Predicting Price Trends in Cryptocurrencies using News Headlines

Daniel Deng, Zain Hemani, and Reid Vender

March 7, 2018

Proposal

Computer Science 4411B

Databases II

Submitted in partial fulfillment of the requirements for the final project

Contents

1	Intr	Introduction		
2	Motivation			2
3	Pro	ject D	etails	4
	3.1	Archit	ecture and Environment	4
		3.1.1	Available Data	4
		3.1.2	Database Structure	5
	3.2	Delive	rables	7
		3.2.1	Features and Finished Product	7
		3.2.2	Project Contribution	8
4	Cor	clusio	n	8

1 Introduction

Unlike stocks, mutual funds, commodities and other speculative instruments, cryptocurrencies are highly volatile. This means price predictions are highly variable and can depend on a number of factors, including the public's perceptions of the perceived value of cryptocurrencies. This perceived value is determined by the market, mainly from the media and various news headlines.

Despite few similar applications that use news headlines to predict price trends, many of these applications do not store their collected information in a database. While these applications do many of their calculations in real-time, which may provide useful guidance for investors, storing this information in a database will also allow for regressions and other analyses to be conducted on the datasets to discover long term trends. Our solution involves the use of the NoSQL database MongoDB, R and Python to store and access our data, explore relationships between stocks and cryptocurrencies, and perform MapReduce to analyze keywords which positively or negatively affect cryptocurrency prices.

2 Motivation

When it comes to pricing predictions on cryptocurrencies, we see all sorts of claims: some people say bitcoin will reach \$16000, while other says it may skyrocket to \$100000 [1]. Unlike stocks, mutual funds, commodities, and other speculative instruments, cryptocurrencies are not regulated and are highly volatile, making such enormous differences in price predictions. This volatility is largely derived from the public's perception of the cryptocurrency and its value. This denotes a clear distinction between the price and value of a cryptocurrency. The price of any cryptocurrency is simply the monetary cost of purchasing it, whereas the value of a cryptocurrency is its perceived benefits and usefulness. The price of a cryptocurrency is not tied to its value but instead, its perceived value. This perceived value is determined in the market, mainly from the media and various news headlines. In many cases, news outlets are the only source of information for many investors in cryptocurrency. This is why correlation with news headlines and the prices of cryptocurrency is an important set of information that can be used to algorithmically and efficiently trade in this ever-growing market.

There is over 20 billion dollars traded each day in various cryptocurrencies around the world. This is a staggering amount which is enormously influenced by news outlets, such as the New York Times, twitter accounts such as John McAfee, and even YouTube videos by large investors. Studies show that tweets by individuals just as McAfee were reliably correlated with price spikes that sent cryptocurrencies worth pennies — even fractions of a penny — temporarily shooting upwards in value anywhere from 50 to 350 percent. For example, on December 15, McAfee tweeted that a coin called SAFEX constituted "the majority of [his] holdings," [2]. Minutes after McAfee tweeted about SAFEX, the price of SAFEX spiked 92 percent, according to the site CoinMarketCap, from \$0.014 USD to a high of \$0.027, before settling down for a more gradual rise up to its current price of \$0.03 [3]. Another example would be the news headlines related to Korean regulations on cryptocurrency exchanges. Many indexes simply took Korean exchange prices off their representations which caused numerous new outlets, such as the New York Times and Bloomberg to release articles [4, 5]. Although no regulation had actually taken place at the time, due to these articles the prices for cryptocurrencies fell significantly — in some cases by more than 75% — leading to a large market fluctuation [6]. As suggested by the above examples, news headlines provide an important insight into future prices for cryptocurrencies and give investors reliable information when looking to long and short their positions.

There are a few similar applications that use news headlines to predict price trends, however many of these applications do not store their collected information in a database. This is mainly done through natural language processing in tokenizing different segments of a headline or tweet and understanding whether the fluctuation will be positive or negative depending on the statement [7]. These applications do many of their calculations in real-time, at the moment the tweet or news headline is received. This does provide useful guidance for investors, however storing this information in a database, as our application plans to do, will not only allow for realtime inferences but also allow for regressions and other analyses to be conducted on the datasets over a longer period of time. This will help determine long term trends and also differentiate between the extent of effect each news or social media outlet has on prices. Our solution has potential benefits for both individual investors and larger institutional investors. News and social media correlations are already widely used in the industry, but looking at this information from a macro lens will allow investors to correctly ascertain which source is more effective and the potential price fluctuation that may occur from the release. This will certainly provide a useful tool to conduct market analysis.

3 Project Details

3.1 Architecture and Environment

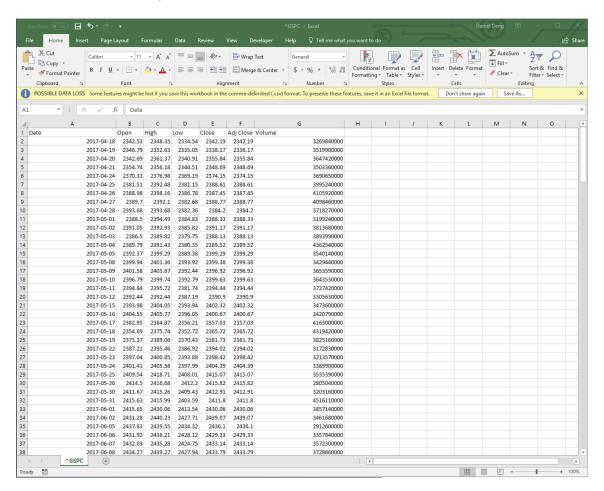
Our project will require datasets from Kaggle and Yahoo Finance. The datasets will come in the form of Excel spreadsheets and they will require cleaning using a python script or built-in functions in Excel. We will be using the NoSQL database MongoDB to store and access our data. Excel and potentially R and Python will be used for our data analysis. Some useful libraries in R include caret, ggplot, dplyr and lattice. PyMongo, NumPy, SciPy, Pandas and Plotly will be essential libraries if we choose to use Python. The hardware that will be used for this project is a standard Windows Laptop with 8Gb of Ram and an Intel i7-7500U processor.

3.1.1 Available Data

The cryptocurrency dataset from Yahoo Finance contains stock opening, high, low and close prices, as well as optional volume, adjusted close and market cap available for dates from 2013 to 2017. "Daily News for Stock Market Prediction" is a dataset from Kaggle which records the top 25 news headlines from /r/worldnews (ranked by Reddit users' votes) from 2008 to 2016 [8]. Together, with over 11,500 stock and cryptocurrency prices and 76,600 news headlines spanning 3,064 days, querying the database could return over 35,000,000 data points.

3.1.2 Database Structure

The figure below shows an example of raw data prices for the S&P (a well-known index that tracks overall market performance). As we can see, there are many different fields but *Date*, *Open* and *Close* will be the most important ones to focus on when we do our analysis.



In order to store the previously described available data, we will first create a single database. Next, we will populate a single collection with an embedded structure:

```
date: <date>,
  headlines: <array>,
  prices: [
    {
      name: <string>,
      open: <double>,
      high: <double>,
      low: <double>,
      close: <double>,
      volume: <long>,
      adj_close: <double>,
      market_cap: <long>
    },
    . . .
 ]
}
```

3.2 Deliverables

3.2.1 Features and Finished Product

The project will produce a report on our findings. Details of this report are outlined in the deliverables below. Our project will consist of the following:

Deliverable	Details
Deliverable 1	Find appropriate datasets and clean the data so that it is
	consistent and usable.
Deliverable 2	Import the dataset into a NoSQL database and perform relevant
	queries to extract data for analysis.
Deliverable 3	Perform MapReduce on the world news data to extract and count
	keywords present in the news headlines.
Deliverable 4	Analyze data for:
	• Correlations in stock and cryptocurrency prices over time
	Keywords in world news headlines which affect stock and
	cryptocurrency prices
	• The stock market is open between 9:30-4pm whereas the
	cryptocurrency market is open $24/7$; the time when the crypto-
	currency market is open but the stock market is not is an area
	of interest
Deliverable 5	(Optional) Build a predictive model using a decision tree like
	random forest or xgb to forecast future direction/prices based on
	keywords in news headlines.
Deliverable 6	(Optional) Build a front-end to showcase our findings interactively
	for the user.
Deliverable 7	Write a report on the findings of our data analysis and its relevance
	to the industry, and any future improvements that could be made.

The overarching goal of our project is to determine the relationships, if any, between cryptocurrency prices and a benchmark index like the S&P, and keywords present in news headlines. Finding a correlation between these datasets can open the doors to a myriad of new trading strategies. Not only are we striving for monetary gain with this project, we are hoping our research will help the average investor improve their investment decisions. Some personal achievements undertaking such a project include getting hands-on experience with setting up and accessing data from a NoSQL database and learning basic data analytics skills.

3.2.2 Project Contribution

As indicated, some of these deliverables are optional but this project is of great interest to us. As investors ourselves, this project can help us explore trading strategies we have never tried before while improving our technical skills. Our unique, database-oriented solution will also allow for regressions and other analyses to be conducted on the datasets to discover long term trends which we hope will be useful guidance for investors. We will be continuing to work on this assignment in the future and seeing it through to the end.

4 Conclusion

Enormous differences in price predictions are, in large part, due to the high volatility of cryptocurrencies. This volatility is largely derived from the public's perception of the cryptocurrency and its value. This distinction between the price and value of a cryptocurrency may be due to news headlines which shape the public's perception. In order to explore this, we plan to design a NoSQL database using MongoDB to store and access our data, R and Python to analyze relationships between stocks and cryptocurrencies, and perform MapReduce to analyze keywords which positively or negatively affect cryptocurrency prices. This is a different approach from similar applications that use news headlines to predict price trends in real-time without the use of databases. This useful tool to conduct market analysis will help us explore trading strategies we have never tried before while improving our technical skills and has potential benefits for both individual investors and larger institutional investors.

References

- [1] A. "Bitcoin Here's Morris, prediction 2018: where investors THIS YEAR," February, 2018. see bitcoin going [Online]. Availhttps://www.express.co.uk/finance/city/913707/Bitcoin-prediction-2018price-crash-bubble-Ripple-ethereum-cryptocurrency-blockchain. [Accessed Feb. 23, 2018].
- [2] J. McAfee, December, 2017. [Online]. Available: https://twitter.com/officialmcafee
- [3] CoinMarketCap, "Safe Exchange Coin," December, 2017. [Online]. Available: https://coinmarketcap.com/currencies/safe-exchange-coin/
- [4] K. Cho, "Why the Cryptocurrency World Is Watching South Korea," February, 2018. [Online]. Available: https://www.bloomberg.com/news/articles/2018-02-04/why-the-cryptocurrency-world-is-watching-south-korea-quicktake. [Accessed Feb. 23, 2018].
- [5] REUTERS, "South Korea Keeps Investors Guessing on Cryptocurrency Regulation," February, 2018. [Online]. Available: https://www.nytimes.com/reuters/2018/02/27/business/27reuters-cryptocurrencies-southkorea.html. [Accessed Feb. 27, 2018].
- "South Flags [6] REUTERS. Korea Regulator Better Deal for Industry," February, 2018. [Online]. Avail-Cryptocurrency https://www.nytimes.com/reuters/2018/02/20/business/20reuterssouthkorea-cryptocurrencies.html. [Accessed Feb. 23, 2018].
- [7] E. R. Connor Lamon, Eric Nielsen, "Cryptocurrency Price Prediction Using News and Social Media Sentiment." [Online]. Available: http://cs229.stanford.edu/proj2017/final-posters/5138197.pdf. [Accessed Feb. 23, 2018].
- [8] Aaron7sun, "Daily news for stock market prediction," March, 2016. [Online]. Available: https://www.kaggle.com/aaron7sun/stocknews. [Accessed Feb. 25, 2018].