

2017

Sentiment Analysis of Tweets Including Emoji Data

Travis LeCompte

Follow this and additional works at: https://repository.lsu.edu/honors_etd



Part of the [Social and Behavioral Sciences Commons](#)

Recommended Citation

LeCompte, Travis, "Sentiment Analysis of Tweets Including Emoji Data" (2017). *Honors Theses*. 836.
https://repository.lsu.edu/honors_etd/836

This Thesis is brought to you for free and open access by the Ogden Honors College at LSU Scholarly Repository. It has been accepted for inclusion in Honors Theses by an authorized administrator of LSU Scholarly Repository. For more information, please contact ir@lsu.edu.

Sentiment Analysis of Tweets Including Emoji Data

Travis LeCompte

Honors Thesis Advisor: Dr. Jianhua Chen

Table of Contents

1 Abstract	2
2 Introduction	3
3 Background	4
3.1 Sentiment Analysis	4
3.2 Classification	4
3.3 Related Works	6
4 Methodology	8
4.1 Data Sources	8
4.2 Emoji Data	9
4.3 Data Labeling	10
4.4 Representation	11
4.5 Experimentation	12
5 Results and Discussion	14
6 Conclusion	18
7 Future Work	19
8 Bibliography	20

1 Abstract

Sentiment analysis of text is a valuable tool used to identify and classify bodies of text for various purposes, including spam detection and threat assessment. In this paper, we examine the effectiveness of using Multinomial Naïve Bayes (MNB) and Support Vector Machine (SVM) classification methods to classify tweets from Twitter that includes nonstandard “Emoji” data. We compare the accuracy of these classification methods with and without the Emoji data included over varying vocabulary sizes. We find that MNB outperforms SVM on the data for large vocabulary sizes, and both classifiers perform slightly better with the Emoji data included.

2 Introduction

Natural language processing (NLP) is a field within artificial intelligence that overlaps heavily with linguistics and focuses on understanding data from raw text. It relies heavily on statistical measures to evaluate sources of information and understand text. NLP is employed throughout software today for various purposes, including syntax and spelling evaluation for word processors, human-computer interaction and automatic language translation. With the growing use of the Internet, a largely text based communication media, NLP has grown in interest for security and management solutions, such as spam detection of emails. Some areas of NLP also work on speech, where emotion and inflection can be used to more accurately understand the information relayed.

Sentiment analysis is a subset of natural language processing that attempts to identify emotion through text. It is a challenging task for machines and can even challenge humans occasionally. The structure of the English language, combined with sarcastic tones and ambiguous grammar, can often be misleading for those without an intuitive understanding of the language. Interest in understanding emotion however has grown alongside the field of artificial intelligence, particularly with the increase in emotional communication over the Internet. The advent of social media in the past decade has led many more individuals than before to communicate about personal and emotional matters through both private and public messaging platforms. Social media owners have become more interested in sentiment analysis to identify emotions on these platforms that can be deemed dangerous, such as threatening or abusive communication. Traditionally this type of behavior is managed by having users report an abusive user, though this can be slow and inaccurate. An automated system can allow for better control and more productive discussion.

3 Background

3.1 Sentiment Analysis

The task of sentiment analysis typically follows a certain framework: collect the raw data, preprocess the data to purify the database, label the data somehow for training purposes, represent the textual data in some vector form for computation, and then feed the data to a certain classifier. Preprocessing is necessary due to the ambiguity inherent in language and human tendency to deviate from standard grammatical structures. This results in very noisy data that is difficult to classify. The cleaning process removes some of this noise, resulting in data that is easier to work with. This can involve removing certain punctuation, trimming words of their various endings to result in only root words (stemming), or filtering the data points by keywords or length.

Most classification relies on supervised learning, where some sort of training is required to teach the system how to recognize the various classes. This in turn requires data to be labeled with classes in advance to teach the system. The labeling process can be done by hand by a human to ensure maximum accuracy, though this is slow and subjective. Other methods involve identifying the most prevalent keywords in text and labeling them, though these automatic methods must then be verified by a human to be accurate. If the dataset itself is not labeled accurately, any classification will also be inaccurate.

Lastly the data must be represented in some way that is easily accessible for classification. This is typically done by representing the text as a vector of features, though it can depend on the classification method. The most common feature descriptors are based on the frequency of words throughout the text. The larger the vocabulary of the data, the larger these feature descriptors will be, which can result in computational limitations. Thus, the features are typically compressed, either by stemming the vocabulary or restricting the feature descriptors to focus only on the most important words.

3.2 Classification

There are several types of classifiers that have been used for text, the two most common being variations on the Naïve Bayes classifier and the Support Vector Machine (SVM). There are other classifiers that can be used with textual data as well, including neural networks and decision trees, though these do not see as much use. Both rely on supervised learning, being trained on one dataset with the elements labeled with proper classes and then tested on a disjoint dataset with the classes missing. Both classifiers learn the features that define each class and then attempt to find which class fits the unlabeled data of the testing dataset.

One of the most commonly used text classifiers is the Naïve Bayes family of classifiers. These classifiers are based around Bayes' theorem of probability. This theorem states that the conditional probability of an event **A** given another event **B** can be defined by the following equation.

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

Figure 1: Bayes' Theorem

This can be extended to a model for classification by defining each data point as some vector of features **x** with an associated class **C_k**. The probability of a given vector pertaining to a specific class can then be calculated as

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

Figure 2: Modified Bayes' theorem for feature vector

However, if we assume that each individual feature **x_i** of the feature vector is conditionally independent of the other features (not always an accurate assumption for text classification, thus the "naïve" title), this calculation can be simplified to

$$p(C_k | x_1, \dots, x_n) = \frac{1}{Z} p(C_k) \prod_{i=1}^n p(x_i | C_k)$$

Figure 3: Bayes' assumption on feature conditional independence

where $p(x_i | C_k)$ is the probability of feature **x_i** occurring in class **C_k** and **Z** is a constant equal to $p(\mathbf{x})$. One can then classify any given feature vector by computing this probability for every possible class and classifying the vector as the class with the highest probability. There are many variations on the model, such as the Multinomial Naïve Bayes, the Gaussian Naïve Bayes and the Bernoulli Naïve Bayes, all of which operate on the probability of the occurrence of words in the text.

By comparison, SVM is a non-probabilistic classifier. The classifier operates by taking a set of points (like the feature vector from Naïve Bayes) and classes associated with these points. The goal is to plot these points and find hyperplanes that best separate the points of each given class, as shown in **Fig. 4**. When testing SVM on an unclassified vector, classification boils down to plotting the feature vector in the required space and observing which area of the space the point falls in. To use **Fig. 4** as an example, if an incoming vector were to fall above the hyperplane, it would be classified as the blue class.

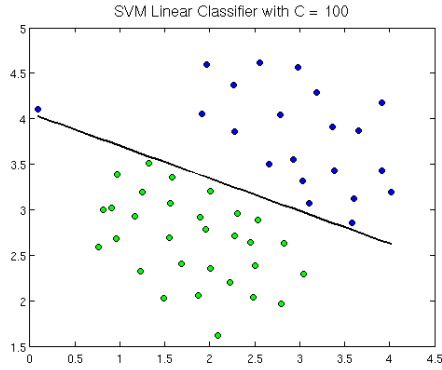


Figure 4: Sample SVM linear classifier output

However, there are many potential hyperplanes that can separate the classes of points shown in the figure above. SVM attempts to find the best fit hyperplane by calculating the distance between the closest points of each class and then selecting a hyperplane between these points that maximizes the margin around the plane to each point (see **Fig. 5**). This ensures the greatest margin of error for the hyperplane and thus the lowest misclassification rate.

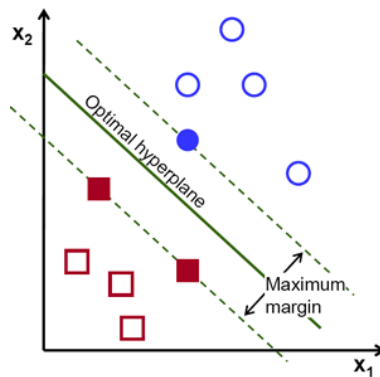


Figure 5: SVM margin representation

There are many variations on SVM, like the variations on Naïve Bayes. Non-linear SVM classifiers can be used to classify datasets that are not linearly separable, while other variations of SVM map the data points to higher dimensions and then find the hyperplanes in these higher dimensions where the points may be more easily separable.

3.3 Related Works

Sentiment analysis in general is a very well-studied topic. The process has been applied to text at many granularities, from classifying entire documents [12] to partitioning documents into sections that can be identified by their opinions, allowing for observing how sentiment changes throughout a document [6], or classifying individual phrases by their positivity and negativity [11]. In this way, sentiment analysis can be used for summarizing reviews for

products in a comparable manner, detecting whether a given user review is positive or negative for a product. Other work has been done to identify the subjectivity of documents, and whether they should be classified as fact or opinion [5].

In the last decade research has increased involving social media, resulting in many experiments on mining political and subjective opinions from blogs and news articles [4, 8]. There has also been other research into the use of Twitter as a source of data [7, 10], and how sentiment analysis can be executed to accurately identify emotion through these short statements [1, 3]. This is like our work, though we include Emoji data in the tweets that is typically unaddressed in these works.

4 Methodology

4.1 Data sources

The first step of sentiment analysis is acquiring a large corpus of textual data for training and testing. Therefore, the initial problem for researchers is the acquisition of high quality, representative samples of data. Traditionally this could be gathered from books or journals and other bodies of printed text. But thanks to the Internet, there is now an ever-growing body of readily available text data to be accessed. This is commonly done via HTML scraping, where one acquires and cleans the HTML source files to access textual data within. This method can quickly gather data and clean it for labeling from a wide variety of websites, be they scientific sources or emotional blogging.

With the advent of social media in the past decade however, there has been a new source of highly emotional human text. Many individuals use Facebook, Twitter, or other social media sites every day and post emotional information or opinions. Studies show that upwards of 80% of individuals in the United States with internet access have Facebook accounts, while between 25-30% of these individuals have accounts on various other media websites (Pinterest, Twitter and Instagram for example) [2]. These sites act as an aggregator for this emotional information that researchers can then access and use for experimentation. Many of them offer some sort of API to developers to access this data with reasonable requests.

For our experimentation, we collected data purely from social media. Specifically, we collected “tweets” from Twitter, which are usually short (less than 140 characters) messages posted by individuals that contain various opinions and emotions. Twitter provides an API for programmers to access their data, requiring a Twitter developer account. We used a package for Python called “twitter” [9] to utilize the API Twitter offers to access their databases and collect the raw tweet data. In this data, all personal identifying information is removed.

To collect specific tweets, we chose to follow the model of Wang et al [10], who tested the viability of tagged tweets as auto-classified textual data. On Twitter, many people use these “hashtags” to mark their tweet as pertaining to a certain topic. For instance, during the 2016 presidential election, many people tagged tweets in favor of candidate Hillary Clinton with “#ImWithHer.” This allows for tweets to be grouped and viewed by their tags to quickly access information, and allows for easy classification of tweets. Our tweets were collected from seven broad emotional categories, described more in Section 4.3.

This very human textual data is extremely raw and dirty, without much structure or regulation. Thus, data must be pre-processed for classification. Our collected data was processed to remove retweets, shorter than four words, or tweets containing URLs. This is an

effort to purify the dataset and keep only those tweets that are relevant for their emotional content. Overall, we collected approximately 700000 tweets from the 7 classes. After cleaning, around 54000 tweets remained as viable for testing. The distribution of the tweets is shown below. While this seems like an extremely small percentage, this is just the nature of raw textual data. This process discards many potential candidates, and cleans the remaining – removing punctuation, trimming references to people or websites, or deleting non-standard characters.

Class	Number of Viable Tweets
Sad	23121
Angry	7672
Happy	7710
Scared	4791
Thankful	4525
Surprised	2167
Love	4075
TOTAL	54061

Table 1: Distribution of Tweets

4.2 Emoji Data

The latter however is of interest due to the recent inclusion of “Emoji” characters in social media. These characters display as certain faces, as seen in **Fig 6**. Their primary use is to display various emotions such as happiness, anger, sadness and love. If the individual in question uses them sans sarcasm, these Emoji characters provide very direct and strong suggestions of the emotion included in the text, even though they are commonly thrown away. This is like the case for punctuation, which is often trimmed in many applications except for sentiment analysis, where question marks and exclamation points can provide insight into the emotions of the writer. The purpose of our experimentation here is thus to investigate the influence that including Emoji data within text can have on the accuracy of emotion classifiers.



Figure 6: Sample of Emoji characters

These Emoji characters are represented in the text as various Unicode values. Web browsers then map these Unicode sequences into the associated images to display. The same process can also be used when collecting the purely textual data – collect the Unicode values and map them to the correct Emoji faces. An example of the Emoji encoding scheme is shown below in **Fig 7**. For example, the Unicode sequence U+1F600 results in a smiling face, while U+1F622 results in a crying face. Instead of mapping the Unicode back into an image, since our plan is to work solely with textual data, we can map the Unicode sequences into specific textual sequences, such as **<ANGRYFACE>** for a mad face and so forth. This allows us to capture the Emoji data in a human recognizable form in text.

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
U+1F5Ex		🔪		👤					💬							💭
U+1F5Fx				📦							📖	🗿	🗼	🗺	🗻	🗺
U+1F60x	😊	😄	😆	😁	😃	😂	😅	😇	😈	😊	😊	😊	😊	😊	😊	😊
U+1F61x	😐	😑	😒	😓	😔	😕	😖	😗	😘	😊	😊	😊	😊	😊	😊	😊
U+1F62x	😡	😠	😞	😟	😟	😟	😟	😟	😟	😟	😟	😟	😟	😟	😟	😟

Figure 7: Emoji Unicode encoding

4.3 Data labeling

After collecting the raw data and cleaning it, the sources must then be labeled with classes to allow for training the classifier. To define the classes of emotion, we followed a similar model to Wang et al [10] where emotion is separated into seven broad classes – *sadness, anger, happiness, fear, thankfulness, surprise, and love*. In collecting the tweets, we then searched for tweets that included one of these seven “hashtags.” Typically, one would have to label each tweet in the data with a class by hand to ensure correct labeling. However, due to the presence of hashtags in these tweets, we can remove the hashtag from the tweet

and classify it as the emotion associated with that hashtag. As shown in Janssens et al [3], the tags can be used for automatically classifying the associated tweets with emotions with high accuracy. For example, a tweet containing the text “#happy” can be assumed to have a happy message.

To ensure the accuracy of the classification process, we need to verify that the automatic labeling of the tweets via their hashtags accurately represents what a human would gather from the tweets. To achieve this, we sampled 100 tweets from the dataset and attempted to label them by hand, comparing our labels to the automatic labels generated by the hashtags. Assuming a uniform distribution of the classes, we would achieve a 1/7 correct rate if we simply guessed the class of each tweet. Our personal ratings however achieved a 70% accuracy rate in agreeing with the automatic labeling process. Therefore, we believe the automatic labeling does a reasonable job of approximating the true classification of each tweet.

4.4 Representation

At this point the data is cleaned to contain only relevant information and labeled for the correct emotional classes. To continue with the classification process, the tweets must be transformed into some numerical representation to feed to a classifier. The most common representations focus on the frequency of the words or characters present in the text. For our experimentation, we utilize both the standard bag of words representation and a bigram representation, both of which rely on word frequency.

Bag of words converts each instance of text into a vector. The standard implementation has each vector of length n , where n is the number of unique words that occur in the entire dataset, also called the vocabulary of the set. The value at each position in the vector is the number of times that specific word occurs in the piece of text in question. This method is simple but effective; its weakness however is that the context and order of the words is lost, only retaining frequency information.

Sentence	Bag of Words	Frequency Vector
I like to play some play sports	{I, like, to, play, some, sports}	[1, 1, 1, 2, 1, 1]
I like to play sports	{I, like, to, play, sports}	[1, 1, 1, 1, 0, 1]

Figure 8: Bag of words representation example

The above figure shows examples of performing the bag of words representation on sample sentences. The resulting vectors can then be used for the classification process. The simplicity allows it to be computed relatively quickly while still retaining high accuracy.

A modification to bag of words is the n-gram approach. Instead of counting the frequency of each individual word, we count the number of times a set of n words occurs in a specific order. Bag of words by default is essentially a unigram approach, considering a set of only one word. By increasing this set size to two or three words, we can retain some contextual and order information within the text while retaining the general frequency information. If we use the same examples from before but apply the bigram approach, we get the following results.

Sentence	Bag of Words	Frequency Vector
I like to play some play sports	{{l, like}, {like, to}, {to, play}, {play, some}, {some, play}, {play, sports}}	[1,1,1,1,1,1]
I like to play sports	{{l, like}, {like, to}, {to, play}, {play, sports}}	[1,1,1,0,0,1]

Figure 9: Bigram representation example

The pairs of words contain more information than simply counting each individual word, though increasing the n-gram size too much would simply result in using the entire sentence as one gram. Typically, bigram or trigram is used for classification. We can then apply the same frequency counting to the bigram representation of the tweet, resulting in a numerical vector that can be used for classification.

4.5 Experimentation

For each of the tweets that pass the previously described cleaning and labeling process, we represent them using both a bag of words and bigram approach. For the bag of words approach, we counted the occurrence of every word present in the dataset. This resulted in a total vocabulary of approximately 10000 words. To reduce computation time, we then restricted the vocabulary to subsets of only the most frequent words. The subset sizes we used were 100, 200, 500, 1000, 2500, and 5000 top words.

For the bigram approach, we again relied on only the most frequent 1000 words. However, for each tweet, instead of measuring the frequency of each individual word, we measure the frequency of each pair of words from the top 1000 words. Overall this results in 1 million possible bigrams (1000^2) to consider. To reduce computation time, we again restricted the set to only the top 10000 and 20000 bigrams, rather than consider all 1 million. Therefore, every tweet is represented as a vector of length 10000 (or 20000) that can be fed to the classifier.

Within each approach, we considered the same tweet dataset in two versions, one including the Emoji data in the tweets, and the other with the Emoji data cleaned out. Our goal is to compare both versions of the dataset with both the bag of words and bigram approach on both the Naïve Bayes and SVM classifier. For the purposes of testing and training, we have varied the testing/training split to 30/70, 50/50 and 70/30. All trials have been run on the same machine under the same conditions to provide consistency between the various trials.

5 Results and Discussion

For each of the experimental cases, we measured the classification accuracy of the system using the given vector size by running each trial ten times and averaging the accuracy to remove variation.

30/70 Split	1%	2%	5%	10%	25%	50%
BOW NE MNB	0.5289	0.55164	0.57752	0.59634	0.61518	0.62516
BOW R MNB	0.52976	0.54962	0.57786	0.59818	0.62116	0.63028
BOW NE SVM	0.54104	0.56486	0.59174	0.5964	0.58244	0.57886
BOW R SVM	0.53994	0.56496	0.59192	0.5992	0.58826	0.58492
BG NE MNB	0.56446	0.573	NA	NA	NA	NA
BG R MNB	0.56414	0.5708	NA	NA	NA	NA

Table 2: Experimental measurements of trials (30/70 split)

In **Table 1** all the raw values of the trials are displayed. The trials are labeled to describe their design (BOW = bag of words, BG = bigram, NE = no emoji, R = emoji, MNB = Multinomial Naïve Bayes, SVM = Support Vector Machine), and the size columns denote the size of the representation vectors (for bag of words it is the % of 10000, bigram the % of one million). The higher percentages for bigram have not been computed as the computation time became too excessive.

For the sizes that have been calculated for bigram, the bigram representation outperforms the bag of words representation. This is expected as the bigram retains contextual information. However, the accuracy increase from a vocabulary size of 100 to 200 words is smaller for bigram than bag of words, and it seems possible that the bigram would fall below the bag of words representation rather than retaining greater accuracy rates. This is not entirely surprising, as the short length of tweets means many bigrams are not present, resulting in very sparse vectors. These sparse vectors in turn are very similar to each other and can be difficult to classify.

The results for the bag of words representation are shown in **Fig 10**. Comparing the two classifiers, SVM outperforms MNB until a vocabulary of 1000, at which point MNB becomes more accurate. Both classifiers however follow a similar linear increase in accuracy up until a vocabulary size of 500 words. At this point, SVM levels off and proceeds to become less accurate, while MNB continues becoming more accurate until beginning to plateau at a vocabulary size of 2500. This is possibly due to the difficulty for SVM in finding the hyperplane at such a high dimension, while the MNB algorithm is not as greatly affected by an increase in feature descriptor size.

Within each classifier, we must also compare the effect of the Emoji data on the accuracy. Neither classifier shows much difference in accuracy until a vocabulary size of 1000. At this point, trials including the Emoji data begin to become more accurate than their counterparts without the included Emoji data. While this is what we hoped for, the margin is very small, with possibly approximately a 0.5% difference in accuracy or less. We believe this is due to the Emoji data largely agreeing with the included text, and thus not providing much more insight than the text itself.

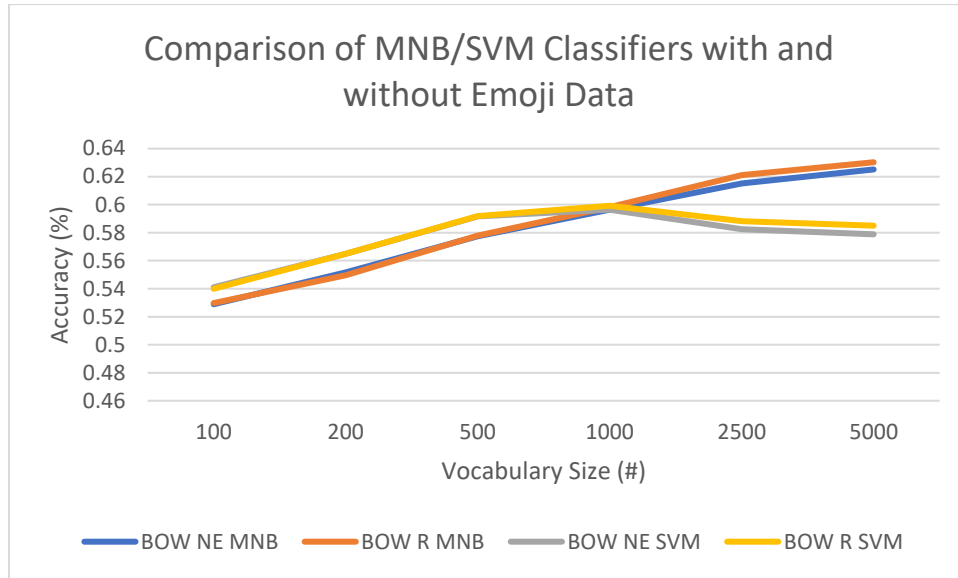


Figure 10: Graph of classification accuracy with varying vocabulary sizes

The results are similar when varying the training/testing split sizes. Though the absolute accuracy rates generally increase, the general accuracy pattern is the same. The tables of the 50/50 split and 70/30 split are each shown below.

50/50 Split	100	200	500	1000	2500	5000
BOW NE MNB	0.52923	0.55543	0.57864	0.60159	0.62349	0.63074
BOW R MNB	0.52999	0.55357	0.58357	0.60499	0.63125	0.63785
BOW NE SVM	0.54126	0.56861	0.59669	0.61577	0.60822	0.59339
BOW R SVM	0.53959	0.56544	0.59847	0.61523	0.60879	0.59505

Table 3: Experimental Measurements of Trials (50/50 Split)

70/30 Split	100	200	500	1000	2500	5000
BOW NE MNB	0.53311	0.55142	0.58325	0.60872	0.62802	0.63583
BOW R MNB	0.53173	0.55733	0.58459	0.6106	0.63246	0.64275
BOW NE SVM	0.54026	0.56657	0.60181	0.62309	0.62533	0.61084
BOW R SVM	0.54211	0.56582	0.60237	0.62553	0.62786	0.60609

Table 4: Experimental Measurements of Trial (70/30 Split)

Each of the classifiers displays approximately a 1% accuracy increase in the 70/30 split over the performance of the 30/70 split, while retaining the same general shape. Again, the best performing classifier is the MNB classifier with Emoji data included. To better compare across the training and testing sizes, we selected this classifier to compare across the varying sizes. This is shown in the graph below.

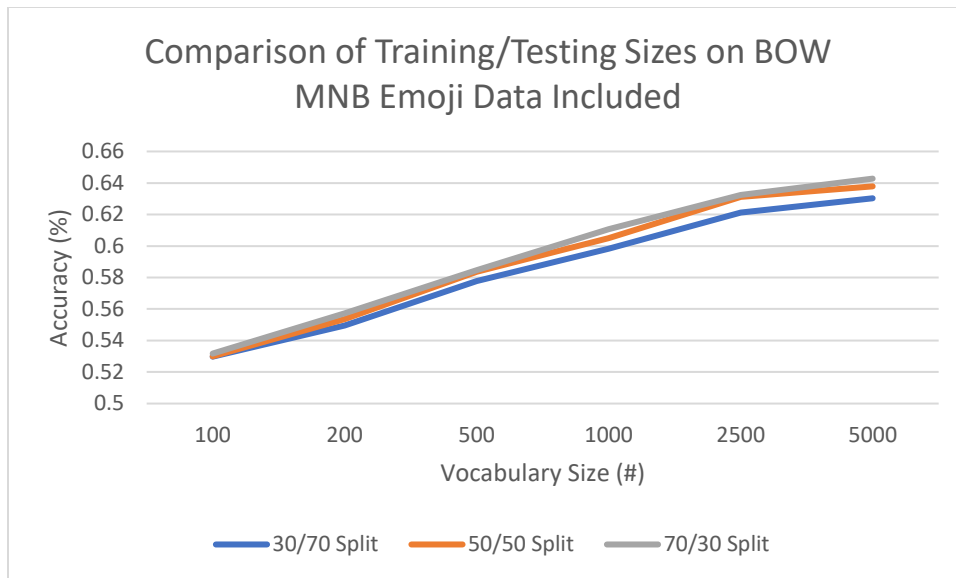


Figure 11: Classification Accuracy of BOW R MNB Over Varying Training/Testing Sizes

There is an accuracy increase by increasing the training data set size. However, the increase magnitude is relatively low, and changing the size does not greatly affect the shape of the graphs whatsoever. They each become more accurate as the vocabulary size is increased, and the training size change has greater affect at higher vocabulary sizes. Due to the continued lack of improvement with the Emoji data in various trials, we created a confusion matrix for the best classifier recorded (BOW MNB emoji data included, 30/70 split with 5000 vocabulary size). This matrix maps the actual classes of tweets to the classes predicted by the classifier. The diagonal of the matrix displays the correct classification rates, while all non-diagonal entries display various misclassification rates.

	SAD	ANGRY	HAPPY	SCARED	THANKFUL	SURPRISED	LOVE
SAD	0.865566038	0.033461	0.038473	0.030808	0.016951651	0.00117925	0.013561
ANGRY	0.605189991	0.272475	0.039852	0.047266	0.02502317	0.00092678	0.009268
HAPPY	0.246793454	0.025652	0.555506	0.029633	0.090667846	0.00044228	0.051305
SCARED	0.60244898	0.025306	0.067755	0.267755	0.024489796	0.00081633	0.011429
THANKFUL	0.302494802	0.012474	0.168399	0.022869	0.467775468	0	0.025988
SURPRISED	0.810810811	0.054054	0.067568	0.027027	0.027027027	0	0.013514
LOVE	0.294653015	0.017065	0.242321	0.028441	0.083048919	0	0.334471

Table 5: Confusion Matrix

The diagonal elements (correct classifications) have been bolded in the confusion matrix, while any occurrence of misclassification greater than correct classification has been highlighted red. Most of the correct classification rates are moderately successful – *sad*, *happy* and *thankful* are very commonly classified correctly. *Sad* dominates the classification accuracy due to its massive size relative to the other classes. *Love* is commonly misclassified as either *sad* or *happy*, though it is still correctly classified as the most common outcome.

The large misclassification rates are prevalent in the remaining classes – *angry*, *scared*, and *surprised*. These are misclassified as *sad* more often than they are classified correctly. The worst case is *surprised*, being misclassified as *sad* over 80% of the time. This is likely due to the relatively small size of *surprised* to the large size of *sad*, which also explains why the other classes are almost never misclassified as *surprised*. *Angry* and *scared* perform better, though they are still misclassified more often as *sad* than their correct classes.

Overall it seems that the classifier struggles to classify correctly due to the relative sizes of the classes. This also affects the accuracy rates of the Emoji data, which are not as prevalent in our experiments due to issues in the data sampling. Since some of the emotion classes show up much more frequently than others (*happy* and *sad* are much more prevalent than *thankful*) which skews the classification process towards predicting *sad* more than any other class. Additionally, our data was collected throughout the 2016 United States presidential election, a social phenomenon that resulted in many individuals taking to Twitter to express very positive and negative emotions regarding the event. We believe this additionally skewed the dataset, making prediction more difficult. For instance, candidate Donald Trump has a mannerism of calling things “sad,” which led to many people tweeting with the hashtag “#sad” while discussing things not related to the emotion, leading to much confusion.

6 Conclusion

Our experiments have shown the effects of including Emoji data in text for sentiment analysis. Although minor, there is an improvement in classification accuracy for including the Emoji data without any substantial computational overhead. We additionally compared both MNB and SVM classifiers on the data, with both bag of words and bigram representations. SVM outperforms MNB at small vocabulary sizes, though this is opposite at large vocabulary sizes. Bag of words also seems to outperform the more advanced bigram approach while saving computation time, as the bigram vectors are very sparse due to the short tweets. Overall this suggests that one should use a bag of words representation with an MNB classifier for large vocabulary sizes including the Emoji data to achieve maximum classification accuracy

7 Future Work

We hope to perform more experimentation on varying the training/testing ratio to produce learning curves for both the MNB and SVM implementations. Additionally, we want to experiment with other possible representation formats apart from bag of words and bigram. Lastly, we would advise ensuring more evenly distributed data for training and testing if possible, as this can skew the classification process and lead to high misclassification rates.

8 Bibliography

- [1] Agarwal, Apoorv, et al. "Sentiment analysis of twitter data." *Proceedings of the workshop on languages in social media*. Association for Computational Linguistics, 2011.
- [2] Greenwood, S., Perrin, A., Duggan, M, "Social Media Update 2016," *PEW Research Center*, 11 Sept 2016. Web. Accessed 5 April 2017.
<http://www.pewinternet.org/2016/11/11/social-media-update-2016/>
- [3] Janssens, O. et al. "Real-time Emotion Classification of Tweets," *IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, 2013.
- [4] Liu, B. "Sentiment Analysis and Opinion Mining," *Synthesis Lectures on Human Language Technology*," May 2012.
- [5] Liu, Bing. "Sentiment analysis and subjectivity." *Handbook of Natural Language Processing, Second Edition*. Chapman and Hall/CRC, 2010. 627-666.
- [6] Nasukawa, T. and Yi, J. "Sentiment analysis: capturing favorability using natural language processing," *K-CAP Proceedings on the 2nd International Conference on Knowledge Capture*, 23 Oct 2003.
- [7] Pak, A. and Paroubek, P. "Twitter as a Corpus for Sentiment Analysis and Opinion Mining," *LREc*. Vol. 10. No. 2010. 2010.
- [8] Pang, B. and Lee, L. "Opinion Mining and Sentiment Analysis," *Foundations and Trends in Information Retrieval*, 7 July 2008.
- [9] Python Twitter Package. Accessed 5 April 2017. <https://github.com/bear/python-twitter>
- [10] Wang, W. et al. "Harnessing Twitter 'Big Data' for Automatic Emotion Identification," *ASE/IEEE International Conference on Social Computing*, 2012.
- [11] Wilson, T., Weibe, J. and Hoffman, P. "Recognizing contextual polarity in phrase-level sentiment analysis", *HLT Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing*, 6 Oct 2005.
- [12] Zhang, C. et al. "Sentiment analysis of Chinese documents: From sentence to document level," *Journal of the American Society for Information Science and Technology*. 2 Sept 2009.