# WHAT MAKES TWITTER INFLUENTIAL?

An analysis of Twitter content, user behavior, and correlation to Bitcoin market value.

NEIL OLIVER

THE NEW SCHOOL | PARSONS SCHOOL OF DESIGN
Submitted in partial fulfillment of the requirements for the degree of Master Science in Data Visualization at Parsons School of Design.

# TABLE OF CONTENTS

# 1  ABSTRACT

Bitcoin's volatile price fluctuations can be linked to its perceived future value, potential use cases & concerns over longevity and legislation. Public discussions on social media can serve as a platform to gain information on Bitcoin from a wide audience, from amateurs to professionals, therefore influencing public opinion and Bitcoin trading activity.

This project dissects the user behavior and content of Twitter to investigate which elements have the strongest correlation to changes in Bitcoin market value. Evaluating a body of 100,000 tweets over a period of one week, volume data of both tweet generation and user interactions is combined with sentiment analysis to create a series of interactive exploratory visualizations.

The visualizations present comparisons of different methods of user interactions, including retweets, likes and comments and allows filtering to remove data that may be deemed uninfluential before aggregated sentiment analysis takes place. During sentiment analysis, weighted models allow for the balance to be shifted to give preference to the opinions of users with a high follower count, or messages with a high level of interaction. As a final stage, time delay between online discussions and Bitcoin trading activity is explored and accounted.

Rather than provide definitive answers, the project invites the user to explore the vast array of filtering, weighting and offset models to form their own conclusions and observations.

# 2  INTRODUCTION

In 2019 US digital advertising revenue surpassed $100 billion (Winterberry Group 2020). In the same year, nearly one in five digital marketing agencies spent over half of their annual budget on social media influencers (SocialPubli 2019). Celebrity endorsements can cost upwards of $1 million dollars per post, but companies can expect a five-fold average return on investment (Influencer Marketing Hub 2020) including Twitter, where users report a 5.2X increase in purchase intent when exposed to influencer marketing (Annalect 2016).

Influencer marketing is not restricted to consumer products and its successful use can be observed across many industries including finance. During 2017, previous antivirus inventor and CEO, John McAfee was famed for the 'McAfee Effect' in relation to the influence of his tweets appeared to have over cryptocurrency (a new form of digital currency that is traded in a similar manner to traditional stocks) value (Toshi Times 2018). While the effect did not sustain over time, the short-term positive influence can be witnessed on several cryptocurrencies with MacAfee's backing. Twitter based support from McAfee could reportedly be purchased for the cost of 25 bitcoins (Twitter 2017); a USD value at the time of more than $200,000.

The high cost of a single cryptocurrency related tweet by an influencer is understandable when put into the context of the cryptocurrency market. With over 5000 different cryptocurrencies, and a total market capitalization of more than $200 billion (reaching a peak of over $795 billion in January 2018), cryptocurrency is a popular and growing technology. The leading (and first ever) cryptocurrency, Bitcoin, has a market share of over 60% of all combined cryptocurrency values. It's huge popularity and market dominance has resulted in all other cryptocurrencies being known as altcoins. Bitcoin's dominance also has effects on the general cryptocurrency market trend, where the direction and movement of Bitcoin value can be observed in the value of the other altcoins.

Existing studies have shown a clear relationship between social media sentiment and financial markets (Phillips and Gorse 2018) including the successful implementation of trading algorithms (Garcia and Schweitzer 2015). This project builds on the aforementioned studies to investigate which specific elements of social media content, and user behavior has the greatest correlation with cryptocurrency value. This is achieved through a focused case study of one social media platform and one cryptocurrency coin.

While more in depth online conversations surrounding cryptocurrencies happen on social media sites such as Reddit (Phillips and Gorse 2018), Twitters additional features of verified status', retweets, likes and comments offer a wider variety of variables that may affect influence. In addition to these features, the ease, speed, volume and length of messages (33 characters average, 280-character limit) reverse chronological order, and public profile focused approach makes Twitter the ideal planform to observe rapid sentiment changes surround cryptocurrency.

This project tracks the messages of Twitter uses discussing Bitcoin (selected as the leading cryptocurrency coin), using hashtags and key terms over a period of one month. Multiple sentiment analysis methods are applied, taking into consideration multiple filtering methods, to only consider users of a particular user groups (for example only those users with verified status) and consider weighted models where different user groups contribute more to the overall sentiment score.

The sentiment analysis is then compared to the changes in Bitcoin value over the same time period, calculating a correlation coefficient and considering time lag to account for any delay between sentiment and the market value.

The resulting visualization demonstrates the effect of each element of Twitter in relation to a filtered or weighted sentiment analysis model. The variations in sentiment score are displayed in relation to Bitcoin market value to show which combination of filters and weights have the greatest correlation with Bitcoin market trend. Each stage of the visualization is interactive to allow the user to explore the full range of possibilities. Due to the high volatility of Bitcoin price, the visualization is not designed to be used as a

prediction tool but is an in-depth exploration of the elements of Twitter have the greatest impact on influencing user behavior.

# 3  BACKGROUND

## 3.1  THE BOOM OF SOCIAL MEDIA INFLUENCERS

Since 2017, influencer marketing has doubled in size to over $6 Billion per year (Influencer Marketing Hub 2020). The market has shifted from product placements from celebrities and public figures who found popularity from outside of social media, to now include people who have grown a follower base for the purpose of generating income through online promotion. The shift is so much that audiences are more receptive of social media influencers than traditional celebrities (Mediakix 2019). One of the most successful influencers, Kylie Jenner, has over 165 million followers on Instagram. The image centric social media platform allows users a preview of her lavish lifestyle, luxury products and exotic holiday destinations. For anyone that dreams of sharing any piece of that lifestyle, the regular additional of paid product placement guide the way to find it.

The use of social media influencers has a strength above traditional digital advertisement; the advertisements are not totally unsolicited; for the majority of cases the user has made the conscious effort to follow or view the influencers profile. As people look to Jenner for lifestyle aspirations, other social media users look for guidance and advice on a wide variety topic. These topics could range from inspirational quotes to health advice or financial planning. Many of these topics do not need to visual centered focus of platforms such as Instagram and are more suited to text based platforms that allow the addition of links and attached media (something that Instagram does not allow). Twitter, with 330 million active monthly users (Report 2019) is a clear example that there is still a desire for such a primary text-based platform.

## 3.2  DIFFERENT PLATFORMS WITH DIFFERENT AFFORDANCES

In comparison to Facebook's 2.45 Billion monthly active users, Twitter may at first appear to be largely insignificant in relation to digital marketing, however the two platforms have a different focus.

While Facebook promotes user privacy (now a default setting) and sharing between friends and family, only 13% of Twitter users have their profile set to private. Twitter has a greater public focus, where the majority of users share posts openly to the any interested party. Links between Twitter users can also be a one-way link. Unlike Facebook's 'friendship' system, a Twitter user can follow any public profile without any requirement for the other user's prior agreement and without the other user reciprocating. This results in celebrities such as Donald Trump gaining over 70 million followers. In comparison to Facebook, the links between users of Twitter have a much broader reach,

with the average twitter user having 707 followers (Kick Factory 2016), vs the average Facebook users 338 friends (Smith 2014).

Another element of Twitter that makes it desirable to marketing and influencers is that the information that users see is not restricted to the individual posts (tweets) from users they follow. The process of 'retweeting' a message is to post another users message onto your profile for your followers to see, therefore spreading a message to a larger audience (discussed further in section 5.4.2). Messages that resonate with a large enough audience can quickly turn 'viral' and receive millions of retweets. Quantifying what will result in a viral message is difficult (Cha et al. 2010), however they can spawn from any account regardless of the number of followers and can resonate for positive or negative reasons. Retweets allow for semi-unsolicited (the user has chosen to follow the user who retweeted but may not have directly followed the original creator of the post) messages to reach a user's profile. This aspect of Twitters platform makes it possible for users to be influenced by messages from users they do not directly follow.

## 3.3   THE VOICE OF A MILLION TWEETS

During the 2016 presidential election campaign, the affordances of the Twitter platform made it ideal for manipulation in the hopes to influence users voting intentions.

*"Highly automated accounts—the accounts that tweeted 450 or more times with a related hashtag and user mention during the data collection period— generated close to 18 percent of all Twitter traffic about the Presidential election"*

(KOLLANYI, HOWARD, AND WOOLLEY 2016)

With an average of 1300 tweets per account per day and messages with pro Trump hashtags reaching over 150,000 each day, automated ('bot') accounts were used heavily to spread messages throughout the social media platform. Without speculating on the impact these messages, this suggests that influence through social media may not only come from a few influential users with high follower numbers, but also from a large body of accounts discussing a similar topic.

There are many other aspects to Twitters platform that can impact the level of influence. User interactions with tweets can take several forms, including retweets, comments, and favoriting. The profile of users not only vary in the number of followers, but also in the number of users they follow back, the frequency that they tweet or if they have been granted a verified status by Twitter (all discussed in detail in section 0). It is not a clear case that user with the most followers have the most influence.

*"...popular users who have high indegree are not necessarily influential in terms of spawning retweets or mentions."*

(CHA ET AL. 2010)

## 3.4 MEASURING INFLUENCE

With the large returns on product placements by influencers and the extreme effort and money spent on digital political campaigns, there is little doubt that Twitter as a body of users can have a large amount of influence, not only over individual buying or voting patterns, but on a scale where it can impact larger markets and events.

When evaluating the sentiment of Twitter discussions surround the stock market, specifically the Dow Jones Industrial Average (DJIA), there was an 86% accuracy in predicting the daily rise and falls in value (Bollen, Mao, and Zeng 2011). Twitter sentiment analysis has also been used to create an algorithmic cryptocurrency trading bot where they confirm "the long-standing hypothesis that trading-based social media sentiment has the potential to yield positive returns on investment" (Garcia and Schweitzer 2015).

When considering the influence of social media in these studies, both the DJIA and cryptocurrency have additional factors that affect their price fluctuations. Each company in the DJIA has physical assets and balance sheets to be taken into consideration for company valuation; cryptocurrency in comparison is not tied to a physical entity. Cryptocurrency is however affected by the current energy price due to the large power consumption that underpins how the technology functions, and can also be affected by discussion of legislation.

In March 2017 the Security and Exchanges Commission (SEC) disapproved a proposal for a Bitcoin (the largest valued cryptocurrency coin) ETF (investment vehicle), citing price volatility and the price being driven by speculation as two concerns (Security and Exchanges Commission 2017). With no inherent value, the largest cause of price changes in cryptocurrency is market supply and demand, driven by speculation and discussions, that happen in part in the public domain via social media websites. This makes cryptocurrency a topic where the effects of social media influence can be witnessed and have a significant impact.

# 4 BITCOIN

## 4.1 BACKGROUND

Originating in 2008, Bitcoin was the first cryptocurrency; a digital asset that was designed as a medium of exchange that is not tied to a central authority. The currency, built on a technology called blockchain, gained popularity in its early stages due to its ability to keep the users anonymous, as payments are tied to a digital address (free for anyone to obtain) and do not require any human identification process. Bitcoin works on a public distributed leger (everyone can view every transaction ever made) and transactions are verified by bitcoin miners (to stop potential fraud or double spending), who are rewarded for their work with new bitcoins. As bitcoin is not a physical object, such as a coin or bank note, readers who are unfamiliar with its concept can consider each coin to be similar to the serial number on every banknote, where everyone can see who (which digital address) owns each number.

Already this description is an oversimplification containing cryptocurrency specific technical language. While the technical details of the implementation of Bitcoin and other cryptocurrencies are not important within the scope of this project, we must understand some principals that affect the value of Bitcoin (and all other cryptocurrencies) in order to consider the affect that social media plays in influencing a change in value.

## 4.2 ELEMENTS AFFECTING VALUE

### 4.2.1 SCARCITY

How scarce an item is can affect its value. A clear example of this is the value of precious metals and gems such as silver, gold and diamonds. While these materials are still being mined and new gems are constantly coming into the market, the flow of new material is slow, and the items are still considered rare. This helps keep their value high due to supply and demand.

This same process is true of Bitcoin. As the popularity of Bitcoin grows, more bitcoins are released as part of the mining process. While this increase in volume of available coins should lower the price, the introduction of new coins is very slow and on a decreasing rate (due to a technical implementation known as block reward halving). Meanwhile, the popularity and adoption of Bitcoin is increasing at a growing rate.

### 4.2.2 MINING & ENERGY PRICES

The process of Bitcoin mining can be compared to guessing a very complex number (a process called proof of work). The complexity of this number is based on previous transactions, and how many transactions there currently are to be verified. The current probability of guessing the correct number is around 1 in 15 trillion. The current level of computation needed is very high and takes dedicated computer setups which use a large amount of energy. This high energy consumption is a factor that must be calculated into

a Bitcoin miner's profit margin, therefore changes in energy price can attract or detract people from mining, which is an underpinning element of the cryptocurrency ecosystem.

The amount of energy used is also a concern on a global scale, in 2018 the energy consumption of the Bitcoin network (34.86 TWh) was more than the whole of Denmark (33TWh) (Digiconomist 2020). This has caused questions over Bitcoins ongoing viability and potential regulation, which is also an area of contention that causes fluctuations in Bitcoin value (Congressional Research Service 2019).

### 4.2.3 ADOPTION & USE CASES

2011 saw the introduction of alternative cryptocurrencies (known as altcoins) that attempted to address some of the limitations and pitfalls of Bitcoin. Other altcoins were created with specific use cases, such as Ripple, a currency designed to improve the speed of interbank monetary transfers. As of 2020 there are over 5000 different cryptocurrencies (CoinMarketCap 2020).

The wide variety of possible use cases for the blockchain and cryptocurrency technology resulted in some investment firms holding cryptocurrency as part of their investment portfolio, in the hope that the technology would find a mainstream, regulated use.

This interest was not only limited to investment firms. The ease of purchase and lack of regulation also saw a growing interest from the general public. Bitcoin's lead on the competition due to being first to market and the most easily accessible through multiple markets and exchanges (discussed in section 4.4.1) made it a popular choice for first time cryptocurrency investors. Bitcoin has remained the cryptocurrency leader, responsible for over 50% of market capitalization (market cap) (CoinMarketCap 2020). By 2017 interest had grown substantially where due to supply and demand the value of a single coin grew from $900 to over $20,000.

## 4.3  HISTORICAL EVENTS

Many of the aspects that make cryptocurrency unique, are the same aspects that can be a cause for concern when considering its regulation and use. The SEC's 2017 denial of regulation due to concerns of speculation driving Bitcoin value can be clearly observed through Bitcoin's price volatility. An example of this is October 25[th] 2019, where the Bitcoin value grew by over 40% where only two days earlier it reached a 5 month low.

### 4.3.1  2018 CRASH

The largest example of Bitcoin price volatility is January 2018 where between January 6[th] and February 6[th], the value of Bitcoin dropped 65%. This was not an unpredicted event, with many analysts predicting the fall due to rise of cryptocurrency popularity being similar to the dot com bubble. With no intrinsic value and the Chicago Board Options Exchange allowing the ability to 'short' (a financial bet against a market) Bitcoin only a month prior; the demand for bitcoin dropped dramatically starting a chain reaction. The

Bitcoin bubble was driven by emotion, fear of missing out and the hope of making a fast and substantial profit.

As part of the huge growth in popularity of Bitcoin, new terminology was born including 'to the moon' and HODL, originally presumed to be a misspelling of hold, now widely accepted to mean 'hold on for dear life' due to cryptocurrency price volatility.

### 4.3.2 COVID-19

In March 2020 the worldwide pandemic COVID-19 caused the stock market to drop over 25%. In the same timeframe Bitcoin value dropped 55%. While this is another clear example of the volatility of cryptocurrency value, it is also an important demonstration of the link between bullish and bearish stock market positions and their effect on cryptocurrency value.

## 4.4  TRADING

The buying and selling of cryptocurrency can be broken down further into the initial purchase of a cryptocurrency from a fiat currency (a currency where the value is backed by a government) and trading between various cryptocurrencies. In order for either of these events to happen, a cryptocurrency exchange must be used in a similar process to using foreign exchange for fiat currencies.

### 4.4.1 EXCHANGES

Multiple cryptocurrency exchange websites are available, each offering different availability to buy cryptocurrency from different fiat currencies and the ability to exchange between different cryptocurrency pairs. Not all cryptocurrency pairs can be directly traded, especially between altcoins, and Bitcoin is often used as an intermediary currency. The price of each cryptocurrency will vary slightly between exchanges but not by a large amount. It is important to understand the need and role of the exchange and variations between exchange for this project, as a single exchange (Binance) is used to provide the valuation data for Bitcoin. Binance has been chosen for its popularity, free and well documented API and its ability to offer expansion to a wide variety of altcoins in the future if the project is expanded.

Additional services such as CoinMarketCap are available that can offer average pricing across multiple exchanges, however, this is a paid services and does not offer a significant benefit to the project considering the low variation in price between exchanges.

### 4.4.2 EXPRESSING VALUE

While Bitcoin price can be expressed as a single monetary value (usually in USD) at any single moment in time, there are several values that can be used to express its current market position, strength and volatility. Many of these values are similar to trading traditional stocks and they will all play a role in determining correlation between Bitcoin value and Twitter influence within this project.

### 4.4.2.1 Candlestick Values

When using a trading platform, a user will view the changes in values between a pair of currencies. For this project we will use the BTC/USD pair (Bitcoin value in relation to the United States Dollar). The user is presented with time series data called candlesticks. A candlestick is similar to a box and whisker diagram in appearance but do not represent the same values. In a candlestick chart, each candlestick represents a user specified time interval. The selectable values are usually between 1 minute, through to 1 day, with intermediate points of 30 minutes & 1 hour. The top and bottom of the lines on the candlestick represent the highest and lowest points (respectively) the currency value reached during the time interval. The top and bottom of the 'box' represents the value of the currency at the start and close of the time interval, however the order depends on a value increase or decrease. If the value increased the bottom of the box presents the open value, if the price decreased the bottom of the box represents the close value. While this may seem complicated for the user to understand quickly, a green / red color coding is used on the box (or a fill / stroke in black and white publications) to show an increase or decrease.

### 4.4.2.2 Volume

In addition to the candlestick values, an additional bar chart is displayed below the main chart (on the same X axis) to show the relative volume of purchases and sales of the currency. This does not show direction on its own but can be used as a strong variable within this project to measure against Twitter volume (discussed in section 4.4.2.2).

### 4.4.2.3 Market Capitalization & Market Share

In addition to the value of a single coin, another valuation of the strength of Bitcoin is its market capitalization (market cap). This is the total value of all bitcoins in existence at the current market value. This is important for direct comparisons between different cryptocurrencies. While direct comparisons between Bitcoin and altcoins is not the focus of this project, Bitcoin's relative market share against all other coins is an indication of its growing strength or decline. This is a strong indicator that can be used to determine correlation to Twitter influence.


# 5  TWITTER

In order to move beyond existing studies that have previously compared the relationship between social media and cryptocurrency, we must break own Twitter into individual components. Each component will be evaluated on the role that is plays within the Twitter ecosystem and what potential it has to impact the level of influence, including a justification for its inclusion or dismissal in this project. This is a crucial step in order to not overlook any specific detail that may potentially, through its inclusion in the project, end up becoming a key in identifying a user behavior that has the strongest link with influencing Bitcoin value.

## 5.1 OVERVIEW

Twitter, similar to most social media platforms, can be used in two main different way, either to consume information or to create and share information. Often users interchange and combine these uses within one session. These different uses each have their own dedicated sections within the platform, with overlap when viewing other users' interactions on your own content.

## 5.2 THE FEED

When acting as a consumer, the main view (and default view when logging into the platform) is the Twitter feed which focuses on the content of other users of the website. The feed is a reverse-chronological display of all of the tweets created by other users that the current user has chosen to follow. In 2016, a significant change took place that affected the way that users receive their content. On the initial load of the page (after a period of not viewing the website), the feed no longer displays tweets in a reverse-chronological order, instead displaying tweets according to Twitters 'relevance model' that shows 'top-ranked' tweets first, with subsequent tweets being in reverse-chronological order (Twitter 2016). The thought process behind this implantation was to drive user engagement through identifying relevant tweets they may have missed (Koumchatzky 2017). In 2018 Twitter users were given the option to control the order of their feed, switching between the two ordering approaches (Twitter 2016).

### 5.2.1 FOLLOWERS

A list of who is following a user is publicly available via the profile page. Every time a user tweets, each of these followers will receive the tweet via their feed (discussed in section 5.2). The likelihood of these tweets being seen by all of their followers is determined by the amount of time each user spends on viewing their feed and the number of users they have chosen to follow. With the average user following 707 people (MacCarthy 2016), a large volume of tweets are queued in each users twitter feed at any moment in time. Combined with the different approaches to feed ordering as discussed in section 5.2, it is difficult to determine the exact reach and consumption of each individual tweet using public information, however this information is available through the engagement API (section 5.4.4). While the use of the engagement API is beyond the scope of this project, the number of followers of each user is still be a relevant and useful indicator for assessing the potential level of influence.

### 5.2.2 FOLLOWING

While the volume of followers should be a clear indication of the number of users interested in the tweets of a user, this is not always the case. The desire that many users have to feel popular has resulted in two user behaviors. The first behavior is that some users may try to manipulate their follower volume through the purchasing of followers from external agencies (Nicholas Confessore 2018). These followers are usually blank

accounts or bot-based accounts that can be purchased in thousands at a time. The second is a phenomenon known as like for like (or in relation to Twitter, follow for follow) and academically is discussed as part of the million followers fallacy (Cha et al. 2010). The premise of the act is that a user will follow another user in the hope or agreement that the other user will reciprocate.

In the same respect as the list of followers, the list of everyone that a user follows is also publicly available. While this does not have an effect on who sees a user's tweets, it can be used in conjunction with the number of followers to create a ratio of followers to following. A higher ratio may be an indication that the user has not gained a higher follower count from the follow for follow phenomenon. Further uses for the following to follower ratio are discussed in section 6.2.6.

### 5.2.3 TWEET VOLUME

Irregular or new Twitter uses can be a red flag when calculating the sentiment of a body of tweets, however the same can be observed in accounts with a high volume of tweets. Twitter accounts are free and quick to setup and tweets can be placed through a number of different methods, including the web interface, mobile application, desktop applications, SMS and programmatically through the use of the Twitter API. This allows for tweets to be posted easily, from a wide range of unrestrictive locations on a regular basis.

A high daily number of tweets can be a sign of an enthusiastic user, but can also be an indicator that the messages are from a computer generated (bot) source (Kollanyi, Howard, and Woolley 2016). Both the total number of tweets and the age of the account are available on the user profile page. These two elements can be combined to create a daily average which can be used as an indicator to the different types of users and accounts (discussed further in section 6.2.1).

### 5.2.4 VERIFIED USER STATUS

One method of identifying users that are not bots and whose popularity is likely not to be due to manipulation either through purchased followers or 'follow for follow' is through the verified user status. This status is given to accounts of celebrities, public figures and influencers through an application directly to Twitter. The verification is not automatic and is investigated and verified through a human process. While these accounts may avoid some of the previous mentioned concerns, less than 1% of active monthly users have a verified status, so are not a representative sample of the twitter population. They may however be useful in a weighted model when investigating influence as these users have usually gained their verified status due to their position in society.

## 5.3 TWEETS

In addition to the sentiment of the written text in each tweet, there are other characteristics that may have an impact on the level of influence it has.

### 5.3.1 LINKS

Links can be included in the body of a tweet to an external source (internal sources are dealt with through a number of different interactions discussed in section 5.4). Their inclusion may be an indication of spam or links to promote external sources. As an argument against this, a link may indicate an especially credible message that is supported by sources links. The evaluation of the credibility of individual links is beyond the scope of this project, however the inclusion of links as a filter will be included.

### 5.3.2 MEDIA

Additional media items such as images, animations and videos can also be included in tweets (originally as linked content and now embedded within the tweet). In a similar respect to links, evaluating the quality of the included media would be difficult within the scope of this project, however the inclusion (or absence) of media items can be used as a filter property.

## 5.4 INTERACTIONS

Interactions refer to any way in which a user can interact with another user's tweet. The process of interaction demonstrates user engagement with the content, therefore displaying some level of influence over the user.

### 5.4.1 LIKES

A like is a single click interaction. The process of liking a tweet shows an appreciation for the post by incrementing a counter displayed at the bottom of the tweet (denoted by a heart symbol). Originally intended (and still referenced within the API) as a favorite button, liked tweets can easily be revisited by the user (and viewed by others) via a link on the users own profile page.

### 5.4.2 RETWEETS

A retweet is another single click interaction. Retweeting places a copy of the original tweet on the users own timeline and in the feed of the people who follow them. The tweet is displayed in the same format as if the user has followed the original content creator. Retweets are important within this project as it immediately increasing the potential reach (and therefore potential influence) of the message. This feature allows the potential for a snowball effect and to create a 'viral' message if the message resonates with other users and they continue to retweet.

### 5.4.3 COMMENTS

The most time consuming of interactions is to reply, comment or quote (no differentiation is given to these terms within the Twitter platform) a Tweet. Due to the additional time cost of this interaction it is less popular than the other methods, however it shows a higher level of engagement (compared to likes or retweets) with the original tweet. Unlike retweets or likes, comments do not necessarily signify an agreement with a tweets content and may in fact be a form of rebuttal or argument against the content.

### 5.4.4 ENGAGEMENTS

While the total number of interactions discussed above are displayed at the footer of each tweet, additional engagements (a term used in Twitters API) are also possible. Each tweet can be clicked on to view the comments and who liked and retweeted the tweet. In addition to this, Twitter also stores how many times the tweet has been viewed (known as impressions), and how many clicks on all media (including links, hashtags and video plays). While these engagements would certainly provide valid variables to consider in influencing outcomes, this information is only available via the enterprise API at a large financial cost and not available through the public API. There are therefore deemed outside of the scope of this project.

# 6 SENTIMENT ANALYSIS

The process of sentiment analysis is to classify text into a series of emotional states, usually for the purpose of analyzing a body of responses. Sentiment analysis does not have fixed classes and is dependent on the subject of the information, however it will usually create a spectrum of results between negative, neutral and positive. Demonstrating correlation between Bitcoin value changes and accurately detected sentiment from Bitcoin related tweets will give the strongest evidence of twitters influence not only over Bitcoin change in price, but also the direction of the price movement.

## 6.1 METHODS

### 6.1.1 NATURAL LANGUAGE PROCESSING

Natural language processing (NLP) is a branch of Artificial Intelligence that focuses on making sense of human language. There are many different approaches to NLP with many employing supervised machine learning techniques. Understanding the specific details of these techniques is not required for the scope of this project however a basic understanding a few principles are needed to understand NLP's relevance and suitability for the project.

Supervised machine learning requires a training set of data. This is data that has been labeled in a particular way usually by a human driven process. In the context of this project, a training set would be a collection of thousands (or hundreds of thousands) of tweets that have been labeled by a human with their judgement of the sentiment of the message. This could be a Boolean (true or false) label or a numeric value. This labeled data then has features extracted from it that the computer uses to 'learn' what positive or negative sentiment looks like. A feature can be anything from the total number of words in each tweet, to the frequency of particular words, the use of punctuation, the list is very long.

Pretrained NLP applications are available and are very successful in accurately predicting sentiment, however the accuracy of these models is dependent on the strength and breadth of the training set.

### 6.1.2 LEXICON BASED

The most basic form on sentiment analysis is the use of a sentiment lexicon or sentiment dictionary. These can both be complied by humans or as a result of machine learning and natural language processing. They are lists of key value pairs that link a word with a sentiment or sentiment value. In its most basic form, a sentiment lexicon could be represented as two lists, one list of positive words and one list of negative words, or with each word being given a value or either +1 for positive or -1 for negative. To use such a list, each word in a body of text would be compared to the list and the score totaled.

In reality the process is more complex, with the words in the lexicon usually being assigned a wider range of values (than simply -1 or 1) as some words can be perceived to be more positive or negative than others. Summation of values is not always the best method of creating an overall score, so a comparative score is usually provided that creates an average sentiment value for each word in the body of text. To ensure an accurate score, the text must first be cleaned including the removal of stop words and lemmatization (discussed in section 6.2).

There are a wide variety of pre-existing sentiment lexicons that are free to use; the most popular being the AFINN-165 library that is available in multiple languages and contains 3382 words.

### 6.1.3 METHOD SELECTION

The performance of machine learning based NLP models on average is more accurate than lexicon-based approaches. It's downfall with relation to this project is the specific language used surround cryptocurrency (and to a larger extent, financial terminology in general). The sentiment of words used within general conversation may have a different sentiment (or strength of sentiment) when used in relation to discussing cryptocurrency. As an example, a term discussed earlier of 'to the moon' may be judged to have no specific (neutral) sentiment, however when discussing Bitcoin this has a strong positive sentiment. Other terms such as bullish or bearish will be rare to be seen in general conversation but have important connotations within financial conversations. Training a machine learning model on a training set of labelled tweets discussing cryptocurrency would provide the most accurate results. This however would require manually labelling a high volume of tweets, which is beyond the scope of this project.

As a compromise, combining both lexicon and NLP approaches has been successfully explored to combine the strengths and weaknesses of both methods. One popular existing system with a specific focus on social media sentiment is VADER (Valence Aware Dictionary and sEntiment Reasoner). While the results of VADER are strong, Lexicon based approaches may still have the ability to provide more accurate results by using

finance specific lexicons such as the one created in 2011 by Loughran and McDonald. This popular financial lexicon was further improved upon with a specific focus on social media analysis with the NTUSD-Fin Sentiment Dictionary (Chen, Huang, and Chen 2018). These different methods and models will be explored within the context of this project to find the which translates the best to identifying sentiment within cryptocurrency specific tweets.

## 6.2 DATA CLEANING

The selection of Twitter as the source for this project is due to its high data volume, userbase that spans a wide demographic (that matches closely to the same demographic that invests in Bitcoin) and ability to represent current trends. The broad scope of Twitter can also be problem area in contrast to other platforms such as Reddit that have separate subreddits for specific subject discussions. In order to successful analyses the sentiment of Twitter in relation to Bitcoin, the information must be filtered and cleaned to remove unrelated or programmatic (bot) content.

### 6.2.1 DAILY TWEET RATE

Detection of high-volume accounts which is an indication of bot accounts. This does not have to be a restrictive cut off as bot accounts usually have very high tweet volume. A limit of 50 tweets a day should ensure that a low number (if any) genuine twitter accounts are filtered out.

### 6.2.2 EXTERNAL LINKS

Links are not processed during sentiment analysis and therefore would not directly impact a sentiment score; however, they may be an indication of tweets that's are promoting external sources and not general discussion around the topic of bitcoin.

### 6.2.3 HASHTAGS

Hashtags are used within the project as one method (combined with a general text-based search) to find Bitcoin related tweets. Hashtags can also be used as a filter method to discard tweets that contain a high level of hashtags (showing no specific topic or an attempt to gain additional attention and promotion) or hashtags that related to

### 6.2.4 LEMMATIZATION

Lemmatization is the process of combining groups of inflections of words so they can be processed together as one. This is an important step in the use of sentiment lexicons to ensure that words are correctly identified and scored.

### 6.2.5 STOP WORDS

Stop words are commonly used words that do not add to the sentiment of the message. If they are not removed before lexicon-based sentiment analysis, then comparative (average based) scores will lower.

### 6.2.6 Follower to Following Ratio

The details of this metric have been discussed in section 5.2.2 and are hard to regulate as the ratio does not always indicate a potential issue. However, in a similar respect to the daily tweet limit, a high filter cut off of 1:10 (1 follower for every 10 accounts followed), would likely filter out accounts whose primary focus is to garner unwarranted attention.

## 6.3 Uncertainty

Comparative scores allow the comparison of sentiment between each individual tweets, however, they also create a single value for text that may contain mixed sentiment. The additional information output by most sentiment analysis tools that allow the full range of sentiment to be displayed should not be dismissed. An important element of this project will be to represent the uncertainty in the sentiment analysis and the spectrum of sentiment in the tweets. In addition to the range of sentiment within a single tweet, it is also important to represent the range of results when aggregating tweets into a time interval when they are compared as time series data against Bitcoin value.

# 7 Correlation

In order to detect if Twitter is having an influence over the Bitcoin value, a correlation must be observed in a quantifiable way. What differentiates this project from previous studies who have already observed a correlation between cryptocurrency price and social media signals is number of different variables being examined and their effect on the strength of the correlation.

## 7.1 Volume

As a baseline variable, both Twitter and Bitcoin can be quantified as a volume amount for a given period of time. Evaluating volume does not give an indication of movement (either price change direction or change in sentiment) but can indicate a level of general interest. Observing correlation by comparing volume is a useful baseline as the only additional variables are the candle size interval (i.e. 1 minute) and the selected period of time being observed (i.e. from January 1 2020 to January 31 2020). With no additional variables, a strong correlation can be used as an indicator to the strength and effectiveness of the other sentiment related correlations and potential patterns in time lagged correlation (discussed later).

## 7.2 Sentiment

The specific analysis and scoring of the sentiment of individual tweets is discussed in section 5, however additional options are available when evaluating a aggregation of tweets. For a correlation to be calculated, both the Bitcoin data and Twitter data must

be examined for the same interval, to achieve this multiple tweets and multiple sales must eb grouped together. This process is commonplace in regard to Bitcoin data and is provided in this format as candlestick data. Individual Tweets and their sentiment values must be aggregated manually. The simplest approach is to take average sentiment scores of each tweet and keeping track of the highest and lowest values to demonstrate a range of results. The resulting sentiment scores from Twitter aggregation can however change dramatically when considering the question *'is everyone equal and relevant?'*

## 7.2.1 FILTERING

In section 0 many different individual elements were discussed and their potential role in affecting the level of influence on Bitcoin price. Allowing live filtering by user input will enable the user of the project to see the effect that each of these has on both the overall Sentiment and how it impacts the correlation value. Filtering can be applied to number of followers, retweets, likes, comments and verified status. Filtering can also be extended to exclude tweets with a close to neutral sentiment value.

Filtering out some of the lowest values should help to remove additional noise that is masking the more influential messages. If we have already accounted for cleaning data that may be affected by social systems, such as follow for follow then it is reasonable to assume that accounts with a large follower counts are popular due to user choice meaning that other users wish to view their content, read their opinions and facts. A large number of followers may not directly correlate to a higher influence, but in combination with the number of retweets, will result in a larger reach than accounts with a smaller follower group and a lower number of retweets.

## 7.2.2 WEIGHTING

In addition (or as a replacement) to filtering, a weighted method can be employed. A weighted method allows for the potential that every person has the ability to influence, however it is likely not that everyone has the same level of influence.

It is easy to understand that an account with 1 million followers has a larger chance of influencing the Bitcoin value when compared against someone with 1 follower. A closer comparison in reality would be to compare the single user with 1 million followers, against 1 million accounts that only have a single follower each. The potential social reach is the same, but is the level of influence?

A weighted model adds a multiplier to the value of each tweet's sentiment value. This multiplier can be any of the numeric variables discussed in section 0. This approach allows the inclusion of all tweets (but still allows filtering if required) with tweets with the highest multiplier value (for instance the highest number of followers) to contribute to the overall score more.

In the same format as filtering, the user will be able to change the weightings and combine weightings to visualize the effect on the sentiment score and correlation to Bitcoin value. For this project we will use the weighted arithmetic mean formula below.

$$\bar{x} = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i},$$

[1]

## 7.3  CALCULATING CORRELATION

A calculated correlation value will be the basis of identifying if the selected body of Tweets (as a collection) are influencing the Bitcoin value. A positive correlation in relation to this project would mean that when the sentiment value is positive, the value of Bitcoin increases and vice versa. A negative correlation would mean that when the sentiment is negative, the price of Bitcoin increases (the variables move in opposite directions).

### 7.3.1  PEARSON CORRELATION COEFFICIENT

The Pearson Correlation Coefficient (Pearson's R) is a measure of linear correlation between two variables. The output of the Pearson's R formula is a value ranging from -1 (negative correlation) to 1 (positive correlation) with 0 representing no correlation.

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

[2]

Pearson's R is a popular method for calculating linear and directional correlation in time series data. The issue with its use in this project is that even if there is a perfect correlation (value of 1), it is unlikely that the increase and decrease in both sentiment and Bitcoin value will happen at the same time. For Twitter to have had an influence, people must react to the sentiment and act by buying or selling Bitcoin. The delay between changes will result in Pearson's R giving a low or null correlation result.

[1] https://en.wikipedia.org/wiki/Weighted_arithmetic_mean

[2] https://en.wikipedia.org/wiki/Pearson_correlation_coefficient

### 7.3.2 TIME LAGGED CROSS CORRELATION

A time lagged cross correlation function (CCF) allows for a correlation to be calculated at different lag (offset) amounts. When normalized, CCF will provide a Pearson's Correlation Coefficient and will provide a lag amount.

$$\rho_{XX}(t_1, t_2) = \frac{K_{XX}(t_1, t_2)}{\sigma_X(t_1)\sigma_X(t_2)} = \frac{E[(X_{t_1} - \mu_{t_1})\overline{(X_{t_2} - \mu_{t_2})}]}{\sigma_X(t_1)\sigma_X(t_2)}$$ [3]

A negative lag and positive coefficient would mean that a positive correlation is being observed where the direction of the Twitter Sentiment is being followed by the same direction in Bitcoin value. A positive lag value would mean that Twitter is reacting to Bitcoin value changes, which would not infer any influence.

### 7.3.3 VARIABLE SELECTION & MANIPULATION

Correlation can be calculated using the baseline data for both sentiment (a value between -1 and 1) and Bitcoin value (either open or close price for each interval). The issue with this approach is that each sentiment value (an aggregation a sentiment values over a set interval) is not related to the previous sentiment value, where the Bitcoin price value is a continuation from the previous candlestick.

If a percentage change is calculated for each candlestick, then it is no longer a continuation of price and can help to identify correlation between price and sentiment direction.

$$\text{Percentage change} = \frac{\Delta V}{V_1} = \frac{V_2 - V_1}{V_1} \times 100.$$ [4]

As an additional step both the sentiment and Bitcoin value can be converted in Binary or Categorical datatypes to represent a similar positive or negative sentiment and a rise or fall in value. This would not show if the strength of sentiment is related to the amount that Bitcoin has changed in value, however this is not necessary to show a correlation and infer a level of influence.

---

[3] https://en.wikipedia.org/wiki/Cross-correlation

[4] https://en.wikipedia.org/wiki/Relative_change_and_difference

# 8  IMPLEMENTATION

## 8.1  TARGET AUDIENCE

While the technical implementation and methodology behind the project are designed to be serious in nature and have sound scientific grounding, the project has not been designed to only cater for users from a technical background. In order to achieve the widest adoption of the project, it has been designed to be accessible by the widest possible audience, including those with little to no prior knowledge of either Bitcoin or Twitter, with an assumption that many will have not read this paper. In order to achieve this, background information will need to be presented on both Twitter and Bitcoin as well as explanations of key areas of investigation, such as weighted analysis and time lag.

## 8.2  DESIGN LANGUAGE

### 8.2.1  TITLE VARIATION

The majority of users will have more experience or knowledge of the role and use of social media and Twitter in comparison to their level of knowledge regarding Bitcoin. Introducing the topic from the standpoint of Twitter's influential powers requires an in-depth explanation of the suitability of Bitcoin as a platform to witness this influence, and therefore requires a detailed explanation of how Bitcoin functions in relation to trading, price variation & additional influencing factors. Changing the focus to initially discuss Bitcoin and the potential factors that impact market value, allows for less background information to be required for both Bitcoin and why Twitter is a suitable influencing factor to focus an investigation on.

As a result, the title of the project varies from the title of this written project, posing the question, 'Is Twitter Feeding Bitcoin?'. The variation aims to quickly introduce the Bitcoin topic, with the opening statement and initial explanation section both having a predominant focus on Bitcoin. This reversal of focus allows for the project to be introduced while assuming only a basic level of knowledge about social media platforms before explaining in detail the specific areas of Twitter that will be investigated and their potential role in affecting influence.

### 8.2.2  WRITTEN LANGUAGE

An informal language is used throughout the project, designed to speak directly to the reader. Within the opening statement, the user is prompted to consider when they may have already encountered Twitter and Bitcoin. As the project progresses, small questions are continually asked as prompts to reflect on current area of investigation and any findings the user may have discovered.

Technical topics such as sentiment analysis and weighted averages are not avoided, rather they are simplified to focus on the details that the user needs to understand to be able to gain insight into the visualizations. At each stage the user is always has an awareness of

if the explanations have been an oversimplification and to what degree they need to understand the technical manipulation and analysis. Each area of investigation can require several explanations. Each explanation is broken up between visualizations, with playful subheading to keep the user from feeling overwhelmed by the technical detail. Visual presentations are including to assist with explanations where possible.

### 8.2.3 VISUAL LANGUAGE

Visual representations of explanations are in place to support or replace where appropriate. With the large number of technical aspects that need explaining to the user, the supporting visuals aim to allow the user to gain an overview of the topics without a detailed read. This is especially important for users that do not have the time to be able to go through each section but would rather focus on one specific section. The visuals follow the same page positioning and formatting to create a consistent feel that is easily recognizable to the user.

#### 8.2.3.1 Color Schemes

In addition to the supporting visuals, two different color palettes have been used identify key features in the project. Using the colors from their respective logos, Blue is used to represent Twitter content and Yellow to represent Bitcoin. These colors are exclusively reserved to present these two assets and are introduced immediately within the subheading of the project through a highlighting affect to draw the users attention to the use of these colors. This negates the need to explain the color usage in each visualization. Some variations on the Blue color are used within the weighted sentiment analysis visualization to differentiate the different weighted models, however they still remain within a blue / purple color palette and still provide a clear separation between the two different assets.

In addition to the blue / yellow palette, a red / green palette is used when visualizing or discussing sentiment analysis. Using these colors to represent a positive or negative issue is familiar to the user and is explained visually through a highlighting effect of positive and negative words.

## 8.3 DATA

Data was collected via a Node.js script using the CCWS and Twit packages to access the Binance & Twitter APIs. The information was processed in accordance to the cleaning details discussed in both section 5 & section 6.2 before being saved to a MongoDB Atlas hosted database. Live stream data from both API sources is gathered via Socket.IO and constantly updating the database using a PM2 instance running on an AWS Cloud 9 service. Bitcoin data is stored in records of 1-minute candle data and Twitter data is stored in records of individual Tweets, removing any information that is not used within the visualizations. The ID of each tweet is stored to allow recall at a later date if additional information is needed.

### 8.3.1 AGGREGATION

Data is requested by the application in two formats, either individual tweets for the volume and filtering based visualizations, and in aggregate form based on a user set time interval of either minutes or hours. The individual tweet calls are limited to a small time period, requesting between $1000 - 2000$ tweets dependent on the user selected start date. Aggregation data does not have the same time restrictions and the user is able to request data over several days, aggregating hundreds of thousands of tweets.

The aggregation is completed using the MongoDB aggregation pipeline and filters the tweets based on the user specified filters on all tweet features (discussed in section 5) before aggregating to give maximum, minimum & average values for each tweet element per time interval (minute or hour). Weighted averages are also calculated within the aggregation call and return within the same results. The aggregation results not allow only the displaying of the changes in sentiment, followers, likes, comments & retweets, but also allows for the uncertainty of the data to be able to be displayed through area charts, such as in the sentiment visualizations.

## 8.4 USER INTERACTION & EXPLORATION

The application is designed to guide the user through the different methods of analyzing and representing a body of Tweets and help the user consider each of the details discussed in section 5 and how they may affect the aggregate values and correlation to Bitcoin. Each stage builds on the tools and considerations of the previous visualization and introduces new elements and filters that the user can manipulate. The variations in the outcomes of each of these elements are displayed using a small multiples approach to allow easy comparisons between tweet elements such as followers, likes an retweets, or when comparing different weighted average approaches.

At each area of investigation, the user is invited to interact with the visualization in some form, either changing the view or scales to help identify different patterns or see the effects of filtering, offset or changing the selected 'weight' in a weighted model. All interactive elements are separated out using a bordered box to keep the user aware of when they can interact. After the initial introduction of the filters, they are displayed in a minimized state to not overwhelm the user. As the visualizations and narrative progress, more filters are added, so that no confusing is caused in the effects of each filter setting before an explanation and introduction has been given.
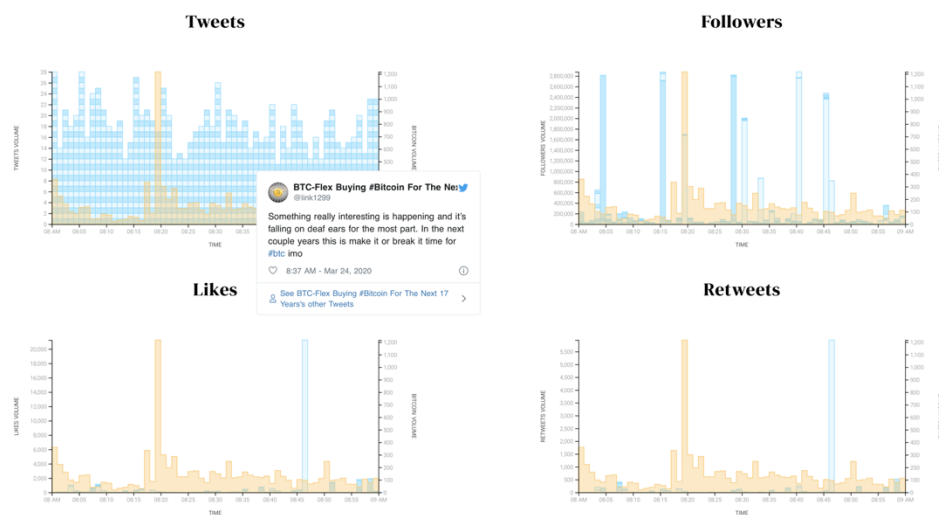
## 8.5 AREAS OF INVESTIGATION

After the initial introduction and explanation of the key terms and elements of twitter that will be investigated, the project is broken down into individual areas of investigation, each with its own visualization style and approach, keeping the visual language discussed in section 8.2.3.

### 8.5.1 VOLUME

The user is first introduced to the most basic form of comparison between Twitter and Bitcoin; a visualization that compares volume of Bitcoin activity (both buying and selling) against various quantifications of Twitter. In its initial and most basic form, the visualization is an area graph with two colored areas of yellow and blue representing the Bitcoin activity and volume of Tweets respectively. The totals for each value are totals per hour, for a period of 8 hours starting at 8am, representing a transitional working day. The initial date can be changed by the user prior to reaching this visualization.

The interactive elements of the visualization allow for the view to be changed from the smooth curves of the area chart into a stacked bar chart. Each block in the stack represents a single tweet, linking the user back to the content that is being aggregated into totals. Each block can be hovered over to display a tooltip of the original Tweet using the official Twitter widget to provide the correct styling and format.

Small multiples are used to compare the differences in totals between tweet volume, retweets, likes, & followers. The blocks for each of the stacks change in size to represent to amount of retweets, likes or followers. This view allows both easy comparison of the reach of the messages and the level of user interaction, but also identifies quickly popular tweets or tweets that have been created by high profile Twitter users.
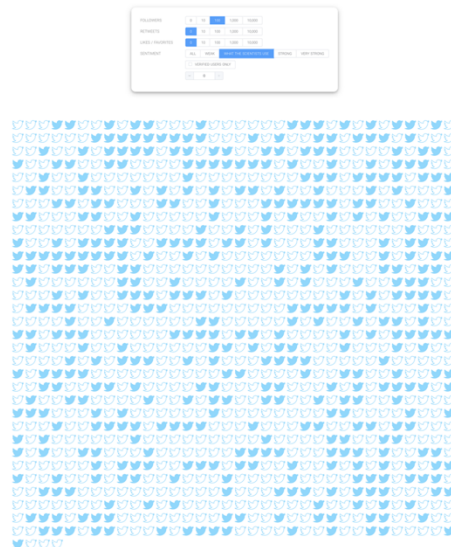
## 8.5.2 FILTERING

The filtering visualization uses the same selection of tweets as the volume visualization to display a grid of twitter icons that each represent a single tweet. Each icon has a tooltip available on hover over in the same manner as the volume visualization.

The purpose of the visualization is to continue to strengthen the users understanding that all of the data being representing throughout the project are genuine tweets created by real users. They are represented the grid format with each icon having no visual bias or difference in appearance as this is how other visualizations and analysis on the links between cryptocurrency and social media content have approached the task, treating each tweet equally.

When the user changes the each of the filters (as detailed in section 5) the tweets that do not match the minimum filter values change appearance to have no fill color. They remain within the visualization so to constantly visualize the total number of tweets that occurred within that time interval. The user is able to quickly visualize the distribution of how many tweets have certain characteristics. Even a seemingly low filter value of 10 likes or retweets will often exclude well over 50% of the tweets clearly showing that while a large volume of tweets may occur in a small time period, only a small proportion of these tweets are interacted with by users. The same observations can be made through filters on the number of followers and verified users. Each of these discoveries are important for the user to understand as they move forward into the upcoming visualizations as these are the elements that will make the differences when comparing different models and approaches in the sentiment and correlation visualizations.
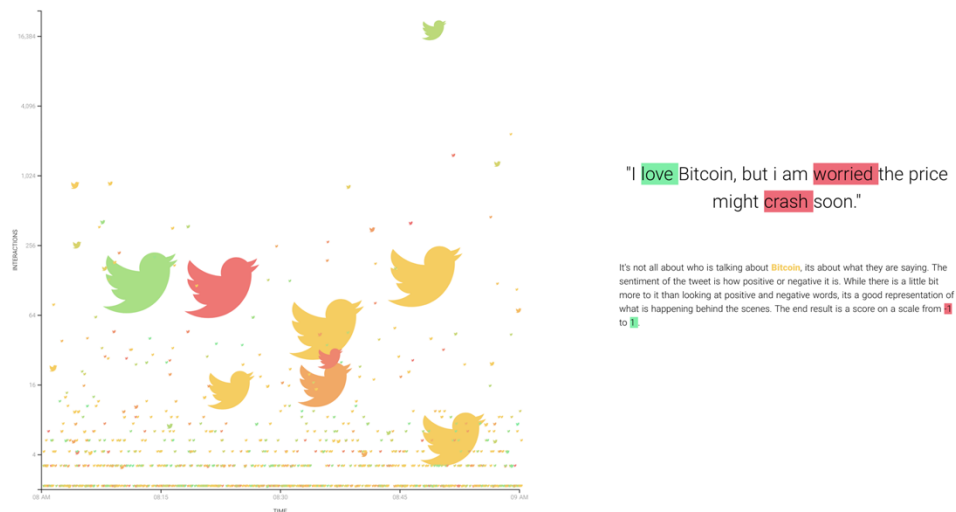
### 8.5.3  TWITTER SENTIMENT OVERVIEW

After a short written and visual explanation of the process of sentiment analysis, the same twitter icon shapes as used in the previous visualization are used color coded to show the sentiment of each of the individual tweets. The layout of each of the tweets moves away from the previous grid system and is displayed on axis with X representing time and Y representing the number of interactions. Finally, the size of each icon varies dependent on the number of followers the user that created the tweet has.

The visualization gives a strong representation of the distribution of the tweet characteristics without any initial or further interaction from the user. The filter values set from the previous visualization apply also to this visualization and the user can change back to the grid format (displaying in the exact format as the previous visualization) to provide object constancy between the two visualizations.

The variation in the size of the icons, combined with a large volume of tweets can create overlap and individual tweets may become indistinguishable. While this may initially appear to be a limitation, it can create a color palette affect where users can get a general feel for the overall sentiment of all of the tweets, or the sentiment of tweets that have a high number of interactions or that happened within a certain section of the displayed time period. Gaining this overview of the sentiment and how certain tweets can contribute to, or dominate an image when considering the combination of all of the tweets is important as the following visualization will only display information based on aggregate results.

## 8.5.4 SENTIMENT VARIATION OVER TIME

As the user reaches this visualization, the time period that is explored is expanded to 7 days with only aggregate results displayed. Each of the user set filter values remain from the previous user explorations and are also available to the user at multiple checkpoints throughout the project so that they can alter the results to see how their changes affect each of the visualization results.
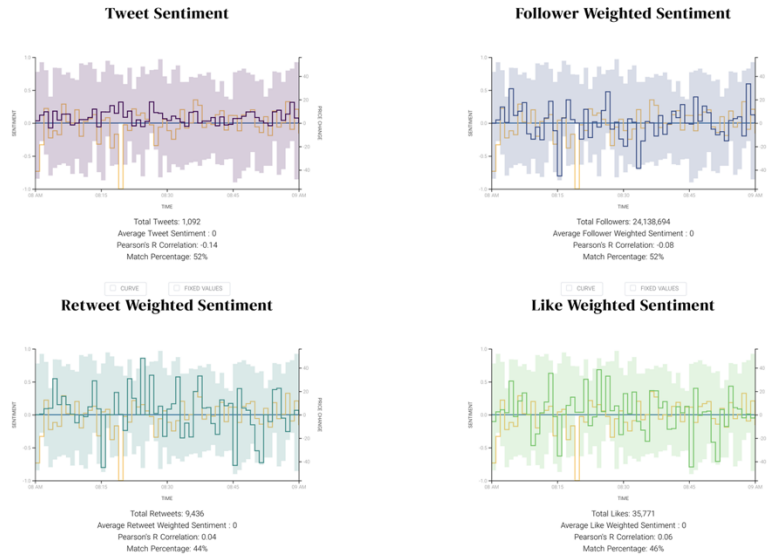
This visualization reintroduces direct comparisons between Twitter and Bitcoin, with the X axis displaying time and the Y axis being split in the middle at a 0 value and representing positive and negatives changes in values. The Twitter value is changes between positive and negative sentiment and the Bitcoin value is increases or decreases in Bitcoin price per coin. The data is displayed as averaged aggregate values in hour increments; however, the user can change this to minute increments.

For both Twitter and bitcoin an area chart is displayed behind the line chart to display both the highest and lowest sentiment and price values for each time interval. This introduces the idea of uncertainty in the data to the user which is explored in more detail in the correlation visualization.

Two interactive elements can change the appearance of the data to help the user in identifying patterns.  Firstly, the user can change the curve type, switching between a smooth curve showing transitions between the values and a hard step, which is more representative of the values as each data point is has a fixed start and end point for the inclusion of data for the aggregation. In addition to the type, the user can choose to exclude the amount of change in the sentiment or Bitcoin value and simply display a positive sentiment or increase in Bitcoin value above the central line and a negative sentiment or price decrease below the central line.

### 8.5.4.1 Weighted Analysis

The topic of weighted averages is displayed through small multiples of the sentiment variation over time visualization described in section 0. Each visualization has the same interaction abilities as the original visualization with the Twitter data varying slightly in color though a blue / purple color scheme to differentiate the differences in each graph. Totals for each weight are displayed under each graph along with Pearson R correlation value and a match value. The match value is explained to the user and ids a focus of the final visualization.

**Tweet Sentiment**

Total Tweets: 1,092
Average Tweet Sentiment : 0
Pearson's R Correlation: -0.14
Match Percentage: 52%

**Follower Weighted Sentiment**

Total Followers: 24,138,694
Average Follower Weighted Sentiment : 0
Pearson's R Correlation: -0.08
Match Percentage: 52%

**Retweet Weighted Sentiment**

Total Retweets: 9,436
Average Retweet Weighted Sentiment : 0
Pearson's R Correlation: 0.04
Match Percentage: 44%

**Like Weighted Sentiment**

Total Likes: 35,771
Average Like Weighted Sentiment : 0
Pearson's R Correlation: 0.06
Match Percentage: 46%

## 8.5.5 TIME LAG & CORRELATION

The final visualization combines all of the concepts from the previous visualizations and aims to provide clear evidence of if the user selected filtering and weighted models are able to show influence from Twitter on Bitcoin price.
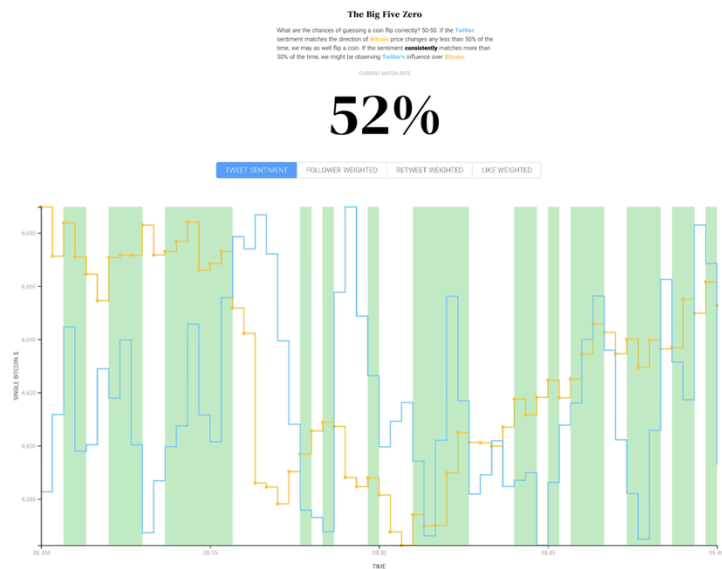
The visualization is displayed full screen, using the same time series as the sentiment variation visualization. Bitcoin value is no longer displayed as price difference per time interval, instead tracking the price change over time, in a similar to format to a traditional trading platform. Twitter sentiment data is transformed to create a continuous value that follows a similar format to the Bitcoin data. While this data does not conform to a scale that has any logical value, it does allow for an easier comparison to be made by the user between the two data sources.

The area chart for both Bitcoin and Twitter that display the range in results is removed and only a single aggregate average value is displayed. The weighted model to be used to create the average value for Twitter can be selected by the user at the top of the visualization. The range of information that was displayed in the area chart is moved into a tooltip and is displayed as an overlapping area chart of a single time interval (either an hour or minute) to allow the user to explore the uncertainty in the data without the additional information cluttering the main visualization and making comparisons harder. The details and how to read the uncertainty tooltip has its own introduction and explanation below the visualization.

A 'match' value calculates the percentage of instances where the Bitcoin value and the Twitter sentiment move in the same direction. Either a positive sentiment and Bitcoin value increasing, or a negative sentiment and Bitcoin value decreasing is deemed a match. The match score not only provides a clear quantification of if there is a correlation between the values but is an easier topic to grasp by a wider range of users. The match

score is visualized through a large total percentage value, with associated explanation, including considerations to how the match score may be the same as a lucky coincidence or the chance of flipping a coin. It is also visualized on the graph through large green bars across the time interval where the two values match. This allows for a very quick judgement and overview of if the changes in filter values or selected weighted model are having a positive or negative effect.

In addition to the tools introduced within earlier visualizations, a final 'offset' control is introduced to the user. The user changeable setting will offset the Bitcoin value by a multiple of the current time interval (either minute or hour). All match, Pearson's R and visualizations are updated to display the offset values to see if influence can be observed after a delayed interval of time. This negates the need to separately calculate a time lagged cross correlation as discussed in section 0 as it produces the same result while also providing a time lagged match value.



# 9   FURTHER ENQUIRY

## 9.1   COMPLETE DATA

While the project does provide a useful insight into the effects of different weighted models and filtering's on correlation strength between Twitter user activity and Bitcoin market value, the data used for the project is restricted. The data has been obtained from the standard (free) Twitter API stream and the 7-day historical tweet search. Investigation and trials were explored with the paid Twitter service but the cost to gain the large volume of tweets required to benefit the project beyond using the free service could not be justified. The time restriction of the 7-day search is being overcome on the project by constantly updating the database using the Twitter stream so that more

information is available to be explored by future users. This data is immediately available without any additional intervention or update.

The main limitation in the standard service data is that some important details for the project such as comment and reply data are not included. If this information were available, then additional comparisons could be made. With comments being one of the most user intensive tasks on the platform, it is not unreasonable to assume that these may be a good indication of influential tweets.

In addition to the missing fields within the provided data, obtaining data from the live stream means that the tweets do not yet have any likes, comments or retweets on them as they have just been created. Retweet events are provided in the stream and an effort has been made to use these retweet events to update the data of existing tweets within the database.

As time progresses the database will improve due to the constant updates, but if a pre-existing dataset could be obtained, especially over key periods of time such as the January 2018 Bitcoin crash, this may provide stronger correlation results. Alternatively, funding would need to be gained in order to purchase the required data directly from Twitter.

In order to accommodate a larger amount of Twitter data over a longer period of time, an additional time interval of a day (along with the existing minute and hour intervals) is required. This may require increased capacity of the MongoDB database to be able to complete the aggregations fast-enough for the 'on the fly' filtering of the user.

## 9.2   ADDITIONAL VOLUME EXPLORATION

Due to the narrative structure of the project that introduces a single tool and topic at a time with an individual topic, the volume visualization does not currently change due to the effects of the user defined filters as the filters had not been introduced. This is within the current capability of the system and could be reintroduced between the filtering and time series visualizations. This may provide a more accurate representation of the Twitter volume data by remove unnecessary noise such as tweets that only reached a very small number of users of that no users interacted with.

In its current form, the volume visualizations do not give any indication to twitter sentiment of direction of Bitcoin value. As an additional level of granularity, the volume visualizations could be split to represent the volume of tweets with a positive sentiment against purchase orders for Bitcoin and an additional set of volume visualizations to display the volume of tweets with a negative sentiment against Bitcoin sell orders. This would provide a more accurate presentation of the information and if correlations were observed then it would indicate a potential influence between Tweets and Bitcoin value.

## 9.3  REALTIME TRADING

The current project was presented as an introduction to the topic to reach a wide target audience. If the target audience is shifted to processional cryptocurrency traders, each of the areas of investigation in the project could be amalgamated into a single platform and tool. Using discoveries from exploration of the current project, user defined settings can be applied to quickly display information of the most up to date Twitter and Bitcoin information. Each of the tools could remain as individual elements as they each evaluate different elements of the Twitter data. Each tool could provide an indication of correlation and could be combined as an ensemble model to provide an overall prediction. Without additional analysis of the Bitcoin trading behavior and external factors effecting Bitcoin's value, the tool would not provide enough information as a stand along system to provide trading advice. It would however be a useful additional source of information to inform traders and their current trading indicators. In addition to providing an indication for manual trades, the information from the tool could be fed into existing automated trading systems as an additional indicator.

# 10 CONCLUSION

The project aim was to form an investigation into the correlation between user behavior on Twitter and Bitcoin. The expectations were not to provide any conclusive results, due to the multitude of external factors at play that also affect Bitcoin price fluctuations. Instead, the project focused on creating a series of exploratory tools that would invite the user to consider which elements of behavior on social media have the largest influence over our behavior.

While no strong correlations were observed during user testing a few key observations were made. Low level filtering of interaction, so that only tweets that had gained even the smallest number of retweets, likes or comments, did improve the match result on the majority of time periods tested. Higher levels of filtering for user interactions did not consistently produce stronger results. This is potentially due to sharp drop off in the number of available tweets when filters are set to a minimum filter value of 100. Higher filter values did show positive effects when used with follower count as on average Twitter users will have a multitude more followers than their average tweet interaction count.

Weighted models produced the strongest observed results. A match result of 50% is deemed insignificant as this would be the same result as blindly guessing the direction of the Bitcoin price change, however follower weighted average combined with low level interaction filtering produced results as high as 60%. It is unlikely that a result this high in hour intervals over a period of 7 days would be produced by chance. Weighted models based on user interaction were expected to produce stronger results based on previous studies that showed user interaction having a stronger correlation to influence. This was

observed in some time periods explored but the most consistent results were follower weighted sentiment average.

An interesting observation that had previously not been considered (although makes logical sense) is that applying a filter to a tweet element and then using the weighted model of that element, does not improve result. In many cases this weakens the result. An example would be applying a filter to the number of retweets and then using the retweet weighted model. This is counter intuitive as the weighted model is already designed to give a lower preference and weight to tweets with a lower retweet count.

Designing the user interface for a novice audience directed the project to be broken down to as more granular level, creating multiple visualizations instead of a smaller set of more complex exploratory tools with a wider range of user interaction. This allowed for each area of exploration to be examined independently, identifying its place in effecting the overall correlation results. The project successfully explains each topic and displays results in an easy to consume manner, allowing the user to explore the vast array of filtering, weighting and offset models to form their own conclusions and observations, whilst maintaining a scientific grounding that can be built on in future developments.

# 11 BIBLIOGRAPHY

Bollen, Johan, Huina Mao, and Xiaojun Zeng. 2011. "Twitter Mood Predicts the Stock Market." *Journal of Computational Science.* https://doi.org/10.1016/j.jocs.2010.12.007.

Cha, Meeyoung, Hamed Haddadi, Fabrício Benevenuto, and Krishna P. Gummadi. 2010. "Measuring User Influence in Twitter: The Million Follower Fallacy." In *ICWSM 2010 - Proceedings of the 4th International AAAI Conference on Weblogs and Social Media.*

Chen, Chung-chi, Hen-hsen Huang, and Hsin-hsi Chen. 2018. "NTUSD-Fin: A Market Sentiment Dictionary for Financial Social Media Data Applications." *Mimeo.*

Garcia, David, and Frank Schweitzer. 2015. "Social Signals and Algorithmic Trading of Bitcoin." *Royal Society Open Science.* https://doi.org/10.1098/rsos.150288.

Kollanyi, Bence, Philip N Howard, and Samuel C Woolley. 2016. "Bots and Automation over Twitter during the First U.S. Presidential Debate." *COMPROP DATA MEMO 2016.1.*

Phillips, Ross C., and Denise Gorse. 2018. "Predicting Cryptocurrency Price Bubbles Using Social Media Data and Epidemic Modelling." In *2017 IEEE Symposium Series on Computational Intelligence, SSCI 2017 - Proceedings.* https://doi.org/10.1109/SSCI.2017.8280809.

Annalect. 2016. *blog.twitter.com.* 05 10. Accessed 03 27, 2020. https://blog.twitter.com/en_us/a/2016/new-research-the-value-of-influencers-on-twitter.html.

CoinMarketCap. 2020. *coinmarketcap.com.* 03 27. Accessed 03 27, 2020. https://coinmarketcap.com/charts/.

—. 2020. *Coinmarketcap.com.* 03 27. Accessed 03 27, 2020. https://coinmarketcap.com/all/views/all/.

Congressional Research Service. 2019. "Congressional Research Service." *https://crsreports.congress.gov.* 08 09. Accessed 03 27, 2020. https://crsreports.congress.gov/product/pdf/R/R45863.

Digiconomist. 2020. *Digiconomist.* 03 27. Accessed 03 27, 2020. https://digiconomist.net/bitcoin-energy-consumption.

Influencer Marketing Hub. 2020. *Influencer Marketing Hub.* 01. Accessed 03 27, 2020. https://influencermarketinghub.com/influencer-marketing-benchmark-report-2020/.

Kick Factory. 2016. *Kick Factory Blog.* 06 23. Accessed 03 27, 2020. https://kickfactory.com/blog/average-twitter-followers-updated-2016/.

Koumchatzky, Nicolas & Andryeyev, Anton. 2017. *Twitter.* 05 07. Accessed 03 27, 2020. https://blog.twitter.com/engineering/en_us/topics/insights/2017/using-deep-learning-at-scale-in-twitters-timelines.html.

MacCarthy, Ryan. 2016. *KickFactory.* 6 23. Accessed 03 23, 2020. https://kickfactory.com/blog/average-twitter-followers-updated-2016/.

Mediakix. 2019. *Mediakix.* 04 18. Accessed 03 27, 2020. https://mediakix.com/blog/power-of-social-media-influencers-trendsetters/.

Nicholas Confessore, Gabriel J.X. Dance, Richard Harris and Mark Hansen. 2018. *The New York Times.* 01 27. Accessed 03 27, 2020. https://www.nytimes.com/interactive/2018/01/27/technology/social-media-bots.html.

Report, Q3 2019 Shareholder. 2019. *Twitter Q3 2019 Shareholder Report.* 1 1. Accessed 03 26, 2020. https://s22.q4cdn.com/826641620/files/doc_financials/2019/q1/Q1-2019-Slide-Presentation.pdf.

Security and Exchanges Commission . 2017. "Release No. 34-80206; File No. SR-BatsBZX-2016-30." *Security and Exchanges Commission .* 03 10. Accessed 03 27, 2020. https://www.sec.gov/rules/sro/batsbzx/2017/34-80206.pdf.

Smith, Aaron. 2014. *Pew Research Center.* 02 03. Accessed 03 27, 2020. https://www.pewresearch.org/fact-tank/2014/02/03/what-people-like-dislike-about-facebook/.

SocialPubli. 2019. *https://socialpubli.com.* 02 04. Accessed 03 27, 2020. http://news.socialpubli.com/2019-influencer-marketing-report-a-marketers-perspective/.

Toshi Times. 2018. *Toshi Times.* 04 08. Accessed 03 27, 2020. https://toshitimes.com/the-mcafee-effect/.

Twitter. 2017. 12 03. Accessed 03 27, 2020. https://twitter.com/btcWhaleclub/status/944643681507737601/photo/1.

—. 2016. *Twitter.* Accessed 03 27, 2020. https://help.twitter.com/en/using-twitter/twitter-timeline.

Winterberry Group. 2020. *Winterberry Group.* 01. Accessed 03 27, 2020. https://www.winterberrygroup.com/our-insights.