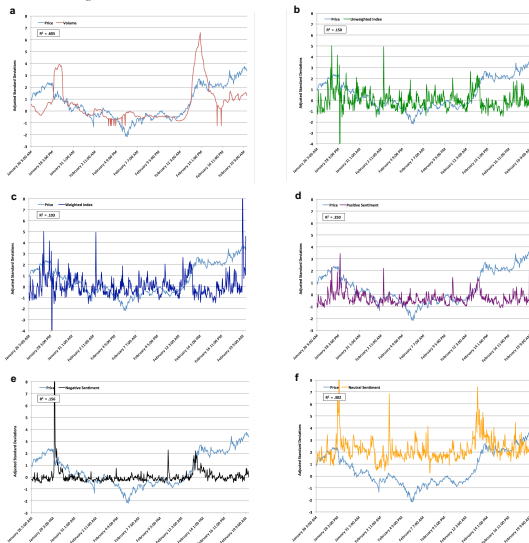


Figure 4: Shown are the price fluctuations versus our chosen six features, respectively, over the time period under consideration. (a) Price vs. Volume, (b) Price vs. Unweighted Index, (c) Price vs. Weighted Index, (d) Price vs. Pos. Sentiment, (e) Price vs. Neg. Sentiment, and (f) Price vs. Neutral Sentiment.

	Correlation Coefficients
Trading Volume	0.605
Neutral	0.302
Positive	0.250
Negative	0.156
Unweighted Index	0.150
Weighted Index	0.103

Figure 5: 5-factor correlation coefficients between the chosen feature and the price data, respectively.

## RESULTS CONTINUED

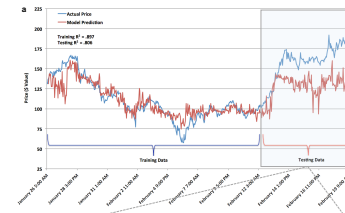


Figure 6: Comparison of model prediction and actual price data. Panel (a) plots the fitted price curve obtained from the training price data and the predicted price curve with respect to the testing data. Panel (b) details the model prediction price data as compared to the testing real price data.

	Positive	Negative	Neutral	Total
Average Hourly Tweets in Training Data	73.7	18.8	75.2	167.7
Percentage of Total	43.9%	11.2%	44.8%	-
Average Hourly Tweets in Testing Data	82.4	23.5	95.5	201.4
Percentage of Total	40.9%	11.7%	47.4%	-

Figure 7: Discrepancies of Twitter sentiments between testing and training data.

## CONCLUSION

In conclusion, our results suggest that by analyzing Twitter sentiment and trading volume, an Extreme Gradient Boosting Regression Tree Model serves as a viable means of predicting price fluctuations within the ZClassic cryptocurrency market. As such, given the complete lack of research within this academic sphere, our model serves as a proof of concept that social media platforms such as twitter can be used to indicate positive sentiment, and that this sentiment is an early signal to future price fluctuations in alternative cryptocurrencies. Of particular interest is seeing whether this approach produces similarly strong results when applied to other alternative cryptocurrencies such as ZCash and Bitcoin Private. However, this discovery sheds light to the possibility of arbitrage opportunities that utilize social media platform sentiment to predict future cryptocurrency prices.

Our pricing model could be further improved by factoring in other social media platforms of data, such as Google Search results, Facebook posts, and Reddit Posts. Moreover, the dictionary that we have used in our model could also be added by adding more sentiment-specific terms that indicate positive and negative sentiment such as "bull" and "bear" respectively. As seen from our manual vs. algorithm cross-analysis, the algorithm's greatest weakness is in classifying tweets that should otherwise be characterized as "negative" as "positive." After careful review it is evident that such inaccurate characterizations are due to the algorithm's inability to detect sarcasm, a pervasive language schema in popular social media platforms. As such, further research to enhance our algorithm to detect sarcasm would increase the reliability of the sentiment analysis, and thereby potentially improve the accuracy of our prediction to retail driven price changes.

## REFERENCES

[1] Goodhart, C. Cryptocurrency. <https://www.bankofengland.co.uk/quarterly-bulletin/2017/04/cryptocurrency>  
[2] Bitcoin.org. <https://www.bitcoin.org/>  
[3] Bitcoin.org. <https://www.bitcoin.org/>  
[4] Bitcoin.org. <https://www.bitcoin.org/>  
[5] Bitcoin.org. <https://www.bitcoin.org/>  
[6] Bitcoin.org. <https://www.bitcoin.org/>  
[7] Bitcoin.org. <https://www.bitcoin.org/>  
[8] Bitcoin.org. <https://www.bitcoin.org/>  
[9] Bitcoin.org. <https://www.bitcoin.org/>  
[10] Bitcoin.org. <https://www.bitcoin.org/>  
[11] Bitcoin.org. <https://www.bitcoin.org/>  
[12] Bitcoin.org. <https://www.bitcoin.org/>  
[13] Bitcoin.org. <https://www.bitcoin.org/>  
[14] Bitcoin.org. <https://www.bitcoin.org/>  
[15] Bitcoin.org. <https://www.bitcoin.org/>  
[16] Bitcoin.org. <https://www.bitcoin.org/>  
[17] Bitcoin.org. <https://www.bitcoin.org/>  
[18] Bitcoin.org. <https://www.bitcoin.org/>  
[19] Bitcoin.org. <https://www.bitcoin.org/>  
[20] Bitcoin.org. <https://www.bitcoin.org/>  
[21] Bitcoin.org. <https://www.bitcoin.org/>