

Predicting How People Vote From How They Tweet

A Senior Project submitted to
The Division of Science, Mathematics, and Computing
of
Bard College

by
Rao Vinnakota

Annandale-on-Hudson, New York
May, 2019

Abstract

In 2016 Donald Trump stunned the nation and not a single pollster predicted the outcome. For the last few decades, pollsters have relied on phone banking as their main source of information. There is reason to believe that this method does not present the complete picture it once did due to several factors—less landline usage, a younger and more active electorate, and the rise of social media. Social media specifically has grown in prominence and become a forum for political debate. This project quantitatively analyzes political twitter data and leverages machine learning techniques such as Naive-Bayes to model election results. Early results are promising, and a true evaluation of the model will come from testing in future elections.

Contents

Abstract	iii
Dedication	vii
Acknowledgments	ix
1 Introduction	1
1.1 Political Polling - A Short Overview	2
1.2 Social Media, Twitter, and Politics	3
1.3 Big Data	3
1.4 Project Roadmap	3
2 Building a Data Pipeline	5
2.1 Streaming Tweets	7
2.2 Databases	8
2.2.1 Overview	8
2.2.2 MongoDB	9
2.3 Data Collection	11
3 Sentiment Analysis	13
3.1 Lexicons	14
3.2 VADER	14
3.2.1 Construction	15
3.2.2 Effectiveness	16
3.2.3 Using VADER	17
3.3 Building a Sentiment Profile	18
4 Exploring the Dataset	21
4.1 Races Collected	21
4.2 Comparing Sentiment	23

4.2.1	Winners vs Losers	23
4.2.2	Democrats vs Republicans	27
4.3	National Trends	30
4.3.1	Exploring the Issues	30
4.3.2	Evaluating Twitter's Response	31
5	Modeling Elections	35
5.1	Approach	35
5.2	A Machine Learning Overview	36
5.3	Understanding the Naive Bayes Classifier	36
5.3.1	Bayes Theorem	37
5.3.2	Naive Bayes Classifier	37
5.3.3	Why Naive Bayes?	38
5.4	Implementation Naive Bayes in Python	39
5.5	Evaluating the Model	40
5.5.1	Criteria	41
5.5.2	Performance	42
6	Conclusion	45
6.1	Project Achievements	45
6.2	Future Work	45
	Appendices	47
A	The How-To Guide	47
A.1	Installation	47
A.2	Gathering Tweets	48
A.3	Building Datasets and Models	49
B	Code	51
	Bibliography	60

Dedication

I dedicate this project to my family. My younger brother Abhi for his unwavering faith, my mother Neelima for never letting me quit, and my father Bapi for being there for me, even when I didn't know I needed it.

Acknowledgments

There are several people without whom this project wouldn't have been possible.

To the Bard Math Department. Being a math major feels like being party of a large, crazy family. All of your endless support even if I wasn't your direct advisee can never be understated. I'm proud to be a member of this incredible community. Thank you.

To my baseball coaches and teammates. This senior season has been a wild ride and I will never forget walking to the gym at 6 am together.

To all my friends both new and old. I wouldn't be the person I am today without you. Thank you.

To my advisor Stefan Mendez-Diez. Thank you for giving me the chance to explore an incredible topic and helping me every step of the way. From computing space, to meetings, to frantic revision, I couldn't have done this project without you. Thank you.

To my family. I wouldn't be here without your endless love and support. Abhi, Mom, and Dad, I somehow made it to graduation. Words just aren't enough. Thank You.

1

Introduction

Donald Trump stunned the world in 2016. Mr. Trump ran an unorthodox and polarizing campaign. Despite suffering multiple public relations setbacks, Mr. Trump still won on election night.

A significant part of Mr. Trump's campaign was his use of social media. Trump used accounts on Twitter and Facebook as a voice to attack candidates, disseminate campaign information, and engage with voters. Mr. Trump's tweets got significant engagement with respondents voicing both support and opposition.

The large disconnect in the 2016 Presidential election was that public polling in 2016 didn't predict the upset. In fact, it didn't come particularly close. The latest tracking polls before election day showed Clinton in the lead with a margin of anywhere from +1 to +6 [7]. Most pollsters believed that Trump had no shot, and would concede early in the night. Obviously, that didn't happen

Tweets following the election showed Trump supporters crowing that victory was inevitable, and that Trump's social media interaction demonstrated it. Was Trump's potential victory present on social media for all to see? If it was, could pollsters have leveraged it to make a more accurate forecast?

Those two questions inspired this project. The goal is to quantitatively analyze social media to model election outcomes. In layman's terms, can you predict how people vote from how they tweet? To do this, tweets were collected, analyzed for opinion, and then used to model with machine learning techniques.

1.1 Political Polling - A Short Overview

The earliest iteration of polling was a straw poll conducted in Harrisburg, Pennsylvania which showed Andrew Jackson leading John Quincy Adams in the 1824 Presidential Election. Journalists from The *Harrisburg Pennsylvanian* asked local voters which they preferred. Jackson winning both the polls and the election created legitimized the poll. Polls continued to be local until the early 20th century when Literary Digest sent out postcards to reader asking for their pick of the presidential election. The Literary Digest correctly predicted 5 straight elections and opinion polls reached the national stage.

In the advertising boom of post-World War II, most polling shifted to telecommunications. With telephones, there was an ability to reach a wider audience, ask deeper questions. Pollsters could now build a profile of a voter, and forecast more accurate results. Accompanying the advance in polling was a focus on sampling. Pollsters couldn't call every single house in the United States. Instead, groups of people or "samples" that would accurately represent the country were called. Responses from the sample were extrapolated to cover the whole.

Today, the digital age and social media have made polling's barrier of entry low. A short scroll through Twitter will show hundreds of "opinion polls". Youtube videos' likes and dislikes function as an informal poll. The Algorithm will often promote videos that are more "well-liked". But, professional political polling is still done along the same lines. A population is sampled, called, and profiled. Profiling technique has improved. Today, forecasters like 538's Nate Silver, will build complex profiles of voters in various regions. But, the source of the information and the method of obtaining the information remains the same.

1.2 Social Media, Twitter, and Politics

The internet's first foray into politics was in 2008. Barack Obama established that it was acceptable for a candidate to leverage the internet as a resource, using the internet as a tool for both message and money. The website MyBarackObama.com helped Obama set records in grassroots donations and mobilizations.

Social websites from the forums of the early 2000s to social media platforms that are popular today are the most visited websites on the internet. These forums encourage conversation, including the discussion of politics. Especially for younger voters who don't own any landlines, these forums are the few places they express political opinions.

This project will specifically look at political opinions expressed on the microblogging platform Twitter. What sets it apart from blogs or other political forums is their ease of use. The longform opinions that are on forums and blogs are short and un-edited in Tweets. Twitter's mandatory message format - less than two hundred and eighty characters - makes sure posts are similar in style, and are far easier to analyze.

1.3 Big Data

A decade ago, attempting to gather, clean, and analyze millions of tweets would have been a fool's errand. It's with big data techniques and packages that this project is now possible. This presents my experiment with several advantages. Unlike traditional pollsters, the data collection process is entirely automated and stored easily.

1.4 Project Roadmap

Chapter 2 covers the data collection process. This goes through the automated tweet gathering, the storage and organization, as well as the optimization of the process.

Chapter 3 introduced and explains the idea of sentiment analysis. It also covers the specific algorithm the project uses: how the algorithm was created, its effectiveness, and its implementation.

Chapter 4 explores the data set that was covered. It breaks down sentiment and other statistics in both the House and the Senate. It covers the gap between winners and losers as well as the gap between candidates from both parties.

Chapter 5 defines the model used and demonstrates the underlying mathematics. It also covers the implementation of the model in Python, and evaluates the model's performance.

Chapter 6 reflects on the project and contemplates future work.

Below represents a high level overview of the project from collecting tweets to making predictions:

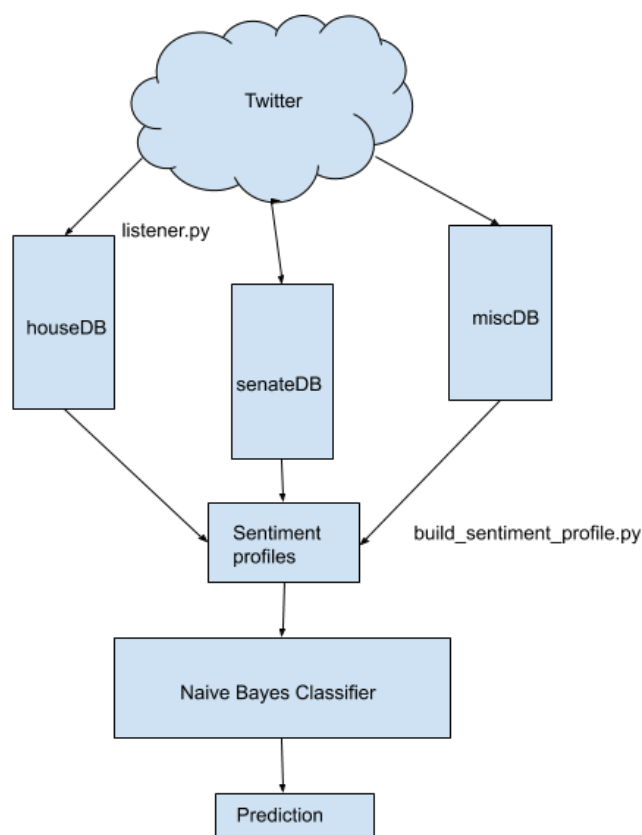


Figure 1.4.1. A high level overview of the project from gathering tweets to building sentiment profiles to making predictions from modeling

2

Building a Data Pipeline

Mathematical models like the one we intend to build run on large amounts of data. Just under three million tweets were collected for this experiment. Obviously that number is too large to be done manually. A script automated a pipeline that collected tweets matching keywords and saved them to be analyzed.

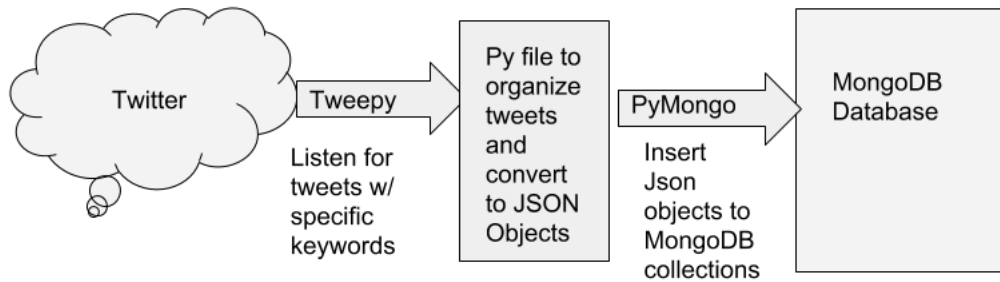


Figure 2.0.1. A high level overview of the tweet collection pipeline

The cloud in Figure 2.0.1 represents all of the tweets on Twitter at any given moment. The arrow represents the streamer. The streamer is a data structure that will monitor Twitter and download all tweets that match certain keywords. Those tweets that the streamer pulls need to be stored. To do that they need to be formatted into JSON or JavaScript Object Notation. When they are, they will be inserted into a database using the management software MongoDB.

To complete all those actions in python, we'll need to use an Application Package Interface (API). APIs are how scripts and other pieces of code can interact with software such as Twitter or MongoDB. Tweepy is the Twitter API and Pymongo is the Python API.

To summarize, the pipeline process can be distilled into three separate tasks.

1. Building a tweet streamer to “listen” to tweets based on keywords
2. Preprocessing the tweets into JSON notation
3. Storing processed tweets into a database running MongoDB

The python script that does this can be found in Appendix B. The explicit instructions on operating the script can be found in Appendix A.

2.1 Streaming Tweets

There are two ways to gather tweets using the Python Twitter API. One is to query a profile or a search term [15]. This method is not preferred for several reasons. First, Twitter has put a hard cap on the number of tweets each query can return at one time. To query every single tweet about a subject will need repeated queries of the maximum amount [15]. Considering that the actual number of tweets is unknown, this is a time-consuming loop that has no estimate-able end. Next, Twitter's search function naturally promotes tweets that have more engagement (likes, comments, and retweets). Collecting only those tweets means we put together an incomplete and biased data set.

The other way to gather tweets is to stream tweets [15]. A Tweepy stream monitors live twitter content and returns all tweets that contain a given set of keywords while the stream is active. Streaming has the disadvantage of gathering only live tweets - it does not search tweets from before it started running. , streaming has the advantage of gathering all tweets that meet the keyword requirement irrelevant of engagement level.

Tweepy's Streamer is an abstracted object [15]. That means that it provides the outline for how a streamer should work, but the specifics are left up to the user. The streamer has functions or methods that need to be defined so that the streamer understands what to do with the tweets it gathers.

The methods that need to be defined here are `on_connect`, `on_error`, and `on_data`. The first two are simple. `On_connect` is used to connect to the internet, and it confirms when the stream is live. If the stream crashes - which can happen for several reasons - `on_error` informs the user why the stream stopped. This can happen for many reasons which include: exceeding the rate limit (tweets streamed per second), a faulty internet connection, and more.

The last method `on_data` is where the streamed tweets are received and processed. Streamed tweets are returned by a streamer in a python dictionary format with the keys as attributes (user, text, etc.). As mentioned earlier, MongoDB only takes in objects in JSON format. The dictionaries need to be converted - which is what happens in the preprocessing. This is a process that happens simultaneously so that there is no need to intermediately store the tweets before they're ready to be stored in the database.

2.2 Databases

The `on_data` method from section 2.1 is inserting tweets into a database. In this section, we'll explore what databases are and how they're specifically applied in this project.

2.2.1 Overview

The term database refers to any structure that electronically stores data. In extension, the Database Management System is the program that allows users to interact with the database. In this paper, when the term database is used, it refers to both the structure and the system.

There are several types of databases. These types are classified by the way the data is organized inside the structure. The relational database model uses tables with rows and columns to store information. Relational databases use Structured Query Language or SQL to organize, extract, and insert data. This has become the industry default which leads us to refer to non-relational databases colloquially to as "NoSQL". Non-relational use an object-oriented approach with different query languages. Each individual entry is an object all by itself and searches are done through a tagging system [4].

Databases offer several advantages that other storage options like Excel or CSVs will not. For one, data integrity is higher. An example of better integrity is that one field can't have entries of different types across all objects [4]. In Excel, a user field could have both a string or an integer and the program would except it without fault. In a database, the error would be spotted and left to be changed at the user's discretion.

Another advantage is the ability to scale a database. They are built to handle gigabytes and even terabytes of data. Excel by contrast is limited to around one million rows [4]. The ability to quickly search and sort through millions of tweets as well as the flexibility to manipulate said data is something that isn't available to Excel.

The question that remains is which database to use here? The database will have to deal with millions of tweets, be able to search terms quickly, and handle inconsistencies in object size. After all, some Tweets have more information than others (location, images, etc.) It's been mentioned several times earlier in this chapter, but a non-relational database running MongoDB was chosen. The next section will give an overview of what MongoDB is and how it's the best choice for this project.

2.2.2 *MongoDB*

MongoDB is a popular open source non-relational database management system. It has an efficient query language as well as a well-defined API Pymongo. In a MongoDB database, each tweet is an object in JSON format. These objects are stored in collections and the collections are stored in a database [10]:

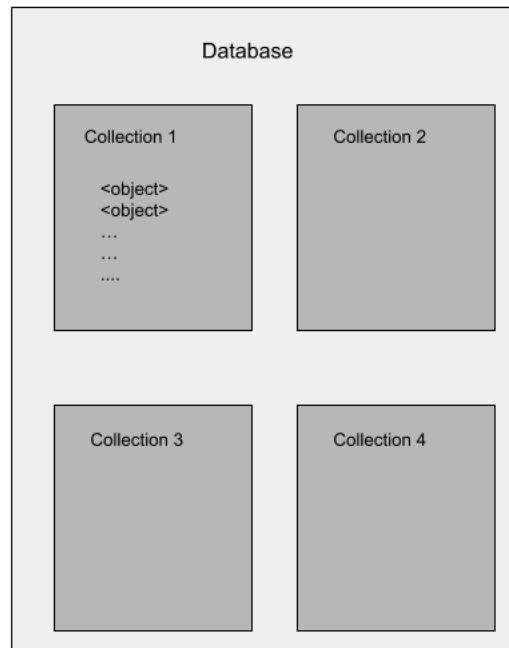


Figure 2.2.1. MongoDB’s storage organization. Tweets are stored individually in collections which are stored overall in a database.

These JSON objects have fields, populated with values. MongoDB allows a user to query collections using field title as well as field values. The queries can be exact matches, relational operators for integers, etc. For example, you can query tweets by date, or location (by exact match) or query all tweets that have more than a hundred retweets(relational). You can also query sub-matches. That means conditions within a field and not an exact match. Sub-matches might be searching for all tweets that contain the words “Donald Trump” or all tweets for all day rather than an exact time down to the second [10].

MongoDB has a mongo shell through which a user can interact with the databases. This is a convenient way to access quick information about a database. Information on using the mongo shell is in Appendix A. For longer and larger actions, we use the Python API Pymongo to write python scripts. Pymongo allows us to insert, edit, and query databases and manipulate the output.

For more specifics on Pymongo and setting up MongoDB, refer to Appendix A or the Mongo documentation [10].

2.3 Data Collection

The first step of implementing the pipeline outlined at the start of the chapter is building the streamer. The first step of streaming requires search terms to stream. The search terms were taken from house and senate races, as well as general terms that refer to the 2018 midterms. Search terms for Senate and House Races were tweets about the candidates as well as a hashtag referring to the race.

To direct the tweets streamed into a database, the `on_data` method is defined. This method is designed for the streamer to deal with the tweets it selects.

```
def on_data(self, data):
    try:

        client = MongoClient(MONGO_HOST)

        # Use twitterdb database. If it doesn't exist, it will be created.
        db = client.miscdb

        # Decode the JSON from Twitter
        datajson = json.loads(data)

        #grab the 'created_at' data from the Tweet to use for display
        created_at = datajson['created_at']

        #print out a message to the screen that we have collected a tweet
        print("Tweet collected at " + str(created_at) + " " + datajson['text'])
```

```
#insert the data into the mongoDB into a collection called twitter_search  
#if twitter_search doesn't exist, it will be created.  
db.misc.insert(datajson)  
  
except Exception as e:  
    print(e)
```

The tweet is converted into a JSON object and inserted into a database and collection of the users choice. This is a barebones definition. Sorting, editing, and adding information to tweets would all be defined here.

The timing of the election in November of 2018 was a hard deadline for streaming tweets. Streaming tweets was a higher priority than having them sorted into neat collections regarding races. Therefore, all the tweets were shoved into one collection without being sorted. This meant that the sorting process was delayed until data exploration and model building.

The data collection process ran through three separate scripts, each with a different set of search terms. Those three scripts fed into three different collections in one over-arching database. Together they assemble a diverse dataset that contains just under three million tweets. Chapter 4 explores the assembled data set.

3

Sentiment Analysis

Sentiment analysis is the computational approach to evaluating text, which includes determining both sentiment (positive/negative) and intensity [9]. There has always been an interest in the idea of sentiment analysis or objectively evaluating text for polarity. But, the field has seen a burst of activity in recent years due to a few key factors - the availability of text through the world wide web, the advance of big data, and the improvement in machine learning techniques.

The question of sentiment analysis's real effectiveness is often brought up. The general consensus is that a ceiling does exist for sentiment analysis, much like a ceiling exists for human analysis. There are statements where humans disagree on the sentiment. The ceiling is statements that have near unanimous agreements. If most humans can agree on a statements, it's not a stretch to assume that an algorithm could come to a similar conclusion.

This paper aims to use sentiment analysis to "poll" twitter users. The model built will make the trade-off specific targeted questions for a large sample size. To make use of collected data, it's important that the sentiment analysis algorithm used is effective and fast; after all, it has to evaluate three million tweets.

3.1 Lexicons

The basic idea of sentiment analysis is really simple. The algorithm recognizes positive and negative words, and then finds a method to average these words and determine an overall sentiment and intensity for the statement. But how does it know which words are positive and which words are negative? It uses a lexicon.

Lexicons are dictionaries. Large dictionaries which contain a list of positive and negative words, each with an intensity expressed. They allow algorithms to quickly ascertain the meaning of a word without having to re-define the word each time. Several well-received lexicons include the LIWC, Bing Lui, and the Harvard Inquirer. The lexicon that we'll use is known as VADER, which will be thoroughly examined later in the chapter.

Lexicons can be broken into two types - semantic orientation or polarity based, and sentiment intensity or valence based [12]. These general classifications cast a wide net and almost all lexicons fit in one or the other.

A well-known polarity-based lexicon is the Linguistic Inquiry and Word Count or LIWC (pronounced "Luke") [14]. LIWC organized words into seventy six categories and has about 900 words that are organized into two categories that are associated with emotion. It's the results of efforts by experts in psychology, sociology and linguistics. It's seen as arguably the best lexicon available, but does not recognize acronyms, abbreviations, or slang [9, 14]. These are all large parts of analyzing social text, and LIWC falls short.

The original valence-based lexicon is The Affective Norms for English Words or ANEW. It created ratings for over one thousand words of the english dictionary on a scale from one to nine, with a neutral midpoint at five. Words with a value greater than five will have a positive and pleasant meaning and negative and unpleasant otherwise [3].

3.2 VADER

The algorithm used to evaluate the data set is known as the Valence Aware Dictionary for sEntiment Analysis or VADER. It was developed by a pair of researchers CJ Hutto and Eric

Gilbert at the Georgia Institute of Technology for the express use of evaluating valence in social media [9].

Social media and microblogging in particular poses a challenge to sentiment analysis algorithms. This is the result of several problems. Earlier lexicons and algorithms were developed before social media's usage rates grew exponentially. They simply cannot handle the rate and volume that social media content is produced. Next, microblogging content is short. Twitter is limited to just two hundred and eighty characters. The average tweet is no more than a sentence or two. This limits context clues for traditional algorithms to pick up. Lastly, social media is famous for its abbreviations. Phrases like "LOL", "IRL", and "TBH" are commonplace on Twitter, but not commonplace on lexicons developed for traditional text. These are problems that VADER solves.

3.2.1 Construction

Hutto and Gilbert had four goals when they constructed the lexicon and the algorithm:

1. The algorithm performs effectively on social media text, but can still be generalized to all media sources.
2. The algorithm should have no training required from the user, but should still be based on a human created and cured lexicon.
3. The algorithm should be fast enough to use online.
4. There's no need to make a speed-performance tradeoff.

A basic overview of the VADER construction process is as follows. Hutto and Gilbert took the LIWC lexicon and made every single term in the lexicon a candidate for a new lexicon. Then, they proceeded to add social media specific terms which include abbreviations, phrases, and emojis as new candidates. These terms were trimmed and weighted in several stages until the final lexicon was produced.

The first stage of trimming used Wisdom-of-the-Crowd. WotC uses a large number of people to evaluate words in the lexicon. The large numbers combats the inexperience of the human raters themselves. Terms with a non-zero valence score were eliminated. Terms with a standard deviation greater than 2.5 which indicates that crowd couldn't reach a reasonable consensus were also eliminated. The terms were left had both a sentiment and an intensity. Terms such as "good" and "okay" while both positive, had a difference in score distinguishing their different meanings.

Their next stage used professional raters to identify parts of texts that can modify intensity despite not necessarily having an intensity themselves. These heuristics includes punctuation, verbal modifiers, and capitalization. Other heuristics identified were terms like "but" that signals a shift in sentiment and negated sentiment.

The last stage used professional raters to distinguish specific intensity modification of heuristics. This was done by taking a tweet, modifying either the syntax or the punctuation, and then randomly inserting these modifications into the dataset for the raters to evaluate. That way the algorithm could put a score on the change in intensity.

3.2.2 *Effectiveness*

The algorithm's effectiveness can only be evaluated by comparing it to the ground truth. The research team obtained ground truth statements for four domains - social media, movie reviews, product reviews, and editorials.

VADER's lexicon's competence was compared to several other well known lexicons such as LIWC, Bing Liu, Harvard Inquirer's, and more. This comparison made sure to use the same rule-based model for all of the lexicons. The three factors to evaluate the lexicons were:

- The correlation of computed intensity to the ground truth intensity.
- Precision - the number of items correctly classified divided by the number of total items.
- Recall - the number of items correctly classified divided by the number of items in the particular class

- F1 Score - the harmonic mean of the precision and the recall.

VADER outperformed some of the best lexicons in the industry currently. The results show that VADER's correlation was consistently higher, and it had the best F1 score for tweets. Hutto and Gilbert concluded that the reason was incorporating human heuristics into their model, essentially balancing a quantitative and qualitative approach.

3.2.3 Using VADER

VADER was written in Python, and is available to be implemented in Python. Below is how to use VADER to evaluate a single tweet.

```
from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer

analyzer = SentimentIntensityAnalyzer()
analyzer.polarity_scores("I'm so happy!")
```

This piece of code will return the following output:

```
{'neg': 0.0, 'neu': 0.318, 'pos': 0.682, 'compound': 0.648}
```

The package returns a series of scores, breaking down how the various components of the model's final score. Positive tweets are tweets that score greater than zero and negative tweets score less than. The intensity is determined the absolute value of the score.

A few interesting observations to note:

- Sarcasm is impossible for the algorithm to understand. To be fair, many humans also struggle with recognizing sarcasm in text.
- Any tweet that's an observation tends to be neutral, even if that observation has a positive or negative slant.

3.3 Building a Sentiment Profile

Each candidate in a given race has tweets about them. Using VADER, we can build a sentiment profile for them. The sentiment profile contains five statistics:

- Total Number of tweets
- Number of Neutral Tweets (Compound = 0)
- Number of Positive Tweets (Compound > 0)
- Number of Negative Tweets (Compound < 0)
- Average Compound Score

The sentiment profile will be the base of the features we'll use to model. Below is a snippet of code that calculates the profile for a given user.

```
sentiment = analyzer.polarity_scores(tweet)

if (sentiment['compound'] == 0):
    neu_tweets += 1

if (sentiment['compound'] > 0):
    pos_tweets += 1

if (sentiment['compound'] < 0):
    neg_tweets += 1

ave_score += sentiment['compound']

count += 1
```

This snippet exists in a for loop that iterates through all the tweets found about a candidate. For the full script, refer to the Appendix.

Here's an example of a sentiment profile:

Candidate	Race	Count	Pos Tweets	Neg Tweets	Neu Tweets	Ave Compound
Ted Cruz	Texas Senate	111459	39642	33701	38116	0.0147602033
Beto O'Rourke	Texas Senate	34926	11445	3830	19651	0.1143955134

Table 3.3.1. A sample sentiment profile for the two main candidates from the race for the Texas Senate. It contains all the parameters outline above

These profiles exist for every single candidate that participated in the election. You can find the sentiment profiles for house and senate candidates in Appendix C. The profiles are the data that will feed into the models that we'll build.

4

Exploring the Dataset

At the end of Chapter 3 we discussed building sentiment profiles for a candidate. These profiles are inspired by the profiles that many polling places create as they ask questions on their candidate. This chapter dives into the dataset that we've created and examines potential patterns and trends.

4.1 Races Collected

In the data gathering process, tweets were collected from thirty-three senate races, twenty nine house races, and a series of miscellaneous hashtags designed to capture national sentiment about the election.

The datasets built in Chapter 3 meant that we could create a large scale look at the House and Senate elections. Displayed in Table 4.1.1 and Table 4.1.2 is a leaderboard of sorts displaying the top three races/candidates in several categories.

The most popular races consists of the races with the most tweets. The most popular candidates is the same, except with candidates instead of races. The Best Liked goes to the candidate that had the highest percentage of positive tweets. Least liked and Most Neutral work the same way, except for negative and neutral tweets respectively. There was a five hundred tweet minimum to qualify for the leaderboard.

Rank	Most Popular Race	Most Popular Candidate	Best Liked	Least Liked	Most Boring
1	Texas Senate	Ted Cruz	Angela Green	Josh Hawley	Tony Campbell
2	Florida Senate	Rick Scott	Lawrence Zupan	Angus King	Kevin Cramer
3	Indiana Senate	John James	John Barasso	Bill Nelson	Matt Rosendale

Table 4.1.1. The tweet statistical overview for the Senate

Rank	Most Popular Race	Most Popular Candidate	Best Liked	Least Liked	Most Boring
1	NY-22	Troy Balderson	Xochitl Small	Josh Hawley	Lizzie Fletcher
2	VA-07	Steve Knight	Harley Rouda	Angus King	John Culberson
3	KY-06	Mike Bishop	Ross Spano	Bill Nelson	Gil Cisneros

Table 4.1.2. The tweet statistical overview for the House

One interesting point is that those who lead the various sentiment-based categories were nowhere close to the leaderboard in the most tweets category. Generally it felt that if a candidate was tweeted about more, their average sentiment score tended to drift towards neutral with the positive and negative tweets balancing out. Don't confuse this neutral score with the Most Boring category. The Most Boring category is constructed using candidates with the percentage of Neutral tweets.

Here is the tweet count, aggregate sentiment score relationship:

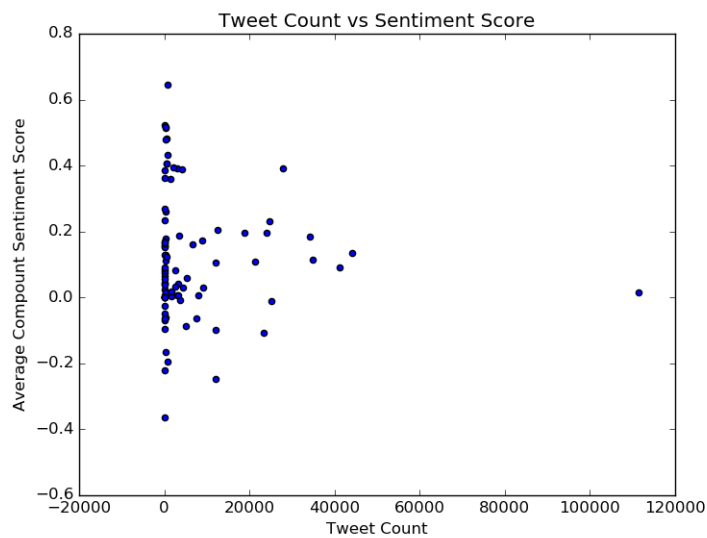


Figure 4.1.1. A scatterplot of tweet count and average sentiment score for all candidates

As the count increases, the range in sentiment scores decrease. The samples with tweet counts greater than twenty thousand all stay within a range of 0.4 on either side of neutral. Conversely, the range for sample sizes less than twenty thousand is from -0.4 all the way to 0.7. Essentially, the larger sample sizes are more balanced with a number of positive and negative tweets. That lends support to the idea that our sample isn't overly biased in one ideological direction.

4.2 Comparing Sentiment

4.2.1 *Winners vs Losers*

The point of an election is to win. The question everyone wants to answer is how? What separates the winners from the losers, and how is that reflected on Twitter? Let's take a look at all the features from the sentiment profiles we created in Chapter 3. We'll begin with the senate.

The graph below shows the average compound sentiment score for winning and losing Senate candidates.

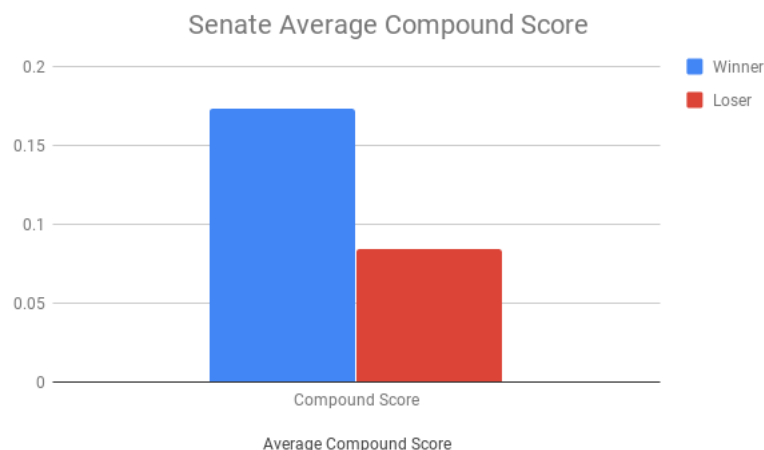


Figure 4.2.1. Average sentiment score for winners and losers from the senate. The winners are in blue and the losers in red

The winning candidates are heads and shoulders above the losing candidate in average sentiment score. That sounds right. Candidates that win elections are favored and liked by the voters and that should be reflected by the average sentiment score

Compound score is a combination of positive and negative tweets. If winners have a higher compound score, it stands to reason that losers will have more negative tweets than winners. The next figure breaks down the count by type of tweet - positive/negative/neutral.

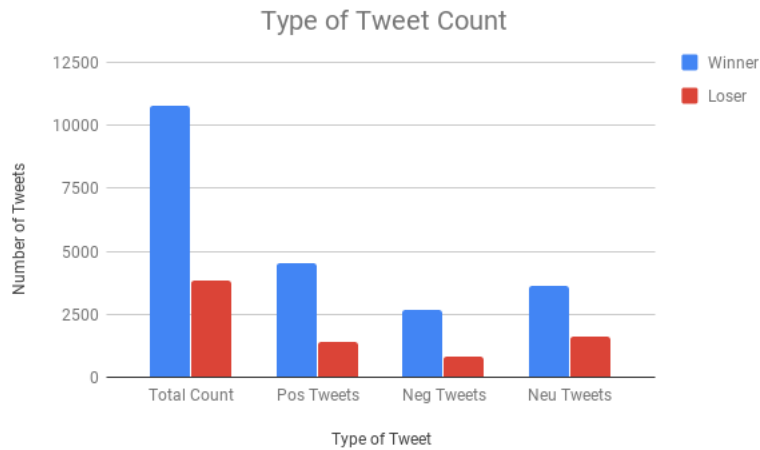


Figure 4.2.2. Averages for total count, positive, negative, and neutral tweets from the senate. Winners are in blue. Losers are in Red

The election losers don't actually have fewer negative tweets. But, what they do have is far fewer tweets than winners. To see whether or not winners have a lower negative score, we should look at the percentages each type of tweets takes from the total count. This is a much more meaningful comparison as it takes into account the number of tweets winners and losers have.

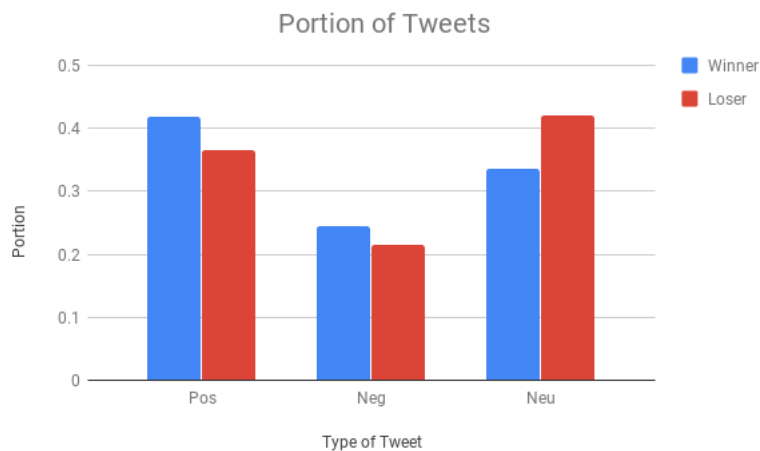


Figure 4.2.3. An expansion of Figure 4.2.2, taking the percentage of positive, negative, and neutral tweets for winners and losers.

A closer look shows that the real difference lies in the fact that losers have a greater portion of neutral tweets than winners do. In terms of positive and negative tweets, the portions differ by less than five percentage points. Neutral tweets have a sentiment score of zero which would dilute and lower the average compound scores. Winners are not only tweeted about more than losers, but they're more likely to have tweets that show *some* sentiment.

But do these patterns hold true for the house? The average compound scores for winners and losers are displayed below:

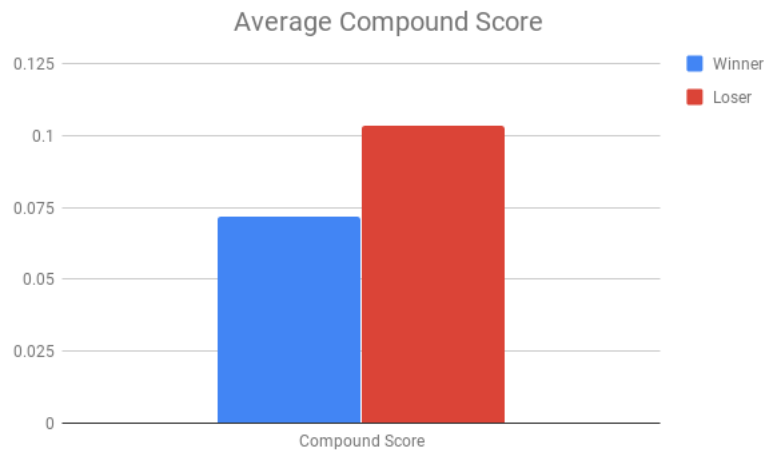


Figure 4.2.4. Average sentiment score for winners and losers from the House. The winners are in blue and the losers in red

Unlike the senate, the losers had a higher average sentiment score than the winners. That's unexpected. Let's take a deeper look into how those numbers break down.

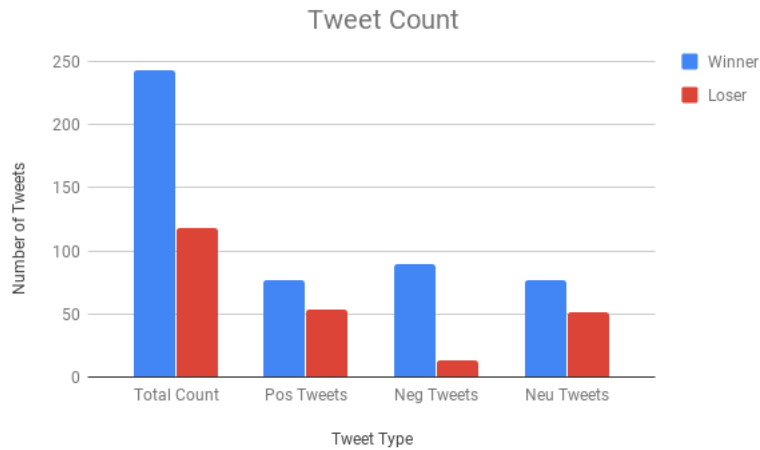


Figure 4.2.5. Averages for total count, positive, negative, and neutral tweets from the House. Winners are in blue. Losers are in Red

Winners do have more positive tweets, but the large difference in the number of negative tweets is what ultimately influences the lower average compound score seen. Unlike the senate, the house races show an even-ness when it comes to neutral tweets. These numbers are hardly conclusive, partially because the large difference in tweet count is inflating these numbers. Here's a breakdown of the tweet percentages:

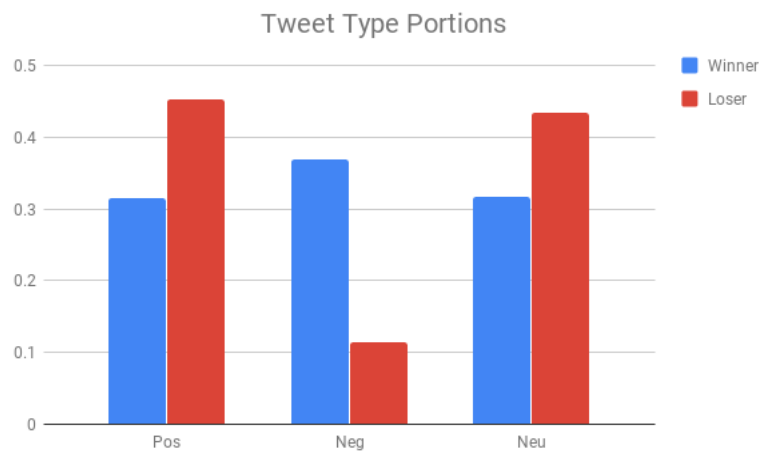


Figure 4.2.6. An expansion of Figure 4.2.5, taking the percentage of positive, negative, and neutral tweets for winners and losers.

The house and the senate winners both roughly have 30% of their tweets designated neutral. The real difference is that negative tweets occupy a far higher proportion of tweets. It doesn't feel that positive and negative tweets have a distinguishable effect on winning.

A major trend that held across both house and senate winners is the fact that winners are tweeted about more. The house had a ratio of about 2-to-1 while the senate had 3-to-1. Candidates that are being talked about more are likely to be on people's minds, both good and bad. Bottom line, positive tweets are preferred, but overall engagement is key.

Senate and house races are more local and don't often engage all eligible voters in the United States. It would be interesting to see if the trend of "more tweets implying more likely to win" holds for a Presidential Election, which is closely followed with lots of voter engagement.

4.2.2 *Democrats vs Republicans*

The two party system is a large part of our election process. The parties command massive influence and many voters choose a candidate based entirely on party affiliation. This means that a party's national platform and performance can influence their performance in more local elections like the house and the senate.

In the house races tracked, the Democrats flipped twenty-one seats which undoubtedly played a large role in their retaking of the House. In the Senate, the Republicans held the majority thanks to the uneven split in seats up for election. In fact, they managed to flip three Democratic seats, allowing them to increase their majority.

Incumbents	House	Senate
Democrats	1	24
Republicans	28	9

Table 4.2.1. Incumbents split by party and seat

Winners	House	Senate
Democrats	22	21
Republicans	7	12

Table 4.2.2. Victors split by party and seat

In the previous section, it showed that overall sentiment and tweet counts were two features that winners of elections had. If those traits hold, it should show that Democrats in the house should have a higher average sentiment and more tweets and the same for Republicans in the Senate.

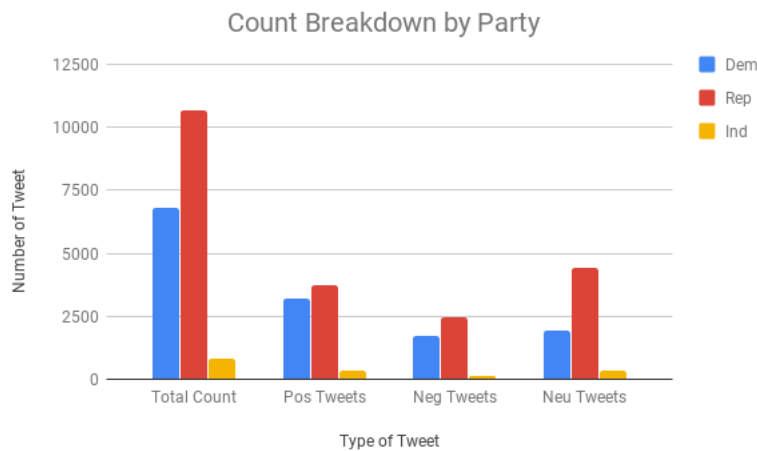


Figure 4.2.7. Averages of total count, positive, negative, and neutral tweets from senate races split by party. Blue is Democrats, red is Republicans, and yellow is Independent

The Republican dominated Senate did indeed carry the higher total count. They also led in every single other category. This will have interesting effects on the actual sentiment comparison.

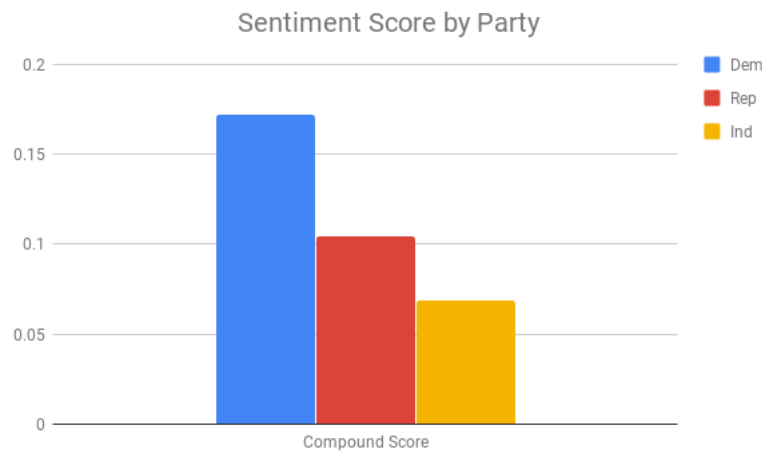


Figure 4.2.8. Average sentiment score from senate races split by party. Blue is Democrats, red is Republicans, and yellow is Independent

Republicans leading both in positive and negative tweets meant that they suffered in the overall sentiment tally. The winning party held the tweet count lead, but not overall sentiment. If tweet count is truly the major feature to follow, the trend should hold in the house.

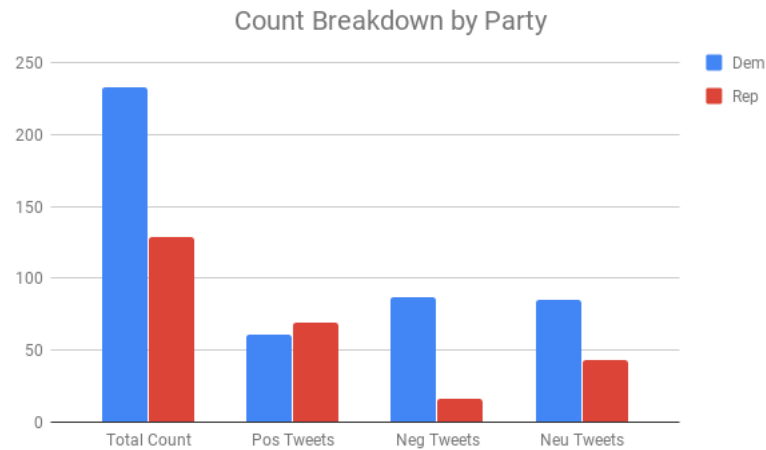


Figure 4.2.9. Averages of total count, positive, negative, and neutral tweets from house races split by party. Blue is Democrat and Red is Republican

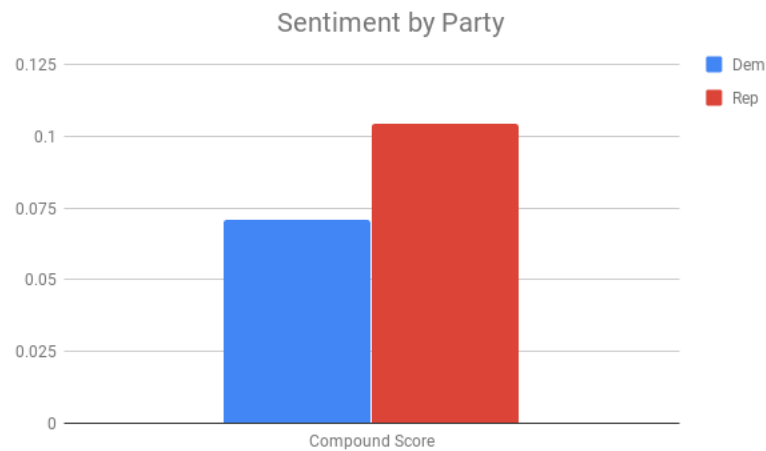


Figure 4.2.10. Average sentiment score from house races split by party. Blue is Democrats, red is Republicans

The main takeaway here is that the tweet count lead continues to hold. The sentiment score favors the Republicans - the losing party - but that correlates with the fact that losing candidates in the house had a better sentiment score as well.

4.3 National Trends

This data set presents the opportunity to explore how national trends are reflected on twitter. Not only that, it'll give a better idea on the kind of sample collected and how it reflects to the samples pollsters take when they observe overall trends they expect to influence elections.

Eight search terms were used to create a database of miscellaneous tweets. These search terms were largely national trends:

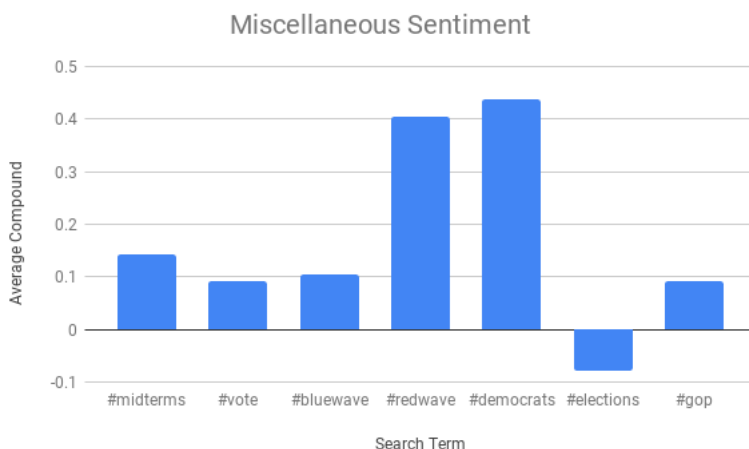


Figure 4.3.1. Average sentiment score from the keywords that made up the miscellaneous database

The sentiment of partisan search terms are about the same. Non-partisan terms have sentiments within a similar range with the peculiar exception of “#elections”. The purpose of this exercise was to show that this database is mostly neutral in contrast to the house leaning democrat and the senate republican.

The three major trends we’re examining are sourced from an ABC news article published on November 6th. The article identified health care, border security, and Donald Trump as three major factors heading into the midterm election [2].

4.3.1 Exploring the Issues

The main motivation for this project *was* Donald Trump. Mr. Trump was a controversial candidate and is now a controversial president. The midterms are often considered a referendum

against the president. Areas which view Mr. Trump favorably should lean republican while those who don't favor him should lean democrat. Since taking office, Mr. Trump has averaged a favorability rating of forty percent while never being higher than forty-five [7].

Health care has had an interesting history in US Politics. In 2010, President Barack Obama passed the Affordable Care Act allowing access for millions of uninsured. Despite many Republican efforts, including a late push by congress in 2018, the ACA still stands. Several moderate Republicans avoided voting against it due to its popularity in their district/state.

According to the Kaiser Family Foundation, the ACA continues to remain a hot topic among voters [11]. A different Pew Research poll showed that sixty percent of voters believe that health care is the government's responsibility [6].

The last major issue is the border wall. Mr. Trump's first campaign promise was the wall and his continued attacks on illegal immigration gained him significant support. After taking office, Mr. Trump has wielded the Immigrations and Customs Enforcement Division (ICE pronounced "ice") as a baton to round up and deport illegal immigrants. A gallup poll had just twenty percent of respondents view Illegal Immigration as "not a problem".

These three issues are front and center, and the data set should have views about them. Having Twitter's numbers agree in some form with that of traditional pollsters shows that the data set is a fair representation of the populace.

4.3.2 Evaluating Twitter's Response

There were six search terms taken from the three issues. Using modified code that can be found in Appendix B, a sentiment profile was put together. Each term was searched independently in the senate, house, and miscellaneous database. The goal was to have the three databases capture a different type of population.

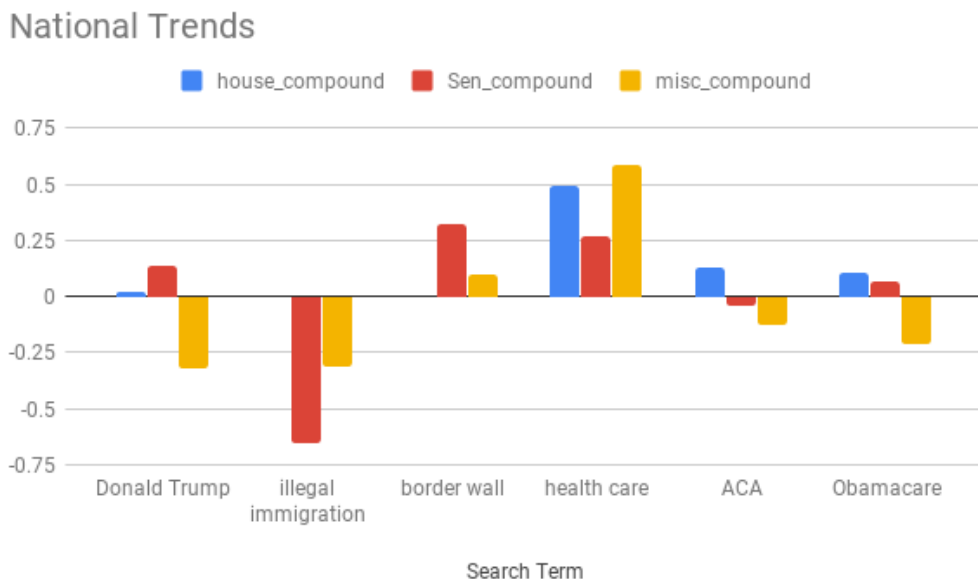


Figure 4.3.2. Sentiment scores of search terms that describe important national trends

Looking at Figure 4.3.2 that assumption holds up. Mr. Trump is viewed favorably in the Republican dominated Senate, while less so in the house and in general. His house sentiment may be a result of a majority of house elections had republican incumbents. His intense negative sentiment with the general populace correlates with that of his general favorability - which is below average.

On illegal immigration, it's interesting to see that Democrats took no stance. That could be because if they did, they would be in the minority. Sentiment is decidedly against illegal immigration and very much for the border wall. This is not only in the republican-dominated Senate but in the more neutral miscellaneous database. It's possible that Democrats saw a path through success by simply not addressing illegal immigration and focusing on other topics that were in their favor like health care.

In section 4.3.1 a Pew research poll was referenced regarding health care and the government. Specifically voters believed that health care was a governmental responsibility. They've echoed that sentiment here, as big fans of health care as an idea in general, but not of specifically the ACA. Notice that the ACA and Obamacare reference the same piece of legislation, but

Obamacare's sentiment is noticeably lower in all three categories. That may be because of the vitriol Obama faces by Republicans and even the nation. It's clear that Obama's controversial presidency is still showing effects. So much so, that just the name on the paper makes all the difference.

5

Modeling Elections

The goal of this project is to build a model that will predict elections. Data has been collected, cleaned, and organized. Tweets have been evaluated with VADER. All that's left is to build a model. This chapter will cover that process and evaluate the resulting model.

5.1 Approach

Part of planning and designing a model is an understanding of the model's purpose. What is the model going to do? Will it make a decision? A recommendation? These are questions that will determine the type and effectiveness of the model that's built.

The challenge with elections is that they are a zero-sum game. When a candidate wins an election, the other candidate loses. The approach taken is to remove the game entirely. Rather than focus on an election the model will focus on a candidate, and act as if they exist in a vacuum. It will predict a candidate's chances of winning their election independent of the other candidates' own probability. That means that the sum of the winning probabilities for all candidates in a given election is not always 1.

This approach will challenge the model constructing portion. This is a shift away from using the tweets as express data input. It's instead a higher-level approach. We will use the summary statistics per candidate as input instead.

5.2 A Machine Learning Overview

The layman's definition of machine learning is to "program a computer to accomplish a task that you don't know how to do". Most coding is to explicitly define a task. Machine learning offers the ability to teach an algorithm the rules of the task, and then allow the algorithm to find its own method to accomplish it. There are many subfields, but the type of learning we'll focus on is supervised learning.

Supervised learning is a type of learning where an algorithm is given input, output pairs and uses those pairs to infer the function. To put it more simply, we give the algorithm both the question and the answer, and it finds out the method to solve it [13]. When creating a supervised learning model, it will need to be trained and tested. A full set of data is split into training and testing sets. The model will use the training set to create the best map it can. The testing set evaluates the map's effectiveness.

A clean, extensive, and unbiased data set is crucial to a successful supervised learning model. The maps that the model creates is reliant entirely on data. A biased dataset will result in a biased model. A few examples of supervised learning include linear and logistic regression, perceptron learning, and naive bayes.

The specific problem in question is a binomial classification problem. Each input has two potential results - winning or losing. With binomial classification, there are a variety of models available. Due to a variety of reasons that will be thoroughly examined in the next section, we've chosen to proceed with the Naive Bayes Classifier, a machine learning technique based on Bayes Theorem of Conditional Probability.

5.3 Understanding the Naive Bayes Classifier

The Naive Bayes classifier is a Bayesian classifier. All Bayesian classifiers are decision making systems that are based off of Bayes Theorem of Conditional Probability. Bayes Theorem makes a prediction based on observation.

5.3.1 Bayes Theorem

Theorem 5.3.1 (Bayes Theorem of Conditional Probability). *Given that you have an object with feature A and a class B , the probability that the object belongs to class B based on the feature A is:*

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (5.3.1)$$

A Bayesian classifier uses that theorem and extends it to make a prediction. Imagine that there are classes A and B and an object with feature x .

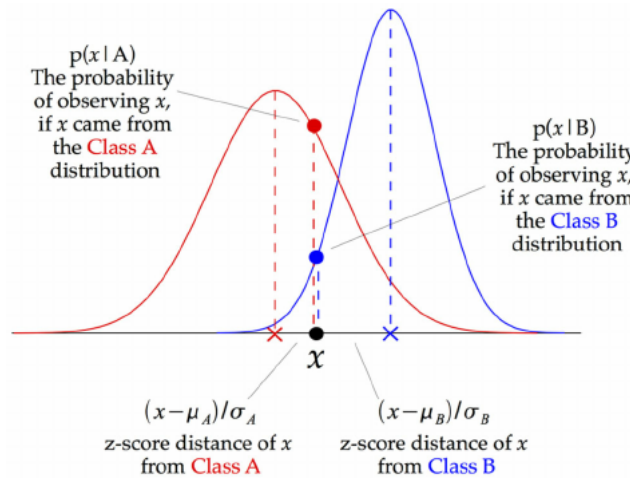


Figure 5.3.1. A demonstration of how Gaussian Naive Bayes operates [1]

The classifier uses Bayes Theorem, calculates $P(x|A)$ and $P(x|B)$ before making a decision based on the resulting probabilities. There are multiple methods to calculate conditional probability and Figure 5.3.1 assumes that the probabilities for each event are distributed normally.

Bayesian classification seems straight forward when you have just one feature to classify. What happens in the case where you have more than one feature? That's the problem the Naive Bayes classifier solves.

5.3.2 Naive Bayes Classifier

The Naive Bayes classifier classifies objects with multiple features into classes. The classifier assumes that each feature is independent from the other. Given a set of features $d =$

$\{d_1, d_2, \dots, d_n\}$ and a class c_j , the probability that d came from c_j is:

$$P(d|c_j) = P(d_1|c_j) \cdot P(d_2|c_j) \cdot \dots \cdot P(d_n|c_j) \quad (5.3.2)$$

which can be simplified to:

$$P(d|c_j) = \prod_{i=1}^n P(d_i|c_j) \quad (5.3.3)$$

This process repeats for all possible classes and the classifier chooses the class with the largest conditional probability. Thus we obtain the following classification rule to determine the class \hat{c} for an object with a set of features d :

$$P(c_j|d_1, \dots, d_n) \propto P(c_j) \prod_{i=1}^n P(d_i|c_j) \implies \hat{c} = \arg \max_{c_j} P(c_j) \prod_{i=1}^n P(d_i|c_j) \quad (5.3.4)$$

In this project, we have two defined classes. Class 1 is losing and class 2 is winning. Every candidate will be classified as a winner or a loser based on their sentiment profile.

Before we can move on to implement a classifier, we have to choose exactly which classifier we need. There are several types of Naive Bayes classifiers and the one we will use is Gaussian Naive Bayes. Figure 5.3.1 is a visual example of a Gaussian Naive Bayes classifier.

Gaussian Naive Bayes is designed for models with continuous parameters. To account for that will require changing the method of calculating conditional probability. With continuous input, we assume the probability to lie on a normal distribution. Keeping the same terminology as equation 5.3.2 we find $P(d_i|c_j)$ as such:

$$P(d_i|c_j) = \frac{1}{\sqrt{2\pi\sigma_{c_j}^2}} \exp\left(-\frac{(d_i - \mu_{c_j})^2}{2\sigma_{c_j}^2}\right) \quad (5.3.5)$$

The parameters σ_{c_j} and μ_{c_j} are estimated using maximum likelihood.

5.3.3 Why Naive Bayes?

The Naive Bayes is the optimal solution given the constraints of the problem. First, the classification problem means that the optimal solution is either a neural network, logistic regression, or Naive Bayes. The chosen approach has limited the size of our dataset. Instead of the three million tweets, the inputs are the hundred and fifty candidates. For smaller data sets, the Naive

Bayes has been shown to outperform other models [8, 16]. Both regression and a neural network in crude terms weight their parameters to produce a final result. To ensure accurate weights, the training data needs to be extensive, something not available to us currently.

The Naive Bayes as explained in section 5.3.2 works off of the probability distribution in the sentiment profiles. Those profiles are created from over two million tweets. The large sample size means that it's likely that the profiles are distributed normally which will improve the accuracy and effectiveness of a Naive Bayes classifier.

5.4 Implementation Naive Bayes in Python

The python library Sci-Kit Learn contains a variety of machine learning models. The machine learning can be imported from Sci-Kit, fitted with training data, and then tested using their predict method. Importing the library is done through the import statement.

```
from sklearn.naive_bayes import GaussianNB
```

The data to train and test the algorithm needs to be imported, split, and organized before any fitting can happen. The dataset has been stored as a CSV, and the pandas package will make importing and organization easy. They are stored in a dataframe, which is a python object that is similar to an excel table.

```
import pandas as pd
```

```
pd.read_csv(FILE_NAME)
```

The dataframe is sliced into data and labels. Data is the parameters that the algorithm will read and weight. These parameters include Tweet Count, Positive Tweets, Average Compound Sentiment Score, etc. Labels are the answers for the algorithm to check itself against. Did the candidate win or lose. To allow the algorithm to process the label, winning is a 1 and losing is a 0.

Both the training and testing sets have a data and label variable. To ensure an unbiased model, we'll randomly select the training and the testing set.

```
test_delim = int(0.75 * senate.shape[0])  
train_data = np.asarray(senate.loc[:test_delim, 'Count':'Party'])  
train_labels = np.asarray(senate.loc[:test_delim, 'Winner'])
```

You'll notice that the inputs to the algorithm are all arrays of integers. Even when dealing with text, data needs to be tokenized or shaped into integers that can be organized into arrays for the algorithm to process. Shaping the information and then understanding the output is on the user.

Fitting the algorithm is easy. The `gaussianNB` object has a `fit` function that takes the training data and labels as input.

```
gnb = GaussianNB()  
gnb.fit(train.data, train.labels)
```

The algorithm is now trained. To evaluate, use the `predict` function on the testing set. The `predict` function will output the prediction, and not the conditional probabilities. To see those, use the `predict_proba` function.

```
output = gnb.predict(test_data)  
probs = gnb.predict_proba(test_data)
```

With the ability to predict, we can now evaluate the effectiveness of Naive Bayes.

5.5 Evaluating the Model

Evaluating a model with just accuracy doesn't capture the full picture of how effective a machine learning model is. It is important to examine subsets of the results to understand how the model interacts with specific classifications. To do this, we use the confusion matrix [5].

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Figure 5.5.1. A sample confusion matrix that outlines the four possible outcomes from a prediction

The confusion matrix outlines the four potential outcomes for any prediction - true/false positives and true/false negatives. These have been used in the medical field for decades when assessing false positives, and gained popularity in machine learning in the 1980's [5]. **True positives** and **true negatives** are when both the predicted output and the actual result are positive and negative respectfully. **False positives** are when the model predicts positive while the actual result is negative. **False negatives** are when the algorithm predicts negative while the actual result is positive.

5.5.1 Criteria

The confusion matrix provides the base for four criteria to evaluate the model. These are precision, accuracy, recall, and f1 [5, 9].

Precision is the ratio of true positives to the sum of true positives and false positives. This is a metric of how often the model is right given that it predicts a candidate to win. It's an important metric because it evaluates the power of the model in regards to predicting a winner.

$$PRECISION = \frac{TP}{TP + FP} \quad (5.5.1)$$

Recall is the ratio of true positives to the sum of true positives and false negatives. It's a good measure of how well the algorithm covers actual positive values.

$$RECALL = \frac{TP}{TP + FN} \quad (5.5.2)$$

Accuracy is the ratio of true positives to total items in the test set, It's an overall measure of how good the model is. On it's own it's an incomplete metric, but it fits in quite well with the other four.

$$ACCURACY = \frac{TP + TN}{TP + TN + FP + FN} \quad (5.5.3)$$

Last, the **f1 score** is the weighted average of the precision and recall. It takes false positives and false negatives into account.

$$F1 = 2 \cdot \frac{PR \cdot RE}{PR + RE} \quad (5.5.4)$$

5.5.2 Performance

To test the model, the data is split into training and testing. First the rows are shuffled and then 75% of the data goes to training while the other 25% goes to testing. Here is the resulting confusion matrix from one run of the model:

Predicted/Actual	Positive	Negative
Positive	3	8
Negative	1	26

Table 5.5.1. The confusion matrix that resulted from running the model on the testing data

From those, we can calculate the metrics outlined in Section 5.5.1. But, to ensure that the results are unbiased we'll run the model ten times, randomly splitting the test and the training set each time. The metrics listed below are averages from those runs.

Metric	Score
Precision	0.75
Accuracy	0.763
Recall	0.31
F1	0.458

Table 5.5.2. Resulting metrics using the values from Table 5.5.1

It's important to keep in mind that this isn't a model predicting any one specific election but rather any election. Randomly splitting the data into training and testing set means that

results we get are results that cover the entire dataset. What these results mean is that when this model predicts any election, it is right seventy-five percent of the time.

Additionally, the precision score is encouraging. It shows that the when the model does choose to predict a positive result, it's often right. Not to mention, it's often right in general. The low recall and f1 score are worrying, but potentially misleading.

Since only one candidate can win an election and some elections have more than two candidates, there is a class imbalance. That means that it is likely that there are more false negatives than true positives which will lower the recall and consequently the f1 score. But, the negative prediction rate (ratio of true negatives to the sum of true and false negatives) is 0.76 which batches the precision/positive prediction rate. What this means is that if the test set were perfectly balanced, it's likely that the recall and the f1 are significantly larger.

6

Conclusion

6.1 Project Achievements

This project used data from twitter to create a predictive model for elections. The model was reasonably successful, predicting elections at a 76% rate. Attempting to model elections also created a deeper understanding of data analysis and modeling. The initial approach of a complex neural network to pull out patterns in the data was discarded when it was clear a simpler approach using rules-based classification was more effective.

The project also leveraged a number of technologies such as MongoDB, the Twitter API, Python and Sci-Kit Learn, and more. It was a great experiment in understanding and immediately applying new technologies, and has laid the groundwork for expansion come future elections.

6.2 Future Work

This project was challenged by the lack of data. Over three million tweets were collected, but those tweets weren't used as individual data points. This project will be improved with more races and candidates collected. The time constraints in November limited the number of races that were followed, but now the data collection pipeline is in place, and all four hundred

and fifty races can be followed in 2020. As more races and data are added, the accuracy and effectiveness of the model should increase. The larger data set also provides more options as far as model selection such as logistic regression, neural networks, and more.

This project was a first investigatory effort into the problem of election prediction. There are a variety of factors that haven't been taken into account - such as bots, attacks, and separation of constituencies. More complex models with more data could be potential for further research.

Appendix A

The How-To Guide

The following chapter is an instruction guide to apply the scripts from Appendix B and recreate various parts of the project. If you're looking for full scripts, those can be found in Appendix B. A quick disclaimer. All the instructions that follow are assuming that you are working on a UNIX system such as MacOS or Ubuntu.

A.1 Installation

There are several programs and packages that need to be installed to gather tweets as well as running a model:

- Python3 - If you're working on a UNIX system, this should come pre-installed. If it doesn't, install it using your system's package installation manager.
- MongoDB - A database management system. Find the distribution for your operating system on the website.
- Text Editor - Choose one that best meets your needs.

There are several python libraries that need to be installed for the scripts in Appendix B to work. They're listed below:

- Tweepy - The Twitter Python API
- Pymongo - The MongoDB Python API
- vaderSentiment - The python library that contains the VADER algorithm discussed in Chapter 3
- Sci-Kit - A Machine Learning Library used to build machine learning. models
- ConfigParser - A python package used to read and parse through configuration files

Python libraries are installed using pip. From the command line, execute the command:

```
pip3 install <library>
```

A.2 Gathering Tweets

Once everything from section A.1 is installed on your machine, use the following steps to stream tweets into a database:

1. Create a developer account on Twitter and create an Twitter application through it. The instructions to do so are provided on <https://developer.twitter.com/>. Once the application is created, it will provide you with an API Key/Secret Pair. It will also provide you with a token key/secret pair. These four randomly generated strings are how the script will interact with twitter.
2. Copy and paste the Listener script from Appendix B. Replace the ckey, csecret, atoken, asecret values with your two key/secret pairs. If you'd like to change the name of the database, rename the db variable in the on_data method.
3. Before we can run the script, mongod has to be ready to receive the tweets. Start the mongod daemon from the command line using the command:

```
$sudo mongod
```

Now, open another tab on the command line and run the script using the command:

```
$python3 misc_listener.py
```

4. To check if the stream is working do the following:

(a) Open a new tab in the command line and start the mongo shell with the command:

```
$mongo
```

(b) In the mongo shell use the command:

```
>show dbs
```

That will output a list of the databases on the server currently. The db name from the script should be listed.

(c) To see the individual tweets in the database, use the following two commands in sequence. Replace db_name with the name of your database in the use command :

```
>use <db_name>
```

```
>db.stats()
```

The stats command will give you an overview of the database. One of the metrics is the number of objects. While the streamer is running, refresh the db.stats() command and observe the object count rising.

A.3 Building Datasets and Models

For this portion, there's little step by step guidance in terms of setup. In Appendix B instructions are included in regards to installation or set up. For any script that's building a sentiment profile or a dataset, the mongod daemon needs to be running. For scripts that are modeling such as naive-bayes.py, the csv needs to be in the same directory as the script is running from. Those instructions are re-iterated in Appendix B.

Appendix B

Code

This project required a number of python scripts to be written. These scripts covered collecting data, organizing and cleaning collected data, as well as building predictive models. Below, you'll find every script written along with its intended purpose.

Name: Listener

Purpose: This script is a filled in template of the code written to stream tweets into a database. The “misc” in the title refers to the fact that this version of the script is streaming tweets from the list of miscellaneous search terms. It has two very similar scripts for the senate and the house.

Execution: It requires a server with the mongod daemon running as well as a Twitter Application with an API key/secret pair. Python packages used are Tweepy and Pymongo. Can be run directly from the command line or through an IDE. Explicit instructions are in Appendix A.2

Expected Output: Will print tweets as they are streamed. Tweets themselves will be stored in the targeted database.

```
1 from tweepy import Stream
2 from tweepy import OAuthHandler
3 from tweepy.streaming import StreamListener
4 from pymongo import MongoClient
```

```

5 import json
6 import sys
7
8
9
10 MONGO_HOST = 'mongodb://localhost/miscdb' # assuming you have mongoDB installed locally
11                                           # and a database called 'twitterdb'
12 ckey = "pF5W6BaHFteEYdkq38cgZ755g"
13 csecret = "AwnBHQtmvVr6cMSYKm89vQOXMO3sPyoGvhXMaow3zlg2G74vLq"
14 atoken = "2646419088-gtzEn525GUUtTlwpegdZGd8dGLWl8m2IGq9jk9y"
15 asecret = "iQbg8ZfJ8KkI06n0lwKCVVprtkzsIp6h1RMXP5LN3Z6o2"
16
17 class StreamListener(StreamListener):
18     #This is a class provided by tweepy to access the Twitter Streaming API.
19
20     def on_connect(self):
21         # Called initially to connect to the Streaming API
22         print("You are now connected to the streaming API.")
23
24     def on_error(self, status_code):
25         # On error - if an error occurs, display the error / status code
26         print('An Error has occurred: ' + repr(status_code))
27         return False
28
29     def on_data(self, data):
30         #This is the meat of the script...it connects to your mongoDB and stores the
31         #tweet
32         try:
33             client = MongoClient(MONGO_HOST)
34
35             # Use twitterdb database. If it doesn't exist, it will be created.
36             db = client.miscdb
37
38             # Decode the JSON from Twitter
39             datajson = json.loads(data)
40
41             #grab the 'created_at' data from the Tweet to use for display
42             created_at = datajson['created_at']
43
44             #print out a message to the screen that we have collected a tweet
45             print("Tweet collected at " + str(created_at) + " " + datajson['text'])
46
47             #insert the data into the mongoDB into a collection called twitter_search
48             #if twitter_search doesn't exist, it will be created.
49             db.misc.insert(datajson)
50         except Exception as e:
51             print(e)
52
53 auth = OAuthHandler(ckey, csecret)
54 auth.set_access_token(atoken, asecret)
55 #input = input("Which districts would you like to observe? Please separate districts by a
56 #comma")
57 #keywords = list(input.split(','))
58 keywords = ["Antonio Delgado", "John Faso", "NY19", "@DelgadoforNY19", "@RepJohnFaso"]
59 keywords = ["#midterms", "#vote", "#bluewave", "#redwave", "#democrats", "#elections", "#
60             gop", "#usa"]
61 twitterStream = Stream(auth, StreamListener())
62 twitterStream.filter(track=keywords)
63 '''

```

Name: Sentiment Profile

Purpose: This code creates a sentiment profile for a given search term and a given database.

The sentiment profile is the same profile that's discussed at the end of Chapter 3

Execution: This code was written under the assumption that a mongodb database is set up with data looking for. You'll also need vaderSentiment and Pymongo installed through pip. The script is executed from the command line using the following template:

```
$python3 sentiment_profile.py <politician> <database>
```

Expected Output: The output will be a sentiment profile printed to the console:

```
Search Term: Claire McCaskill Race: Missouri Senate Count: 27940 \newline
```

```
Pos Tweets: 23663 Neg Tweets: 2328 Neu Tweets: 1949 Average Compound: 0.392
```

```
1 import ast
2 import sys
3 from configparser import ConfigParser
4 from pymongo import MongoClient
5 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
6
7 MONGO_HOST = 'mongodb://localhost/'
8 POLITICIAN = sys.argv[1]
9 DB = sys.argv[2]
10
11 def gather_tweets(phrased, collection):
12     cursor = collection.find({"text": {"$regex": phrased, "$options": "Si"}})
13     entries = [entry for entry in cursor]
14     return(entries)
15
16 def strip_tweet(text):
17     text = text.replace('\n', '')
18     return(text)
19
20 if __name__ == "__main__":
21     config = ConfigParser()
22     client = MongoClient(MONGO_HOST)
23     analyzer = SentimentIntensityAnalyzer()
24     if sys.argv[2] == 'house':
25         db = client.housedb
26         presort = db.presort
27     if sys.argv[2] == 'senate':
28         db = client.senatedb
29         presort = db.presort
30     if sys.argv[2] == 'misc':
31         db = client.miscdb
32         presort = db.misc
33     entries = gather_tweets(POLITICIAN, presort)
34     ave_score = 0
35     pos_tweets = 0
36     neg_tweets = 0
37     neu_tweets = 0
38     count = 0
39     for i in entries:
```

```

40     tweet = strip_tweet(i['text'])
41     sentiment = analyzer.polarity_scores(tweet)
42     if (sentiment['compound'] == 0):
43         neu_tweets += 1
44     if (sentiment['compound'] > 0):
45         pos_tweets += 1
46     if (sentiment['compound'] < 0):
47         neg_tweets += 1
48     ave_score += sentiment['compound']
49     count += 1
50 if count == 0:
51     print("No Tweets for this Politician exist in this database")
52 else:
53     print("Average score:", ave_score/count, "Positive Tweets:", pos_tweets,
54           "Negative Tweets:", neg_tweets, "Neutral Tweets:", neu_tweets)

```

Name: Build Search Profile

Purpose: This script takes a list of all the races and search terms and creates a sentiment profile to each one and then writes that sentiment profile to a file. That allows use that file for analysis any which way we like.

Execution: This script needs an accompanying configuration file. That means that you'll also need the Python library ConfigParser to read that configuration file. The other requirements are the same as those required by the Sentiment Profile script. To run the script, use the command:

```
python3 build_search_profile.py senate
```

To build the profile for the house, use house as the input instead of senate with the same command format as above.

Expected Output: The file prints whichever race it's currently working on to the console.

The final output will be a .csv file it creates in the working directory.

```

1 import ast
2 import sys
3 from configparser import ConfigParser
4 from pymongo import MongoClient
5 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
6
7 MONGO_HOST = 'mongodb://localhost/'
8 FILE_NAME = 'races.ini'
9 keywords = ["#midterms", "#vote", "#bluewave", "#redwave", "#democrats", "#elections", "#gop", "#usa"]
10 #POLITICIAN = sys.argv[1]
11 DB = sys.argv[1]
12
13 def section_to_dict(section, parser):
14     #change path to your ini file if running locally
15     parser.read(FILE_NAME)
16     out_dict = {}
17     for key in parser[section]:
18         temp = ast.literal_eval(parser[section][key])

```



```

19     out_dict[temp[-1]] = temp[:-1]
20     return(out_dict)
21
22 def gather_tweets(phrase, collection):
23     cursor = collection.find({"text": {"$regex": phrase, "$options": "$i"}})
24     entries = [entry for entry in cursor]
25     return(entries)
26
27 def strip_tweet(text):
28     text = text.replace('\n', '')
29     return(text)
30
31 def build_sentiment_profile(politician, db_choice):
32     client = MongoClient(MONGO_HOST)
33     analyzer = SentimentIntensityAnalyzer()
34     if db_choice == 'house':
35         db = client.housedb
36         presort = db.presort
37     if db_choice == 'senate':
38         db = client.senatedb
39         presort = db.presort
40     if db_choice == 'misc':
41         db = client.miscdb
42         presort = db.misc
43     entries = gather_tweets(politician, presort)
44     ave_score = 0
45     pos_tweets = 0
46     neg_tweets = 0
47     neu_tweets = 0
48     count = 0
49     for i in entries:
50         tweet = strip_tweet(i['text'])
51         sentiment = analyzer.polarity_scores(tweet)
52         if (sentiment['compound'] == 0):
53             neu_tweets += 1
54         if (sentiment['compound'] > 0):
55             pos_tweets += 1
56         if (sentiment['compound'] < 0):
57             neg_tweets += 1
58         ave_score += sentiment['compound']
59         count += 1
60     return (ave_score, pos_tweets, neg_tweets, neu_tweets, count)
61
62 if __name__ == "__main__":
63     config = ConfigParser()
64     if DB == 'senate':
65         senate = section_to_dict('SENATE', config)
66         f = open('senate_profile.csv', 'w+')
67         f.write('Search Term,Race,Count,Pos Tweets,Neg Tweets,Neu Tweets,Average Compound\n')
68         for i in senate:
69             print(i)
70             for term in senate[i]:
71                 ave_score, pos_tweets, neg_tweets, neu_tweets, count =
72                 build_sentiment_profile(term, 'senate')
73                 if (count > 0):
74                     line = str(term) + ',' + str(i) + ',' + str(count) + ',' + str(
75                         pos_tweets) + ',' + str(neg_tweets) + ',' + str(neu_tweets) + ',' + str(ave_score/
76                         count) + '\n'
77                 else:
78                     line = str(term) + ',' + str(i) + ",0,0,0,0,0\n"
79                 f.write(line)
80     if DB == 'house':
81         senate = section_to_dict('HOUSE', config)

```

```

79     f = open('house_profile.csv', 'w+')
80     f.write('Search Term,Race,Count,Pos Tweets,Neg Tweets,Neu Tweets,Average Compound
    \n')
81     for i in senate:
82         print(i)
83         for term in senate[i]:
84             ave_score, pos_tweets, neg_tweets, neu_tweets, count =
build_sentiment_profile(term, 'house')
85             if (count > 0):
86                 line = str(term) + ',' + str(i) + ',' + str(count) + ',' + str(
pos_tweets) + ',' + str(neg_tweets) + ',' + str(neu_tweets) + ',' + str(ave_score/
count) + '\n'
87             else:
88                 line = str(term) + ',' + str(i) + ",0,0,0,0,0\n"
89             f.write(line)
90     if DB == 'misc':
91         f = open('misc_profile.csv', 'w+')
92         f.write('Search Term,Count,Pos Tweets,Neg Tweets,Neu Tweets,Average Compound\n')
93         for i in keywords:
94             print(i)
95             ave_score, pos_tweets, neg_tweets, neu_tweets, count =
build_sentiment_profile(i, 'senate')
96             line = str(i) + ',' + str(count) + ',' + str(pos_tweets) + ',' + str(
neg_tweets) + ',' + str(neu_tweets) + ',' + str(ave_score/count) + '\n'
97             f.write(line)
98         f.close()
99     f.close()

```

Name: Naive-Bayes Modelling

Purpose: The script models a candidate's probability of winning given their sentiment profile.

It reads in the data from a given csv, splits the data into training and testing sets, and then trains the algorithm. The script then predicts a series of outputs given the testing set, and then calculates the error.

Execution: Running this will need the python libraries Sci-Kit Learn or sklearn, pandas, and numpy. It will also need the sentiment profiles that were created from the previous scripts

Expected Output: The script will print out the accuracy of the model to the console. It can be toggled with to print the output as well as the probabilities for the output (bayesian modeling returns a probability of being 1 or 0).

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.linear_model import LogisticRegression
5
6 FILE_NAME = "senate_profile.csv"
7 FILE_2 = 'house_profile.csv'
8
9 def score(output, answers):
10     true_pos = 0
11     true_neg = 0
12     false_pos = 0
13     false_neg = 0

```

```

14     print("Predicted , Actual , Result")
15     for i in range(len(answers)):
16         if (output[i] == 0 and answers[i] == 0):
17             true_neg += 1
18         if (output[i] == 1 and answers[i] == 1):
19             true_pos += 1
20         if (output[i] == 1 and answers[i] == 0):
21             false_pos += 1
22         if (output[i] == 0 and answers[i] == 1):
23             false_neg += 1
24     return(true_pos , true_neg , false_pos , false_neg)
25
26 senate = pd.read_csv(FILE_NAME)
27 house = pd.read_csv(FILE_2)
28 house = house.drop(columns=['Unnamed: 9' , 'Unnamed: 10'])
29 house = house.dropna()
30 senate = pd.concat([house, senate])
31 senate = senate.sample(frac=1).reset_index(drop=True)
32 senate = senate.replace('Dem', 0)
33 senate = senate.replace('Rep', 1)
34 senate = senate.replace('Ind', 2)
35 senate = senate.replace('Misc', 2)
36 #test_data = senate.sample(frac=1).reset_index(drop=True)
37 test_delim = int(0.75 * senate.shape[0])
38 train_data = np.asarray(senate.loc[:test_delim, 'Count': 'Party'])
39 train_labels = np.asarray(senate.loc[:test_delim, 'Winner'])
40 test_data = np.asarray(senate.loc[test_delim:, 'Count': 'Party'])
41 test_labels = np.asarray(senate.loc[test_delim:, 'Winner'])
42
43 gnb = GaussianNB()
44 gnb.fit(train_data , train_labels)
45 output = gnb.predict(test_data)
46 print(output)
47 print(test_labels)
48 probs = gnb.predict_proba(test_data)
49 true_pos , true_neg , false_pos , false_neg = score(output , test_labels)
50 print(score(output , test_labels))
51 print("Accuracy" , gnb.score(test_data , test_labels))

```


Bibliography

- [1] *Illustration of gaussian naive bayes*, ResearchGate, Jul 2013.
- [2] *5 major issues dominating the midterm elections: Analysis*, Nov 2018.
- [3] Margaret M Bradley and Peter J Lang, *Affective norms for english words (anew): Instruction manual and affective ratings*, Tech. report, Citeseer, 1999.
- [4] Thomas Connolly and Carolyn Begg, *Database systems: A practical approach to design, implementation, and management: Global edition*, Pearson Education, 2014.
- [5] Tom Fawcett, *An introduction to roc analysis*, Pattern recognition letters **27** (2006), no. 8, 861–874.
- [6] Inc Gallup, *Immigration*, Apr 2019.
- [7] ———, *Presidential approval ratings – donald trump*, Apr 2019.
- [8] David J Hand and Keming Yu, *Idiot’s bayes—not so stupid after all?*, International statistical review **69** (2001), no. 3, 385–398.
- [9] Clayton J Hutto and Eric Gilbert, *Vader: A parsimonious rule-based model for sentiment analysis of social media text*, Eighth international AAAI conference on weblogs and social media, 2014.
- [10] MongoDB Inc., *The mongodb 4.0 manual*¶.
- [11] Ashley Kirzinger, Cailey Muñana, and Mollyann Brodie, *6 charts about public opinion on the affordable care act*, Dec 2018.
- [12] Anna Korhonen, *Automatic lexical classification-balancing between machine learning and linguistics*, Proceedings of the 23rd Pacific Asia Conference on Language, Information and Computation, Volume 1, vol. 1, 2009.
- [13] Andrew Ng, *Cs229 lecture notes*, CS229 Lecture notes **1** (2000), no. 1, 1–3.

- [14] James W Pennebaker, Martha E Francis, and Roger J Booth, *Linguistic inquiry and word count: Liwc 2001*, Mahway: Lawrence Erlbaum Associates **71** (2001), no. 2001, 2001.
- [15] Twitter, *Tweepy documentation*¶.
- [16] Harry Zhang, *The optimality of naive bayes*, AA **1** (2004), no. 2, 3.