

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/317272393>

# Sentiment based Analysis of Tweets during the US Presidential Elections

Conference Paper · June 2017

DOI: 10.1145/3085228.3085285

CITATIONS

26

READS

1,209

4 authors, including:



**Ussama Yaqub**

Rutgers Business School

9 PUBLICATIONS 106 CITATIONS

[SEE PROFILE](#)



**Soon Ae Chun**

CUNY CSI & Graduate Center

161 PUBLICATIONS 2,105 CITATIONS

[SEE PROFILE](#)



**Vijayalakshmi Atluri**

Rutgers, The State University of New Jersey

175 PUBLICATIONS 4,097 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Social Media sentiment and data analytics [View project](#)



dg.o 2020: TRACK 3. Smart Cities: Intelligent Innovation and Transformation [View project](#)

# Sentiment based Analysis of Tweets during the US Presidential Elections

Ussama Yaqub  
Rutgers Business School  
1 Washington Park  
Newark, NJ, 07102  
[ussama.yaqub@rutgers.edu](mailto:ussama.yaqub@rutgers.edu)

Vijayalakshmi Atluri  
Rutgers Business School  
1 Washington Park  
Newark, NJ, 07102  
[atluri@rutgers.edu](mailto:atluri@rutgers.edu)

Soon Ae Chun  
IS & Informatics, CSI  
City University of New York  
Staten Island, NY, USA  
[soon.chun@csi.cuny.edu](mailto:soon.chun@csi.cuny.edu)

Jaideep Vaidya  
Rutgers Business School  
1 Washington Park  
Newark, NJ, 07102  
[jsvaidya@rutgers.edu](mailto:jsvaidya@rutgers.edu)

## ABSTRACT

In a relatively short period of time, social media has gained significant importance as a mass communication and public engagement tool for political and governance purposes. Rapid dissemination of information through social media platforms such as Twitter, provides politicians and campaigners with the ability to broadcast their message to a wide audience instantly and directly while bypassing the traditional media channels. In this paper, we investigate the nature and characteristics of the political discourse that took place on Twitter during the American Presidential elections of November 2016. The goal of this study is to perform exploratory sentiment based analysis of Twitter data that was gathered both before and after the Election Day. Our objective is to identify the nature and sentiment of discussions along with understanding the behavior of users with respect to their Twitter profile and associated attributes of their tweets. We also aim to inspect popular Twitter discussion topics and their relation with important news and events occurring simultaneously.

## CCS CONCEPTS

• Human-centered computing → Collaborative and social computing → Collaborative and social computing theory, concepts and paradigms → **Social media**

## KEYWORDS

Sentiment analysis, social media, behavior analysis, elections.

## ACM Reference format:

U. Yaqub, S.A. Chun, V. Atluri, and J. Vaidya. 2017. Sentiment based Analysis of Tweets during the US Presidential Elections. In *Proceedings*

*of ACM Digital Government Research, Staten Island, NY, USA, June 2017 (dg.o'17)*, 10 pages.

DOI: 10.1145/3085228.3085285

## 1. INTRODUCTION

Use of social media for communication has greatly increased over the last few years. Due to the availability of online social networks such as Facebook and Twitter, people are now increasingly relying on social media to not only interact with one another but also to read and share news, discuss important events and engage in political discussions. Additionally, proliferation of smart phones has further facilitated the use of this medium, allowing users to communicate without any limitation on location.

This opportunity presented by social media has been recognized by politicians and political parties globally [19]. The potential role social media can play in political events was first highlighted during the US Presidential elections of 2008. Twitter played an important part in the campaign of Barack Obama. The Obama campaign made an effective use of Twitter to post campaign updates and inform followers with opportunities to volunteer with the campaign [33]. During the 17 months of the election campaign starting from April 2007 to Election Day November 5th 2008, the Obama campaign posted 262 tweets and gained approximately 118,000 new followers [34]. In light of this successful Twitter campaign, all major candidates and political parties now have some form of presence on social media.

This growth of Twitter usage during elections by politicians, campaigners and public has led to ever increasing research in the areas of social media sentiment analysis and data analytics. Twitter analytics with regards to election prediction and candidate popularity have grown tremendously over the years [14]. Some studies have even gone so far as to state that sentiment analysis of tweets can be used as a substitute for traditional polls monitoring consumer confidence and political approval ratings [10].

Twitter averaged 313 million monthly active users as of 2nd quarter of 2016 [7]. This massive number of people utilizing the service for information gathering and expressing their views has provided political actors with the opportunity to put their message across quickly and cheaply without going through the traditional media briefings and news conferences [19].

Popularity of Twitter provides a unique opportunity for e-government initiatives, specifically with regards to communication between government institutions and citizens

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [Permissions@acm.org](mailto:Permissions@acm.org).  
dg.o '17, June 07-09, 2017, Staten Island, NY, USA  
© 2017 Association for Computing Machinery.  
ACM ISBN 978-1-4503-5317-5/17/06...\$15.00  
<http://dx.doi.org/10.1145/3085228.3085285>

[13]. Over the years, governments around the world have worked to increase openness and transparency into the working of their institutions [12]. Social media in general and Twitter in particular can play a constructive role in this regard. Hence from civic services to police department, information sharing through Twitter can lead to greater transparency and more confidence of citizens on their state and local institutions [11]. One aspect of this public engagement is through political discourse on social media. In this paper, we perform Twitter data analysis to understand public sentiment and dialogue concerning general elections. Techniques discussed here however can be applied in broad areas of public policy ranging from local governance to bi-directional communication with state institutions.

In the recently concluded US Presidential Elections of 2016, Twitter played a very important role in the dissemination of information regarding various policy points for both major candidates. Both candidates had millions of followers on Twitter and had their tweets closely monitored by general public and by the mainstream media. Indeed, even today, the debate regarding the circulation of fake news on social media and its effect on the elections still rages on [8]. Although it is hard to quantify the role Twitter played in the elections 2016, all are in agreement that it was nonetheless significant. This means that political players cannot ignore the role of social media as a communication channel. Overall, social media presents an exciting avenue of opportunity for politicians, campaigners and political activists to not only broadcast their message but also to engage in dialogue with proponents of competing political ideas and ideologies.

In this paper, we attempt to understand the political discourse that took place on Twitter during the US Presidential Elections of 2016. This is an exploratory study where we analyze Twitter data to identify user behavior on Twitter along with the nature of conversations that took place. For this purpose, we gathered over 3 million tweets for 21 days daily, starting from 29th of October up until 18th of November. Metadata of these tweets was then extracted and the text was analyzed for sentiment. We used the publicly available SentiStrength [1] software to assign tweets with positive and negative sentiment scores.

The exploratory analysis also includes measuring the popularity of both major presidential candidates (Donald Trump and Hillary Clinton) on Twitter. Daily average sentiment score of tweets containing each candidate is calculated. This daily average is then compared with the overall average daily sentiment of all tweets present in the database. A sentiment trend is also computed for both candidates. This trend is compared and contrasted with the public polls conducted during this time period.

Additionally, to better understand Twitter conversations in terms of popular topics and trends, message sentiment, message propagation, content creation versus reuse of content, we further refine our analysis based on numerous user characteristics. User attributes such as number of followers and friends, account creation date, total number of tweets, number of favorites etc., are available as metadata. These attributes are utilized along with tweet text to better understand the sentiment displayed and the content generated or shared by the users.

The rest of the paper is structured as follows. In Section 2, we review the related work in the literature. In section 3 we propose our user behavioral model of social media user behavior in the context of the 2016 Elections. In this section we will also provide an overview of our data and its preparation for analysis. Section 4 provides the data characteristics and the actual content based

analytical methods. Section 5 presents the results of our data analysis. In section 6 we discuss these results while in section 7, we conclude the paper.

## 2. RELATED WORK

The role of Internet and communication technologies (ICT) in modern society cannot be understated. Individuals and institutions around the world are trying to increase public engagement by utilizing Web 2.0 [12, 13]. This provides a quick and cost effective platform to political actors and state institutions to communicate quickly and directly with public [11]. For example, Twitter has now been used by city governments to benefit their populations by raising information awareness in a simple, low cost fashion. The idea is to enhance the responsiveness of different branches of local governments that deal primarily in performing tasks on behalf of the citizens [13].

Along with governance, Twitter sentiment analysis is also being used in vast array of areas related with governance and public trust ranging from predicting resentment against government policies to predicting general election results [2, 14]. Various models have been developed that try to understand the user behavior and retweeting on Twitter [28]. The emerging field of techno-social systems aims to comprehend and predict this behavior. Although this area of study is still evolving and generating a lot of enthusiasm, nonetheless, a debate on the efficacy of using Twitter sentiment analysis to predict elections and other real world events still continues [15, 16]. Important questions such as how representative Twitter users are of general population remain to be answered. These issues become acute when these analyses are conducted on data obtained from developing countries where a relatively small percentage of population has access to internet.

Another aspect is the varying levels of user activity. Some users are far more active than others online thus having a greater 'weight' to their opinions when compared with low activity users. At the same time, there exists much noise on Twitter in the form of automated activity and spam, which exploit trending topics to advertise various unrelated products or content. Different solutions have been proposed to differentiate between human activity and that generated by bots [17].

Other studies have looked at how the information is diffused on the social networks and what role sentiment plays in this information diffusion [4]. Most agree that sentiment does play an important role in information diffusion on twitter. Some have gone as far as saying that there exists a positivity bias in information spread and that positive tweets are retweeted more and reach a wider audience than negative tweets [3, 4].

Studies have also been conducted to discover correlation between tweets sentiment and public opinion polls. A high correlation of 80% was claimed by one study between the Index of Consumer Sentiment (ICS) conducted by Reuters and Twitter sentiment [10]. The study also found high correlation between Gallup's daily tracking poll for job approval rating of President Barack Obama and Twitter sentiment over the course of 2009. According to the authors this high correlation between Twitter sentiment analysis and public survey data indicated potential of tweets as a substitute for the traditional polls [10].

With regards to the use of Twitter in politics, researchers have examined the ways in which Twitter influence communications of mainstream news and journalism. Recent research shows that

social media and Twitter in particular are playing an increased role in mainstream media as a news source. This can be in form of a quote, background information or policy issues outlined through twitter messages by politicians or other political actors such as news commentators and observers [29]. Twitter is now ever more used as news agenda building tool for mainstream media [30, 31]. This was a very commonly observed phenomenon during the recently concluded US elections of 2016.

Utilization of Twitter by politicians and their campaigns is a popular subject of study. Usage of Twitter during the campaign cycle of 2008 in USA by Barack Obama generated interest in understanding Twitter's role in political campaigns [33, 34]. Similar research was also conducted in analyzing Twitter activity of US Congress members during their election campaigns. Studies showed that congress members frequently posted information on Twitter regarding their political positions on various issues along with posting information related with their constituencies [35, 36]. Social media user behavior analysis has also been conducted from the knowledge creation and sharing perspective in e-government context. Studies have been conducted on evaluating knowledge creation and sharing behaviors depending on the level of activity of individual Twitter users [9]. User tweeting behavior in-terms of reusing existing content vis-à-vis new content creation can be dependent upon how frequently they tweet.

### 3. METHODOLOGY

United States of America has the highest number of Twitter users in the world [7]. As of May 2016, there are approximately 67.5 million active users of the microblogging site in the country [7]. This large user base combined with a significant event such as the elections makes Twitter data an ideal case study of social media usage in political discourse. Furthermore, both major presidential candidates and their respective political parties made extensive use of social media for campaign purposes, and used Twitter for various purposes: from outlining policies and issues such as security, economy, healthcare, gender equality to promoting slogans. Twitter was one of the favorite communications tool for Donald Trump, who utilized it very regularly to react to news concerning his candidacy and other issues of importance as they rose before and after a hotly contested election.

In light of this important role that social media played in the general elections, our aim is to utilize user generated data to understand their behavior and important issues discussed and highlighted online. Our approach relies on data gathered from social media platform (i.e., Twitter) and to perform data analytics to understand the nature of discussions and user behavior on the microblogging site. Our methodology comprised of the following steps:

1. Use of search terms "Trump", "Clinton" and "Election2016" to gather Twitter data for our period of interest which is the Presidential elections of 2016 in USA.
2. Data cleaning and extraction.
3. Sentiment tagging and classification of gathered tweets.
4. Importing data in MySQL database to perform exploratory analysis of data.
5. Development of user behavioral model, formulate hypotheses and derive findings proving or disproving the hypotheses.

We utilized Python to gather Twitter data using streaming API. Python was also used for data cleaning and extraction of each

tweet's associated metadata. As mentioned earlier, SentiStrength [1] is used to identify message sentiment while data analysis is performed using the MySQL database.

#### 3.1 Data Set

We utilized Twitter data for our study of US Presidential Elections 2016 that was gathered using the public Twitter streaming API [18]. The streaming API allows near real time access to global stream of Twitter data. We have collected tweets for a total of 21 days, starting from 29th October 2016 and ending on 18th November 2016. We have gathered 10 days of data prior to the elections day and 10 days of post-election data. In total, 3,108,058 tweets are utilized for this study.

Scientific studies using Twitter messages either employ hashtags or specific keywords to collect relevant data. Both approaches follow the same principle that hashtag or keywords indicate a message's relevance to a given topic [32]. In our case, we use the keywords 'Trump', 'Clinton' and 'Election2016' to download tweets. Here we have the two major candidate names along with an impartial term to capture neutral sentiment. The reason for using keywords instead of hashtags was primarily motivated by the reasoning that hashtags are utilized by users who are somewhat familiar with the concept of trends on Twitter and are hence more experienced than a novice user. By downloading data utilizing keywords, we have attempted to make our data sample more inclusive of novice Twitter users. For each tweet, we extracted metadata details such as tweet time and date, its id, creator id and user name, location, etc.

#### 3.2 Sentiment Analysis

With data cleaned and relevant tweet fields extracted, the text messages are then analyzed and tagged with sentiments using SentiStrength [1]. SentiStrength is a freely available software that has been used to perform sentiment analysis in various studies utilizing Twitter data [2,3,4]. One of the advantages of using SentiStrength is that the tool has been specifically developed to capture sentiment of short, informal texts [20]. Studies conducted on short texts have shown this tool to be able to capture positive sentiment with 60.6% accuracy and negative sentiment with 72.8% accuracy [23]. Hence, this makes SentiStrength an ideal tool for our exploratory analysis.

SentiStrength operates by assigning two scores to each text message it analyzes. It assigns a negative and a positive score, with the scores ranging between [-1, -5] and [1, 5], respectively. A score of -1 or 1 indicates a somewhat neutral text sentiment while a score of -5 or 5 indicates a very high negative or positive sentiment respectively. In-order to classify a tweet as overall positive or negative, we assigned a total sentiment score to each tweet. To do this, we have added both the positive and negative sentiment scores for each tweet as a total sentiment.

$$\text{Sent (total)} = \text{Sent (positive)} + \text{Sent (negative)}$$

Thus, a total sentiment score of 4 (or -4) indicates a strong positive (or negative) sentiment for the tweet respectively while if the total score adds to 0 then the tweet can be classified as neutral. Finally, the sentiment tagged text messages and the associated metadata are stored in MySQL database for further analysis.

### 3.3 Removing Noise from Data

To increase accuracy of our analysis, the next step is to remove noise from our dataset. Presence of spam on Twitter is a well-known phenomenon. Although Twitter tries hard to identify and remove automated accounts, not all bots are easily identifiable as social bots are designed specifically to impersonate human behavior. Much research has been conducted to identify automated non-human activity on Twitter. It has been found that up to 10.5% of Twitter accounts might be bots [17]. Studies have also concluded that as high as 9% of tweets are generated by automated accounts [37].

In order to identify and remove spam present in our dataset in the form of automated activity, we have removed tweets belonging to accounts having abnormally high tweet rates. Through literature review of various studies, we have discovered that users tweeting over 150 times a day can be safely classified as bots [5]. We have found this to be a safe assumption and decided to remove all tweets from our dataset where the associated account had an average of over 150 tweets a day. We have also removed from our dataset, tweets associated with accounts having names such as "iPhone giveaways" etc. as these tweets are not intended to add towards the political discussion but are rather promotional in nature.

A total of 209,370 tweets are resultantly identified as having generated by abnormally high activity accounts or by accounts that had names or descriptions that can be classified as spam. After excluding the spam tweets, we are left with 2,898,688 tweets for 21 days generated by 1,131,232 unique users. It is interesting to note here that while the average tweet per user in our data set was 2.56, the top 9% of the users are responsible for almost 52% of the tweets while around 69% of users have only 1 associated tweet in the dataset. [Table 1](#) below has the dataset details.

**Table 1: User related dataset statistics.**

Tweet per user	Total tweets	Total Users	Percent of all users	Percent of all tweets
1	781,575	781,575	69.09%	26.96%
2	309,675	154,964	13.70%	10.68%
3	178,293	59,487	5.26%	6.15%
4	127,424	31,896	2.82%	4.40%
5 or more	1,501,726	103,310	9.13%	51.81%
<b>Total</b>	<b>2,898,693</b>	<b>1,131,232</b>	<b>100%</b>	<b>100%</b>

## 4. SENTIMENT AND USER BEHAVIORAL MODEL

As discussed earlier, the aim of the proposed model is to understand the characteristics of communication and behavior of users that took place on Twitter conversations during the US Presidential Elections 2016. In this section, we discuss the hypothesis development to be tested using this data.

### 4.1 Analysis of User Behavior

Once data is sanitized, we then proceed in employing data mining and analysis techniques to perform data analytics and find useful information. While performing exploratory data analysis we try to find evidence supporting some of our beliefs that we develop by reviewing literature review of similar studies analyzing Twitter

data for insights into user behavior and tweeting patterns. Following are some of our suppositions that we test on this data, looking for evidence either corroborating or refuting them in the context of US Elections 2016.

**4.1.1 Content Creation:** Content creation on social media remains a very interesting subject of study. Researchers have looked at the question of why some content becomes popular and is retweeted thousands of times while many other tweets are never retweeted [14]. The relationship between social ties and the similar types of content that users create and share online along with the motivation to create new content is also an important issue to understand in this regard [24]. Furthermore, for political online social media content, researchers have observed a high rate of reusability [22]. In light of these research questions, we test whether most of the users in our dataset are speaking their mind and participating actively in online discussions or whether they are mostly passive and reusing content or thoughts other users created by simply retweeting them. To test this proposition, we propose the following hypothesis:

**Hypothesis 1 (H1):** *Majority of users commenting on the elections are not creating new content but are rather finding interesting and useful information which they are retweeting with other people in their network.*

**4.1.2 Twitter Activity Based on Following and Usage:** Research has been conducted on the message framing behavior of users on Twitter as a function of various characteristics including number of followers and the level of activity. Adding hashtags (#) preceding to keywords in messages allows users to make them searchable by other users. Burns et al. [25] have conducted an entire study on the use of hashtags in the context of coverage and discussion of news. The study utilizes large scale quantitative research on the use of hashtags in Twitter discussions allowing users to become a part of them.

The use of hashtag allows users to become part of Twitter trends and also enables them to reach a large audience by making their tweets searchable. Studies have claimed that infrequent users of Twitter are not as concerned about getting their tweets searched as compared to heavy users who are more motivated to make their tweets to reach a bigger audience allowing them to further increase their reach and following [9].

In light of this discussion, we believe that a similar trend will be discovered in our election dataset. Users with large following and heavy usage will be more concerned about making their tweets searchable than those having fewer followers and lesser number of tweets. By framing their messages using hashtags, heavy users are able to reach a broader audience, making their messages more searchable. Thus we will test the following hypothesis:

**Hypothesis 2 (H2):** *Users with large following and heavy usage are more concerned with making their tweets searchable and hence use hashtag (#) significantly more often in their messages compared with users who have few followers and are light users, who are indifferent towards reaching a greater audience in context of discussions related with elections.*

**4.1.3 Message Targeting:** Twitter allows users to broadcast their message to multiple people with one single tweet. However, Twitter also lets its users interact one-to-one by addressing a



person directly. This enables them to respond to other user's tweets paving way for a dialogue. Various studies regarding conversations on Twitter during elections have stated that people do not only use Twitter to post their political opinions but also engage in interactive discussions [14]. Nonetheless, direct messaging also creates complexities for users in having to handle multiplicity and one-to-one conversations at the same time [26]. Management of audience especially becomes challenging as the number of followers of a user grows.

Based on the above discussion, we assume similar behavior amongst the users of our dataset and believe that there will be a high number of one-to-one messages indicating interactive political dialogue. Hence we will test the following hypothesis:

**Hypothesis 3 (H3):** *Users in context of elections do not use Twitter only to voice their opinions but also use the platform to interact with other users in one-to-one discussions on political issues.*

**4.1.4 Reflection of current news:** Due to the instant nature of communication on Twitter, it can be used as a real time latest news identification tool. Several studies have been conducted in this regard, which attempt to identify real world events by analyzing Twitter streaming data [21]. Studies have claimed that based on trending topics based on active time period of tweets showed that as many as 85% of topics are headlines or persistent real world news [27]. Studies have also claimed that Twitter allows users to engage in real-time discussion of live televised broadcasting. Hence during major sports, entertainment and political events, Twitter is used to provide running commentary of real-time world events as they unfold on live television [35]. Thus it can be stated that analysis of daily tweets can provide us with the current news events taking place in the real world. In light of the above discussion, we will test this hypothesis on our dataset. We will analyze our tweets to search for high frequency terms in-order to test the following hypothesis.

**Hypothesis 4 (H4):** *Frequency of popular terms in Twitter discussions can be utilized to identify significant real world events and news taking place related with the elections.*

## 4.2 Method

In order to verify our hypothesis, we perform quantitative analysis of our data set. For these purposes we utilize following proxies:

- Tweets that are not original and are retweeted by the user contain RT string at the beginning of the message. Original tweets do not contain this string. Furthermore, in our dataset, retweets have their original creation date and tweet id. These fields are NULL for new tweets.
- Hashtags (#) enable users to make their messages searchable, allowing them to become part of trends.
- When a user tweets directly to another Twitter user, the message begins with "@" character. Hence tweets beginning without "@" are broadcast intended to all audiences while tweets starting with "@" are direct messages.

## 5. RESULTS

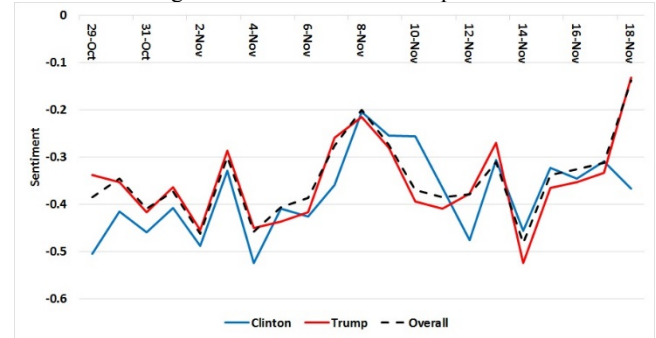
In this section, we present our findings of the data analysis. Some preliminary findings have helped us understand the overall sentiment of data along with the number of positive and negative tweets. We have also performed deeper data analysis to

corroborate the hypotheses that we have developed in the previous section and test them through our insights.

### 5.1 Sentiment Analysis

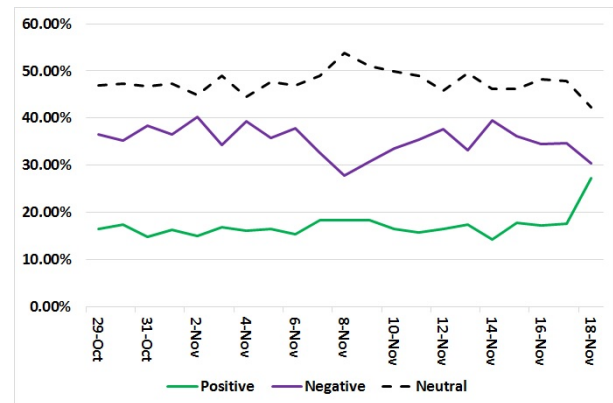
**5.1.1 Dataset and Candidate Sentiment:** Our preliminary data analysis involves analyzing the overall sentiment of entire dataset and of both candidates individually. With all tweets tagged with sentiment scores, we calculate average daily sentiment of the entire dataset along with tweets mentioning only Trump or Clinton in-order to create a comparison amongst them. The purpose of this analysis is to identify the overall sentiment of the Twitter conversations related with elections 2016 along with the sentiment of discussions involving both presidential candidates.

We discover that the average daily sentiment is negative for all 21 days of messages. Not only is it negative overall, but also for both candidates. Fig. 1 shows the average daily sentiment of all tweets in the database along with average daily sentiment of tweets containing the terms Clinton or Trump.



**Figure 1: Average daily sentiment overall and for tweets containing "Trump" and "Clinton".**

**5.1.2 Positive, Negative and Neutral Sentiment Tweets:** While the daily average sentiment is negative for all days, the number of neutral tweets in the database is higher than the negative and positive tweets. However neutral tweets have 0 sentiment score and thus have little effect on the daily average sentiment score. Fig. 2 exhibits the percentage of neutral, negative and positive tweets in the dataset for each day.



**Figure 2: Daily percentage of negative, positive and neutral tweets.**

We have also observed that there are more tweets with negative sentiment than positive. This finding is different from other studies that have been conducted using SentiStrength to perform sentiment analysis of Twitter messages [2, 3, 4]. These overall negative scores might indicate the bitter nature of the political campaign associated with Elections 2016. The negative score might also be due to the fact that almost 90% of tweets in our dataset contained either or both candidate names (Clinton or Trump) and the negative sentiment can thus indicate strong negative feeling exhibited towards these two candidates by their opponents.

**5.1.3 Candidate Popularity:** Prior to the elections, most polls conducted by various organizations showed Hillary Clinton leading Donald Trump [39]. Table 2 shows the results of most well-known polls conducted during 29<sup>th</sup> Oct to 7<sup>th</sup> Nov. We want to contrast this with the sentiment from our Twitter dataset. By creating daily sentiment average of tweets associated with terms “Clinton” and “Trump”, we want to see which candidate has the better sentiment score and thus favorable opinion among Twitter users.

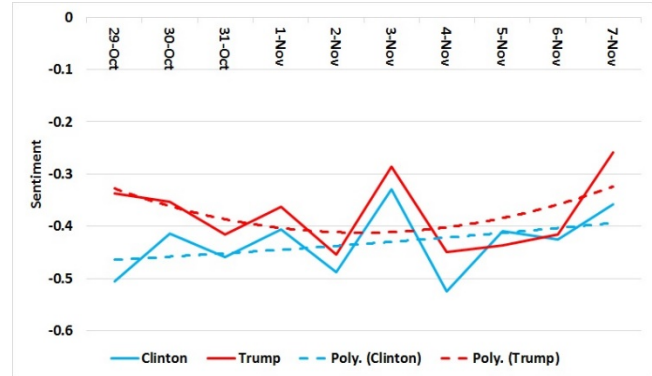
**Table 2: Polls results for elections. [38]**

DATE	POLL	CLINTON	TRUMP
6-Nov	Economist/YouGov	49	45
5-Nov	Fox News	48	44
5-Nov	Bloomberg	46	43
5-Nov	ABC News/Wash Post	49	46
5-Nov	NBC News/SM	51	44
5-Nov	CBS News	47	43
5-Nov	IBD/TIPP	43	42
5-Nov	Monmouth	50	44
5-Nov	LA Times/USC	43	48
4-Nov	NBC News/WSJ	48	43
3-Nov	Reuters/Ipsos	44	40
2-Nov	Fox News	46	45
2-Nov	McClatchy/Marist	46	44
31-Oct	CBS News/NYT	47	44
31-Oct	Economist/YouGov	48	45
30-Oct	Gravis	50	50
29-Oct	NBC News/SM	51	44

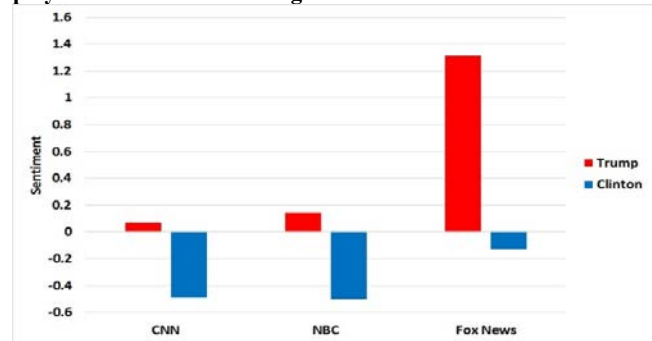
In the subsequent sentiment analysis of our data, we discover Donald Trump leading Hillary Clinton. Using only first 10 days of data, from Oct 29<sup>th</sup> to Nov 7<sup>th</sup>, we observe that tweets containing “Trump” only have a lower average negative sentiment than tweets mentioning only “Clinton”. Fig. 3, depicts this finding. As the sentiment is fluctuating for both candidates during these ten days, we have created a polynomial trend line of degree 2 to make this sentiment difference more observable. We can see here that sentiment trend associated with Donald Trump is consistently less negative than Hillary Clinton. This finding is contrary to majority of the pre-election polls predicting a Hillary Clinton victory.

Although sentiment remains negative for both candidates overall across various measures, we are nonetheless able to discover a loyal base for Donald Trump on Twitter that showed a positive sentiment towards him regardless of the topic or news source. Fig. 4 shows the sentiment of this base towards both candidates when sharing news from 3 of the most popular news organizations mentioned in our database.

**5.1.4 Sentiment based on account reputation:** A study by Chu et al. [17] on detection of automation on Twitter created a parameter of user reputation to discriminate between human and automated activity. The study examined 500,000 Twitter users for this purpose. It stated among other findings that a measure of account reputation can be used to distinguish human activity from bot.



**Figure 3: Sentiment of tweets for both candidates along with polynomial trend-line of degree 2.**



**Figure 4: Sentiment of tweets for Trump loyal base mentioning popular news organizations.**

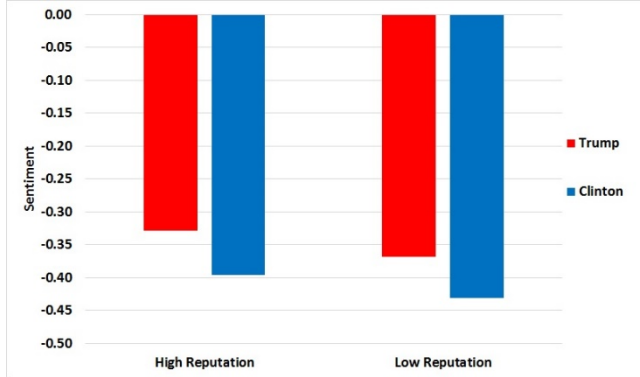
In order to do this, they formulated that humans have a near equal number of followers and friends. This is not true for bots as in-order for bots to create a bigger network, the common strategy is to follow a large number of users in hope that some of these users will follow them back. Thus an account reputation is defined as following equation:

$$\text{Account Reputation} = \frac{\text{num of followers}}{\text{num of followers} + \text{num of friends}}$$

It is noted that for celebrity accounts this measure comes close to 1. This can be observed for many popular accounts such as Katy Perry, who as of January 2017, has 95.3 million followers but is only following 186 other Twitter users.

The study concluded that around 60 percent of bots had a reputation count of less than 0.5 [17]. We have decided to use account reputation as a factor as well in our study to observe the sentiment for both terms “Trump” and “Clinton”. All users are divided into high reputation and low reputation accounts. Accounts having a reputation score of greater than 0.5 are defined as high reputation while those with a reputation of less than or

equal to 0.5 are defined as low reputation. Fig. 5 below shows sentiment for both terms based on account the tweet account reputation.



**Figure 5: Sentiment of tweets for both candidates based on user reputation.**

Although both account types have overall negative sentiment for both terms, in general the negative sentiment is less for “Trump”. The account reputation measure is further utilized in our data analysis as well to better understand user behavior.

## 5.2 Hypotheses Testing

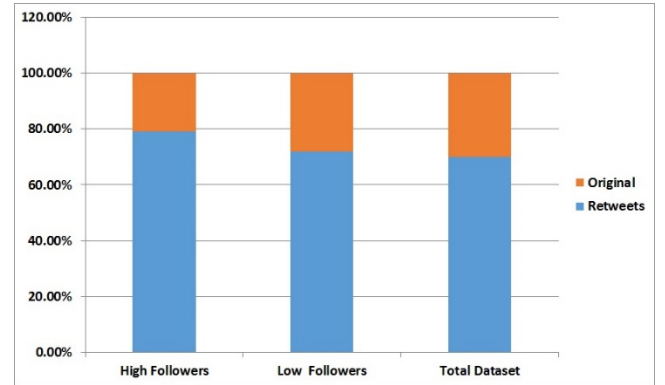
We now present the results obtained through analysis of our data that can help us to refute or verify our initially developed hypotheses presented in the last section.

**5.2.1 Content Creation:** Studies have shown that number of retweets on Twitter usually ranged from 1.44% to 19.1% of all messages [14]. Majority of the tweets created on Twitter are never retweeted. However, this is contrary to our finding where in our dataset number of retweets is as high as 70%. Majority of the users in our dataset, commenting on the elections are not creating any new content but rather reusing the information already present. Furthermore, 100 most retweeted tweets appear 7081 times in our dataset at an average of almost 71 times each. This behavior of high retweets is present amongst all user groups regardless of the number of followers they have or the frequency with which they tweet. For top 100 users with followers and friends greater than 10,000, 79% of their tweets are retweets while for the bottom 100 users with 1 follower and 1 friend this number is 72%. The statistics are shown in Fig. 6.

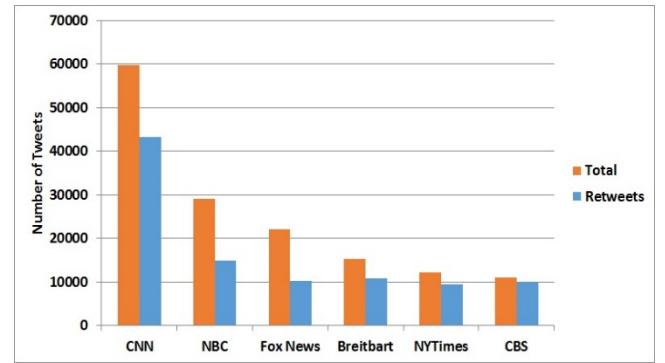
Hence H1 is supported by our dataset: majority of users are not creating new content but retweeting information amongst their network. This reusing of the information corroborates the study conducted by McKenna et al. [22] who discovered that 87% of political bloggers provide links to news articles and other blogs in their blog posts. Fig. 7 shows the most popular news outlets mentioned in all tweets and retweets. However, there is one small group of users who remain an exception to this rule. Users with a following of greater than 50,000 and a reputation of over 0.8, have substantially fewer retweets as a percentage of their total tweets. Fig. 8 depicts this high following, high reputation user’s content creation pattern.

**5.2.2 Twitter Activity Based on Following:** In terms of message framing by utilizing hashtag (#) to reach a broader

audience, we discover similar tweeting behaviors for both heavy and light users. We have defined heavy and light users based on their following and activity on Twitter.



**Figure 6: Tweeting behavior of users.**



**Figure 7: Most popular news outlets discussed in tweets.**

Heavy users and light users are defined as those with large number of followers and high tweeting versus fewer followers and low tweeting respectively [8]. In our dataset, heavy users use hashtag only marginally more than light users. The top 100 users in terms of followers in our dataset have hashtag in almost 21% of their tweets while bottom 100 users in terms of following have a hashtag in only 15% of their tweets. Table 3 shows the breakdown of these two user types.

Both of these user types show somewhat similar hashtag usage behavior when compared with the overall percentage of tweets containing hashtag. Hence H2 is not supported by our dataset: users with large following and heavy tweeting use hashtag (#) only marginally more than users with few followers and tweets. Thus in context of elections 2016, heavy users have not used framing significantly more than light users.

**Table 3: Hashtag (#) use by heavy and light users.**

	Followers on Avg.	Friends on Avg.	Tweets per user	Tweets with #	Percent with #
Heavy users	273,627	37,020	4.24	2.56	21.23%
Light user	1	55	1.75	0.58	15.43%



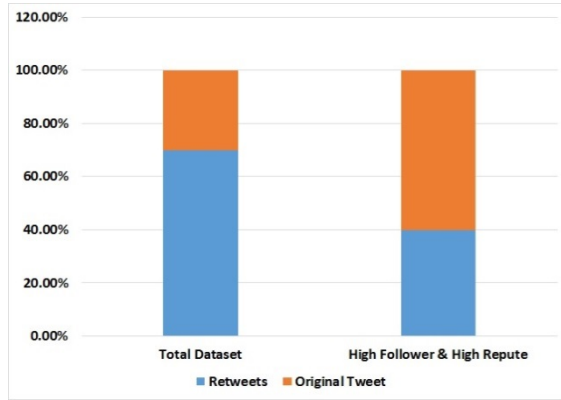


Figure 8: Tweeting behavior of users with high number of followers and high reputation.

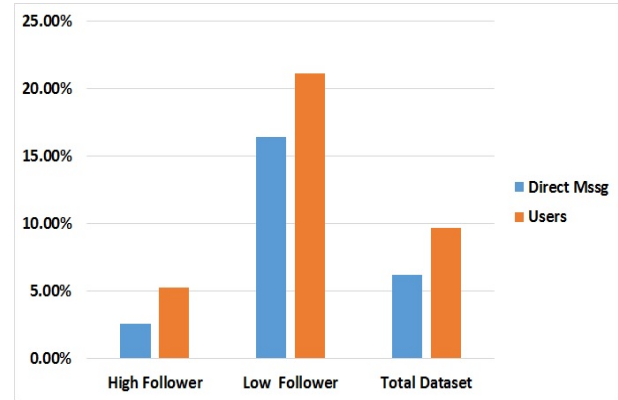


Figure 9: Comparison of direct messages for users with high and low following.

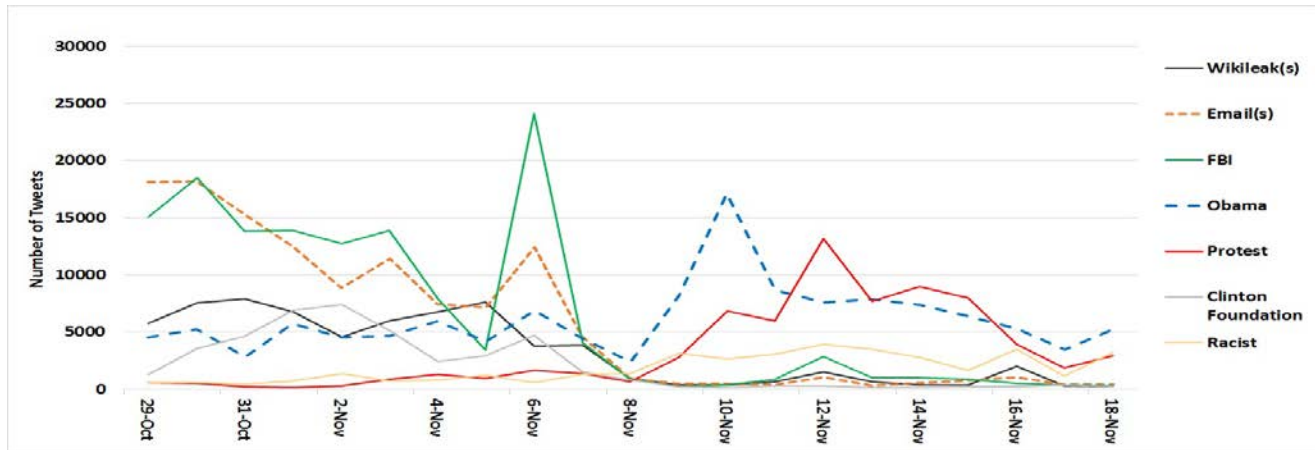


Figure 10: Frequency diagram of popular discussion topics.

**5.2.3 Message Targeting:** In order to test message targeting hypothesis, we look at all the tweets in our database that start with the expression “@”. Other studies have also utilized a similar approach to gauge interactive discussions on Twitter [10]. We discover that few tweets are targeted to other users directly and this low one-to-one interaction with other users is significantly less for users who have a large number of followers. For the entire dataset, users engaging in direct messaging are only 9.66% while direct messages account for only 6.17% of all tweets. This number is lower than that of 10% claimed by other studies analyzing political discussions on Twitter in context of elections [10]. Fig. 9 shows the percentage of direct messaging in total dataset and according to number of user followers.

Hence, we can state that H3 is not supported by our dataset: users on Twitter do not engage in one-to-one discussions with other users and primarily use the platform to simply post their opinions.

**5.2.4 Reflection of current news events:** We have created a word cloud of the most popular terms used in our dataset. We find that different terms are popular before and after the elections. Figure 10 shows popular terms and their occurrence on each day. Here we can observe that some terms are popular prior to the

Election Day and became relatively obscure post elections. Hence for example *WikiLeaks*, *Emails* and *FBI* are popular discussion topics before 8<sup>th</sup> November but became irrelevant later on as the candidate associated with them lost the elections. On the other hand, term such as *Protest* is infrequent prior to the Election Day but becomes popular later on due to the street protests that ensued post elections. Finally, terms such as *Obama* remained relatively frequent during this entire period. This is due to President Obama’s presence in the news for campaigning before elections and helping president elect in transition post elections.

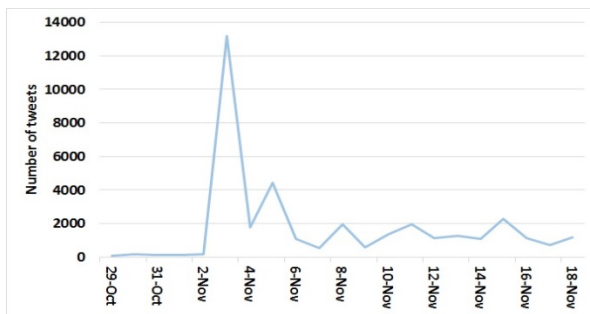
These popular terms indicate real world events, discussions and news appearing real time in Twitter conversations. We can also detect this by closely observing the peaks of frequent terms in Fig. 10. For example, *FBI* and *Email* both peak sharply on 6<sup>th</sup> November. This is the very same day on which FBI made the announcement that they have completed their review of emails and do not recommend any action against Hillary Clinton.

Similarly, we can also spot the term *Obama* abruptly peaking on 10<sup>th</sup> November, which is the day President Barack Obama met President-Elect Donald Trump in the white house, 2 days after his election victory. This was the most important news item for that day and we can see it reflected on the Twitter conversations occurring that day.

**Table 4: Three most frequent daily terms.**

Most Popular terms			
	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>
29-Oct	Email(s)	FBI	WikiLeaks
30-Oct	FBI	Email(s)	WikiLeaks
31-Oct	Email(s)	FBI	WikiLeaks
1-Nov	FBI	Email(s)	Clinton Foundation
2-Nov	FBI	Email(s)	Clinton Foundation
3-Nov	FBI	Email(s)	WikiLeaks
4-Nov	FBI	Email(s)	WikiLeaks
5-Nov	WikiLeaks	Email(s)	Obama
6-Nov	FBI	Email(s)	Obama
7-Nov	Email(s)	Obama	FBI
8-Nov	Obama	Email(s)	WikiLeaks
9-Nov	Obama	Racist	Protest
10-Nov	Obama	Protest	Racist
11-Nov	Obama	Protest	Racist
12-Nov	Protest	Obama	Racist
13-Nov	Obama	Protest	Racist
14-Nov	Protest	Obama	Racist
15-Nov	Protest	Obama	Racist
16-Nov	Obama	Protest	Racist
17-Nov	Obama	Protest	Racist
18-Nov	Obama	Racist	Protest

Finally, *Protest* peaks on 12<sup>th</sup> November. By this time protests are being held against the election results in many major cities of the United States and remained daily headline news item. [Table 4](#) displays the top 3 most popular terms for each day from 29<sup>th</sup> October to 18<sup>th</sup> of November. We can further see a proof of this in trend for the term *Melania*. Melania Trump gave her first major campaign speech in Pennsylvania on 3<sup>rd</sup> November 2016. She was one of the most mentioned terms on Twitter that day, remaining relatively obscure before and after that event. [Fig. 11](#) below plots this trend. Hence we can state that H4 is supported by our data analysis and we can state that frequency of popular terms in Twitter discussions can indeed be utilized to identify significant real world events and news that are taking place in relation with the elections 2016.

**Figure 11: Daily frequency chart for tweets mentioning "Melania".**

## 6. DISCUSSION

### 6.1 Candidate Sentiment on Twitter

Although both candidates had negative sentiment overall, tweets mentioning "Trump" did comparatively better than those mentioning "Clinton". This remained consistent in our analysis

regardless of measures such as overall sentiment based on user reputation or news organization mentioned in tweets.

This finding is in contrast with most polls conducted during this time period showing Hillary Clinton leading Donald Trump, predicting a Hillary Clinton win. [Table 2](#) displays the data from these polls. [Fig. 3](#) above displays the average daily sentiment for both candidates along with polynomial trend lines. It indicates a better daily sentiment trend for Donald Trump 10 days of data prior to the Election Day. [Fig. 4](#) shows the sentiment of a loyal base of Twitter users who displayed positive sentiment towards Donald Trump regardless of the news source.

There has been much debate on the usability of Twitter sentiment to gauge public opinion and various studies have claimed that Twitter sentiment analysis is indeed a good measure of public opinion [10]. We believe that in case of US Elections 2016, Twitter sentiment proves to be an accurate indicator of real world public sentiment.

### 6.2 User Behavior on Twitter during Elections 2016

According to our exploratory data analysis, we believe that although Twitter is a popular tool for political discussions and debate, a very small number of users dominate this platform. [Table 1](#) displays this dominance, where almost 52% of all tweets in our dataset originate from around 9% of the users while over 69% users accounted for only 26% of total tweets. Support for hypothesis 1 in our dataset, further solidifies this conclusion where 70% of all tweets in our dataset are retweets. Hence, majority of users are simply following trends and discussions through retweets. Most users are passive and do not actively participate in conversations by expressing their original thoughts. One exception to this rule however is the small group of users who have high reputation and have large number of followers. From [Fig. 6](#), we can see that these users retweet only 40% of the times while 60% of time they are posting original tweets. These users actively participate in discussions by voicing their own thoughts and opinions through Twitter.

Rejection of Hypothesis 3 indicates that in context of US elections 2016, Twitter was primarily used to spread political opinion and not to discuss these opinions with other users. Only 6.17% of messages in our sample are direct messages. This finding of using Twitter for broadcasting rather than engagement for political conversations is in contrast with other Twitter studies conducted during elections that claim people use Twitter to engage in interactive discussions [14].

### 6.3 Twitter for Live Coverage of Important Election Events

Several studies have looked at Twitter streaming data as a source for identifying current news and real world events [22, 28]. They have concluded that Twitter trends are usually the most important events of the day and can be used to predict headline news. We believe that this assumption holds true in context of US Elections 2016. As discussed in result of hypothesis 4 in previous section, the most popular daily terms in our sample set are usually the most important news stories for the day. Hence in context of Elections 2016, Twitter remained a good proxy to identify the most significant daily events that are taking place.

## 7. FUTURE WORK & CONCLUSIONS

For this data study, we have analyzed approximately 3 million Twitter messages associated with US Presidential elections of 2016. These tweets were collected over a period of 21 days, before and after the elections that were held on November 8th. The tweets were filtered based on their messages mentioning either the two major Presidential candidates (Hillary Clinton and Donald Trump) or the term Election 2016.

The primary purpose of this study was to understand the nature of political discourse that took place on Twitter during the elections in terms of sentiment, frequently mentioned terms, and popular tweets / retweets. User attributes such as number of followers and friends, duration of activity on Twitter, number of messages to date and reputation were also utilized in this regard. A user behavior model was developed for this purpose and various hypotheses were tested to understand characteristics of political discussions on Twitter. It was observed that the overall sentiment of the tweets was negative. This was true for both candidates and can be taken as indication to the bitter nature of the recently concluded elections. It was also observed that little original content was created by users during discussions and most were rather retweeting.

In future we would like to expand this study into other regions and countries of the world where Twitter is gaining popularity as a political campaigning tool and where politicians and public are turning towards social media for political broadcasts and information. We would also expand this study into areas other than general elections and from politicians to state institutions. The lessons learned can be used to gauge the sentiment of public discussion regarding not only political statements but also towards broadcasts or information provided by state institutions and functionaries to communicate with public.

## REFERENCES

- [1] SentiStrength. <http://sentistrength.wlv.ac.uk/>
- [2] Calderon, Nadya A., et al. "Mixed-initiative social media analytics at the World Bank: Observations of citizen sentiment in Twitter data to explore" trust" of political actors and state institutions and its relationship to social protest." In *Big Data (Big Data)*, 2015 IEEE International Conference on, pp. 1678-1687. IEEE, 2015.
- [3] Ferrara, Emilio, and Zeyao Yang. "Measuring emotional contagion in social media." *PloS one* 10, no. 11 (2015): e0142390.
- [4] Ferrara, Emilio, and Zeyao Yang. "Quantifying the effect of sentiment on information diffusion in social media." *PeerJ Computer Science* 1 (2015): e26.
- [5] An In-Depth Look at the Most Active Twitter User Data. © 2007-2017 SYSOMOS. <https://sysomos.com/inside-twitter/most-active-twitter-user-data>
- [6] First Monday. © *First Monday*, 1995-2017. ISSN 1396-0466. <http://journals.uic.edu/ojs/index.php/fm/article/view/2793/2431#p4>
- [7] Statista. <https://www.statista.com/statistics/282087/number-of-monthly-active-twitter-users/>, 2016
- [8] Paul Mozur and Mark Scott. The New York Times. (2016, Nov 17). [http://www.nytimes.com/2016/11/18/technology/fake-news-on-facebook-in-foreign-elections-thats-not-new.html?\\_r=0](http://www.nytimes.com/2016/11/18/technology/fake-news-on-facebook-in-foreign-elections-thats-not-new.html?_r=0)
- [9] Schwartz-Asher, Daphna, Soon Ae Chun, and Nabil R. Adam. "Social Media User Behavior Analysis in E-Government Context." In *Proceedings of the 17th International Digital Government Research Conference on Digital Government Research*, pp. 39-48. ACM, 2016.
- [10] O'Connor, Brendan, Ramnath Balasubramanyan, Bryan R. Routledge, and Noah A. Smith. "From tweets to polls: Linking text sentiment to public opinion time series." *ICWSM* 11, no. 122-129 (2010): 1-2.
- [11] Heverin, Thomas, and Lisl Zach. "Twitter for city police department information sharing." *Proceedings of the American Society for Information Science and Technology* 47, no. 1 (2010): 1-7.
- [12] Bertot, John C., Paul T. Jaeger, and Justin M. Grimes. "Using ICTs to create a culture of transparency: E-government and social media as openness and anti-corruption tools for societies." *Government information quarterly* 27, no. 3 (2010): 264-271.
- [13] Lorenzi, David, et al. "Utilizing social media to improve local government responsiveness." In *Proceedings of the 15th Annual International Conference on Digital Government Research*, pp. 236-244. ACM, 2014.
- [14] Tumasjan, Andranik, Timm Oliver Sprenger, Philipp G. Sandner, and Isabell M. Welpe. "Predicting Elections with Twitter: What 140 Characters Reveal about Political Sentiment." *ICWSM* 10 (2010): 178-185.
- [15] Gayo Avello, Daniel, Panagiotis T. Metaxas, and Eni Mustafaraj. "Limits of electoral predictions using twitter." In *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media. Association for the Advancement of Artificial Intelligence*, 2011.
- [16] Metaxas, Panagiotis T., Eni Mustafaraj, and Dani Gayo-Avello. "How (not) to predict elections." In *Privacy, Security, Risk and Trust (PASSAT) and 2011 IEEE Third International Conference on Social Computing (SocialCom)*, pp. 165-171. IEEE, 2011.
- [17] Chu, Zi, Steven Gianvecchio, Haining Wang, and Sushil Jajodia. "Detecting automation of twitter accounts: Are you a human, bot, or cyborg?." *IEEE Transactions on Dependable and Secure Computing* 9, no. 6 (2012): 811-824.
- [18] <https://dev.twitter.com/streaming/overview>
- [19] Romero, Daniel M., Brendan Meeder, and Jon Kleinberg. "Differences in the mechanics of information diffusion across topics: idioms, political hashtags, and complex contagion on twitter." In *Proceedings of the 20th international conference on World wide web*, pp. 695-704. ACM, 2011.
- [20] Thelwall, Mike, et al. "Sentiment strength detection in short informal text." *Journal of the American Society for Information Science and Technology* 61, no. 12 (2010): 2544-2558. Becker, Hila, Mor Naaman, and Luis Gravano. "Beyond Trending Topics: Real-World Event Identification on Twitter." *ICWSM* 11 (2011): 438-441.
- [21] McKenna, Laura, and Antoinette Pole. "What do bloggers do: an average day on an average political blog." *Public Choice* 134, no. 1-2 (2008): 97-108.
- [22] Thelwall, Mike, Kevan Buckley, and Georgios Paltoglou. "Sentiment in Twitter events." *Journal of the American Society for Information Science and Technology* 62, no. 2 (2011): 406-418.
- [23] Zeng, Xiaohua, and Liyuan Wei. "Social ties and user content generation: Evidence from Flickr." *Information Systems Research* 24, no. 1 (2013): 71-87.
- [24] Bruns, Axel, and Jean Burgess. "Researching news discussion on Twitter: New methodologies." *Journalism Studies* 13, no. 5-6 (2012): 801-814.
- [25] Marwick, Alice E. "I tweet honestly, I tweet passionately: Twitter users, context collapse, and the imagined audience." *New media & society* 13, no. 1 (2011): 114-133.
- [26] Kwak, Haewoon, Changhyun Lee, Hosung Park, and Sue Moon. "What is Twitter, a social network or a news media?." In *Proceedings of the 19th international conference on World wide web*, pp. 591-600. ACM, 2010.
- [27] Cheng, Justin, Lada Adamic, P. Alex Dow, Jon Michael Kleinberg, and Jure Leskovec. "Can cascades be predicted?." In *Proceedings of the 23rd international conference on World wide web*, pp. 925-936. ACM, 2014.
- [28] Broersma, Marcel, and Todd Graham. "Social media as beat: tweets as a news source during the 2010 British and Dutch elections." *Journalism Practice* 6, no. 3 (2012): 403-419.
- [29] Parmelee, John H. "Political journalists and Twitter: Influences on norms and practices." *Journal of Media Practice* 14, no. 4 (2013): 291-305.
- [30] Wallsten, Kevin. "Microblogging and the news: political elites and the ultimate retweet." *Political Campaigning in the Information Age* (2014): 128-147.
- [31] Jungherr, Andreas. "Twitter in politics: a comprehensive literature review." Available at SSRN 2402443 (2014).
- [32] Abroms, Lorien C., and R. Craig Lefebvre. "Obama's wired campaign: Lessons for public health communication." *Journal of health communication* 14, no. 5 (2009): 415-423.
- [33] Baumgartner, Jody C., Jenn Burleson Mackay, Jonathan S. Morris, Eric E. Otenyo, Larry Powell, Melissa M. Smith, Nancy Snow, Frederic I. Solop, and Brandon C. Waite. *Communicator-in-chief: How Barack Obama used new media technology to win the White House*. Edited by John Allen Hendricks, and Robert E. Denton Jr. Lexington Books, 2010.
- [34] Glassman, Matthew, Jacob R. Straus, and Colleen J. Shogan. "Social networking and constituent communication: Member use of Twitter during a two-week period in the 111th Congress." *Congressional Research Service, Library of Congress*, 2009.
- [35] Golbeck, Jennifer, Justin M. Grimes, and Anthony Rogers. "Twitter use by the US Congress." *Journal of the American Society for Information Science and Technology* 61, no. 8 (2010): 1612-1621.
- [36] Highfield, Tim, Stephen Harrington, and Axel Bruns. "Twitter as a technology for audiencing and fandom: The# Eurovision phenomenon." *Information, Communication & Society* 16, no. 3 (2013): 315-339.
- [37] Hausteine, Stefanie, et al. "Tweets as impact indicators: Examining the implications of automated "bot" accounts on Twitter." *Journal of the Association for Information Science and Technology* 67, no. 1 (2016): 232-238.
- [38] BBC News. Copyright © 2017 BBC. Nov 7, 2016. <http://www.bbc.com/news/election-us-2016-37450667>