# Improved Forecasting of Cryptocurrency Price using Social Signals

**Maria Glenski**
Data Science and Analytics
Pacific Northwest National Laboratory
maria.glenski@pnnl.gov

**Tim Weninger**
Computer Science and Engineering
University of Notre Dame
tweninge@nd.edu

**Svitlana Volkova**
Data Science and Analytics
Pacific Northwest National Laboratory
svitlana.volkova@pnnl.gov

## Abstract

Social media signals have been successfully used to develop large-scale predictive and anticipatory analytics. For example, forecasting stock market prices and influenza outbreaks. Recently, social data has been explored to forecast price fluctuations of cryptocurrencies, which are a novel disruptive technology with significant political and economic implications. In this paper we leverage and contrast the predictive power of social signals, specifically user behavior and communication patterns, from multiple social platforms GitHub and Reddit to forecast prices for three cyptocurrencies with high developer and community interest – Bitcoin, Ethereum, and Monero. We evaluate the performance of neural network models that rely on long short-term memory units (LSTMs) trained on historical price data and social data against price-only LSTMs and baseline autoregressive integrated moving average (ARIMA) models, commonly used to predict stock prices. Our results not only demonstrate that social signals reduce error when forecasting daily coin price, but also show that the language used in comments within the official communities on Reddit (r/Bitcoin, r/Ethereum, and r/Monero) are the best predictors overall. We observe that models are more accurate in forecasting price one day ahead for Bitcoin (4% root mean squared percent error) compared to Ethereum (7%) and Monero (8%).

## Introduction

Cryptocurrencies, like Bitcoin, Ethereum, and Monero, are a new and disruptive technology that are often leveraged in highly volatile and fast-evolving environments. As with stocks and other securities, cryptocurrencies are bought, held, and traded. Unlike traditional currencies and stocks, these digital currencies rely on decentralized systems and cryptographic technologies, *e.g.,* blockchain ledgers, rather than a centralized institution, *e.g.,* banks. In this new paradigm, money is moved more quickly, independently, and often anonymously or semi-anonymously. As a result, the wide adoption and historic volatility of cryptocurrencies have significant political and economic implications. Price speculation, where traders buy securities in the hopes that they will quickly rise in price, often occurs in these highly volatile markets. Speculative trading typically occurs with little (or no) regard to the asset's fundamental value, but rather in regard to patterns in the the asset's price movement, rumor, or other suppositious data.

Signals from social media have been extensively used to predict real world events such as election results (Tumasjan et al. 2010; Sang and Bos 2012; Cameron, Barrett, and Stewardson 2016; Dokoohaki et al. 2015; Wang and Lei 2016; Khatua et al. 2015), movie sales (Mishne, Glance, and others 2006; Asur and Huberman 2010; Tang, Yeh, and Lee 2014; Abel et al. 2010), protests (Maharjan et al. 2018), public health events (Volkova et al. 2017; Corley et al. 2010; Lamb, Paul, and Dredze 2013; Paul, Dredze, and Broniatowski 2014; Bodnar and Salathé 2013), and stock market activity (Bollen, Mao, and Zeng 2011; Chen et al. 2014; Makrehchi, Shah, and Liao 2013; Oh and Sheng 2011; Mao et al. 2012; Martin 2013; Porshnev, Redkin, and Shevchenko 2013; Oliveira, Cortez, and Areal 2013; Rao and Srivastava 2012; Zimbra, Chen, and Lusch 2015; Li, Zhou, and Liu 2016; Zhao et al. 2016). Ding et al. (2014) introduced a deep neural network approach to predict the directionality of stock prices and the S&P 500 index using signals from related news events and Tetlock (Tetlock 2007) highlighted the correlation between media pessimism and market prices and volume. Bollen et al. (2011) predict relative differences in the daily Dow Jones industrial average using measures of collective mood states derived from Twitter activity and found the addition of some but not all possible states improved the predictive ability of their proposed models. Similar to Bollen et al. , we analyze the benefit of including or excluding a range of social signals in our proposed models.

Like in stock markets and securities, cryptocurrency market activity has also been predicted using social signals. Previous work by Kim et al. (2016) predicted the direction of price fluctuations for cryptocurrency coins utilizing social signals from online cryptocurrency forums. Other studies focus on predicting relative changes, *i.e.,* the return, in coin prices (Rao and Srivastava 2012; Wang and Vergne 2017). For example, Wang and Verne (2017) proposed a model to
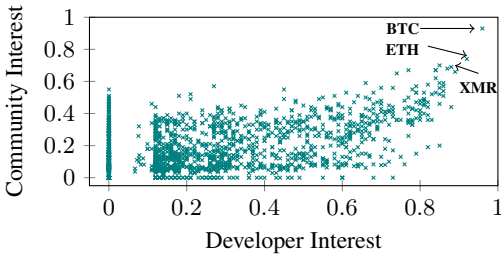
Figure 1: Coins plotted by developer and community interest. The present work will focus on three of the most popular: Bitcoin (BTC), Ethereum (ETH), and Monero (XMR)

predict the return, *i.e.,* the change in price relative to the opening value. Phillips and Gorse (2017) predict the beginning and end of spikes in cryptocurrency prices, which they call "price bubbles", using hidden Markov models that were previously used for the detection of influenza outbreaks. Although these previous works were able to predict certain changes in a cryptocurrency's price, they were not able to predict the actual price of the asset. In the present work we evaluate the benefit of incorporating a variety of social signals into models in order to forecast the *actual daily price high-values* of three popular cryptocurrencies.

With the increasing use and reliance on these digital currencies, price fluctuation forecasting is an interesting yet difficult challenge. We address this problem by leveraging the predictive power of social signals from GitHub and Reddit to forecast immediate and near future prices of three popular cryptocurrencies: Bitcoin, Monero and Ethereum. In summary, our main contributions are:

1. we develop neural network models that incorporate both social and price signals to generate forecasts of coin prices,

2. we present an in-depth analysis of model performance for coin price forecasts up to 3 days in advance for Bitcoin, Ethereum, and Monero when incorporating a variety of social signals and the relative improvement over models that rely solely on price history, and

3. we report average performance of models for forecasts of coin price up to two weeks in advance.

## Why Bitcoin, Ethereum, and Monero?

A preliminary analysis of potential cryptocurrencies identified Bitcoin, Ethereum, and Monero as the top three cryptocurrencies in terms of both developer interest on GitHub and community interest on social media platforms. To determine which coins to focus on in this study, we collected data from CoinGecko[1] for 1,742 cryptocurrencies and performed an initial analysis of coins in terms of the developer interest (a measure of activity in public repositories on GitHub and Bitbucket) and community interest (a measure of discussions and popularity on social media) features released by CoinGecko. The key results of this analysis are illustrated in Figure 1, which shows that Bitcoin, Ethereum, and Monero

have the highest degree of both developer and community interest.

## Social and Financial Data Collection

In this study, we focus on taking advantage of social data to forecast the price for three cryptocurrencies: Bitcoin (BTC), Ethereum (ETH), and Monero (XMR). Bitcoin is the first decentralized cryptocurrency and holds the highest market capitalization (market cap)[2]. Ethereum holds the second highest market cap and relies on the same blockchain technology that underpins Bitcoin. Monero, while still popular, holds a much lower market cap than Bitcoin or Ethereum but focuses on privacy; transactions are private (with the origin, destination, and amounts obfuscated) and untraceable (transactions cannot be linked to a particular cyber- or real-world identity).

Are social signals relevant to cryptocurrency prices? To gain an initial understanding of this question, we illustrate the alignment of price, social interactions, and real-world events related to Bitcoin in Figure 2.

Historical price data (daily high, low, and price at market open and close) for each coin was collected from Crypto-Compare[2]. This resulted in a price history from:

- 2010/07/16 through 2018/05/21 for Bitcoin,

- 2015/08/06 through 2018/05/21 for Ethereum,

- 2015/01/28 through 2018/05/21 for Monero.

In addition to financial data, we collected publicly available data for two social platforms, GitHub and Reddit, from which we extracted social signals for each coin.

**GitHub Dataset**  GitHub is a collaborative software social network primarily used to develop and share novel technologies and software. As of 2017, 67M repositories (repos) used by 24M users and 1.5M organizations were hosted on the site.[3] We collected interactions with the main repo for each coin (bitcoin/bitcoin, ethereum/go-ethereum, and monero-project/monero) from a public subset of the GitHub archive.[4] These interactions can be divided into two main categories of event types: 1) indications of user interest or repo *popularity* such as watching or forking the repo and 2) *direct contributions* when a user reviews or directly contributes to code, reports issues, comments on issues or participates in code reviews.

**Reddit Dataset**  Similarly, we collected all posts and comments submitted to the official subreddits[5] for each of the cryptocurrencies of interest (r/bitcoin, r/ethereum, and r/monero) across the three years (2015 – 2017) for which we collected GitHub events data. Reddit is a popular social

---

[1]https://www.coingecko.com/en/coins/all

[2]https://www.cryptocompare.com/

[3]https://octoverse.github.com/

[4]https://www.githubarchive.org/

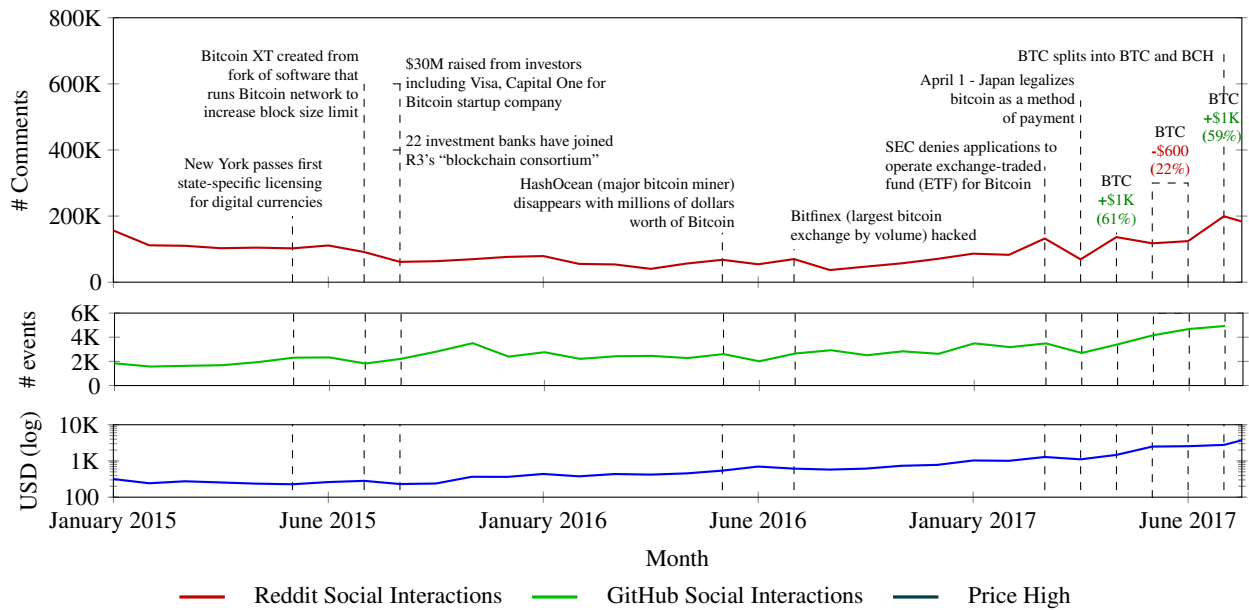[5]Reddit posts and comments are publicly available via an archive hosted at https://files.pushshift.io/reddit/.

Figure 2: Motivation: the alignment of price high (USD) and social interactions on Reddit (r/bitcoin) and GitHub (bitcoin/bitcoin).

news aggregator[6] that allows communities of users to share and discuss information, opinions, and entertainment media. Each of the subreddits has substantial traffic including 3.6K, 500, and 380 comments posted each day, on average, with community-sizes of 913K, 337K, and 137K subscribers (as of August 2018) for r/bitcoin,r/ethereum, and r/monero, respectively.

## Methodology

In this section, we present the methodology used to evaluate the benefit of incorporating signals from social media into models that forecast the future price of a cryptocurrency. We trained and evaluated models that rely on 1) historical price alone, 2) historical price and each social signal, and 3) historical price and combinations of each of the social signals.

### Forecasting Tasks

We define the forecasting tasks as predicting the daily price high of each cryptocurrency of interest $j$ days in advance, $Y_{t_i+j}$, focusing on the immediate and near future $Y_{t_i+1}$, $Y_{t_i+2}$, $Y_{t_i+3}$ (1, 2, or 3 days in the future, respectively) using predictive signals from $k$ days in the immediate past $X_{[t_{i-k},t_i]}$ and evaluate any benefits that arise from incorporating up to two weeks of signal history by varying $k$ from 1 to 14 days. We consider models that incorporate social signals and historical price values versus those that consider only historical pricing as predictive signals to identify the performance gains (if any exist) when social signals are included. Models are trained and evaluated independently for each of the cryptocurrencies of interest.

### Social Signals

Along with the daily price high, we use a variety of signals of popularity, activity, and language used in discussions across two popular social media platforms: GitHub and Reddit. Here, we describe the social signals we extracted from social media activity related to the three coins of interest.

For the GitHub platform, we consider two types of social signals:

- $GH_{Pop}$ – a vector representation of the daily totals for each **popularity** event: the Watch event, where users star a repo in order to receive notifications, and the Fork event, where users create a copy of the repo code. These event types provide a measure of how popular a given repo is among users who may or may not be direct contributors.

- $GH_{All}$ – a vector representation of the **overall activity**, *i.e.,* daily counts for each popularity (Watch and Fork) and direct contribution event types (CommitComment, Issue, IssueComment, PullRequest, PullRequestReview-Comment, and Push).

For the Reddit platform, we consider four types of social signals:

- $R_{Vol}$ – the **volume** of comments posted each day, a signal for the size of discussion within the coin's official subreddit.

- $R_{Lang}$ – the **language** used in comments represented as 10k-dimensional vectors of word-level daily-normalized statistics that focus on the most frequent unigrams.

- $R_{Score}$ – a vector representation of the quartiles of Reddit scores (*i.e.,* # upvotes - # downvotes) for comments posted each day, an indication of the range of **popularity** of comments for the given day.

- $R_{Sent}$ – quartiles for the subjectivity and polarity of comments each day which provides a signal of the distribution of **sentiment** in discussions within the coin's official subreddit.

Before evaluating forecasting models, we explore the relationships between coin-related social signals and their respective price high time-series. To do so, we first examine the correlation of social signals ($x$) and coin-price ($y$) over a sample of $N$ days for each of the three coins of interest using Pearson R correlation to examine the linear relationships between social signals and coin price. Pearson $R$ correlation ranges from -1 (perfectly inversely correlated) to 1 (perfectly correlated), where a score of 0 indicates linearly independent variables, *i.e.,* no correlation. Next, we consider non-linear correlations between signals ($x$) and price ($y$) using distance correlation (Székely and Rizzo 2013; Szekely, Rizzo, and others 2014; Székely et al. 2007):

$$dCorr(x,y) = \left\{ \begin{array}{ll} \frac{dCov(x,y)}{\sqrt{dVar(x)dVar(y)}}, & dVar(x)dVar(y) > 0 \\ 0, & dVar(x)dVar(y) = 0 \end{array} \right\}$$

where distance covariance ($dCov$) and distance variance ($dVar$) are defined as:

$$dCov(x,y) = \sqrt{\frac{\sum_{k,l=1}^{n} A_{kl}B_{kl}}{n^2}}, \ dVar(x) = dCov(x,x)$$

and $A_{kl} = a_{kl} - \frac{1}{n}\sum_{l=1}^{n} a_{kl} - \frac{1}{n}\sum_{k=1}^{n} a_{kl} + \frac{1}{n^2}\sum_{k,l=1}^{n} a_{kl}$.

$B_{kl}$ is defined similarly with $b_{kl}$ in place of $a_{kl}$ where $a_{kl}$ and $b_{kl}$ are Euclidean distance matrices of $x$ and $y$, respectively, defined as

$$a_{kl} = |x_k - x_l|, \quad b_{kl} = |y_k - y_l|.$$

Distance correlation varies from 0 to 1, where a distance correlation of 0 indicates independence of variables. Finally, we also examine the variation of each signal ($x$) by the interquartile range (IQR) and standard deviation $\sigma$ for each of the signals of interest, identifying those that remain consistent over the time period of interest (and thus are potentially less informative).

## Forecasting Models

Neural network models with long short-term memory (LSTM) layers have previously been used to effectively forecast influenza dynamics from a combination of clinical and social media signals, outperforming models that did not incorporate the social signals (Volkova et al. 2017). LSTMs are a type of recurrent neural networks with built in memory cells that store information and can exploit long range context (Hochreiter and Schmidhuber 1997). These networks are surrounded by gating units that allow or prohibit the reset, read, and writing of such information. Inspired by the influenza dynamic models, we propose a neural network model, shown in Figure 3, that also utilizes LSTM layers.

The proposed neural network architecture consists of a 400-dimensional LSTM layer followed by an 800-dimensional LSTM then a dense layer with a single unit.
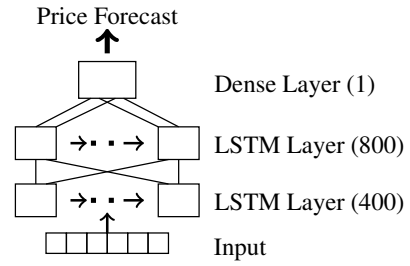


Figure 3: Neural network model architecture.

We do not employ an activation layer or activation function in the final dense output layer because we have defined each forecasting task as a regression task to predict numerical high values so an activation transformation is not needed. The optimizer used by the model is the ADAM optimizer (a parameter specific adaptive learning rate method)(Kingma and Ba 2014) and we optimized performance by minimizing the mean squared error of a 20% subset of the training data which is used as a validation set. To avoid overfitting, we employed early stopping callbacks with a maximum limit of 20 epochs. As a result, although models could have trained for up to 20 epochs, our models typically needed, on average, between five and six epochs only. The described architecture of the proposed model enables a relatively short training time for an expedient and scalable framework.

Signals used as input, e.g., social signals and historical price data, are fed to the network as a single concatenated vector. Before concatenation, each signal value is normalized to range between 0 and 1 using min-max normalization for the given feature across the entire dataset. Target price values ($T$) and signal input vectors ($I$) are defined as:

$$T(Y,j) = \{Yj, Y_{j+1}, ...Y_{|Y|}\}$$

$$I(X,k) = \{< X_{i-k}, ..., X_{i-1}, X_i > for \quad i \in [k, |X|] \}$$

where $X$ is a 2D array of the min-max normalized signals of interest, $Y$ is the historical price high time series, $j$ is the size of the forecasting window (*i.e.,* how many days in advance to predict the price), and $k$ is the number of days of signal history used.

**Parameter Tuning**   To identify our final model configuration, we performed a series of parameter tuning. We varied several of the model parameters: batch sizes of 16, 32, and 64; learning rates of 0.1, 0.01, 0.001, 0.0001, and 0.00001; and combinations of LSTM layers ranging from 10 to 400-dimensional layers followed by 20 to 800-dimensional layers. We found best performance when using a batch size of 16, learning rate of 0.001, and 400 and 800 units for the first and second LSTM layers, respectively.

**Baseline**   The autoregressive integrated moving average (ARIMA) model is a commonly used architecture when forecasting the price or return of stocks (Pai and Lin 2005; Zhang 2003). Essentially ARIMA models treat the future value of a variable (*e.g.,* price) as a linear combination of

past values and errors. We train a variety of ARIMA models and identify the top performing ARIMA model for each of our forecasting windows of interest to use as baselines. To identify these baseline models, we fix the size of the moving average window to 0, the difference order to 1 (to make the time series stationary), and perform autocorrelation analysis to identify a range of appropriate lag parameters. We found that models with a lag order of 0 achieved the best performance overall across coins. $R$ results for these baseline models were at least 0.95 ($p < 0.001$) for Bitcoin, at least 0.92 ($p < 0.001$) for Ethereum, and at least 0.91 ($p < 0.001$) for Monero.

## Evaluation

To ensure performance comparison across identical train and test dates for all combinations of training and forecasting window sizes, the train and test periods were identified using the largest possible window sizes. We first restricted the dataset to the period for which we had data for all signals across all coins. Then, input vectors and target prices were formatted for the largest training (14) and forecasting (3) window sizes and split into 80% for train and 20% for test. This resulted in a training period from 2015/11/11 through 2017/04/18 (525 days) and a test period from 2017/05/04 through 2017/08/31 (120 days) used for all model configurations.

Model performance is evaluated using the following error measurements: root mean squared error (RMSE), mean squared percentage error (MSPE), mean absolute percentage error (MAPE), max absolute percentage error (MaxAPE), and root mean squared percentage error (RMSPE). For a set of $N$ predictions ($\hat{y}$) and true price values ($y$), MAPE, also known as mean absolute percentage deviation, is defined as:

$$MAPE(\hat{y}, y) = \frac{1}{N} \sum_{i=0}^{N} \frac{|\hat{y}_i - y_i|}{y_i}$$

maximum absolute percentage error is defined as:

$$MaxAPE(\hat{y}, y) = \max\left(\frac{|\hat{y}_i - y_i|}{y_i}\right)$$

and mean squared percentage error is defined as:

$$MSPE(\hat{y}, y) = \frac{1}{N} \sum_{i=0}^{N} \frac{(\hat{y}_i - y_i)^2}{y_i}$$

Prices vary widely between coins and, for some, within coins over time. Therefore, we primarily report the performance metrics that consider percentage errors as they allow a relative comparison across the three coins of interest.

## Social Signal Analysis

In this section we explore the relationships between daily price and signals of popularity and direct contributions to GitHub repositories and volume, sentiment, and popularity of cryptocurrency-related discussions on Reddit for our three coins of interest. To identify if an informative relationship exists, we examine person and distance correlation of

Table 1: Pearson $R$ correlation and distance correlation ($DC$) of daily price and social signals. Pearson results are significant ($p < 0.001$) unless indicated with a dash '—' ($p \geq 0.05$).

| Social Signal | | Bitcoin $R$ | Bitcoin $DC$ | Ethereum $R$ | Ethereum $DC$ | Monero $R$ | Monero $DC$ |
|---|---|---|---|---|---|---|---|
| GitHub | Watch | 0.87 | 0.86 | 0.68 | 0.73 | 0.72 | 0.68 |
| | Fork | 0.75 | 0.72 | 0.40 | 0.38 | 0.41 | 0.48 |
| | Issues | 0.05 | 0.09 | 0.22 | 0.05 | 0.00 | 0.36 |
| | IssueComment | 0.13 | 0.14 | 0.36 | 0.29 | 0.25 | 0.54 |
| | Push | 0.06 | 0.09 | 0.06 | 0.17 | 0.07 | 0.11 |
| | CommitComment | 0.06 | 0.09 | 0.04 | 0.12 | 0.08 | 0.08 |
| | PullRequest (PR) | 0.18 | 0.20 | 0.14 | 0.22 | 0.15 | 0.21 |
| | PRReviewComment | 0.39 | 0.37 | 0.29 | 0.36 | 0.22 | 0.44 |
| Reddit | Comment Volume | 0.58 | 0.62 | 0.48 | 0.51 | 0.67 | 0.78 |
| | Subjectivity | — | 0.00 | — | 0.05 | 0.16 | 0.25 |
| | Polarity | — | 0.02 | 0.13 | 0.15 | 0.31 | 0.43 |
| | Score | 0.34 | 0.37 | 0.47 | 0.59 | 0.69 | 0.79 |

Table 2: Standard deviations ($\sigma$) and Inter-Quartile Range (IQR) of daily price and social signals.

| Signal | Bitcoin $\sigma$ | Bitcoin $IQR$ | Ethereum $\sigma$ | Ethereum $IQR$ | Monero $\sigma$ | Monero $IQR$ |
|---|---|---|---|---|---|---|
| Price High | 868.00 | 626.85 | 99.44 | 34.67 | 19.71 | 11.96 |
| Watch | 10.34 | 7.00 | 2.09 | 2.00 | 1.60 | 1.00 |
| Fork | 4.55 | 4.00 | 0.66 | 0.00 | 0.94 | 1.00 |
| Issues | 3.34 | 3.00 | 27.42 | 1.00 | 2.47 | 2.00 |
| IssueComment | 24.98 | 34.00 | 7.93 | 7.00 | 9.51 | 10.00 |
| Push | 3.57 | 5.00 | 1.06 | 0.00 | 3.80 | 1.00 |
| CommitComment | 0.90 | 0.00 | 0.18 | 0.00 | 0.26 | 0.00 |
| PullRequest (PR) | 5.92 | 8.00 | 1.14 | 0.00 | 5.15 | 4.00 |
| PRReviewComment | 13.25 | 15.00 | 2.57 | 0.00 | 5.50 | 1.00 |
| Comment Volume | 1990.43 | 1981.00 | 695.86 | 369.00 | 376.34 | 466.00 |
| Score | 242.97 | 129.00 | 80.00 | 47.00 | 18.34 | 19.00 |
| Subjectivity | 0.00 | 0.00 | 0.00 | 0.00 | 0.13 | 0.00 |
| Polarity | 0.01 | 0.00 | 0.06 | 0.00 | 0.23 | 0.25 |

social signals with coin price and the variance of all features to identify which signals to include in the battery of model configurations we consider in our ablation experiments.

In Table 1, we see that price is correlated ($p < 0.001$) with daily volume for each of the GitHub event types, although to varying degrees within and across coins. We also see that both Fork and Watch events are highly correlated for Bitcoin while Watch events are highly correlated and Fork events are moderately correlated for the other two coins. For the Reddit platform, we see that daily comment volume and comment scores are both correlated with daily price highs across all three coins of interest. While there is no significant linear correlation between subjectivity and polarity of the daily batch of comments and price highs for Bitcoin, these features have varying levels of correlation across the remaining two coins of interest. As we saw with GitHub features, there is some variation across coins of interest in the relationships between the various social signals and price timeseries.

When we consider how the signals vary within themselves in Table 2, we find similar patterns across and within coins. The most highly correlated social signals have larger, rela-

tive to other signals, variation as summarized by the standard deviations and inter-quartile ranges of the signal vectors. We find that the subjectivity and polarity signals from Reddit comments linked to Bitcoin and Ethereum that held no significant correlation with price also show little to no variation. However, we see they vary slightly to moderately for the third coin, Monero. As each of our social signals indicate a relationship with price for at least one of our coins of interest, we include all GitHub and Reddit social signals as well as combinations of signals from both the GitHub and Reddit platforms in our ablation experiments and highlight the results in the following section.

## Forecasting Results

In this section, we describe the performance of models that rely on historical price alone, historical price and each signal from GitHub or Reddit, and historical price and combinations of each GitHub and Reddit signal. In particular, we highlight models which incorporated social signals that achieved high performance relative to the baseline ARIMA models and LSTMs that relied on historical price alone.

First, we explored the benefit of increased signal history with a comparison of model performance using signals from one to fourteen days in the past in Figure 4. We plot MSPE as a function of training window size, *i.e.,* the number of days of signal history to rely on, for neural network (LSTM) models that rely on each of the combinations of predictive signals when forecasting price one to two days in advance. *Interestingly, we see that models with the smallest window size (1) achieved the best performance.* Therefore, we use a signal history window size of 1 day for subsequent model evaluation.

Next we focused on identifying model configurations that achieve the best performance to evaluate the benefit of incorporating social signals from a variety of platforms and the benefit of a variety of combinations of social signals. To do so, we perform an ablation study where we compare models that incorporate each combination of price and social signals with models (LSTM and ARIMA baselines) that relied solely on historical prices. To identify the best, overall model, we then averaged percentage errors across the three coins and ranked social signal-infused models with the baseline ARIMA and neural network models that did not incorporate social signals.

We present a summary of the averaged error for the top performing models in Table 3, ordered by the mean of mean RMSPE over the three forecasting tasks. *The top performing model is the proposed LSTM that relies on price history and $R_{Lang}$, the representation of the language used in comments within the official subreddit.* Here, we see that LSTM models that incorporate social signals outperform the price only baselines when averaged across the three coins and forecasting windows. If we only consider the immediate forecasting window of one day in advance, the proposed LSTM model that relies solely on price history outperforms the others.

We then consider performance for each coin individually in Table 4. Here we see that, in most cases, neural network models that incorporated social signals slightly outperformed both ARIMA baselines and neural network models
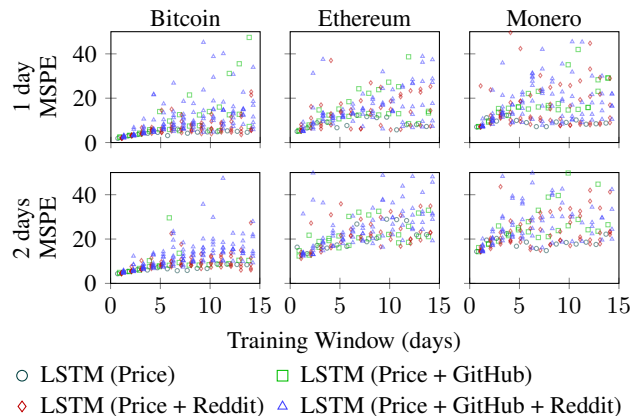


Figure 4: MSPE results for neural network models for each coin plotted by varying training window sizes for the first two forecasting windows, *i.e.,* when predicting one or two days in advance. Jitter has been added to x-axis for enhanced readability.

Table 3: RMSPE averaged over the three coins of interest (Bitcoin, Ethereum, and Monero) when predicting price up to 3 days in advance for ARIMA baseline and top performing neural network (LSTM) models. Lowest values highlighted in bold.

| Model | Signal | Forecasting Window (days) | | | |
|---|---|---|---|---|---|
| | | 1 | 2 | 3 | mean |
| LSTM | $\$ + R_{Lang}$ | 6.70 | **9.88** | **12.06** | **9.55** |
| LSTM | $\$ + GH_{Pop} + R_{Lang}$ | 6.64 | 9.99 | 12.40 | 9.68 |
| LSTM | $\$ + R_{Vol}$ | 6.78 | 9.98 | 12.48 | 9.75 |
| LSTM | $\$$ | **6.60** | 10.46 | 12.32 | 9.79 |
| ARIMA | $\$$ | 7.30 | 10.56 | 13.10 | 10.32 |

that relied solely on price history, but these results are not statistically significant. However, we see that our proposed LSTM neural network architecture, and especially neural network models that incorporate social signals, minimize the maximum absolute percentage error (MaxAPE). That is, in worst-case prediction performance, we see a much lower error rate for our proposed models that rely on social signals alongside price history. We see the best performance is achieved when models forecast the price one day in the future (FW = 1 Day); unsurprisingly, it is easiest to predict the next day's price using signals from the day before.

If we expand the range of forecasting windows beyond the near and immediate future, we see that, again, the LSTM model that incorporates the $R_{Lang}$ social signal alongside historical price has the best performance across the three coins. Table 5 presents the RMSPE of top performing models when forecasting price up to two weeks in advance. Figure 5 illustrates how this top performing model, performs similarly to price-only LSTM and ARIMA models, on average, outperforms the price-only models in respect to the worst-case prediction errors (MaxAPE).

We illustrate the actual predictions of the top performing social signal enhanced neural network model against the true

Table 4: MAPE, RMSPE, and MaxAPE results for baseline ARIMA, neural network models that rely solely on historical price, and the top performing social signal enhanced neural network models, as identified during ablation experiments. Results for neural network models (LSTM) reported used the top performing training window size of 1 day. Lowest error rates are highlighted in bold.

| | Model | Signals | Bitcoin MAPE | Bitcoin RMSPE | Bitcoin MaxAPE | Ethereum MAPE | Ethereum RMSPE | Ethereum MaxAPE | Monero MAPE | Monero RMSPE | Monero MaxAPE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **$ in 1 day** | ARIMA | $ | 3.24 ± 0.65 | 4.83 ± 3.35 | 23.19 | 5.13 ± 1.18 | 8.29 ± 6.50 | 48.19 | 5.84 ± 1.19 | 8.78 ± 6.26 | 42.71 |
| | LSTM | $ | **3.11 ± 0.55** | **4.35 ± 2.69** | 17.16 | **4.71 ± 0.96** | **7.08 ± 4.48** | **26.18** | **5.71 ± 1.12** | 8.39 ± 5.96 | 42.53 |
| | LSTM | $ + $R_{Lang}$ | 3.26 ± 0.55 | 4.45 ± 2.69 | **16.66** | 4.76 ± 0.97 | 7.16 ± 4.54 | 26.53 | 5.94 ± 1.11 | 8.49 ± 5.85 | **41.83** |
| | LSTM | $ + $GH_{Pop}$ + $R_{Lang}$ | 3.19 ± 0.56 | 4.44 ± 2.68 | 16.77 | 4.73 ± 0.97 | 7.10 ± 4.50 | 26.35 | **5.71 ± 1.12** | **8.38 ± 5.97** | 42.70 |
| | LSTM | $ + $R_{Vol}$ | 3.35 ± 0.55 | 4.52 ± 2.71 | 16.49 | 5.06 ± 0.99 | 7.44 ± 4.68 | 27.36 | 5.69 ± 1.12 | 8.38 ± 5.99 | 42.72 |
| **$ in 2 days** | ARIMA | $ | 5.35 ± 0.87 | 7.18 ± 4.13 | 23.35 | 8.44 ± 1.56 | 12.05 ± 7.83 | 48.35 | 8.74 ± 1.61 | 12.44 ± 7.56 | 41.13 |
| | LSTM | $ | 5.14 ± 0.76 | 6.61 ± 3.68 | 22.03 | 9.81 ± 1.49 | 12.78 ± 6.84 | 37.30 | 8.44 ± 1.55 | 11.99 ± 7.40 | 41.35 |
| | LSTM | $ + $R_{Lang}$ | 5.38 ± 0.79 | 6.90 ± 3.82 | 23.79 | **7.67 ± 1.34** | **10.61 ± 6.07** | **33.64** | 8.90 ± 1.50 | 12.12 ± 7.26 | **39.80** |
| | LSTM | $ + $GH_{Pop}$ + $R_{Lang}$ | 5.40 ± 0.78 | 6.89 ± 3.78 | 23.23 | 7.87 ± 1.43 | 11.11 ± 6.33 | 35.01 | **8.32 ± 1.56** | **11.95 ± 7.48** | 41.73 |
| | LSTM | $ + $R_{Vol}$ | **5.11 ± 0.76** | **6.59 ± 3.67** | **21.78** | 8.11 ± 1.44 | 11.34 ± 6.38 | 35.31 | 8.44 ± 1.56 | 12.02 ± 7.48 | 41.80 |
| **$ in 3 days** | ARIMA | $ | 6.80 ± 1.03 | 8.85 ± 4.72 | 23.94 | 10.85 ± 1.83 | 14.81 ± 8.74 | 50.03 | 11.14 ± 1.99 | 15.64 ± 9.17 | 49.89 |
| | LSTM | $ | **6.33 ± 0.92** | **8.08 ± 4.29** | **23.01** | 10.36 ± 1.66 | 13.79 ± 7.35 | 37.17 | **10.77 ± 1.92** | **15.07 ± 8.90** | 50.23 |
| | LSTM | $ + $R_{Lang}$ | 6.37 ± 0.92 | 8.12 ± 4.31 | 23.92 | 9.86 ± 1.53 | 12.96 ± 6.90 | 35.41 | 10.96 ± 1.90 | 15.10 ± 8.85 | **49.36** |
| | LSTM | $ + $GH_{Pop}$ + $R_{Lang}$ | 6.36 ± 0.92 | 8.13 ± 4.30 | 23.13 | **9.84 ± 1.53** | **12.92 ± 6.87** | 35.31 | 11.77 ± 2.02 | 16.16 ± 9.28 | 52.58 |
| | LSTM | $ + $R_{Vol}$ | 6.36 ± 0.92 | 8.13 ± 4.28 | 23.08 | 10.38 ± 1.67 | 13.84 ± 7.37 | 37.31 | 10.97 ± 1.99 | 15.46 ± 9.07 | 51.53 |

Table 5: RMSPE averaged over the three coins of interest (Bitcoin, Ethereum, and Monero) when predicting price up to 14 days in advance for ARIMA baseline and top performing neural network (LSTM) models.

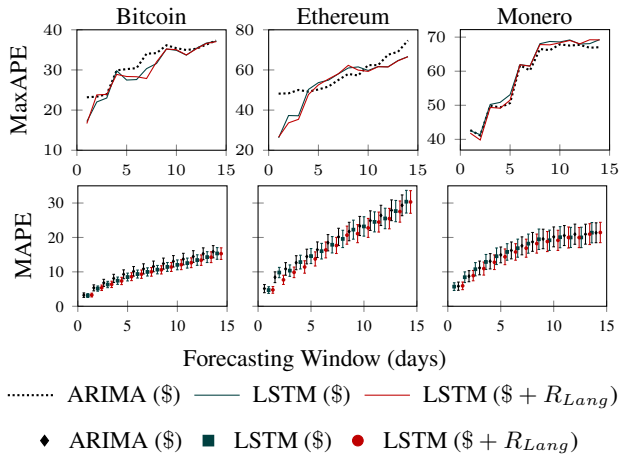| Model | Signals | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LSTM | $ + $R_{Lang}$ | 6.70 | **9.88** | **12.06** | **13.92** | **15.83** | 17.57 | **18.72** | 20.45 | **21.22** | **22.25** | **23.38** | **23.87** | 25.23 | 26.67 | **18.41** |
| LSTM | $ | **6.60** | 10.46 | 12.32 | 14.63 | 16.35 | 17.38 | 18.74 | 20.08 | 21.92 | 22.44 | 23.43 | 23.92 | **24.96** | 26.71 | 18.57 |
| LSTM | $ + $GH_{Pop}$ | 6.80 | 10.27 | 12.76 | 14.50 | 16.19 | **17.36** | 18.75 | **19.97** | 21.40 | 23.13 | 23.76 | 24.26 | 25.66 | 27.32 | 18.72 |
| LSTM | $ + $GH_{Pop}$ + $R_{Lang}$ | 6.64 | 9.99 | 12.40 | 14.19 | 16.17 | 17.66 | 19.14 | 20.56 | 21.62 | 23.16 | 24.01 | 24.90 | 26.11 | 26.67 | 18.80 |
| ARIMA | $ | 7.30 | 10.56 | 13.10 | 15.07 | 16.62 | 18.15 | 19.62 | 20.95 | 22.17 | 23.17 | 23.99 | 24.88 | 25.58 | **26.32** | 19.11 |



Figure 5: MaxAPE (above) and MAPE with 95% confidence intervals (below) for each coin plotted by varying forecasting window sizes (1 to 14 days in advance). Jitter has been added to x-axis of MAPE plots for enhanced readability.

price values over the test period in Fig. 6. This allows us to visualize not only the performance of each model relative to the true price by day but also relative to the trends in true price values across coins. As we saw in Table 2, the range of true price values differs greatly by coin. While Bitcoin ranges from $1000 to $5000, Ethereum ranges from $100 - $400 and Monero ranges between $30 and $50 for the majority of the test period with a spike to between $100 and $150 in the final 10 days of the test period. *The differences in the variance of true price parallel the differences in model performance.* These results indicate that variations in price or the range of price values may heavily affect predictability of coin price; coins with lower price values or variance are more difficult to predict.

## Discussion

In each of our analyses it is apparent that models perform best when forecasting the price for Bitcoin then Ethereum and worst when forecasting for Monero. In Fig. 6, we saw evidence that this may be attributable to variations in true price but we hypothesize it may also be tied to the size of activity on social media or market cap – Monero has both the smallest social activity from which we drew social signals (GitHub event and comment counts noted above when we described data collection) and the lowest market cap (as is shown in Fig. 7). As we saw in Fig. 6, Monero also had the smallest range of price values. This suggests that activity volumes may have a direct effect on cryptocurrency prices, or at the least, performance of models that forecast
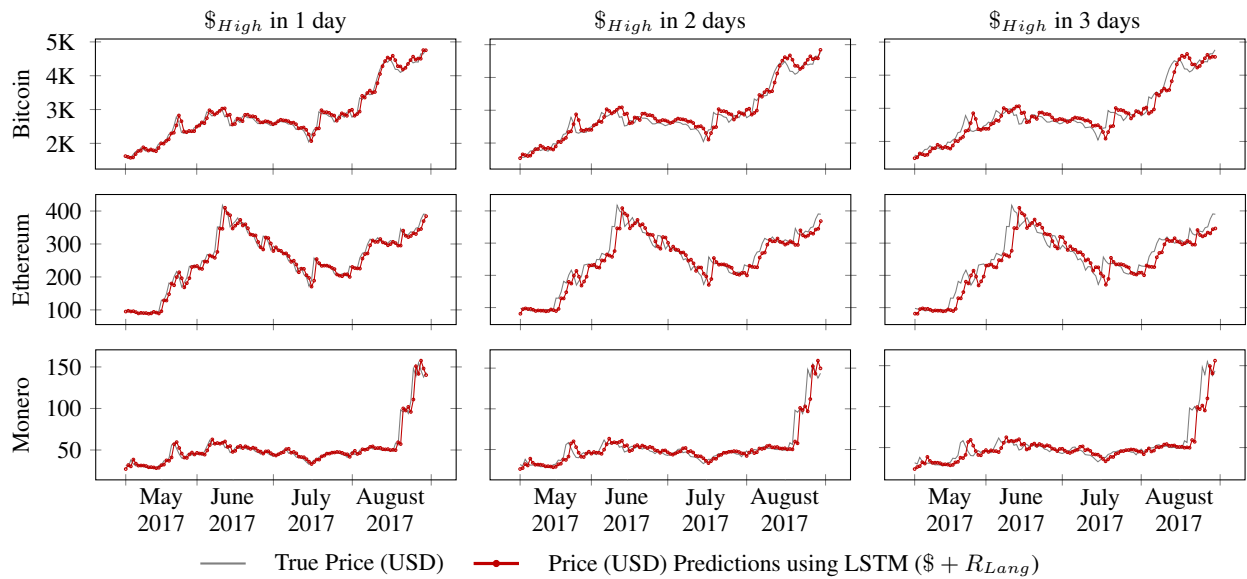
Figure 6: Predictions from the overall top performing neural network model (incorporating price and $R_{Lang}$, language used in comments posted to the official subreddit communities). Solid gray lines represent the true price high values for each day across the test period (2017/05/04-2017/08/31) and colored lines with markers (representing individual predictions) plot predictions.

coin price.

As we saw above, the best model performance is achieved for forecast windows of 1 day in advance, an intuitive result as this is an easier task than longer windows of 2 and 3 days. However, we also find better performance across all forecasting windows for Bitcoin. Bitcoin is both the oldest and most established coin with the largest market cap. Within our three coins of interest, Ethereum has the next highest market cap but Monero has the next longest lifetime. As a result, we plot MAPE as a function of lifetime and as a function of market cap[7] in Fig. 7 to explore whether one of these characteristics may have an effect on price values or model performance. We calculate lifetime as the number of days from when the genesis block was mined (*i.e.,* the beginning of the "ledger" of transactions) to the end of the testing period (2017/08/31).

When we examine Fig. 7, we find that MaxAPE tends to decrease as market cap increases among top performing models but we do not find the same pattern when we plotted MAPE by lifetime. It is important to note that the sample size of coins is small (N=3) and these patterns may not hold among a larger sample of coins or among a set of coins with more variation in lifetimes or market cap. However, these results combined with the increased performance of models that rely on indications of popularity (event counts, comment volume, etc.) suggest that the popularity of a coin affects the performance of predictive models that forecast coin prices. Intuitively, popularity should affect price – a coin that no one knows about or that is less popular probably has a lower price.

There are several potential impacts of identifying models
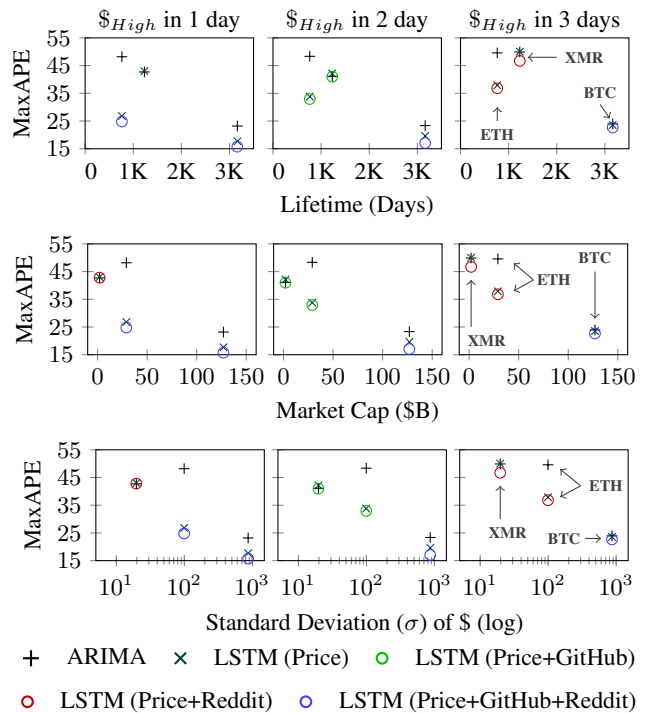


Figure 7: MaxAPE as a function of lifetime (# days between mining of genesis block and end of test period), as a function of market cap (in billion USD), and standard deviations of coin price high (USD) for ARIMA and LSTM price only and top performing price and social signal LSTM models.

---

[7]Market cap for each coin collected from CryptoCompare.com on September 4th 2018.

that reliably forecast the actual prices of cryptocurrencies, beyond the obvious use of identifying which coin to purchase and when. One such use would be as indications of irregular and potentially fraudulent or deceptive activity. For example, a "pump and dump" where a group artificially inflates the price of a coin with excessive transactions and then sells to unsuspecting speculators outside the group at the peak. Although illegal in stocks and securities, this behavior is not yet regulated in cryptocurrencies and is openly advertised (Fuscaldo 2018; Town 2018). Similarly, fabricated or misleading news stories and social media postings have been used to artifically inflate or decrease coin prices. Say there is a model that reliably predicts the price of a cryptocurrency and it begins to fail at an unexpected magnitude, could this be used as an indication of "pump and dump" activity or other manipulation?

## Conclusions and Future Work

In the current work, we have presented models that incorporate social signals for improved forecasts of cryptocurrency price values. Rather than predicting the returns (relative price difference) or the direction of fluctuations (increase versus decrease in price) – both are much easier tasks compared to our task formulation, our models predict the daily price high in USD. We focused on social signals from the language, volume, and sentiment of discussions on Reddit and indications of popularity and direct contribution activity on GitHub. Our analysis of model performance and comparisons to baselines that rely solely on price history have identified the benefit of social signals. Although performance improvements, on average, are not statistically significant, in a task such as price value prediction in a market as volatile as cryptocurrencies, even a modest performance improvement can have a notable impact. More so, the minimization of worst-case error (MaxAPE) when using our proposed social signal-infused model is of significant benefit.

With the speed and volatility of cryptocurrency markets, it may be more valuable to model price in granularities of hours or minutes rather than days. However, a limitation of this approach is the accessibility of such fine-grained pricing data. As the price data available was provided in daily increments, we presented the results of models for predictions of more coarse-grained daily price high-values. The proposed neural network models and framework of social signal extraction could be easily adapted to a more fine-grained approach.

Due to data collection (API) constraints, we performed a similar analysis of models using Twitter social signals for a subset of the time period of interest. Results were inconclusive but future work will also consider such social signals from Twitter and other related platforms alongside GitHub and Reddit. We are also in the process of collecting data to cover not only the time period considered in the current work but also a second dataset that considers the periods of historic spikes and dips in Bitcoin pricing, *i.e.,* the period covering the significant spike in the price of Bitcoin at the end of 2017 and steep drops that occurred over the first few months of 2018.

Our analysis and discussion have also identified several other avenues of future work. One such avenue is the potential to generalize the top performing models for the most popular coin or coins, *e.g.,* Bitcoin, to forecast price values for other coins, such as new or less popular coins. Can we combine the social signals for Bitcoin or cryptocurrencies in general with a coin's historical price features to predict coins that do not have a sufficient volume of social activity? Another direction for future work is the adaption of these models for use in anomaly detection — can models be adapted to identify coins that may be compromised (*e.g.,* through targeted deception in news stories and social media discussions or fraudulent behavior like "pump and dump" schemes)? Finally, a third avenue is an expanded exploration into the trends identified as potential drivers of price predictability: social activity volumes, market share, coin lifetimes, variance, etc. Our analysis has identified several trends and potential explanatory variables however we focused on the three coins with the highest community and developer interest. An expanded analysis that focuses on a larger variety of coins along both axes would provide a more rigorous evaluation.

## References

[Abel et al. 2010] Abel, F.; Diaz-Aviles, E.; Henze, N.; Krause, D.; and Siehndel, P. 2010. Analyzing the blogosphere for predicting the success of music and movie products. In *ASONAM*, 276–280. IEEE/ACM.

[Asur and Huberman 2010] Asur, S., and Huberman, B. A. 2010. Predicting the future with social media. In *IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology-Volume 01*, 492–499. IEEE Computer Society.

[Bodnar and Salathé 2013] Bodnar, T., and Salathé, M. 2013. Validating models for disease detection using twitter. In *International Conference on World Wide Web*, 699–702. ACM.

[Bollen, Mao, and Zeng 2011] Bollen, J.; Mao, H.; and Zeng, X. 2011. Twitter mood predicts the stock market. *Journal of Computational Science* 2(1):1–8.

[Cameron, Barrett, and Stewardson 2016] Cameron, M. P.; Barrett, P.; and Stewardson, B. 2016. Can social media predict election results? evidence from new zealand. *Journal of Political Marketing* 15(4):416–432.

[Chen et al. 2014] Chen, C.; Dongxing, W.; Chunyan, H.; and Xiaojie, Y. 2014. Exploiting social media for stock market prediction with factorization machine. In *IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)*, volume 2, 142–149. IEEE.

[Corley et al. 2010] Corley, C. D.; Cook, D. J.; Mikler, A. R.; and Singh, K. P. 2010. Using web and social media for influenza surveillance. In *Advances in Computational Biology*. Springer. 559–564.

[Ding et al. 2014] Ding, X.; Zhang, Y.; Liu, T.; and Duan, J. 2014. Using structured events to predict stock price movement: An empirical investigation. In *Proceedings of*

*the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1415–1425.

[Dokoohaki et al. 2015] Dokoohaki, N.; Zikou, F.; Gillblad, D.; and Matskin, M. 2015. Predicting swedish elections with twitter: A case for stochastic link structure analysis. In *ASONAM*, 1269–1276. IEEE/ACM.

[Fuscaldo 2018] Fuscaldo, D. 2018. 'pump and dump' hits cryptocurrency market. *Investopedia*.

[Hochreiter and Schmidhuber 1997] Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.

[Khatua et al. 2015] Khatua, A.; Khatua, A.; Ghosh, K.; and Chaki, N. 2015. Can# twitter_trends predict election results? evidence from 2014 indian general election. In *HICSS*, 1676–1685. IEEE.

[Kim et al. 2016] Kim, Y. B.; Kim, J. G.; Kim, W.; Im, J. H.; Kim, T. H.; Kang, S. J.; and Kim, C. H. 2016. Predicting fluctuations in cryptocurrency transactions based on user comments and replies. *PloS one* 11(8):e0161197.

[Kingma and Ba 2014] Kingma, D. P., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

[Lamb, Paul, and Dredze 2013] Lamb, A.; Paul, M. J.; and Dredze, M. 2013. Separating fact from fear: Tracking flu infections on twitter. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 789–795.

[Li, Zhou, and Liu 2016] Li, Q.; Zhou, B.; and Liu, Q. 2016. Can twitter posts predict stock behavior?: A study of stock market with twitter social emotion. In *ICCCBDA*, 359–364. IEEE.

[Maharjan et al. 2018] Maharjan, S.; Shrestha, P.; Porterfield, K.; Arendt, D.; and Volkova, S. 2018. Towards anticipatory analytics: Forecasting instability across countries from dynamic knowledge graphs. In *Proceedings of the 5th Pacific Northwest Regional NLP Workshop (NW-NLP 2018), Redmond, Washington*.

[Makrehchi, Shah, and Liao 2013] Makrehchi, M.; Shah, S.; and Liao, W. 2013. Stock prediction using event-based sentiment analysis. In *IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)-Volume 01*, 337–342. IEEE Computer Society.

[Mao et al. 2012] Mao, Y.; Wei, W.; Wang, B.; and Liu, B. 2012. Correlating s&p 500 stocks with twitter data. In *International workshop on hot topics on interdisciplinary social networks research*, 69–72. ACM.

[Martin 2013] Martin, V. 2013. Predicting the french stock market using social media analysis. In *International Workshop on Semantic and Social Media Adaptation and Personalization*, 3–7. IEEE.

[Mishne, Glance, and others 2006] Mishne, G.; Glance, N. S.; et al. 2006. Predicting movie sales from blogger sentiment. In *AAAI spring symposium: computational approaches to analyzing weblogs*, 155–158.

[Oh and Sheng 2011] Oh, C., and Sheng, O. 2011. Investigating predictive power of stock micro blog sentiment in forecasting future stock price directional movement. In *ICIS*, 1–19. Citeseer.

[Oliveira, Cortez, and Areal 2013] Oliveira, N.; Cortez, P.; and Areal, N. 2013. Some experiments on modeling stock market behavior using investor sentiment analysis and posting volume from twitter. In *International Conference on Web Intelligence, Mining and Semantics*, 31. ACM.

[Pai and Lin 2005] Pai, P.-F., and Lin, C.-S. 2005. A hybrid arima and support vector machines model in stock price forecasting. *Omega* 33(6):497–505.

[Paul, Dredze, and Broniatowski 2014] Paul, M. J.; Dredze, M.; and Broniatowski, D. 2014. Twitter improves influenza forecasting. *PLoS currents* 6.

[Phillips and Gorse 2017] Phillips, R. C., and Gorse, D. 2017. Predicting cryptocurrency price bubbles using social media data and epidemic modelling. In *SSCI*, 1–7. IEEE.

[Porshnev, Redkin, and Shevchenko 2013] Porshnev, A.; Redkin, I.; and Shevchenko, A. 2013. Machine learning in prediction of stock market indicators based on historical data and data from twitter sentiment analysis. In *ICDMW*, 440–444. IEEE.

[Rao and Srivastava 2012] Rao, T., and Srivastava, S. 2012. Analyzing stock market movements using twitter sentiment analysis. In *ASONAM*, 119–123. IEEE/ACM.

[Sang and Bos 2012] Sang, E. T. K., and Bos, J. 2012. Predicting the 2011 dutch senate election results with twitter. In *Workshop on semantic analysis in social media*, 53–60. ACL.

[Székely and Rizzo 2013] Székely, G. J., and Rizzo, M. L. 2013. Energy statistics: A class of statistics based on distances. *Journal of statistical planning and inference* 143(8):1249–1272.

[Székely et al. 2007] Székely, G. J.; Rizzo, M. L.; Bakirov, N. K.; et al. 2007. Measuring and testing dependence by correlation of distances. *The annals of statistics* 35(6):2769–2794.

[Szekely, Rizzo, and others 2014] Szekely, G. J.; Rizzo, M. L.; et al. 2014. Partial distance correlation with methods for dissimilarities. *The Annals of Statistics* 42(6):2382–2412.

[Tang, Yeh, and Lee 2014] Tang, W.-H.; Yeh, M.-Y.; and Lee, A. J. 2014. Information diffusion among users on facebook fan pages over time: Its impact on movie box office. In *International Conference on Data Science and Advanced Analytics (DSAA)*, 340–346. IEEE.

[Tetlock 2007] Tetlock, P. C. 2007. Giving content to investor sentiment: The role of media in the stock market. *The Journal of finance* 62(3):1139–1168.

[Town 2018] Town, S. 2018. How to spot a pump and dump (and avoid it). *Investopedia*.

[Tumasjan et al. 2010] Tumasjan, A.; Sprenger, T. O.; Sandner, P. G.; and Welpe, I. M. 2010. Predicting elections with twitter: What 140 characters reveal about political sentiment. In *ICWSM*, volume 10, 178–185. AAAI.

[Volkova et al. 2017] Volkova, S.; Ayton, E.; Porterfield, K.; and Corley, C. D. 2017. Forecasting influenza-like illness dynamics for military populations using neural networks and social media. *PLOS ONE* 12(12):1–22.

[Wang and Lei 2016] Wang, M.-H., and Lei, C.-L. 2016. Boosting election prediction accuracy by crowd wisdom on social forums. In *CCNC*, 348–353. IEEE.

[Wang and Vergne 2017] Wang, S., and Vergne, J.-P. 2017. Buzz factor or innovation potential: What explains cryptocurrencies returns? *PloS one* 12(1):e0169556.

[Zhang 2003] Zhang, G. P. 2003. Time series forecasting using a hybrid arima and neural network model. *Neurocomputing* 50:159–175.

[Zhao et al. 2016] Zhao, S.; Tong, Y.; Liu, X.; and Tan, S. 2016. Correlating twitter with the stock market through non-gaussian svar. In *ICACI*, 257–264. IEEE.

[Zimbra, Chen, and Lusch 2015] Zimbra, D.; Chen, H.; and Lusch, R. F. 2015. Stakeholder analyses of firm-related web forums: Applications in stock return prediction. *ACM TMIS* 6(1):2.