

Spring 2018

## Don't Take This Personally: Sentiment Analysis for Identification of "Subtweeting" on Twitter

Noah L. Segal-Gould  
*Bard College*

Follow this and additional works at: [https://digitalcommons.bard.edu/senproj\\_s2018](https://digitalcommons.bard.edu/senproj_s2018)



Part of the [Artificial Intelligence and Robotics Commons](#)



This work is licensed under a [Creative Commons Attribution-Noncommercial-No Derivative Works 4.0 License](#).

---

### Recommended Citation

Segal-Gould, Noah L., "Don't Take This Personally: Sentiment Analysis for Identification of "Subtweeting" on Twitter" (2018). *Senior Projects Spring 2018*. 244.  
[https://digitalcommons.bard.edu/senproj\\_s2018/244](https://digitalcommons.bard.edu/senproj_s2018/244)

This Open Access work is protected by copyright and/or related rights. It has been provided to you by Bard College's Stevenson Library with permission from the rights-holder(s). You are free to use this work in any way that is permitted by the copyright and related rights. For other uses you need to obtain permission from the rights-holder(s) directly, unless additional rights are indicated by a Creative Commons license in the record and/or on the work itself. For more information, please contact [digitalcommons@bard.edu](mailto:digitalcommons@bard.edu).

# Don't Take This Personally: Sentiment Analysis for Identification of "Subtweeting" on Twitter

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

by  
Noah Segal-Gould

Annandale-on-Hudson, New York  
May, 2018



# Abstract

The purpose of this project is to identify subtweets. The Oxford English Dictionary defines “subtweet” as a “[Twitter post] that refers to a particular user without directly mentioning them, typically as a form of furtive mockery or criticism.” This paper details a process for gathering a labeled ground truth dataset, training a classifier, and creating a Twitter bot which interacts with subtweets in real time. The Naive Bayes classifier trained in this project classifies tweets as subtweets and non-subtweets with an average  $F_1$  score of 72%.



# Contents

<b>Abstract</b>	<b>iii</b>
<b>Dedication</b>	<b>ix</b>
<b>Acknowledgments</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Methods</b>	<b>5</b>
2.1 Gathering Tweets for the Ground Truth Dataset . . . . .	5
2.1.1 The Twitter API . . . . .	5
2.1.2 Searching for Tweets . . . . .	6
2.1.3 Changes in Data Acquisition . . . . .	8
2.2 Language Analysis & Naive Bayes . . . . .	9
2.2.1 Resources for Language Analysis . . . . .	9
2.2.2 Cleaning and Preparing the Data . . . . .	10
2.2.3 TF & TF-IDF . . . . .	12
2.2.4 Naive Bayes . . . . .	12
2.2.5 Statistical Considerations . . . . .	14
2.2.6 K-Folds Cross-Validation . . . . .	16
2.2.7 Training & Testing Naive Bayes . . . . .	17
2.3 Creating the Twitter Bot . . . . .	18
<b>3 Results</b>	<b>21</b>
3.1 Ground Truth Dataset . . . . .	21
3.1.1 Characteristics of the Ground Truth Dataset . . . . .	23
3.2 Statistics . . . . .	26
3.2.1 The Confusion Matrix . . . . .	26
3.2.2 Precision, Recall, & F1 . . . . .	27
3.3 Known Subtweeters . . . . .	27

3.4	Most Informative Features . . . . .	28
3.5	The Twitter Bot . . . . .	29
3.6	Discussion . . . . .	30
<b>4</b>	<b>Conclusion</b>	<b>33</b>
4.1	Summary of Project Achievements . . . . .	33
4.2	Future Work & Considerations . . . . .	33
	<b>Bibliography</b>	<b>35</b>

# List of Figures

2.2.1 K-Folds Cross-Validation Example (Borovicka et al., 2012) . . . . .	17
3.1.1 Histogram of Ground Truth Tweet Lengths ( $\bar{x}_s = 103.71, s_s = 71.99$ ), ( $\bar{x}_n = 105.63, s_n = 74.95$ ) . . . . .	23
3.1.2 Histogram of Ground Truth Tweet Punctuation ( $\bar{x}_s = 1.79, s_s = 1.54$ ), ( $\bar{x}_n = 2.16, s_n = 1.68$ ) . . . . .	24
3.1.3 Histogram of Ground Truth Unique English Stop Words ( $\bar{x}_s = 7.14, s_s = 1.31$ ), ( $\bar{x}_n = 7.05, s_n = 1.39$ ) . . . . .	25
3.1.4 Histogram of Ground Truth English Unique Words ( $\bar{x}_s = 16.80, s_s = 10.95$ ), ( $\bar{x}_n = 16.25, s_n = 11.15$ ) . . . . .	25
3.2.1 Confusion Matrix on Test Data from Each Fold on the Ground Truth Dataset . .	26
3.3.1 Distribution of Subtweet Probabilities for 11 Known Subtweeters ( $\bar{x} = 0.520, s = 0.127$ ) . . . . .	27
3.3.2 Distribution of Subtweet Probabilities for 3 Known Subtweeters ( $\bar{x}_{min} = 0.430, s_{min} = 0.116$ ), ( $\bar{x}_{med} = 0.523, s_{med} = 0.128$ ), ( $\bar{x}_{max} = 0.561, s_{max} = 0.123$ ) .	28
3.5.1 Example of the Twitter Bot Quoting Users' Tweets . . . . .	30
3.5.2 Example of the Twitter Bot Replying to a User's Tweet . . . . .	30





# Dedication

I dedicate this senior project to @jack, who has willfully made numerous changes to Twitter which inevitably angered millions.



# Acknowledgments

Thank you professors Sven Anderson, Keith O'Hara, and Rebecca Thomas for making this project possible through your combined efforts to teach and advise me. Thank you Benjamin Sernau '17 for enduring through three years of Computer Science courses with me and being a source of unending joy in my life. Thank you to Julia Berry '18, Aaron Krapf '18, and Zoe Terhune '18 for being my very best friends and giving me things worth caring about. Finally, thank you to my parents Tammy Segal and Emily Taylor for your constant support and patience throughout my four years at Bard College.



# 1

## Introduction

There were over 140 million active Twitter users who sent 340 million text-based tweets to the platform every day by March of 2012 (Twitter, 2012). Since Twitter-founder Jack Dorsey sent the first Tweet in March of 2006 (Dorsey, 2006) social scientists, political scientists, and computer scientists have applied machine learning techniques to understand the patterns and structures of the conversations held on the platform. Through sentiment analysis, we are able to use machine learning to identify patterns within natural language which indicate particular feelings both broadly (e.g. positive, neutral, or negative) and toward topics (e.g. politics, terrorism, brands, etc.).

On Twitter, the most common way to publicly communicate with another user is to compose a tweet and place an “@” before the username of that user somewhere in the tweet (e.g. “How are you doing, @NoahSegalGould?”). Through this method, public discussions on Twitter maintain a kind of accountability: even if one were to miss the notification that they were mentioned in a tweet, one’s own dashboard keeps a running list of their most recent mentions.

“Subtweet” was coined in December of 2009 by Twitter user Chelsea Rae (Rae, 2009) and was entered into Urban Dictionary the following August (Urban Dictionary, 2010). In “To tweet or subtweet?: Impacts of social networking post directness and valence on interpersonal impressions” (Edwards and Harris, 2016), Edwards and Harris sought to analyze student participants’

perceptions of known subtweeters. In the news, too, subtweets have garnered attention in *The Atlantic* (Madrigal, 2014), *The Washington Post* (Dewey, 2016), and *Slate* (Hassler, 2016). In news media, subtweets garner attention for their prevalence among government officials as well. Following President Donald Trump’s inauguration, The Washington Post compiled its “A running list of all the possible subtweets of President Trump from government Twitter accounts,” (Ohlheiser, 2017) cementing subtweets as particularly newsworthy.

If an individual sought to disparage or mock another, they could certainly do so directly. But the targeted user would probably notice, and through the search functions of the platform, anyone could see who has mentioned either their own or another’s username. Instead, a phenomenon persists in which users of the platform deliberately insult others in a vague manner by making complaints while omitting the targets of those complaints.

Although the OED’s definition states that a subtweet “...refers to a particular user without directly mentioning them, typically as a form of furtive mockery or criticism,” it is perhaps too restrictive. Some individuals believe subtweets abide by this definition, but others expand it to allow inclusion of others’ real names (especially if that individual does not own a Twitter account), and some do not even require that a particular user be the target of the tweet. Because subtweeting is colloquial in nature, we will expand the definition of subtweet to permit these less restrictive features.

Sentiment analysis on social networking services such as Twitter has garnered attention within seemingly distinct fields of interest. In “Text mining for market prediction: A systematic review,” Nassirtoussi et al. surveyed varied methods for text-mining social media for sentiment analysis of financial markets and approached that problem with both behavioral and economic considerations in mind (Nassirtoussi et al., 2014). Following a terrorist event in Woolwich, London in 2013, Burnap et al. analyzed the immediate Twitter response following the attack to inform statistics on how long it takes for responses from official sources to disseminate during crises (Burnap et al., 2014). Prior research of these kinds utilizes sentiment analysis techniques on tweets, but no known research exists which specifically performs any sentiment analysis on subtweets.

Long before Twitter, psychologist Gordon Allport wrote about “antilocution” in *The Nature of Prejudice* (Allport, 1954). For Allport, antilocution was the first of several degrees of apathy which measure prejudice in a society. It represented the kind of remarks which target a person, group, or community in a public or private setting but do not address the targeted individual directly. Different from both hate speech and subtweeting, antilocution necessitates that an in-group ostracize an unaware out-group through its biases.

The most germane research available focuses on sentiment analysis of figurative language. Determining sentiment based on features of text which are distinctly separate from their literal interpretations presents difficulties for human readers as well as computer programs. In *SemEval*, the International Workshop on Semantic Evaluation, analysis of figurative language on Twitter has been a core task for their competition since 2015 (Ghosh et al., 2015) and returns this year with a specific focus on ironic tweets (Van Hee et al., 2018). In this year’s description for “Task 3: Irony detection in English tweets,” Van Hee et al. touch upon online harassment as a potential point of significance for sentiment analysis of ironic tweets.

This project pursues sentiment analysis of subtweets in order to challenge the hypocrisy of utilizing a service which presents itself as a public forum to speak in distinctly private ways. Toward this end, these are our goals: this project will provide a framework for collecting examples of subtweets, train a classification algorithm using those examples, and finally utilize that classifier in real time to make tweets which were intended to be unseen by specific parties easily accessible to all parties. In presenting covertly hurtful content as obviously hurtful in a public fashion, perhaps it will promote a particular awareness that tweets posted by public accounts are indeed publicly accessible, and that Twitter’s Terms of Service (Twitter, 2016) allows for this kind of monitoring.

Using a machine-learning approach to perform sentiment analysis, syntactic and linguistic features are typically utilized in probabilistic (e.g. Naive Bayes and Maximum Entropy) and linear (e.g. Support Vector Machines and Neural Networks) classification algorithms. The probabilistic approach is sometimes called *generative* because such models generate the probabilities



of sampling particular terms (Medhat et al., 2014). Linear classification utilizes the vectorized feature space of words, sentences, or documents to find a separating hyperplane between multiple classes.

This project approaches the problem of identifying subtweets using the probabilistic Naive Bayes classification algorithm. We demonstrate this method by programmatically identifying subtweets in real time. Using a Twitter bot—an automated program which accesses Twitter—we use this classifier to detect and interact with subtweets as they are posted.

# 2

## Methods

### 2.1 Gathering Tweets for the Ground Truth Dataset

#### *2.1.1 The Twitter API*

Twitter provides a free Application Programming Interface (API) to registered users and has done so since September of 2006 (Stone, 2006). The API allows developers to programmatically access and influence tweets individually or through real time search filters, and also read and write direct messages (Twitter, 2018). The creation of a Twitter application which utilizes the API requires creation and email verification of an account, and developers are also required to agree to the Terms of Service (Twitter, 2016). Creation of an application provides developers with authentication tokens which can then be used to access the API.

To make creation of Twitter applications easier, Tweepy (Roesslein, 2009) is an open source library for the Python programming language which provides methods and classes used to interact with the API and its status objects (Twitter, 2018). A Twitter status object is a dictionary of key and value pairs which contains text, media, and user information associated with particular tweets. There are rate limits for both reading and writing to the API which must be kept in mind when programming for it.

### 2.1.2 Searching for Tweets

For acquisition of a ground truth dataset, we consider **true subtweets** to be tweets to which another Twitter user replied who specifically called it out as a subtweet. We consider **true non-subtweets** to be tweets to which another user replied who specifically did **not** call it out as a subtweet.

```

1 def get_subtweets(max_tweets=2500,
2                   query=("subtweet AND @ since:2018-04-01"
3                         "exclude:retweets filter:replies")):
4     subtweets_ids_list = []
5     subtweets_list = []
6     for potential_subtweet_reply in tweepy.Cursor(api.search, lang="en",
7                                                   tweet_mode="extended",
8                                                   q=query).items(max_tweets):
9         potential_subtweet_original = first_tweet(potential_subtweet_reply)
10        if (not potential_subtweet_original.in_reply_to_status_id_str
11            and potential_subtweet_original.user.lang == "en"):
12            if (potential_subtweet_original.id_str in subtweets_ids_list
13                or "subtweet" in potential_subtweet_original.full_text
14                or "Subtweet" in potential_subtweet_original.full_text
15                or "SUBTWEET" in potential_subtweet_original.full_text):
16                continue
17            else:
18                subtweets_ids_list.append(potential_subtweet_original.id_str)
19                subtweets_list.append({"tweet_data": potential_subtweet_original._json,
20                                     "reply": potential_subtweet_reply._json})
21            with open("../data/other_data/subtweets.json", "w") as outfile:
22                json.dump(subtweets_list, outfile, indent=4)
23    return subtweets_list

```

After the API credentials are loaded, this Python code is used to download tweets with replies which **do** and **do not** call them out as subtweets. Tweepy provides the `Cursor` object which is used to iterate through different sections of the Twitter API. In this example, we use it with Twitter’s search API to find tweets matching our query. This particular version finds **true subtweets** with their accusatory replies, but it requires only modification of its query argument to download **true non-subtweets**. Within the `for` loop, we confirm that the following qualities hold true:

- The alleged subtweet or non-subtweet is not in reply to any other tweet.
- The user who posted it used English as their primary language.
- We do not download duplicate tweets.
- It does not contain the string “subtweet” (with variations in capitalization).

For acquisition of both subtweets and non-subtweets, this function can be modified to accept different queries. Respectively, these are the queries which were utilized in acquiring so-called subtweet-accusations (i.e. these replies typically claim that the tweet to which they reply was a subtweet) and non-subtweet-accusations (i.e. these replies are essentially normal replies to normal tweets):

```
"subtweet AND @ since:2018-04-01 exclude:retweets filter:replies"
```

```
"-subtweet AND @ since:2018-04-01 exclude:retweets filter:replies"
```

The only difference between the two is that the former searches for tweets which contain the string “subtweet” and the latter searches for tweets which exclude it. Both access the API as far back as it will allow (typically one week for free API users) and exclude retweets (i.e. unoriginal tweets which are reposted) while specifically searching for tweets which were in reply to other tweets. Unfortunately, tweet status objects representing replies to tweets do not contain the object data of the tweet to which they reply.

```
1 def first_tweet(tweet_status_object):
2     try:
3         return first_tweet(api.get_status(tweet_status_object.in_reply_to_status_id_str,
4                                           tweet_mode="extended"))
5     except tweepy.TweepError:
6         return tweet_status_object
```

To obtain this object, this code essentially goes up a chain of recursively finding the original tweet to which a reply was made until the API can no longer find another tweet object with an `in_reply_to_status_id_str` attribute, indicating that the tweet is an original **true subtweet** or **true non-subtweet**. The two versions of this program for acquiring **true subtweets** and **true non-subtweets** ran for three weeks between March and April of 2018, acquiring over 20,000 tweets which compose our ground truth dataset.

	True Subtweet Data	True Non-Subtweet Data
Tweet	Talk to him again about “dropping me” and you’ll get your teeth knocked out	That’s been one of my biggest issues here; the onus is on ordinary people who, in their spare time, must campaign for the basic services of a city. This is not how progressive cities should be built
Reply	Thomas don’t subtweet me during work hours	i guess i am not as creative as i thought

This table features a **true subtweet** and its associated reply and a **true non-subtweet** and its associated reply.

### *2.1.3 Changes in Data Acquisition*

This approach for creating a ground truth dataset relies on a particular phenomenon in which Twitter users call-out the subtweets of their peers. The following pattern was observed: a user posts a subtweet which is easily recognized by a peer, and that peer then replies to that tweet in order to complain that the original user was subtweeting or to ask if the tweet was indeed a subtweet. Initially, the program for acquiring the ground truth dataset used the Twitter API’s search functionality to specifically search for replies to tweets which contained some form of the string “subtweet.” It utilized the API’s status object to access the tweet to which it was replying. For 73 days, each day’s alleged subtweets and their associated accusatory replies were saved.

Previously, the classifier was trained using a dataset which was half composed of these alleged subtweets and half composed of tweets randomly selected from a pre-labeled sentiment analyzed tweets dataset (Go et al., 2009). This procedure failed to make the training data representative of **true subtweets** and **true non-subtweets**. The alleged subtweets downloading program was revised and it has been set to download tweets with replies which specifically do **not** contain the string “subtweet.” In both the program which downloads subtweets and the program which downloads non-subtweets, the assumptions about these interactions will not hold true in every case. They are intended as generalizations which make acquiring a ground truth dataset for use in performing binary classification significantly easier and less time-consuming than finding and labeling subtweets and non-subtweets by hand. Indeed, the dataset utilized by Go et al. uses a similar method for acquiring labeled data. In their *Sentiment140* dataset, the labels were acquired according to emoticons present within the tweets instead of through hand-labeling by actual humans.

## 2.2 Language Analysis & Naive Bayes

Before training the classifier, the ground truth dataset is modified such that the important features within it are easily accessible. Changes are made to preserve the characteristics of the text which are relevant to the goal of the analysis and to leave out the ones which are irrelevant. Because we will be using Naive Bayes, we must keep in mind which features in each tweet (e.g. URLs and user names) ought to influence the probabilities that an entire feature-set (i.e. that whole tweet) suits a particular class.

### *2.2.1 Resources for Language Analysis*

#### *2.2.1.1 Regular Expressions*

For text classification through machine learning, it is popular to modify the ground truth dataset to make features which are not important to the classification problem as **generic** as possible. For classification of subtweets, the classifier will treat URLs, mentions of usernames, and English first names generically. In other words, it will keep track of the existence of those features but specifically will not encounter the text contained within them. In identification of subtweets, there exists no syntactic or linguistic significance in the format of a URL or the name a user chooses to associate with themselves or another. However, the existence of those features within the tweet remains important. For this kind of substring searching, pattern matching through regular expressions was used to replace every occurrence of URLs, usernames, and first names with special tokens which were not already in the dataset. The top 100 most common English names for both men and women over the last century were acquired from the United States Department of Social Security.

#### *2.2.1.2 Tokenization*

Instead of training the classifier on entire strings, **tokenization** is necessary in order to extract individual features from the text. The Natural Language Toolkit (Bird and Loper, 2004) provides a tweet tokenizer to achieve this. For some string, the tokenizer splits apart words, usernames, URLs, hashtags, and punctuating characters as individual tokens. NLTK's tweet tokenizer also

appropriately distinguishes between punctuating characters and emoticons composed of punctuating characters.

### 2.2.1.3 *N-Grams*

An **n-gram** is a contiguous sequence of  $n$  tokens in a piece of text. For example, given a string such as “This is a test,” the bigrams ( $n = 2$ ) for this string are “This is,” “is a,” and “a test.” Instead of training the classifier using unigrams ( $n = 1$ ) exclusively, we train it using unigrams, bigrams, and trigrams ( $n = 3$ ). Thus, when the probability that some specific token within a tweet belongs to a specific class is calculated, its neighbors are also considered in combination with it. **n-grams** enable the classifier to treat particular groupings of tokens with some size  $n$  as importantly as it treats the individual tokens, thus identifying particular word groupings most associated with the classes.

### 2.2.1.4 *Stop Words*

A list of stop words typically contains the most common words in a language. For English text, the list is often composed of words such as “the,” “it,” and “of.” Tokens matching stop words are ignored during classifier training because they are too common to help the classifier distinguish subtweets from non-subtweets.

## 2.2.2 *Cleaning and Preparing the Data*

Although **true subtweets** and **true non-subtweets** are identified using characteristics of the Twitter API’s tweet status objects, we utilize Naive Bayes exclusively for text classification on the unicode text contained within each alleged subtweet and non-subtweet status object. Thus, we ignore the replies to these tweets.

```

1 def load_data(filename, threshold=0.1):
2     data = [(urls_pattern.sub("GENERIC_URL",
3         at_mentions_pattern.sub("GENERIC_MENTION",
4         names_pattern.sub("GENERIC_NAME",
5         t["tweet_data"]["full_text"])))
6         .replace("\u2018", "'')
7         .replace("\u2019", "'')
8         .replace("\u201c", "\"")
9         .replace("\u201d", "\"")
10        .replace("&quot;", "\"")
11        .replace("&amp;", "&")
12        .replace("&gt;", ">")
13        .replace("&lt;", "<"))
14    for t in json.load(open(filename))

```

```

15         if t["tweet_data"]["lang"] == "en"
16         and t["reply"]["lang"] == "en"
17         and t["tweet_data"]["user"]["lang"] == "en"
18         and t["reply"]["user"]["lang"] == "en"]
19     new_data = []
20     for tweet in data:
21         tokens = tokenizer.tokenize(tweet)
22         english_tokens = [english_dict.check(token) for token in tokens]
23         percent_english_words = sum(english_tokens)/len(english_tokens)
24         if percent_english_words >= threshold:
25             new_data.append(tweet)
26     return new_data

```

The `load_data` function **genericizes** URLs, mentions of user names, and mentions of English first names. Using regular expressions for pattern matching, these substrings are replaced with special identifiers. The tweets are also cleaned to make HTML characters and unicode characters more consistent. The list comprehension intentionally excludes **non-English** tweets and those which were **not** posted by accounts which list their primary language as English. NLTK's tweet tokenizer is utilized at the end to check for tweets which contain at least 10% English tokens. This language detection is performed on each token in the tweet using the pyEnchant library (Kelly, 2016) which primarily serves as a spell-checker. The resulting dataset of either **true subtweets** or **true non-subtweets** is returned as a list.

After text cleaning, we remove tweets which are present in both the dataset of **true subtweets** and **true non-subtweets**. Duplicates may have appeared because one user thought a tweet was a subtweet but another did not.

```

1     subtweets_data = load_data("../data/other_data/subtweets.json")
2     non_subtweets_data = load_data("../data/other_data/non_subtweets.json")
3     subtweets_data = [tweet for tweet in subtweets_data
4                       if tweet not in non_subtweets_data]
5     non_subtweets_data = [tweet for tweet in non_subtweets_data
6                          if tweet not in subtweets_data]

```

With the duplicates gone, we limit the size of the larger dataset to be the same as the smaller of the two. Thus, both the **true subtweets** and **true non-subtweets** compose the entire final ground truth dataset equally.

```

1     smallest_length = len(min([subtweets_data, non_subtweets_data], key=len))
2
3     subtweets_data = sample(subtweets_data, smallest_length)
4     non_subtweets_data = sample(non_subtweets_data, smallest_length)
5
6     subtweets_data = [(tweet, "subtweet") for tweet in subtweets_data]
7     non_subtweets_data = [(tweet, "non-subtweet") for tweet in non_subtweets_data]
8
9     training_data = subtweets_data + non_subtweets_data

```



In this code, Python’s `random` library provides the `sample` function which randomly returns a list of  $n$  items from a list. After both datasets are made the same length, the lists of strings are made into lists of tuples. For each tuple, the **true subtweet** or **true non-subtweet** is associated with its label (i.e. subtweet and non-subtweet). Finally, these two lists of tuples are put together for use in training the classifier.

### 2.2.3 TF & TF-IDF

Term frequency (TF) is a simple method for vectorizing text in which all terms (e.g. tokens, features, words, etc.) in the corpus are featured in a vector for each document, and the frequency of each term is reflected in the number representing the corresponding term. The bag of words model essentially vectorizes features using this method. Thus,  $tf(t, d) = f_{t,d}$ .

Unfortunately, TF falls short when the corpus of documents contains terms which appear frequently but do not necessarily help inform the classifier on terms that are best associated with a particular class. TF-IDF, or term frequency-inverse document frequency, is the product of the TF for a specific term and the inverse document frequency (IDF) for that same term. The TF is equal to the ratio between the number of occurrences of a term in a document, and the total number of words in that document. IDF, then, is the logarithm of the ratio between the number of documents in the corpus, and the number of documents which contain that term. The product of TF and IDF assigns weights which appropriately value terms which are frequent within a document but rare in the entire corpus of documents. Thus,  $idf = \log \frac{N}{n_t}$  where  $N$  is the total number of documents and  $n_t$  is the number of documents which contain the term. Because the weighted feature vectors calculated using TF-IDF follow a multinomial distribution, our classification algorithm is specifically a Multinomial Naive Bayes classifier.

### 2.2.4 Naive Bayes

Naive Bayes stands out as particularly simple and common for use in text classification. A **bag of words** model typically ignores word positions in favor of keeping track of raw token frequencies, which are weighted to produce TF-IDF feature vectors. Then, Bayes theorem is

utilized to predict the probability that a given feature set (e.g. words, sentences, etc.) belongs to a particular label (i.e. a category or class). Bayes theorem is a means of predicting the **posterior probability** that an event occurs given the observation of another event. It relies on the **conditional probability** that the other event occurs given the observation of the event and the **prior probabilities** that both events occur in general.

Thus,

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

where  $P(A|B)$  is the **posterior probability**,  $P(B|A)$  is the **conditional probability**, and  $P(A)$  and  $P(B)$  are the **prior probabilities**.

The **prior probability**  $P(A)$  is otherwise known as the **class probability** and is equal to the general probability of encountering a particular class. In our case, we have chosen to keep our classes balanced (i.e. there are equal numbers of documents in both), so the **class probability** will always be 50%. The other **prior probability**  $P(B)$  or **evidence**, then, refers to the probability of encountering the features within a document **independent** of the class label. The **conditional probability**  $P(B|A)$  refers to the likelihood of encountering the features within a document given that those features belong to a particular class. The **evidence** is often ignored in the final classification step on the basis that it acts merely as a scaling factor when trying to maximize which class produces the greater **posterior probability**. The **posterior probability** is the probability that a particular document belongs to a class given the observed features within that document. The Naive Bayes algorithm maximizes this probability in order to predict which class best fits a document. In all cases where we must calculate the probability for an entire feature-set, we simply take the product of all the extracted feature probabilities in the document. Consider this classification example:

if  $P(\omega = \text{subtweet} \mid \mathbf{x}) \geq P(\omega = \text{non-subtweet} \mid \mathbf{x})$  classify as subtweet,  
 else classify as non-subtweet.

Dropping the evidence term on the basis that it is constant for both classes, we expand that:

$$P(\omega = \text{subtweet} \mid \mathbf{x}) = P(\mathbf{x} \mid \omega = \text{subtweet}) \cdot P(\text{subtweet})$$

$$P(\omega = \text{non-subtweet} \mid \mathbf{x}) = P(\mathbf{x} \mid \omega = \text{non-subtweet}) \cdot P(\text{non-subtweet})$$

Assuming the **posterior probability** for the former is greater than or equal to the latter, the classifier predicts that the document fits that class. The naive assumption maintains that all features are treated as conditionally independent (i.e. that the presence or omission of a particular feature does not change the likelihood of encountering other features), and although this is frequently violated, Naive Bayes often performs well anyway (Zhang, 2004).

For cases in which the classifier encounters a feature absent from the features which were used to train it, a so-called **zero probability** appears. Because the probability of encountering the feature is 0, **additive smoothing** is often utilized to appropriately weight new features using an extra term  $\alpha$ , so the probability that an entire feature-set fits into a specific class is not 0. In this project, we use **Laplace smoothing** ( $\alpha = 1$ ). This technique smooths categorical data by including the pseudo-count  $\alpha$  into each probability estimate. Thus, the probability of a feature given a particular class becomes:

$$P(x_i | \omega_j) = \frac{N_{x_i, \omega_j} + \alpha}{N_{\omega_j} + \alpha d} \quad (i = (1, \dots, d))$$

where  $N_{x_i, \omega_j}$  is the number of times feature  $x_i$  appears in samples from class  $\omega_j$ ,  $N_{\omega_j}$  is the total count of all features in class  $\omega_j$ ,  $\alpha$  is the parameter for additive smoothing, and  $d$  is the dimensionality of the feature vector  $\mathbf{x} = [x_1, \dots, x_d]$ . Compared to datasets which contain millions of tweets such as *Sentiment140* (Go et al., 2009), our classifier has access to significantly fewer documents. We utilize **Laplace smoothing** because when the classifier is tested on new tweets it will likely encounter never before seen features given the limited size of the ground truth dataset.

### 2.2.5 Statistical Considerations

In the binary classification of subtweets and non-subtweets, we consider true positives (TP) to be **true subtweets** which were correctly labeled as **predicted subtweets**, false positives (FP) to be **true non-subtweets** which were incorrectly labeled as **predicted subtweets**, true negatives (TN) to be **true non-subtweets** which were correctly labeled as **predicted non-subtweets**, and false negatives (FN) to be **true subtweets** which were incorrectly labeled

as **predicted non-subtweets**. As such, there are two ways for the classifier to be wrong: it can produce false negatives and false positives. Using TP, FP, TN, and FN, we can calculate statistical measurements on the performance of our classifier such as **precision**, **recall**,  $F_1$  **score**, and **null accuracy**.

#### 2.2.5.1 Precision

**Precision** refers to the ratio between the true positives, and the true positives and false positives. It is also known as the **positive predictive value**.

$$P = \frac{TP}{TP + FP}$$

#### 2.2.5.2 Recall

**Recall**, then, refers to the ratio between the number of true positives, and the true positives and false negatives. It is also known as the **sensitivity**.

$$R = \frac{TP}{TP + FN}$$

#### 2.2.5.3 Accuracy

The **accuracy** is the ratio between the true positives and the true negatives, and the true positives, true negatives, false positives, and false negatives. **Accuracy** alone is a particularly bad quantifier of how well a classifier performs when working with data which is class-imbalanced (i.e. there are not equal numbers of items in each class). In our ground truth dataset, the classes are balanced so measuring accuracy will still be informative.

$$A = \frac{TP + TN}{TP + TN + FP + FN}$$

#### 2.2.5.4 F1 Score

The  $F_1$  **score** is a weighted average of the **precision** and **recall**. Thus, it takes both false positives and false negatives into account.

$$F1 = \frac{2 * (P * R)}{P + R}$$

#### 2.2.5.5 Null Accuracy

The **null accuracy** is the accuracy which is obtained by always predicting the most frequent class. Because there are two classes and the tweets within the ground truth dataset equally compose both, the **null accuracy** will always be 50%.

#### 2.2.6 K-Folds Cross-Validation

Instead of using the entire ground truth dataset as training data for the Naive Bayes classifier, we can split it into a **training set** and a **test set**. The **training set** is fed to the classifier, and the **test set** is used to observe statistics about its performance. Figure 2.2.1 depicts cross-validation using  $k$ -folds, in which we make random splits in the dataset  $k$  times to create several **training set** and **test set** sections. In  $k$ -folds, we then take statistical measurements of how well the classifier performs on the **test set** for each fold. If we made one single split into training and testing sections of our ground truth dataset and only used that single **test set** to gather statistics on the classifier’s performance, we would not be able to confirm that those statistics were representative of all the data in the entire ground truth dataset. Instead, we perform 10-fold cross validation, choosing 90% of the data to be the **training set**, and 10% to be the **test set** in each fold. **Precision**,  $F_1$  **score**, and **recall** are calculated within each iteration of the 10 folds, thus utilizing 10 different **test sets**. Finally, we use the averages of those statistics across all folds to measure the overall performance of the classifier.

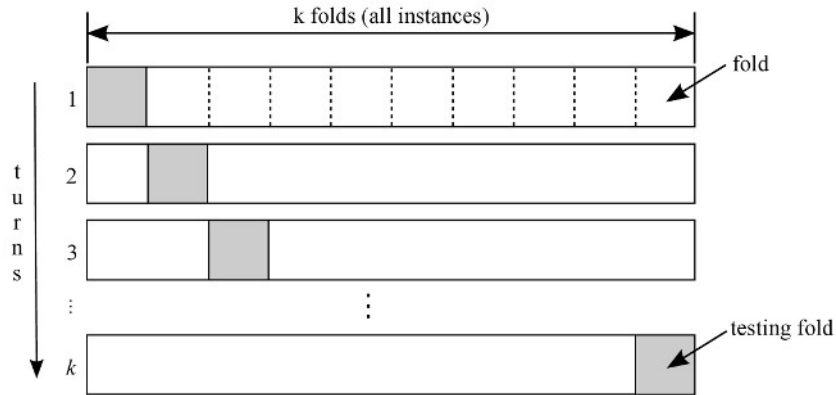


Figure 2.2.1. K-Folds Cross-Validation Example (Borovicka et al., 2012)

### 2.2.7 Training & Testing Naive Bayes

We utilize Scikit Learn’s API for machine learning (Pedregosa et al., 2011) to create a **pipeline**. In Scikit, **pipelines** make managing machine learning algorithms easy by consolidating their parts into one object with configurable attributes.

```

1 sentiment_pipeline = Pipeline([
2     ("vectorizer", TfidfVectorizer(tokenizer=tokenizer.tokenize,
3                                   ngram_range=(1, 3),
4                                   stop_words="english")),
5     ("classifier", MultinomialNB())
6 ])

```

Our **pipeline** contains a vectorizer and a classifier. We change the default arguments for Scikit’s TF-IDF vectorizer to use NLTK’s tweet tokenizer, and specify that we want to calculate our TF-IDF vectors using unigrams, bigrams, and trigrams. Then, we set the vectorizer to use Scikit’s default English language stop words. The Multinomial Naive Bayes Classifier we implement using Scikit includes **Laplace smoothing** ( $\alpha = 1$ ) by default.

Scikit also has a convenient **KFold** object which we utilize to perform cross-validation on our classifier. The goal of  $k$ -folds cross-validation is to train the classifier and acquire statistics on its performance while treating  $k$  different parts of the ground truth dataset as **test sets**. Within a single iteration for 10 iterations, the **test set** will always be 10% of the entire ground truth dataset, however the next iteration will use a different section. We ultimately average those statistics to understand the overall performance of the classifier.

```

1 def confusion_matrices(training_data, num_folds=10):
2     text_training_data = np.array([row[0] for row in training_data])
3     class_training_data = np.array([row[1] for row in training_data])
4     kf = KFold(n_splits=num_folds, random_state=42, shuffle=True)
5     cnf_matrix_test = np.zeros((2, 2), dtype=int)
6     for train_index, test_index in kf.split(text_training_data):
7         text_train, text_test = (text_training_data[train_index],
8                                 text_training_data[test_index])
9         class_train, class_test = (class_training_data[train_index],
10                                  class_training_data[test_index])
11
12         sentiment_pipeline.fit(text_train, class_train)
13         predictions_test = sentiment_pipeline.predict(text_test)
14         cnf_matrix_test += confusion_matrix(class_test, predictions_test)

```

In each iteration of the 10 folds, the above program splits apart a training set and a **test set**.

The classifier is trained on the **training set** using `sentiment_pipeline.fit`, and the classifier's predictions for the **test set** in that fold are used to add toward a confusion matrix which will categorically visualize the performance of the classifier. We also calculate the **precision**, **recall**, and  $F_1$  **score** for each fold's individual **test set**.

## 2.3 Creating the Twitter Bot

We use Tweepy to interact with the Twitter API. It provides a convenient object for streaming Twitter data in real time. The `StreamListener` class can track tweets by searching for specific users, locations, and keywords. For our purposes, it has to be extended to track subtweets.

```

1 class StreamListener(tweepy.StreamListener):
2     def on_status(self, status):
3         id_str = status.id_str
4         screen_name = status.user.screen_name
5         created_at = status.created_at
6         retweeted = status.retweeted
7         in_reply_to = status.in_reply_to_status_id_str
8         text = status.full_text
9
10        # Genericize extra features and clean up the text
11        text = (urls_pattern.sub("GENERIC_URL",
12                                at_mentions_pattern.sub("GENERIC_MENTION",
13                                                        names_pattern.sub("GENERIC_NAME",
14                                                                        text)))
15                .replace("\u2018", "'")
16                .replace("\u2019", "'")
17                .replace("\u201c", "\"")
18                .replace("\u201d", "\"")
19                .replace("&quot;", "\"")
20                .replace("&amp;", "&")
21                .replace("&gt;", ">")
22                .replace("&lt;", "<"))
23
24        tokens = tokenizer.tokenize(text)
25
26        english_tokens = [english_dict.check(token) for token in tokens]

```

```

27     percent_english_words = sum(english_tokens)/float(len(english_tokens))
28
29     # Make sure the tweet contains some English
30     is_english = False
31     if percent_english_words >= 0.1:
32         is_english = True
33
34     # Calculate the probability using the pipeline
35     subtweet_probability = sentiment_pipeline.predict_proba([text]).tolist()[0][1]

```

This part of the program is for live classification of subtweets and gathers information on tweet status objects as they are encountered. We extract data from the tweet object including the ID of that tweet and the text contained within it. We utilize the same techniques for cleaning and genericizing the tweet which we used in preparing our ground truth data for the classifier. The `pipeline` has a `predict_proba` method which takes as its input a list of strings and outputs an array of probabilities for each class. `subtweet_probability`, then, uses that method to predict the probability that the tweet fits the “subtweet” class according to the Naive Bayes classifier and the vectorizer we used in our `pipeline`. We also check that the potential subtweet meets specific requirements before the Twitter bot will interact with it.

```

1  if all([subtweet_probability >= THRESHOLD,
2         "RT" != text[:2], is_english,
3         not retweeted, not in_reply_to]):

```

Included in this conditional statement is a comparison to determine if the probability that the tweet is a subtweet meets a specific threshold. We do not want to call out subtweets unless the probability is high enough. Following this check, we can interact with the tweet in several ways.

```

1  # Quote the tweet
2  api.update_status(("Is this a subtweet? {:.2%} \n" +
3                   "https://twitter.com/{}/status/{}".format(subtweet_probability,
4                                                             screen_name,
5                                                             id_str))
6
7  # Like the tweet
8  api.create_favorite(id_str)
9
10 # Reply to the tweet
11 api.update_status("@{} Is this a subtweet? {:.2%}".format(screen_name,
12                                                             subtweet_probability),
13                    id_str)

```

To quote the potential subtweet means that the tweet being referenced is embedded within our own tweet with a caption above it. Finally, we instantiate our custom `StreamListener` class and use it to filter through tweets in real time.



```
1 stream_listener = StreamListener()
2 stream = tweepy.Stream(auth=api.auth, listener=stream_listener, tweet_mode="extended")
3
4 stream.filter(follow=user_ids, stall_warnings=True, languages=["en"], async=True)
5 print("Streaming has started.")
6 sleep(DURATION)
7 stream.disconnect()
```

The `user_ids` list contains strings of Twitter user IDs. Every time a user whose ID is in the list sends a tweet, the program will classify that tweet and use the predicted probability that it is a subtweet to call it out on Twitter from the account linked to the application's API credentials. We filter the stream asynchronously in order to use the `sleep` function from the `time` Python library, so we can run the Twitter bot for a limited number of seconds.

# 3

## Results

### 3.1 Ground Truth Dataset

No corpus already exists which provides examples of subtweets and non-subtweets. To produce such a corpus, the acquisition of the ground truth dataset relies on assumptions and generalizations of how Twitter users accuse others of subtweeting. For each subtweet and non-subtweet we acquired, we **did not** also acquire every single reply to that tweet, because doing so would violate the API’s rate limits (Twitter, 2018). Instead, we maintain that using Twitter’s search API, a single reply to a tweet is enough to indicate that it **is** or **is not** a subtweet. By ignoring the associated metadata of tweet objects and their replies, we obtain the unicode text contained within subtweets and non-subtweets for our ground truth dataset.

True Subtweets	True Non-Subtweets
I know who I want to take me home	He’s followed the putrid smell of GENERIC_MENTION which has led him to GENERIC_MENTION’s whereabouts. Some odd neighbor boy was watching him as he approached the house. Liam didn’t hesitate to drink the young one dry. “HONEY, I’M HOME!” He calls, kicking in the door.
One bad chapter doesn’t mean your story is over	YG so low? I guess the only group that sold physicals well is Bigbang GENERIC_URL
What do you offer someone who doesn’t like tea, coffee, or hot chocolate? YOU’RE A GUEST, YOU NEED A HOT DRINK, IT IS THE LAW.	is everything ok?? GENERIC_URL
on the flip side, i agree that “big” accounts have an obligation to do their due diligence to not sic their followers on others	The parade of falsehoods about CIA nominee Gina Haspel GENERIC_URL
im bout to drop like 20lbs and then become a rapper.	Very first day of fortnite got 2 second place games and a third place gg not bad for a trash player :p
It’s hilarious that your only “godly” when it’s convenient for you or when you want to put it on twitter.	Rosie the Corpse Flower bloomed at Tucson Botanical Gardens: See pix allery here GENERIC_URL #Tucson GENERIC_MENTION GENERIC_URL
It really bothers me when people try to take advantage or assume I do shoots for free GENERIC_URL	What is Uncle’s GREATEST prediction?
I just don’t get people lol	Didi Gregious is my new favorite player.
y’all kids today with your “waifu this” and “waifu that,” back in my day our waifus had metal jaws and carried high-frequency blades	GENERIC_MENTION we need to get you drawn as Dekamaster Doggy Krueger
I follow this gorgeous dog account on insta and his owner has such an annoying voice. SHUT UP AND LET ME ENJOY YOUR MAGNIFICENT DOG! I’M NOT HERE FOR YOU!	“Huge Caravan” Of Central Americans Is Headed For The U.S. Border In Hopes Of Asylum GENERIC_URL

When the vectorizer runs over these strings in order to tokenize them and extract features for use in the classifier, GENERIC\_URL, GENERIC\_MENTION, and GENERIC\_NAME stand in place of actual URLs, mentions of user names, and English first names in the original tweets. We use these generic tokens because we want the classifier to probabilistically acknowledge the presence of these features while specifically ignoring the content within them.

### 3.1.1 Characteristics of the Ground Truth Dataset

With tens of thousands of tweets contained within the ground truth dataset, understanding the characteristics of these tweets by example alone is insufficient. To improve our overall understanding of these characteristics, we gathered statistics on the distributions of tweet lengths, punctuating characters, stop words, and unique words in both the **true subtweets** and **true non-subtweets** parts of the dataset.

Twitter serves as a micro-blogging framework which limits the lengths of English text posts to 280 characters. Because Twitter is also a multimedia platform which supports embedding URLs, images, and videos, tweets do not necessarily contain this maximum number of characters. In contrast to the popularity of multimedia tweets, subtweets are characterized as exclusively text-based. Between both **true subtweets** and **true non-subtweets**, Figure 3.1.1 illustrates that the ground truth dataset contains more short tweets ( $1 \leq \text{length} \leq 125$ ) than long tweets ( $125 < \text{length} \leq 280$ ). The subscripts  $s$  and  $n$  respectively indicate subtweets and non-subtweets.

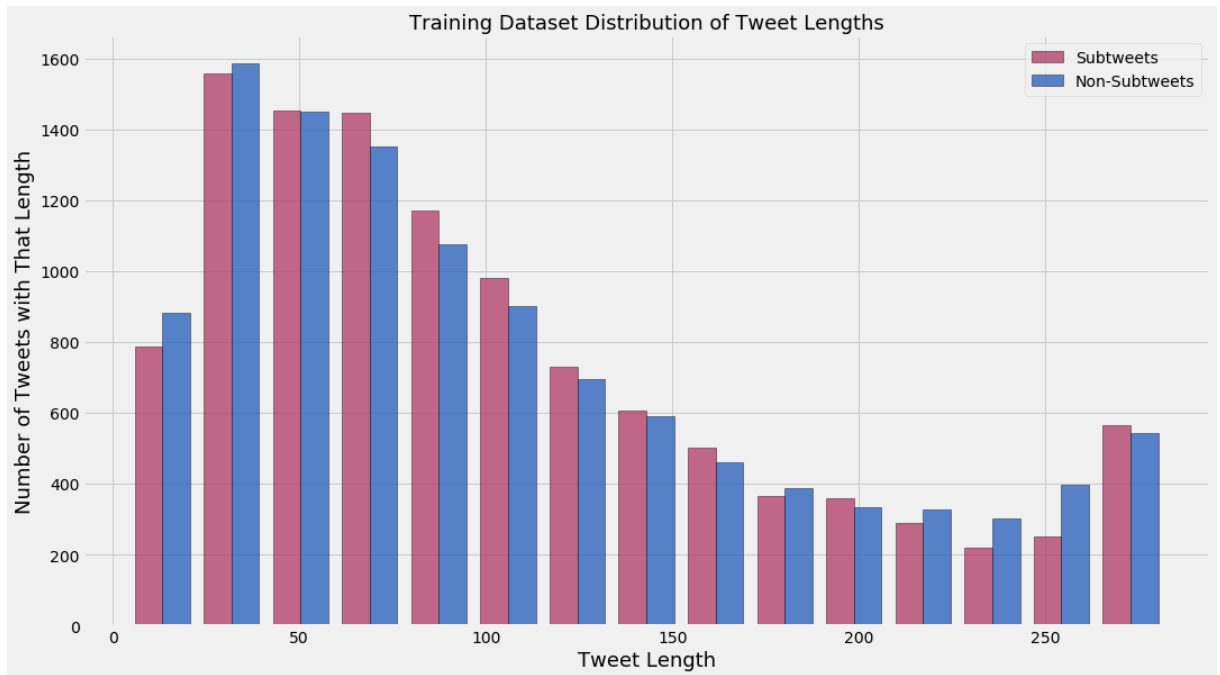


Figure 3.1.1. Histogram of Ground Truth Tweet Lengths  
 $(\bar{x}_s = 103.71, s_s = 71.99), (\bar{x}_n = 105.63, s_n = 74.95)$

Twitter does not enforce any grammar or syntax regulations in tweets, thus punctuating characters are often ignored altogether for use in the stream-of-consciousness style of writing. As shown in Figure 3.1.2, both **true subtweets** and **true non-subtweets** follow this trend in omitting punctuation, on average containing just under 2 punctuating characters (such as quotation marks, periods, commas, and apostrophes).

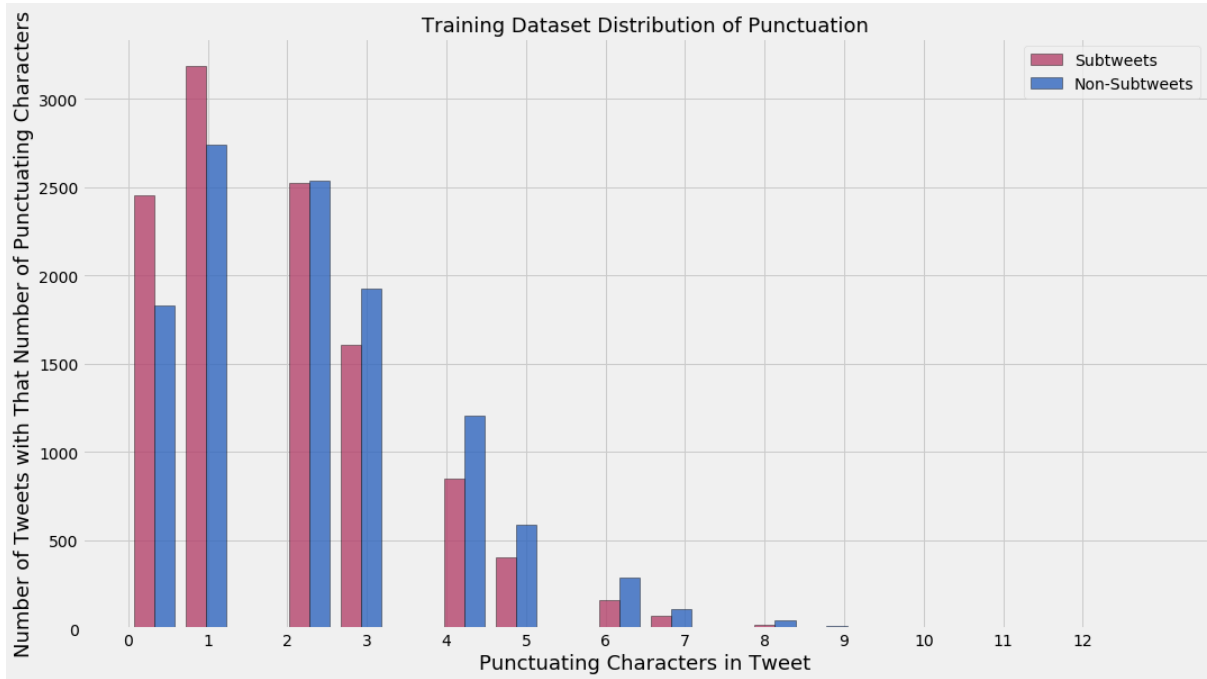


Figure 3.1.2. Histogram of Ground Truth Tweet Punctuation  
 $(\bar{x}_s = 1.79, s_s = 1.54), (\bar{x}_n = 2.16, s_n = 1.68)$

English stop words are ignored by the vectorizer because they are too common to help the classifier distinguish between the appropriate features for the classes. Figure 3.1.3 shows that our ground truth dataset contains tweets which are more likely to contain between 5 and 8 unique stop words than they are to contain fewer. In comparison, Figure 3.1.4 illustrates the distribution of the number of **unique English words** per tweet in the ground truth dataset. As we expect of text-based posts under 280 characters in length, tweets containing fewer unique English words are more common.

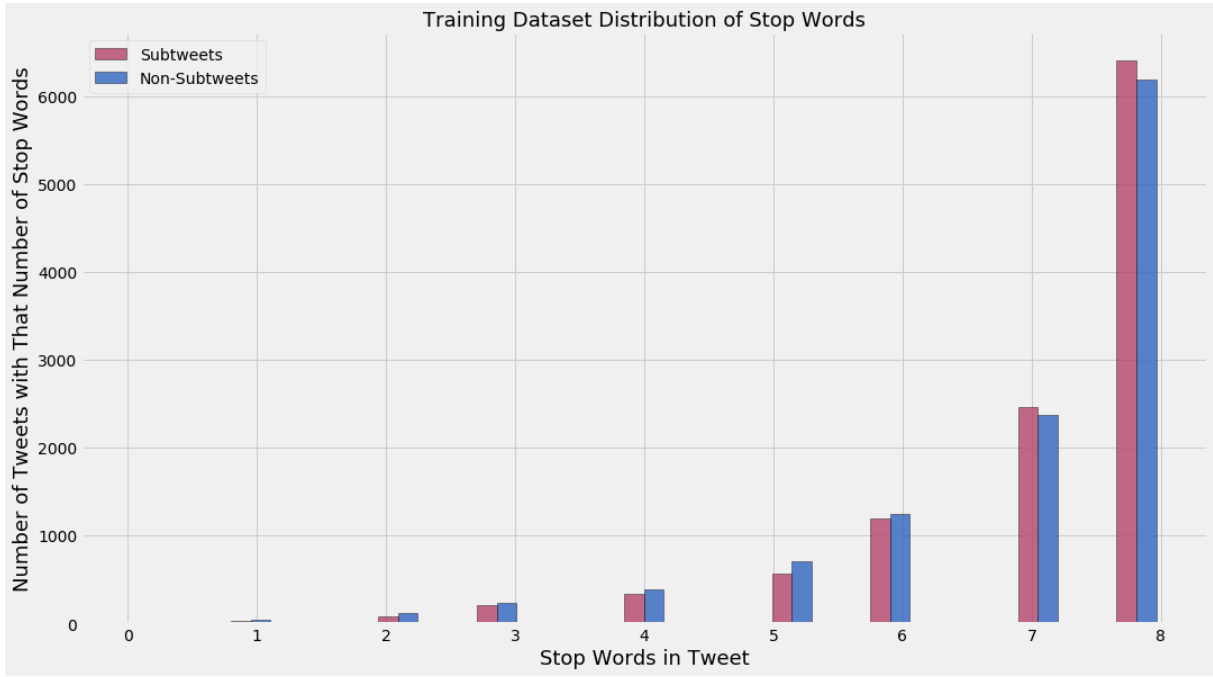


Figure 3.1.3. Histogram of Ground Truth Unique English Stop Words  
 $(\bar{x}_s = 7.14, s_s = 1.31), (\bar{x}_n = 7.05, s_n = 1.39)$

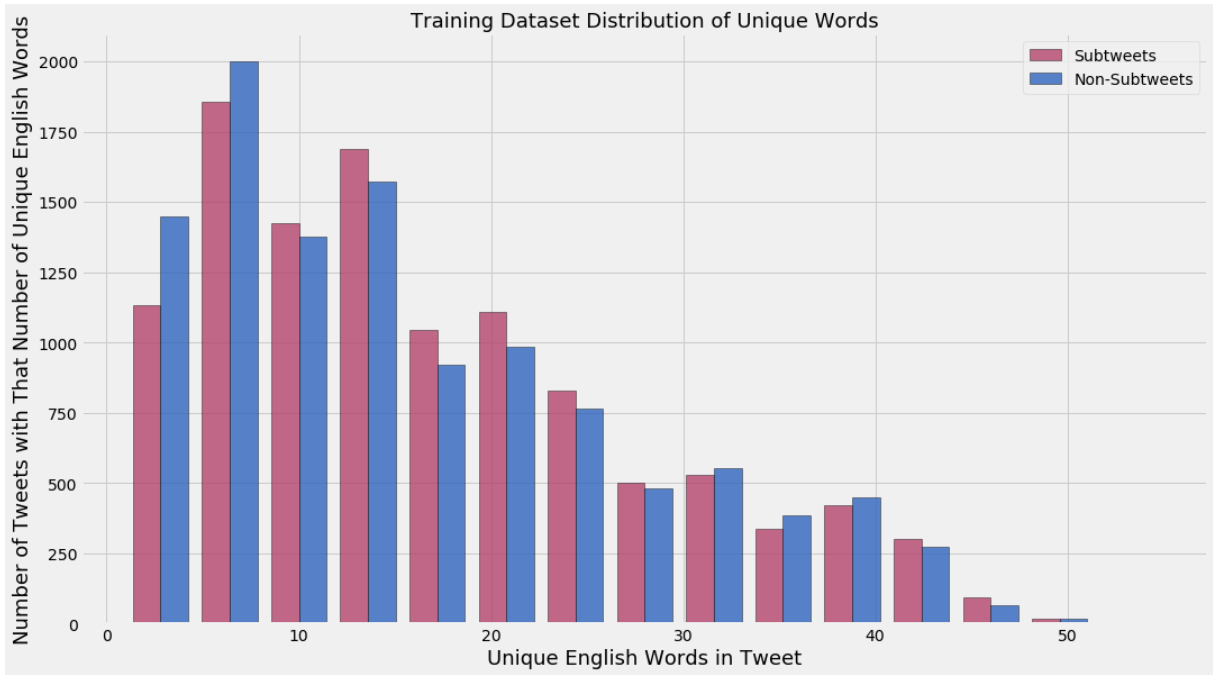


Figure 3.1.4. Histogram of Ground Truth English Unique Words  
 $(\bar{x}_s = 16.80, s_s = 10.95), (\bar{x}_n = 16.25, s_n = 11.15)$

## 3.2 Statistics

### 3.2.1 The Confusion Matrix

We used  $k$ -folds cross-validation to test the performance of our Naive Bayes classifier on each **test set** fold of our ground truth dataset. Each individual **test set** was composed of only 10% of the tweets in the entire dataset. By keeping track of the classifier's predictions on this **test set**, we accumulated a **confusion matrix** of **true positives**, **true negatives**, **false positives**, and **false negatives** for all 10 folds. Figure 3.2.1 illustrates these outcomes in terms of raw counts and normalized over the entire ground truth dataset.

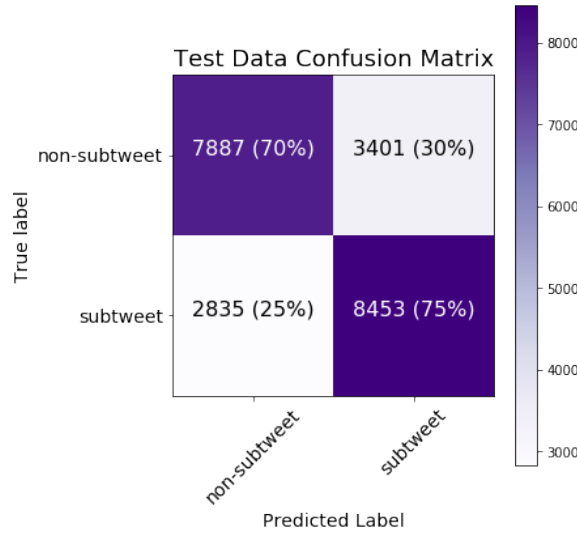


Figure 3.2.1. Confusion Matrix on Test Data from Each Fold on the Ground Truth Dataset

We read the confusion matrix by matching rows with columns to see how well the classifier performs on **true subtweets** and **true non-subtweets** in making classifications to produce **predicted subtweets** and **predicted non-subtweets**. These results indicate that the Naive Bayes classifier performs 5% better when classifying true positives than it does when classifying true negatives. In comparison, it is 5% **worse** at classifying false positives than it is at classifying false negatives.

### 3.2.2 Precision, Recall, & $F_1$

The statistics gathered using  $k$ -folds cross-validation were averaged across all 10 folds, and measure the overall statistical performance of our classifier.

	Precision	Recall	$F_1$ Score
Non-Subtweets	0.7357	0.6988	0.7166
Subtweets	0.7132	0.7490	0.7305

The average  $F_1$  scores for non-subtweets and subtweets are respectively 71.66% and 73.05%.

Thus, our classifier performs similarly well above the 50% **null accuracy** on both.

## 3.3 Known Subtweeters

In order to test how the classifier performed on tweets from known subtweeters, we acquired all publicly available tweets from 11 different accounts from which users had previously subtweeted. Figure 3.3.1 shows the distributions of tweets classified as subtweets under particular subtweet probability thresholds using the classifier on all the tweets from these accounts.

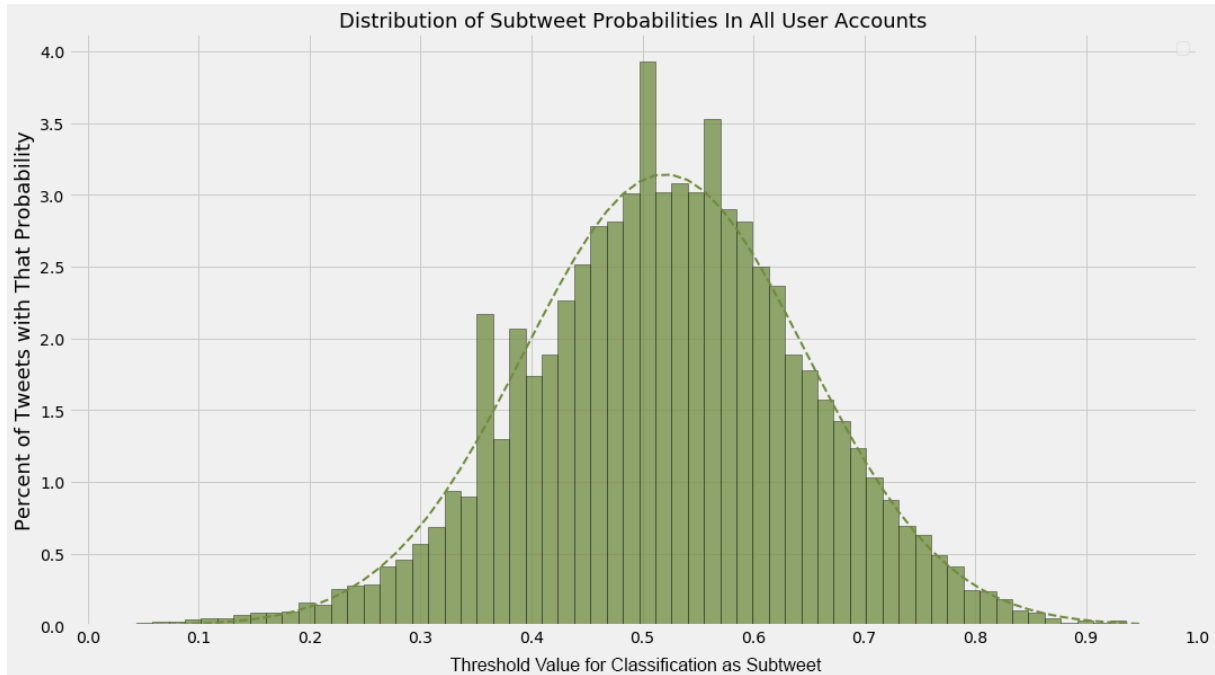


Figure 3.3.1. Distribution of Subtweet Probabilities for 11 Known Subtweeters( $\bar{x} = 0.520$ ,  $s = 0.127$ )



Of these 11 accounts, the Naive Bayes classifier predicts that out of 26,928 tweets, 15,538 were subtweets (if the threshold for classifying a particular tweet as a subtweet was a probability of 50%). In comparison, if we only accept classifications of at least 75% and make that our threshold, then just 924 tweets (3.4%) were predicted to be subtweets. Of these accounts, we selected the users who produced the minimum percentage of subtweets (27.6%), the median percentage of subtweets (58.8%), and the maximum percentage of subtweets (68.4%) using the former (50%) threshold. Figure 3.3.2 shows the distributions of tweets classified as subtweets under particular subtweet probability thresholds using the classifier for all the tweets from these three users' accounts.

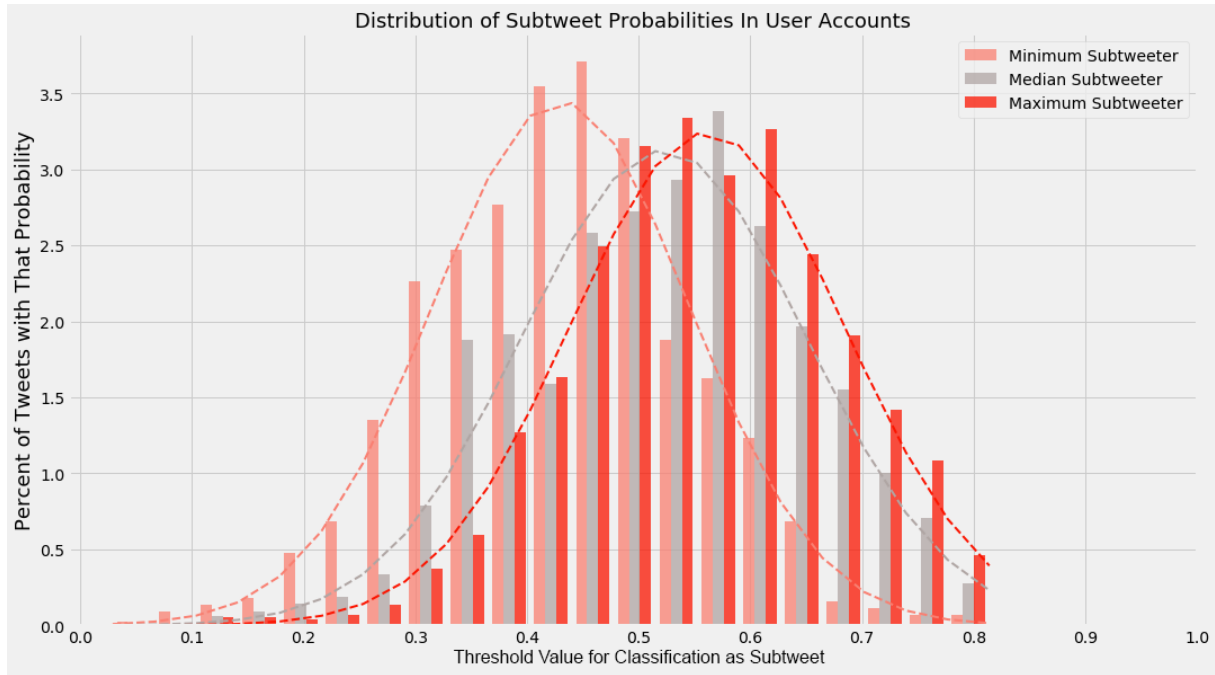


Figure 3.3.2. Distribution of Subtweet Probabilities for 3 Known Subtweeters  
 $(\bar{x}_{min} = 0.430, s_{min} = 0.116), (\bar{x}_{med} = 0.523, s_{med} = 0.128), (\bar{x}_{max} = 0.561, s_{max} = 0.123)$

### 3.4 Most Informative Features

By iterating through the entire vocabulary of features on which the classifier was trained, and sorting those features by the probability that the feature belongs in the subtweet class, we are able to see which features are most informative for classifying tweets as **predicted subtweets**.

Token	$P(T \omega = \textit{subtweet})$
.	-7.655
,	-8.022
GENERIC_URL	-8.049
”	-8.148
people	-8.427
?	-8.552
like	-8.673
don’t	-8.683
just	-8.744
i’m	-8.833
!	-8.864
it’s	-9.080
. GENERIC_URL	-9.100
:	-9.110
know	-9.184
...	-9.204
you’re	-9.210
twitter	-9.284
love	-9.329
*	-9.404
friends	-9.418

The only bigram in the top 20 most informative features is “. GENERIC\_URL,” which can only occur when a URL follows the end of a sentence. In comparison, the rest of these terms are unigrams of which 8 are exclusively punctuating characters (period, comma, quotation, question, exclamation mark, colon, ellipsis, and asterisk). The presence of these most informative features indicates the importance of punctuation in predicting subtweets.

### 3.5 The Twitter Bot

After training and testing our classifier, we utilized it in creation of a Twitter bot which interacts with **predicted subtweets** in real time. By limiting the minimum **predicted subtweet** probability threshold to 75%, we prevent the Twitter bot from interacting with users too much. Figures 3.5.1 and 3.5.2 show (censored) examples of these interactions.



Figure 3.5.1. Example of the Twitter Bot Quoting Users' Tweets

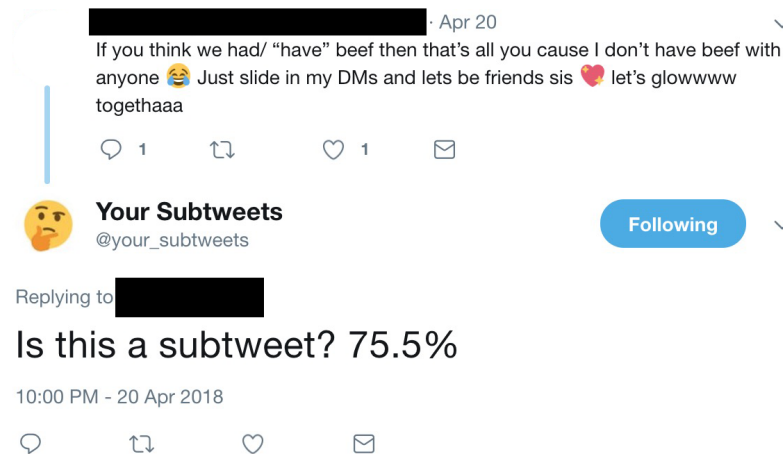


Figure 3.5.2. Example of the Twitter Bot Replying to a User's Tweet

### 3.6 Discussion

Subtweets have garnered attention from news organizations (Madrigal, 2014), social scientists (Edwards and Harris, 2016), and governmental officials (Ohlheiser, 2017), but sentiment analysis of subtweets has been entirely unresearched. Utilizing the Twitter API, we acquired data for training and testing a Naive Bayes classifier, and developed a Twitter bot which actively calls-out subtweets in real time. We obtained 12,169 **true subtweets** and 21,411 **true non-subtweets** to

be used in our ground truth dataset. After confirming they were English and removing duplicates, there were 11,288 **true subtweets** and 19,289 **true non-subtweets** remaining. In order to preserve a balance between the classes, we limited the number of **true non-subtweets** to be the same as the number of **true subtweets**. Thus, our final ground truth dataset contains 22,576 tweets. Our classifier identifies non-subtweets with an  $F_1$  **score** of 71.66% and subtweets with an  $F_1$  **score** of 73.05%, performing significantly better than the **null accuracy** of 50%.

For the 2017 International Workshop on Semantic Evaluation (SemEval), numerous teams submitted approaches for task 4, subtask A: English Twitter polarity classification on a 3-point scale (**positive**, **neutral**, and **negative** sentiment) (Rosenthal et al., 2017). Although this particular classification is ternary instead of binary (such as in our case of classifying subtweets and non-subtweets), the outcomes of the competition on this particular subtask are still informative for us in comparing the performance of our method to others. The team which performed best on this subtask was **DataStories** with an average  $F_1$  **score** of 67.7% (Baziotis et al., 2017). Baziotis et al. utilized Long Short-Term Memory (LSTM) networks to achieve this score. In comparison to our average  $F_1$  **score** of 72.4%, we performed similarly well on a particularly less complicated problem than 3-point general sentiment classification.

Because we utilized our own novel method for data acquisition, we did not acquire a ground truth dataset which is as successful in representing the features of actual subtweets and non-subtweets as those which can be obtained using a service such as Amazon’s Mechanical Turk which hand-labeled crowdsourced data. The Twitter API rate limitations prevented acquisition of tweets which were more than one week old at the time of downloading, and we only ever acquired one single reply to each subtweet and non-subtweet as the distinguishing characteristic of **true subtweets** and **true non-subtweets**. If a Twitter API search query exists which can find all the replies to a particular tweet, it will serve better than our method for finding only one.



# 4

## Conclusion

### 4.1 Summary of Project Achievements

We treated a new colloquial form of online harassment—“subtweeting”—as a text classification problem and used the probabilistic Naive Bayes classification algorithm to identify it. We judged the performance of the algorithm and went beyond identification to engage with subtweets, promoting publicity for content which is deliberately written to be unseen by the targeted party. We utilized Twitter’s API to demonstrate a potential use for sentiment analysis on this kind of text.

All data acquired and all programs developed for this project have been made publicly available on GitHub at <https://github.com/segalgouldn/Senior-Project-Subtweets>.

### 4.2 Future Work & Considerations

Our implementation of Naive Bayes exclusively made classifications using features from the unicode text contained within tweets, but other features related to the metadata contained within tweet objects and their replies will probably prove fruitful in producing a better subtweets classifier. We did not test other classification algorithms in favor of pursuing Naive Bayes singularly, but there exists no reason to **not** utilize others. Using topic modeling via methods such as non-

negative matrix factorization and latent dirichlet allocation, we can analyze the topics about which most users are tweeting, and potentially group these users into subtweet-topic networked communities through clustering algorithms such as  $k$ -means.

Finally, we treated this project as a binary classification problem between subtweets and non-subtweets, but it can be expanded. What would we find if we trained a classifier to distinguish ironic tweets, sarcastic tweets, mocking tweets, and subtweets? What are the features of the figurative language used in each of these examples? These questions motivate further research.

# Bibliography

Allport, Gordon W. 1954. *The nature of prejudice*.

Baziotis, Christos, Nikos Pelekis, and Christos Doukeridis. 2017. *Datastories at SemEval-2017 task 4: Deep LSTM with attention for message-level and topic-based sentiment analysis*, Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017), 747–754.

Bird, Steven and Edward Loper. 2004. *NLTK: the natural language toolkit*, Proceedings of the ACL 2004 on Interactive poster and demonstration sessions, 31.

Borovicka, Tomas, Marcel Jirina Jr, Pavel Kordik, and Marcel Jirina. 2012. *Selecting representative data sets*.

Burnap, Pete, Matthew L Williams, Luke Sloan, Omer Rana, William Housley, Adam Edwards, Vincent Knight, Rob Procter, and Alex Voss. 2014. *Tweeting the terror: modelling the social media reaction to the woolwich terrorist attack*, Social Network Analysis and Mining **4**, no. 1, 206.

Dewey, Caitlin. 2016. *Study confirms what you always knew: People who subtweet are terrible*, available at <https://goo.gl/SeV3mx>.

Dorsey, Jack. 2006. *inviting coworkers*, available at <https://twitter.com/jack/status/29>.

Edwards, Autumn and Christina J Harris. 2016. *To tweet or subtweet?: Impacts of social networking post directness and valence on interpersonal impressions*, Computers in Human Behavior **63**, 304–310.

Ghosh, Aniruddha, Guofu Li, Tony Veale, Paolo Rosso, Ekaterina Shutova, John Barnden, and Antonio Reyes. 2015. *Semeval-2015 task 11: Sentiment analysis of figurative language in twitter*, Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), 470–478.

Go, Alec, Richa Bhayani, and Lei Huang. 2009. *Twitter Sentiment Classification using Distant Supervision*, CS224N Project Report, Stanford, 12.

Hassler, Chelsea. 2016. *Subtweeting looks terrible on you. (you know who you are.)*, available at <https://goo.gl/NCz27z>.



- Kelly, Ryan. 2016. *PyEnchant: a spellchecking library for Python*.
- Madrigal, Alexis C. 2014. *Behind the machine's back: How social media users avoid getting turned into big data*, available at <https://goo.gl/h36jxx>.
- Medhat, Walaa, Ahmed Hassan, and Hoda Korashy. 2014. *Sentiment analysis algorithms and applications: A survey*, Ain Shams Engineering Journal **5**, no. 4, 1093–1113.
- Nassirtoussi, Arman Khadjeh, Saeed Aghabozorgi, Teh Ying Wah, and David Chek Ling Ngo. 2014. *Text mining for market prediction: A systematic review*, Expert Systems with Applications **41**, no. 16, 7653–7670.
- Ohlheiser, Abby. 2017. *A running list of all the possible subtweets of President Trump from government Twitter accounts*, available at <https://goo.gl/hFh81R>.
- Pedregosa, Fabian, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, and Vincent Dubourg. 2011. *Scikit-learn: Machine learning in Python*, Journal of machine learning research **12**, no. Oct, 2825–2830.
- Rae, Chelsea. 2009. *I hate when i see people...*, available at [https://twitter.com/Chelsea\\_x\\_Rae/status/6261479092](https://twitter.com/Chelsea_x_Rae/status/6261479092).
- Roesslein, Joshua. 2009. *tweepy Documentation* **5**, available at <http://docs.tweepy.org/en/v3.5.0/>.
- Rosenthal, Sara, Noura Farra, and Preslav Nakov. 2017. *SemEval-2017 task 4: Sentiment analysis in Twitter*, Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017), 502–518.
- Stone, Biz. 2006. *Introducing the Twitter API*, available at [https://blog.twitter.com/official/en\\_us/a/2006/introducing-the-twitter-api.html](https://blog.twitter.com/official/en_us/a/2006/introducing-the-twitter-api.html).
- Twitter. 2018. *Twitter API Docs*, available at <https://developer.twitter.com/en/docs>.
- . 2016. *Twitter Terms of Service*, available at <https://twitter.com/en/tos>.
- . 2018. *Tweet objects*, available at <https://developer.twitter.com/en/docs/tweets/data-dictionary/overview/tweet-object>.
- . 2012. *Twitter turns six*, available at [https://blog.twitter.com/official/en\\_us/a/2012/twitter-turns-six.html](https://blog.twitter.com/official/en_us/a/2012/twitter-turns-six.html).
- . 2018. *Rate Limiting*, available at <https://developer.twitter.com/en/docs/basics/rate-limiting>.
- Urban Dictionary. 2010. *Subtweet*, available at <https://www.urbandictionary.com/define.php?term=subtweet>.
- Van Hee, Cynthia, Els Lefever, and Vronique Host. 2018. *Semeval-2018 Task 3: Irony detection in English Tweets*, Proceedings of the 12th International Workshop on Semantic Evaluation (SemEval-2018).
- Zhang, Harry. 2004. *The optimality of naive Bayes*, AA **1**, no. 2, 3.