

Identifying and Analyzing Cryptocurrency Manipulations in Social Media

Mehrnoosh Mirtaheri
USC Information Sciences Institute
mehrnoom@usc.edu

Sami Abu-El-Haija
USC Information Sciences Institute
haija@isi.edu

Fred Morstatter
USC Information Sciences Institute
fredmors@isi.edu

Greg Ver Steeg
USC Information Sciences Institute
gregv@isi.edu

Aram Galstyan
USC Information Sciences Institute
galstyan@isi.edu

ABSTRACT

Interest surrounding cryptocurrencies, digital or virtual currencies that are used as a medium for financial transactions, has grown tremendously in recent years. The anonymity surrounding these currencies makes investors particularly susceptible to fraud—such as “pump and dump” scams—where the goal is to artificially inflate the perceived worth of a currency, luring victims into investing before the scammers can sell their holdings. Because of the speed and relative anonymity offered by social platforms such as Twitter and Telegram, social media has become a preferred platform for scammers who wish to spread false hype about the cryptocurrency they are trying to pump. In this work we propose and evaluate a computational approach that can automatically identify pump and dump scams as they unfold by combining information across social media platforms. We also develop a multi-modal approach for predicting whether a particular pump attempt will succeed or not. Finally, we analyze the prevalence of bots in cryptocurrency related tweets, and observe a significant presence of bots during the pump attempts.

KEYWORDS

cryptocurrency, pump and dump, social media data mining, anomaly detection

1 INTRODUCTION

The inception of blockchain technology [14] gave birth to the popular cryptocurrency Bitcoin (symbol BTC). Since then, thousands of cryptocurrencies have emerged, and their hype has caused massive price swings on the trading markets. In December 2017, BTC quadrupled in market value in just over a month, then within a few days started a gradual decline until it reached half of its peak value. These price changes allowed some investors to realize huge profits, contributing to the allure of cryptocurrencies. Even though most investments are made in relatively established cryptocurrencies, including Bitcoin (BTC) and Ethereum (ETH), there are thousands of other smaller cryptocurrencies. These currencies are prime targets for manipulation by *scammers*, as evidenced by the proliferation of pump and dump schemes.

Pump and dump schemes are those in which a security price inflates due to deliberately deceptive activities. Those fraudulent schemes originated in the early days of the stock market and are now growing rapidly in the cryptocurrency market. The fact that

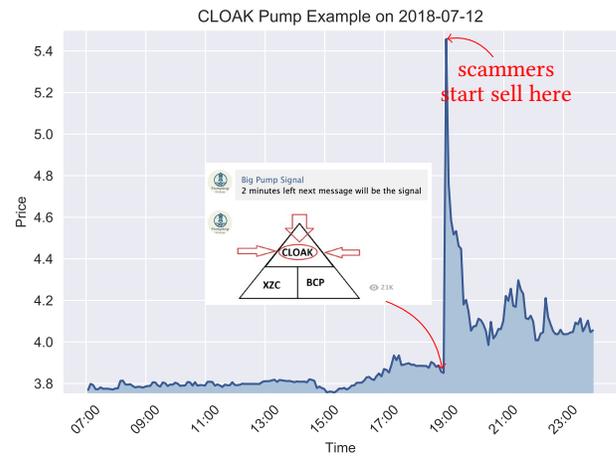


Figure 1: Pump announcement by admin of a Telegram Channel “Big Pump Signal,” overlaid on the market values of \$CLOAK. The announcement precedes the price swing.

the Commodity Futures Trading Commission (CTFC) and U.S. Securities and Exchange Commission (SEC) have issued several warnings [28] against cryptocurrency pump and dump schemes highlights the severity of the threat.

Although in the early days of cryptocurrencies, pump and dump schemes were taking place by marketing teams in ICOs [26] (Initial Coin Offerings), they are taking different forms nowadays. Pump and dump schemes have three major components: (1) a group of scammers; (2) a private or semi-private communication medium where scammers coordinate their illicit activities; and (3) a social media platform for conducting orchestrated campaigns to hype a given coin. In a typical scenario, scammers create groups on platforms such as Telegram or Reddit to coordinate group purchases of a particular cryptocurrency, while creating false hype around it by making public posts (i.e., *pump*) on social media platforms such as Twitter. Normal traders, who only see the rise in the price and are unaware of malicious activity, might buy the coin hoping to anticipate the next trend, thus boosting the price even further. Once a certain price target is met, the scammers start to sell (i.e., *dump*) their holdings, leading to a precipitous drop in the price. We illustrate the process with an example mentioned in a Wall Street

Journal article¹ that we also found in our collected data. In this example, the fraudsters coordinated their activities using a public group “Big Pump Signal” on the Telegram messaging app, which has more than 60,000 members. As shown in Figure 1, after the coin CLOAK posted on the group at 7:00 PM GMT, the price inflated immediately after the pump message and dropped shortly afterwards.

As cryptocurrency trading attracts more public attention, it becomes extremely important to be able to detect such fraudulent activities and inform potentially susceptible people before they become victims of these crimes. Toward this end, we study the extent to which pump and dump groups on online forums, such as Telegram, are accompanied and correlated with suspicious activity on Twitter and cryptocurrency price movements. In addition, we quantify the ability of machine learning models to predict pump and dump schemes from our data sources. In particular, we propose and address the following two predictive tasks:

- (1) Predict an unfolding pump operation/advertisement campaign happening on Telegram.
- (2) Predict whether a detected operation will succeed, i.e., will the target price mentioned in the Telegram message be reached by the market shortly after the announcement?

Although closely monitoring covert communication of fraudulent users on Telegram enables us to detect pump and dump schemes, in many realistic scenarios we might not have access to such communication. For instance, some scammers are using private channels with restricted access and membership fees, or communicate on different platforms altogether. Furthermore, most pump and dump events happen in a very short time after a coin is announced on Telegram. Thus, we examine the possibility of detecting pump and dump schemes *without* relying on availability of the Telegram data. Remarkably, we demonstrate that it is indeed possible to detect pump and dump events by leveraging only publicly available information, such as Twitter and historical market data.

Our main contributions are as follows:

- (1) We propose a multi-modal approach to monitor potentially malicious activities in cryptocurrency trading by combining data from three distinct sources: (i.) Real-time market data on cryptocurrency trading including both price and volume information; (ii.) Twitter data of *cashtag* mentions for cryptocurrencies; and (iii.) Telegram data that contains potential mentions of and instructions for pump and dump activities.
- (2) We identify pump and dump operations on Telegram messages by manually labeling a fraction of those messages as “*pump*” versus “*not-pump*”, and then building a classifier to label the remaining messages. Our approach is efficient for dealing with huge numbers of unlabeled messages mentioning cryptocurrencies and achieves a high precision (Section 3).
- (3) In Section 4, we explore the possibility of forecasting specific pump and dump activities, as quantified on the two classification tasks, based on different combinations of features

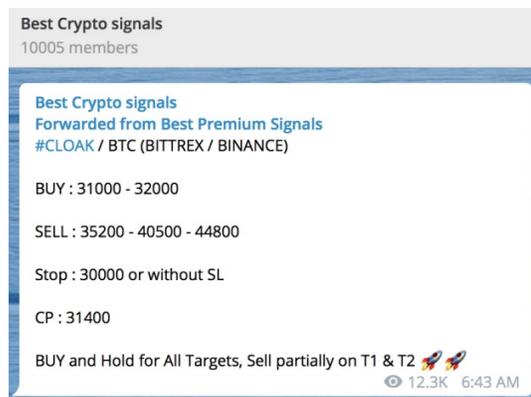


Figure 2: The anatomy of a common pump announcement, containing a “BUY” price and multiple target “SELL” prices. “CP” and “Stop” refer to the current coin price and stop loss respectively.

extracted from the above-listed data sources. Our results indicate that it is indeed possible to forecast such events with reasonable accuracy.

- (4) We study the efficacy of pump and dump operations on cryptocurrency price movements and Twitter activity in Section 3.2 and investigate the prevalence of Twitter bots in cryptocurrency-related tweets, especially during the alleged pump and dump attacks, and observe that the majority of highly active users are bots. (Section 5)
- (5) We release² a comprehensive dataset containing coins, timestamps, pump messages, and indications of whether or not the pumps were successful—together with features we extract from Twitter.

2 DATA DESCRIPTION

In this section we describe the datasets used in our study, explain the data collection process, and provide basic statistics of the data.

2.1 Telegram Data

In the context of the cryptocurrency market, scammers coordinate within groups to inflate the market value for a coin using social media platforms. In particular, the messaging platform Telegram is widely used for sharing cryptocurrency-related information, including pump announcements. The reason for Telegram’s popularity among scammers is that it provides anonymity for the users. For example, a Telegram channel consists of an anonymous admin and a set of members; however, the only person who can post to the channel is the admin who also is the only one who can see the list of members, while his/her identity is anonymous to the members.

We implemented a crawler using the Telegram API³ to collect data from Telegram channels. The crawling process starts with a set of a few initial channel IDs, and then we extend the list by extracting the other channels’ hyperlinks advertised in the seed channels and add those channels to our set. We continue to snowball out as new

¹<https://www.wsj.com/graphics/cryptocurrency-schemes-generate-big-coin/>

² Available at <https://github.com/Mehrnoom/Cryptocurrency-Pump-Dump>

³<https://telethon.readthedocs.io/en/latest/>

channel IDs appear in the Telegram channels. Due to the nature of how private channels are joined (the password to a private channel can be passed via its URL), it is entirely possible that we would crawl private channels if such URLs are posted in a public channel. We make no distinction between the two classes in our experiment.

Table 1 shows the statistics of the Telegram data. We extracted all the messages containing at least one coin from our coin list including all the cryptocurrencies provided by CoinMarketCap.com, which resulted in 195,576 messages. The Telegram channels in the table are categorized by their size (number of members) because channels with more members are more likely to contribute to a pump event.

Table 1: Telegram data statistics. Top: Histogram of number of channels grouped by their size (i.e., members in each channel). Middle: Summary statistics. Bottom: Number of messages containing the coins in our coin list.

Channel Size(S)	Count
$S \in [0, 100]$	51
$S \in (100, 1000]$	159
$S \in (1000, 10000]$	50
$S \in (10000, \infty)$	63
Number of channels	423
Average channel size	6,084
Median channel size	1,005
Total number of messages	195,576

2.2 Twitter Data

One way of promoting finance-related information on social media, especially Twitter, is to use *cashtags*, which are ticker symbols of stocks or cryptocurrencies prefixed by \$, e.g., \$BTC is the appropriate cashtag for Bitcoin. Using the Twitter streaming API,⁴ we implemented a system that tracks all the cryptocurrencies provided by CoinMarketCap.com, including 1,600 cashtags. We began the data collection process on March 15, 2018. However, cryptocurrencies had received a considerable amount of attention prior, toward the end of the 2017, due to growth in the Bitcoin price which paved the path for fraudulent activities. For better coverage of the potentially fraudulent events occurring before March 15, 2018, when we began our data collection, we also purchased tweets from September 1, 2017, to March 31, 2018.

The resulting dataset includes 30,760,831 tweets and 3,708,176 users in total from September 1, 2017, to August 31, 2018. Figure 3 also shows the distributions of the number of users per cashtag. The distribution is heavy-tailed, which suggests that many users are interested in only a few cashtags, while a small number of users tweet about many cashtags.

2.3 Cryptocurrency Market Data

This dataset consists of the time-series of market values for many cryptocurrencies. We developed a crawler to collect data from CoinMarketCap.com. Instead of using end-of-the-day historical

⁴<https://developer.twitter.com/en/docs/tweets/filter-realtime/overview.html>

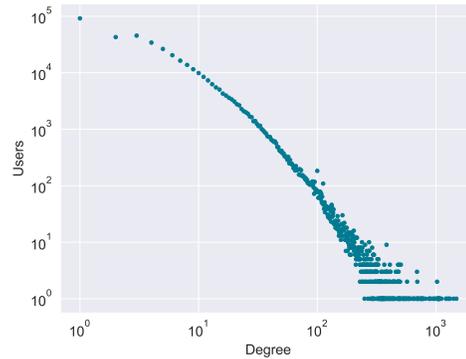


Figure 3: Users’ degree distribution, showing most users are interested in only a few cashtags, while some users discuss a large number of cashtags.

data, we chose to use data with a five-minute granularity because pump and dump schemes usually happen within a very short period. The dataset includes all the coins listed on the website at the time we started the collection process.

3 PUMP AND DUMP ACTIVITIES ON TELEGRAM

As mentioned above, Telegram is a popular choice for scammers to organize and coordinate pump and dump operations. To analyze such activities, let us define the following two notions:

- (1) **Pump attempt:** The act of targeting a coin on Telegram by posting a *pump* message mentioning the coin as a “pump attempt.” In Section 3.1 we describe our approach to detect pump messages.
- (2) **Successful pump attempt:** A pump attempt is *successful* if the actual price approaches the target price within a time window after the first pump message has been posted.

We next describe a simple method for detecting individual pump attempts and assessing whether those attempts were successful. The pump attempts (either successful or not) are used as ground truth for building and evaluating predictive models proposed in Section 4.

3.1 Pump Message Identification on Telegram

Most of the Telegram messages in our dataset are about cryptocurrency-related news, advice, and advertisements that are not relevant in the context of predicting pump and dump activities. Only a small fraction of those messages contain actual pump announcements. As shown in Table 1, the number of cryptocurrency-related messages is large, which would make it prohibitively expensive to manually label them as pump-relevant or not. Fortunately, however, most of the pump announcements follow specific patterns, or redundancies, which we are able to detect with machine learning techniques.

Text is the most common format used by Telegram channels for broadcasting pump events, although some channels embed the coin name in an image to prevent trading bot activities (Figure 1). A

Table 2: Telegram pump detection performance

Base Rate	Accuracy	Precision	Recall	F1
0.603	0.879	0.895	0.908	0.901

Table 3: Pump attempts statistics

Pump Attempts	62,850
Pump Messages	47,992
Channels	209
Coins	543

pump text message includes the name of the coin, the price to buy the coin, and one or more desirable target prices to achieve.

We leverage this specific common pattern to extract pump-related messages using a weakly supervised approach. Toward this goal, we labeled 1,557 messages in total as pump/not-pump, coming from 15 channels. To avoid bias toward a specific coin, we replace all the cryptocurrency symbols based on whether they are OOV (object out of vocabulary). We represent each post as a TF-IDF vector. For a given post, an entry in its TF-IDF vector corresponds to the frequency of a token appearing in the post (TF) divided by the number of posts in which the token appeared (IDF). In general, the size of the vector is equal to the size of the vocabulary of the entire corpus. We can use word n -grams (a sequence of n words) to construct the vocabulary.

We train a linear SVM with an SGD optimizer; this achieves an accuracy of 87% and a precision of 89%. The optimal parameters for linear classifier and TF-IDF tokenizer are obtained by cross-validation. The best result is achieved when we use both unigrams and bigrams, $\max_{DF} = 0.5$, $\min_{DF} = 0.01$, and L_2 penalty. The classifier scores are included in Table 2. Using the trained model, we then label the entire messages as pump/not-pump. From each message, we extract the coins mentioned in the message and the message timestamp. We call the (coin, timestamp) pair a pump attempt. We aggregate multiple timestamps into the earliest one, if they appear within a 3-hour window. We also remove the messages that mention many coins, as a pump message usually targets only a small number of coins. The statistics of the final set of messages and pump attempts are given in Table 3.

3.2 Signatures of Pump and Dump Activities on Market and Social Data

Having established the prevalence of pump and dump operations on Telegram, we now analyze the effectiveness of those operations, by juxtaposing pump messages with cryptocurrency market data and social signals collected from Twitter. Specifically, we focus on the overall effect of pump attempts on cryptocurrency prices and determine whether there is an indication of concurrent fraud activity on Twitter.

Let $S_i = (c_i, t_i)$ be a pump attempt, with t_i the time of the attempt, and c_i the target coin (in the next section, we explain in detail our approach for detecting pump attempts). For every pump attempt $S_i = (c_i, t_i)$ we extract two time series segments: The price

and tweet volume of coin c_i , denoted as $\tilde{P}_{t_i-w:t_i+w}^{(c_i)}$ and $\tilde{V}_{t_i-w:t_i+w}^{(c_i)}$ respectively, where w is a time window equal to three days. Each segment is normalized between 0 and 1 by a minmax normalization, transforming each point x to $\tilde{x} = (x - \min)/(\max - \min)$. We calculate Q^{price} and $Q^{twitter}$ as follows:

$$Q^{price} = \frac{1}{n} \sum_i^n P_t^{(c_i)} \quad \text{and} \quad Q^{twitter} = \frac{1}{n} \sum_i^n V_t^{(c_i)},$$

Note that $i \in [1, n]$ means that we have considered all the pump attempts. For each coin, we also select a set of timestamps uniformly at random, with the size equal to the number of pump attempts targeting the coin. In the same manner, we make two aggregated timeseries for the random timestamps.

Figure 4a depicts an average normalized price of each coin centered around the pump timestamps (Q^{price}) and random timestamps. This figure shows a pattern of spikes occurring within one hour of a pump message, followed by a general downward trend. A significant increase in tweet volume is also observable in Figure 4b, which shows the average tweet volume of each coin around the pump timestamps ($Q^{twitter}$) and random timestamp. The seasonality pattern present in the Telegram timestamps, and not in the random timestamps, suggests that most of the pump attempts happen around a specific time of day.

Successful Pump Attempts. Although Figure 4a indicates a price spike for the aggregated data, it is reasonable to expect that not all the pump messages are actually followed by a spike in the coin price. Furthermore, even if there is a spike, it might fall short of meeting the “target” price, thus resulting in *failed* pump. We used a simple rule-based approach to extract the “buy” and “target” prices from all the messages, and augmented this information with the coin price data to decide whether a given pump attempt was successful or not. Figure 5 shows the percentage of successful pump messages for different thresholds and time windows. Fewer than 5% of the pump messages meet the most strict conditions, meaning that the coin price reaches a higher price than the extracted target price, and within an hour of the pump message.

4 PUMP ATTEMPTS PREDICTION

In this section we study the feasibility of predicting pump and dump events from the social media and market data only, without relying on the availability of Telegram messages⁵. Specifically, we focus on the following two classification tasks:

- Task I: Detect whether there is an unfolding pump operation/advertisement campaign happening on Telegram, by considering social signals *only* from Twitter, and historical market data.
- Task II: Given a pump message on Telegram for a specific coin, predict whether the operation will succeed (i.e., will the target price, set by the scammers, be met within 6 hours of the message is posted on Telegram).

Both of these tasks are cast as binary classification problems. The feature vector for each record is extracted at a specific timestamp:

⁵We use Telegram messages only as ground truth for evaluating our predictive models, but not as input to those models.

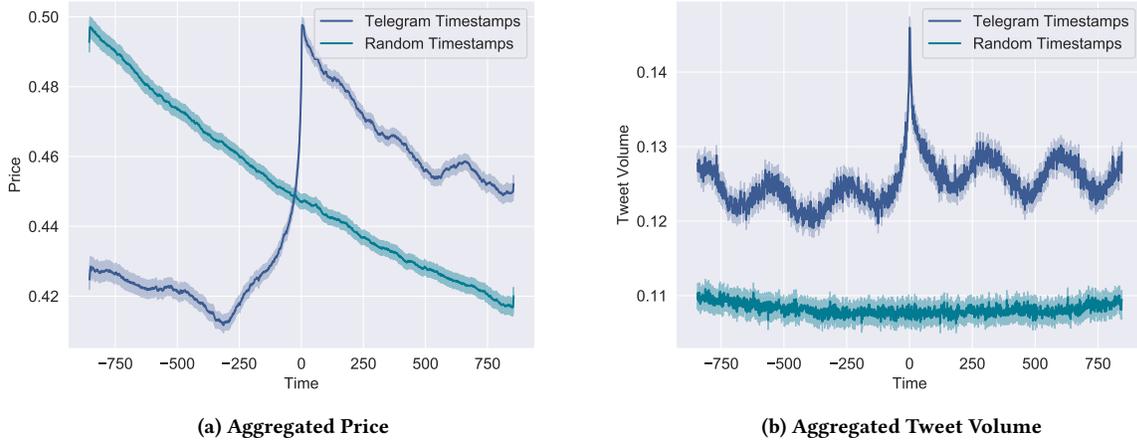


Figure 4: Price (a) and Twitter (b) time series segments of cryptocurrencies separated by whether they are mentioned in a Telegram pump and dump channel. Segments start 3 hours before and end 3 hours after the mention in a Telegram channel. Timestamps are selected uniformly at random.

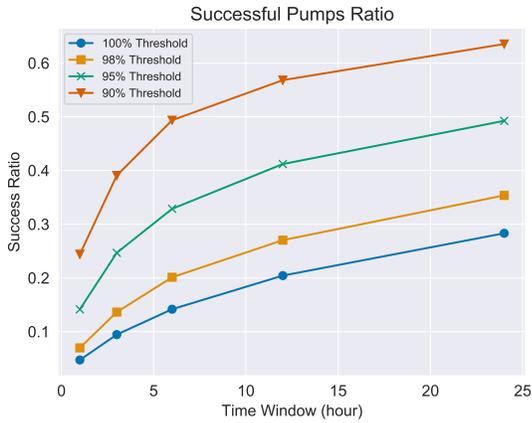


Figure 5: Ratio of successful pumps for different thresholds and time windows. “100% threshold” indicates the coin price beats the extracted target price. “x% threshold” is when the coin price either beats the price or reaches the x% of the extracted price.

specifically, it contains data from 6 hours prior to the timestamp. The features of the timestamps are explained next.

4.1 Features

Table 4 explains all the features extracted from our data sources that are used for the two prediction tasks. Graph features include: (i.) PageRank score; (ii.) CorEx user embeddings; and (iii.) user connected components. These are extracted from the coin-coin network, the pump-user network, and the user-user network, respectively. These networks are undirected and temporal, where the presence (or absence) of an edge depends on the time range in

which the networks are constructed. Each network is explained in its corresponding feature section.

Economic Features. Include (i) Coin Market Cap, (ii) Volume, and (iii) Price, in $h \in \{1, 2, \dots, w\}$ hours before timestamp t . For each of these three, we also include percentage change. For example, price percentage change at hour h is $\frac{\text{price}_{h+1} - \text{price}_h}{\text{price}_h}$.

Target Price Features. This feature will only be used in the second classification task (whether the pump succeeds, introduced in Section 4). As explained in Section 3, the definition of success depends on the target price mentioned in a pump message. For a given pump message with a mentioned target price x , we include $\frac{x - \text{price}_h}{x}$ for $h \in \{1, 2, \dots, w\}$ hours before the pump message timestamp.

Twitter Statistics. For a given timestamp t , the features extracted from time period $[t, t - w]$ include (i.) the number of tweets mentioned the cashtag; (ii.) the number of unique users who mentioned the cashtag in a tweet; and (iii.) the average sentiment of all the tweets mentioning the cashtag calculated using [7].

PageRank Score. We calculate the PageRank score [16] of a coin in the **coin-coin** graph created at time t , in which the nodes correspond to the cashtags, connected by an edge if they are mentioned by the same user, within the period $[t - w, t]$. The edge weights correspond to the number of times two coins are co-mentioned by a user.

CorEx User Embeddings. Consider a bipartite graph containing pump attempts and users as nodes. The edge weight between a pump attempt $S_i = (c_i, t_i)$ and user u_j is equal to the number of times u_j mentions c_i in a tweet within the period $[t - w, t]$. Here we chose $w = 6$ hours. We call this graph **pump-user network** and denote $\mathbf{B} \in \mathbb{R}^{|S| \times |U|}$ as the affiliation matrix of this bipartite graph, where S and U are the set of pump attempts and users respectively. B_{ij} is equal to the weight of the edge connecting S_i

Table 4: Description of the features used in the pump predictions and user clustering.

Feature type	Description
Twitter Features	<ul style="list-style-type: none"> • Number of tweets mentioned the cashtag in the period $[t, t - w]$ • Number of unique users mentioned the cashtag in the period $[t, t - w]$ • Average sentiment of all the tweets mentioning the cashtag in the period $[t, t - w]$ calculated using [7] • PageRank score [16] of a coin in the Coin-Coin graph created at time t • Twitter User Connected Components • CorEx user embedding
Economic Features	<ul style="list-style-type: none"> • Coin market cap and market cap percentage change at h hour before the pump where $h \in \{1, 2, \dots, 12\}$ • Coin volume and volume percentage change at h hour before the pump where $h \in \{1, 2, \dots, 12\}$ • Coin price and price percentage change in BTC unit at h hour before the pump where $h \in \{1, 2, \dots, 12\}$ • Target price percentage difference with coin price at h hour before the pump where $h \in \{1, 2, \dots, 12\}$

and u_j . In fact, matrix \mathbf{B} represents each pump attempt as a $|U|$ -dimensional vector of user activities, meaning that each user is considered as a variable. We want to cluster users if their activity are correlated. Total Correlation Explanation (CorEx) discovers a latent representation of complex data based on optimizing an information-theoretic approach. More specifically, given a set of n -dimensional vectors $\mathbf{X} \in \mathbb{R}^{m \times n}$, CorEx aims to find latent variables $\mathbf{Y} = \mathbf{W}\mathbf{X}$ that best describe the multivariate dependencies of \mathbf{X} by minimizing $TC(\mathbf{X}|\mathbf{Y}) + TC(\mathbf{Y})$. Here, TC is “total correlation” or multivariate mutual information [31].

We apply linear CorEx on \mathbf{B} , and find the best number of latent factors k by plotting the sum of the total correlation for each latent variable against k . This value starts decreasing significantly when $k = 24$. The weight matrix $\mathbf{W} \in \mathbb{R}^{m \times k}$ obtained from applying linear CorEx is used as the embedding for the users. Later in Section 5 we further explore pump-user network and CorEx clusters for analysing bot activities.

User Connected Components. We would like to characterize groups of users that participate in coordinated spreading of “pump” messages on Twitter. Given a coin c and all the pump instances targeting c , $S_i = (c, t_i)$, a bipartite graph is built from the pump instances and users. u_j is connected to S_i if u_j tweeted about c within the time period $[t_i - 6, t_i]$, measured in hours. From this bipartite graph, we then build a **user-user** network, in which users i and j are connected with a weighted edge that corresponds to the time the users co-tweet about c . After this process, the graph is very dense with usually $> 80\%$ of nodes belonging to one connected component. We sparsify the graph by keeping the top- k edges per user. Then we calculate the connected components from the graph, dropping connected components consisting of less than 25 users. Each user is now represented by the ID of its connected component. Users that do not correspond to the connected components are ignored in this feature representation. Manual inspection of the connected components showed that they are indeed meaningful, with the largest connected component usually corresponding to users who are likely involved in pump and dump operations. Given

a coin and a timestamp, we create a feature vector containing the counts of users tweeting about the coin, grouped by the connected components they belong to. Intuitively, these features represent how active a group (connected component) has been within a period from $t - w$ to t . We tried $k \in \{1, 2, 3\}$, and the final result does not change much. However, higher values produced just a handful of sizable connected components.

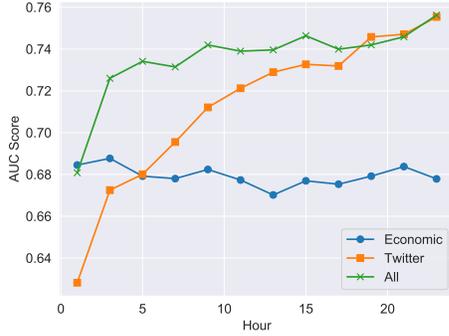
4.2 Classification Tasks

In this section we explain the experimental setting and prediction tasks design.

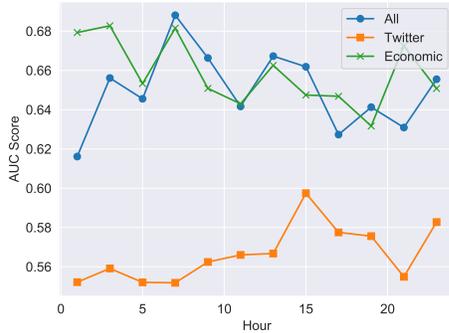
Task I: Predicting Pump Attempts. We propose a binary classification task for predicting Telegram pump messages that will happen in the future, from Twitter activity. In our setup, we use the timestamps of Telegram messages labeled as pump attempts by our classifier as positives. We use an equal number of random timestamps as negatives.

Task II: Will the Pump Succeed? The positives of this task are a subset of the positives in Task I—namely, the pumps that have succeeded. In other words, the target price mentioned in the pump message was successfully met by the market. The negatives of this task include all the negatives of Task I and some of the positives of Task I—specifically, the pump messages that were not successful. In other words, the target price was not reached within 6 hours of the pump message.

For both tasks, we train a binary Random Forest classifier. It might be possible to improve the classification accuracy by using sophisticated approaches, i.e., neural nets, but since our focus is to demonstrate the possibility of classification, we focus on traditional methods. We split the dataset into train/test, such that each sample (c, t) corresponding to coin c at time t is in the training set if $t < d$ and in the test set otherwise. For each coin, d is picked in way to reach a 75%/25% split. Let $Train = \{S_1, S_2, \dots, S_d\}$ and $Test = \{S_{d+1}, S_{k+2}, \dots, S_n\}$, sorted by their timestamps. For each coin, we train a separate classifier in the following way:



(a) Classification Task I



(b) Classification Task II

Figure 6: AUC-ROC score of the model, varying the time period of consideration from 1 to 24 hours.

Algorithm 1

```

Input:
1: Coin  $c$ 
2:  $Train = \{S_i | S_i = (c, t_i)\}_{i=1}^d$ , sorted by  $t_i$ 
3:  $Test = \{S_i | S_i = (c, t_i)\}_{i=d+1}^n$ , sorted by  $t_i$ 
for  $S_i$  in  $Test$  do
     $model = RandomForestClassifier.train(Train)$ 
     $P_i = model.inference(S_i)$ 
    Add  $S_i$  to  $Train$ 
    Remove  $S_i$  from  $Test$ 
    Add  $P_i$  to  $Prob$ 
return  $Prob$ 
    
```

We evaluate the approach using the Area Under the Receiver Operating Characteristic Curve (ROC-AUC), which gives a baseline of 0.5 for random guesses; the higher the metric, the better. The features are chosen as explained in Section 4.1. Figure 6 shows AUC-ROC score by time period in the past. Based on this plot, for the first task, we choose $w = 15$ for economic and Twitter features. For the second task, we choose $w = 15$ for Twitter features and $w = 7$ for economic features.

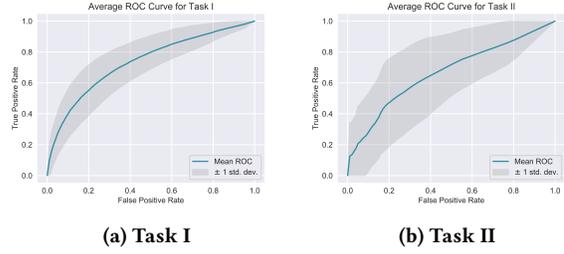


Figure 7: Average ROC plot for the two classification Tasks that we consider. All features are used for the classification.

4.3 Results and Discussion

Table 5a summarizes Task I results. The columns describe the prediction accuracy measured in AUC score when using Twitter-only, economic-only, and combined features. When using a combination of features, the AUC score averaged over all the coins is 0.74, which is significantly better than random baseline. We can see that social media features are more effective for this task, although adding economic features provides a slight increase. Finally, when we take the average over only the top 20 highest volume coins, the prediction accuracy increases slightly. In the table we also show the AUC numbers for the five coins with the highest AUC score when using all the features.

The results for Task II are shown in Table 5b. Overall, we see that the average accuracy is lower for this task compared to Task I, but is still considerably better compared to a random baseline. We also note that compared to Task I, the variance in the AUC scores is considerably higher. Interestingly, we observe that in average, economic features are much more useful for this task (average AUC of 0.7) than Twitter features (average AUC of 0.59). Furthermore, adding Twitter features to the economic features actually deteriorates performance. The average ROC curve when using all features is plotted in Figure 7 for both classification tasks.

We suggest that the Twitter features are more predictive for the first task, due to the social nature of this process. To be more specific, when scammers target a coin (pump attempt), they try to promote it in different social media platforms, resulting in correlated activity on Telegram and Twitter. However, as we observed in Section 3.2, only a small ratio of pump attempts “succeed” based on our definition. Therefore, although a pump operation might generate extra traffic in Twitter, its effect on the coin price depends on many other factors such as market characteristics. This might be the reason that for the second task, using only market/economic features gives us better performance. One possible explanation for performance drop when we add Twitter features in the second task is that the number of positive labels for this task is very low, and adding Twitter features increases the sparsity. Although we removed the coins with less than five positive in their training set, the number of positives is 14 in average for the remaining coins (less than 15% ratio in average, and for some coins as low as 6%). The average dimension of the Twitter features is 52 (it could vary across the coins, because of the number of the connected components) and more than a hundred for some coins.

Table 5: ROC-AUC test performance on the two binary tasks, averaging 10 different train:test partitions. We show evaluations for three feature-sets. From right to left: The right-most is “Both,” utilizing both financial and Twitter. The “Economic” and “Twitter” utilize economic and Twitter features separately. The average AUC score is reported for (i) all the coins (ii) 20 coins with the highest dollar volume. The top five most-predictable coins are also shown for each task.

(a) Classification Task I: Predicting Pump Attempts

	Twitter	Economic	Both
Average AUC (all coins)	0.73 ± 0.07	0.67 ± 0.1	0.74 ± 0.08
Average AUC (top 20 vol.)	0.75 ± 0.06	0.68 ± 0.1	0.75 ± 0.07
\$ADA	0.85	0.84	0.88
\$NCASH	0.84	0.81	0.88
\$DGB	0.79	0.76	0.86
\$RCN	0.82	0.61	0.85
\$TRX	0.83	0.81	0.84

(b) Classification Task II: Will the Pump Succeed?

	Twitter	Economic	Both
Average AUC (all coins)	0.59 ± 0.16	0.70 ± 0.19	0.66 ± 0.17
Average AUC (top 20 vol.)	0.62 ± 0.18	0.76 ± 0.17	0.71 ± 0.16
\$NMR	0.71	0.36	0.94
\$XEM	0.63	0.66	0.90
\$XRP	0.77	0.84	0.90
\$QTUM	0.83	0.80	0.89
\$ARK	0.47	0.77	0.88

We investigated the reasons for obtaining high accuracy scores for some coins but not the others. In particular, we analyzed potential relations between the prediction accuracy and financial indicators of a coin such as market cap, volume, and so on. As shown in Table 5b, the accuracy score is typically higher for coins with higher dollar volume. Our preliminary analysis of other features produced rather ambiguous results, as we did not find meaningful correlations between accuracy and coin features. This may require more thorough investigation in the future.

5 PREVALENCE OF TWITTER BOTS

In this section, we study the presence of bot activity around pump attempts by exploring the pump-user network that we explained in Section 4.1. This bipartite network has 36 connected components, but the largest connected component contains roughly 99% of the users and 99% of pump attempts. Below we discuss the involvement of Twitter bots in those attempts.

First, from our tweet dataset we extract the tweets containing a Telegram invitation link (e.g., <http://t.me/Monsterpumper>). We label the users associated with these tweets as **telegram active users**.

Next, we use two approaches for classifying a user as a *bot*.

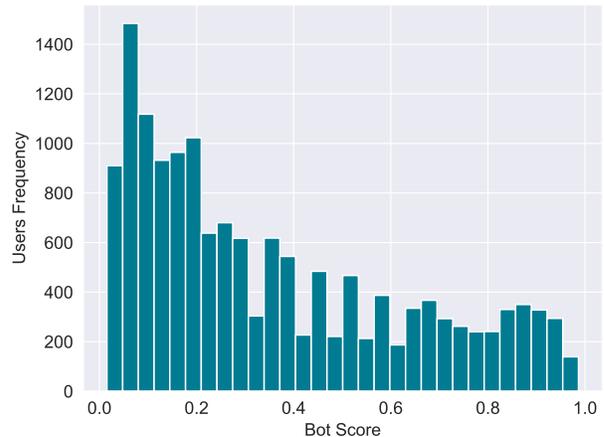


Figure 8: Score distribution obtained using Botometer API. Shows the probability of a user being bot.

- **Twitter Suspended List.** Using the Twitter API, we collected the most recent account status of the users in our dataset and checked whether they are still **active** or **suspended** by Twitter.
- **Botometer Score [29].** Twitter works based on algorithms with high precision but low recall, since they do not want to mistakenly suspend users that are not bots. So we employed the Botometer API ⁶ to detect other potential bots. Given a user id, Botometer returns a probability of that user being a bot. Botometer classifies users using six types of features: friend, network, content, sentiment, temporal, and user. Figure 8 shows the distribution of the classifier score using all six features.

A user was classified as a bot if either the user was suspended or its Botometer score is above 0.55 ⁷.

Table 6 shows an increasing ratio by degree. For a given user, the degree is the sum of the weights of its adjacent edges in the pump-user network and is the number of times the user participated in the pump operations. This suggests that a larger fraction of highly active users are bots, and the ratio increases with the activity level. For example, 84% of the users that participated more than 10K times in the pump activities are either suspended or are bots, according to their Botometer score.

User Clustering. Now we look at different user clusters. Using the weight matrix \mathbf{W} obtained from applying CorEx on \mathbf{B} (explained in more detail in Section 4.1) we cluster the users of the pump-user network by assigning user u_i to $\text{argmax} \mathbf{W}_i$.

Figure 9 shows the ratio of the bots and Telegram-active users in each cluster. The blue bar shows the ratio of the users that are bots **and** Telegram-active. Note that the bot ratio is the number of users that are either suspended or have a botometer score > 0.55 . The first and second cluster have 734 and 1,011 members respectively where

⁶<https://botometer.iuni.iu.edu/#!/api>

⁷We use 0.55 instead of 0.5 suggested by [29] to avoid edge cases

Table 6: Ratio of scammer users based on the number of times a user contributes to the pump attempts. The ratio of Telegram active users and suspended users is higher among highly active users.

Degree	Total # Users	Suspended	Telegram Active	Botometer Score > 0.55
50	20881	0.19	0.27	0.20
100	10969	0.23	0.34	0.20
500	2577	0.40	0.37	0.24
1000	1435	0.46	0.36	0.24
5000	340	0.56	0.30	0.25
10000	179	0.42	0.29	0.36

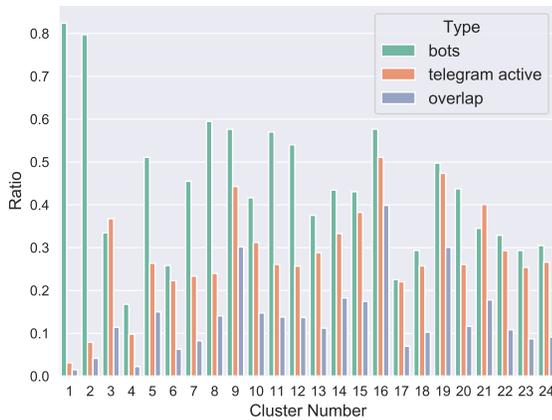


Figure 9: Ratio of bots and Telegram-active users in each cluster

more than 80% of them are bots. Cluster 16 is also interesting in a sense that it has 600 members where 50% of them are Telegram-active, and around 60% are bots. Cluster 17 and 18 have 3K members each with a low ratio of bot members.

6 RELATED WORK

We split the discussion of the related work into three complementary threads. First, we focus on fraud in social media and some of the efforts that have been taken to address it. Next, we discuss other work that analyzes the relation between financial market and social media data. We then describe the effort on studying cryptocurrency activity, with a special focus on work that includes social media in its analysis.

Social media has a long history of fraudulent activity, and some types of fraud appear in our work. First, when scammers attempt to pump a coin by making it look more popular than it actually is, they are engaging in a specific type of misinformation. Misinformation is a major problem on social media [32], and several recent efforts have tried to detect it [11, 19, 24]. Another burgeoning line of work is bot detection [25]. There is a known connection between bots

and misinformation, wherein bots are actively employed to spread misinformation in social networks [5, 12]. We leverage previous literature in these areas in our approach. The bot labeling approach that we use in the bot assessment portion is based on previous work [8]. Additionally, we study the dynamics of users’ reactions to the pump and dump campaigns. This is similar to previous work on social media where similar inputs are used to identify susceptible users [15, 22].

The literature on predicting various financial market properties exploiting signals from different social media data is quite extensive. [1, 2, 30] use sentiment features to predict stock price movements, while [21, 23, 27] employ network features extracted from various social media interactions. However, only a few papers study stock market manipulations on social media. [20] studies pump and dump in OTC (Over The Counter or Penny stocks) and shows that an abnormally high number of messages on Twitter is associated with price increase followed by a price reversal. [3] shows the presence of bot and spam activity on Twitter stock microblogs and compares Twitter activity with financial data. Authors in [4] uncover and introduce *Cashtag Piggybacking*, a type of malicious activity on Twitter in which scammers try to highlight low-value stocks by co-mentioning them in the tweets with high-value stocks.

Several papers study cryptocurrencies, and analysts build models to predict their price movement. An example is [9]. In this paper the authors use cryptocurrency forums to predict the price and volume of cryptocurrencies. In another effort [17] the authors build a model to predict price fluctuations of cryptocurrencies. Specifically, they use epidemic models on social media activity to predict price bubbles of cryptocurrencies. [10] further tests how users discuss cryptocurrencies and how that discussion impacts price. They found that specific topics are likely to be tied to price movements. [18] extends this analysis by using wavelets to predict price movements based on social media data. [6] looks into the dynamics underlying social media and how they correlate with cryptocurrency price. They find that opinion polarization has a significant effect on price, and use this to build a model that predicts the price of the cryptocurrency. In an effort to understand the dynamics of cryptocurrency discussions, [13] performed topic modeling on a popular cryptocurrency discussion forum. They identified several common threads of discussion, such as bitcoin theft. Moreover, they showed that different mining technologies have different patterns of adoption on the forums. Our work stands apart from these methods by moving away from predicting price and volume movements, and instead identifying patterns of malicious behavior.

The work that is most related to ours is [33], where the main goal is predicting which coin will be pumped based on social signals from Telegram. The authors focus on “pre-pump” messages that announce an upcoming pump operation, but do not mention a coin. They developed a model to predict the likelihood of each coin being the target of the subsequent pump operation following the “pre-pump” message. Our work is complementary in that we consider a richer set of prediction problems, we use social signals from Twitter, and we provide a user-centric analysis of such pump attacks.

7 CONCLUSION

In this paper we present a novel computational approach for identifying and characterizing cryptocurrency pump and dump operations that are carried out in social media. Specifically, given financial and Twitter data pertaining to a particular coin, our method is able to detect, with reasonable accuracy, whether there is an unfolding attack on that coin on Telegram, and whether or not the resulting pump operation will succeed in terms of meeting the anticipated price targets. We also analyze activities of users involved in pump operations, and observe a prevalence of Twitter bots in cryptocurrency-related tweets in close proximity to the attack.

In future work, we plan to augment our datasets with other sources (e.g., Reddit posts) to help with the prediction tasks considered here. Also, while our analysis of bot activity relied on suspended accounts, it will be interesting to develop a bot detection tailored to the cryptocurrency domain. Finally, as a practical outcome of the work presented here, we envision building a cryptocurrency monitoring system that will detect impending pump attacks in real-time and warn susceptible users.

8 ACKNOWLEDGEMENTS

We thank Saurabh Birari for developing a crawler for collecting historical market data. We would also like to thank Emilio Ferrara and Pegah Jandaghi for their helpful comments.

REFERENCES

- [1] Francesco Bellini and Nicola Fiore. 2017. Exploring Sentiment on Financial Market Through Social Media Stream Analysis. Springer, Cham, 115–129. https://doi.org/10.1007/978-3-319-49538-5_18
- [2] Wenhao Chen, Yi Cai, Kinkeung Lai, and Haoran Xie. 2016. A topic-based sentiment analysis model to predict stock market price movement using Weibo mood. *Web Intelligence* 14, 4 (11 2016), 287–300. <https://doi.org/10.3233/WEB-160345>
- [3] Jonathan Clarke, Hailiang Chen, Ding Du, and Yu Jeffrey Hu. 2018. Fake News, Investor Attention, and Market Reaction. (2018).
- [4] Stefano Cresci, Fabrizio Lillo, Daniele Regoli, Serena Tardelli, and Maurizio Tesconi. 2019. Cashtag piggybacking: Uncovering spam and bot activity in stock microblogs on Twitter. *ACM Transactions on the Web (TWEB)* 13, 2 (2019), 11.
- [5] Michelle Forelle, Phil Howard, Andrés Monroy-Hernández, and Saiph Savage. 2015. Political bots and the manipulation of public opinion in Venezuela. *arXiv preprint arXiv:1507.07109* (2015).
- [6] David Garcia and Frank Schweitzer. 2015. Social signals and algorithmic trading of Bitcoin. *Royal Society open science* 2, 9 (2015), 150288.
- [7] 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/icwsm14.vader.hutto.pdf>.
- [8] Xia Hu, Jiliang Tang, Huiji Gao, and Huan Liu. 2014. Social spammer detection with sentiment information. In *Data Mining (ICDM), 2014 IEEE International Conference on*. IEEE, 180–189.
- [9] 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 (2016), e0161197.
- [10] Young Bin Kim, Jurim Lee, Nuri Park, Jaegul Choo, Jong-Hyun Kim, and Chang Hun Kim. 2017. When Bitcoin encounters information in an online forum: Using text mining to analyse user opinions and predict value fluctuation. *PLoS one* 12, 5 (2017), e0177630.
- [11] Sejeong Kwon, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. 2013. Prominent features of rumor propagation in online social media. In *2013 IEEE 13th International Conference on Data Mining*. IEEE, 1103–1108.
- [12] David M J Lazer, Matthew A Baum, Yochai Benkler, Adam J Berinsky, Kelly M Greenhill, Filippo Menczer, Miriam J Metzger, Brendan Nyhan, Gordon Pennycook, David Rothschild, et al. 2018. The science of fake news. *Science* 359, 6380 (2018), 1094–1096.
- [13] Marco Linton, Ernie Gin Swee Teo, Elisabeth Bommers, CY Chen, and Wolfgang Karl Härdle. 2017. Dynamic Topic Modelling for Cryptocurrency Community Forums. In *Applied Quantitative Finance*. Springer, 355–372.
- [14] Satoshi Nakamoto. 2008. Bitcoin: A Peer-to-Peer Electronic Cash System. In *bitcoin.org*.
- [15] Pinar Ozturk, Huaye Li, and Yasuaki Sakamoto. 2015. Combating rumor spread on social media: The effectiveness of refutation and warning. In *System Sciences (HICSS), 2015 48th Hawaii International Conference on*. IEEE, 2406–2414.
- [16] Larry Page. 1997. PageRank: Bringing Order to the Web. In *Stanford Digital Library*.
- [17] 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*. IEEE, 1–7.
- [18] Ross C Phillips and Denise Gorse. 2018. Cryptocurrency price drivers: Wavelet coherence analysis revisited. *PLoS one* 13, 4 (2018), e0195200.
- [19] Jacob Ratkiewicz, Michael Conover, Mark R Meiss, Bruno Gonçalves, Alessandro Flammini, and Filippo Menczer. 2011. Detecting and tracking political abuse in social media. *ICWSM* 11 (2011), 297–304.
- [20] Thomas Renault. 2017. Market manipulation and suspicious stock recommendations on social media. (2017).
- [21] Eduardo J. Ruiz, Vagelis Hristidis, Carlos Castillo, Aristides Gionis, and Alejandro Jaimes. 2012. Correlating financial time series with micro-blogging activity. In *Proceedings of the fifth ACM international conference on Web search and data mining - WSDM '12*. ACM Press, New York, New York, USA, 513. <https://doi.org/10.1145/2124295.2124358>
- [22] Justin Sampson, Fred Morstatter, Liang Wu, and Huan Liu. 2016. Leveraging the implicit structure within social media for emergent rumor detection. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. ACM, 2377–2382.
- [23] Sunil Saumya, Jyoti Prakash Singh, and Prabhath Kumar. 2016. Predicting Stock Movements using Social Network. Springer, Cham, 567–572. https://doi.org/10.1007/978-3-319-45234-0_150
- [24] Kate Starbird, Jim Maddock, Mania Orand, Peg Achterman, and Robert M Mason. 2014. Rumors, false flags, and digital vigilantes: Misinformation on twitter after the 2013 boston marathon bombing. *ICConference 2014 Proceedings* (2014).
- [25] V. S. Subrahmanian, A. Azaria, S. Durst, V. Kagan, A. Galstyan, K. Lerman, L. Zhu, E. Ferrara, A. Flammini, and F. Menczer. 2016. The DARPA Twitter Bot Challenge. *Computer* 49, 6 (June 2016), 38–46. <https://doi.org/10.1109/MC.2016.183>
- [26] The U.S. Commodity Futures Trading Commission. 2018. Investor Alert: Public Companies Making ICO-Related Claims. (2018). https://www.cftc.gov/sites/default/files/idc/groups/public/@customerprotection/documents/file/customeradvisory_pumpdump0218.pdf
- [27] Wenting Tu, David W. Cheung, Nikos Mamoulis, Min Yang, and Ziyu Lu. 2016. Investment Recommendation using Investor Opinions in Social Media. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval - SIGIR '16*. ACM Press, New York, New York, USA, 881–884. <https://doi.org/10.1145/2911451.2914699>
- [28] U.S. Securities and Exchange Commission. 2017. Customer Advisory: Beware Virtual Currency Pump-and-Dump Schemes. (2017). <https://www.investor.gov/additional-resources/news-alerts/alerts-bulletins/investor-alert-public-companies-making-ico-related>
- [29] Onur Varol, Emilio Ferrara, Clayton A Davis, Filippo Menczer, and Alessandro Flammini. 2017. Online human-bot interactions: Detection, estimation, and characterization. In *Eleventh international AAAI conference on web and social media*.
- [30] Yinglin Wang. 2017. Stock market forecasting with financial micro-blog based on sentiment and time series analysis. *Journal of Shanghai Jiaotong University (Science)* 22, 2 (4 2017), 173–179. <https://doi.org/10.1007/s12204-017-1818-4>
- [31] Satoshi Watanabe. 1960. Information theoretical analysis of multivariate correlation. *IBM Journal of research and development* 4, 1 (1960), 66–82.
- [32] Liang Wu, Fred Morstatter, Xia Hu, and Huan Liu. 2016. Mining misinformation in social media. *Big Data in Complex and Social Networks* (2016), 123–152.
- [33] Jiahua Xu and Benjamin Livshits. 2018. The anatomy of a cryptocurrency pump-and-dump scheme. *arXiv preprint arXiv:1811.10109* (2018).