

The construction of the list was contextual, as some words that appear frequently may still be useful for sentiment analysis. To get a bet-

ter picture of these common, but useful words, three more word clouds were constructed over the word lists per sentiment (Figures 3, 4, and 5).

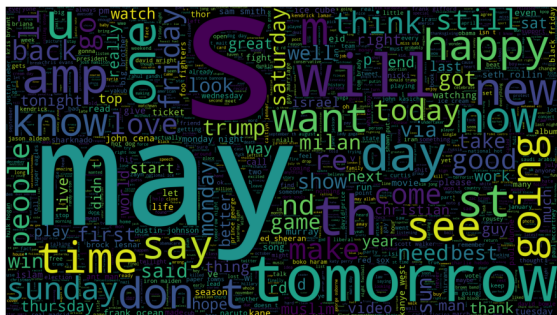


Figure 3: Word cloud over positive training tweets

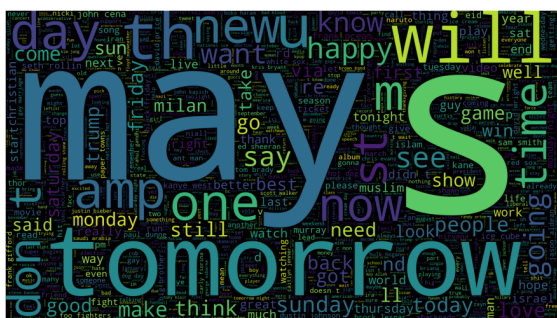


Figure 4: Word cloud over neutral training tweets

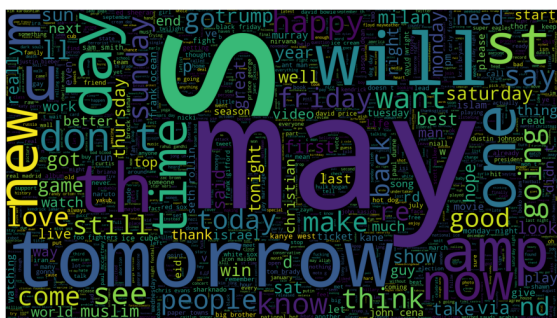


Figure 5: Word cloud over negative training tweets

However this method did not consider stopwords in other languages (since english tweets are an overwhelming majority in the dataset). Since manual stopword list construction was not as exhaustive as the data required, the Python NLTK module's stopwords corpus is used where necessary (Bird et al., 2009).

2.3 Vectorization

Classifiers in the **SciKit-learn** Python library require that all data be vectorized to a real-valued space. The following vectorizers were used.

2.3.1 Count

Vectorising the tweet into term counts can highlight terms which appear more often in tweets for certain sentiments. This is implemented using **SciKit-learn's** **CountVectorizer** (Pedregosa et al., 2011).

2.3.2 TF-IDF

This measure of relative word frequency provides more insight, as it measures words that approach more often in one tweet relative to their overall frequency. This is implemented using **SciKit-learn's** **TfidfVectorizer** (Pedregosa et al., 2011).

2.3.3 Metrics

The final type of vectorization that will occur is in the form of metrics. Certain metrics such as word length may be distributed differently based on the sentiment of the tweet. This is implemented by creating a list of mappings to metrics, then vectorizing with **SciKit-learn's** **DictVectorizer** (Pedregosa et al., 2011).

2.3.4 Why not hashing?

The **SciKit-learn** library suggests another text feature extractor in the form of hashing. This is not used here as the distinct advantage this vectorizer has over others is saving on space and time (Pedregosa et al., 2011). While the dataset contains more than 20000 tweets, using the other vectorizers did not present such issues in practice.

2.4 Features

While the data is given as a raw text format, there are multiple features which can be extracted for the purpose of sentiment analysis.

2.4.1 N-grams

This includes extraction of individual words or characters (1-grams of each) and word pairings (2-grams). If $N > 1$, important orderings/configurations of words can be identified, but at the cost of exponentially increasing the possible number of features. With a vocabulary of w unique words, there can be as many as w^N distinct N-grams. Therefore, the N-grams used for model construction are:

- 1-grams of words,

- 1-grams of characters,
- 2-grams of words.

This method will requires that tweets are cleaned, but may not need stopwords removal.

2.4.2 Stems

Can find the roots for english words by removing stems. It can also be used as a cleaning step. However, this may be less effective with words that aren't in english (as they may use different suffixes), or words which are misspelled.

2.4.3 Lemmas

This method finds the true roots of words. This suffers from the same drawback as stemming, relying the text being in english.

2.4.4 Word Lengths

A tokenization where the tokens generated are the lengths of the words in the tweet.

2.4.5 Character Frequencies

2.4.6 Links

2.4.7 Hashtags

2.4.8 Mentions

2.4.9 Emoticons

2.4.10 Simple Metrics

2.4.11 Phonetic Frequencies

2.4.12 Poetic Phonetics

2.5 Model Selection

2.5.1 Feature Selection

2.5.2 Classifiers

2.5.3 Evaluation Metrics

3 Results

4 Analysis

5 Conclusions

References

- Steven Bird, Edward Loper, and Ewan Klein. 2009. *Natural Language Processing with Python*. O'Reilly Media Inc.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- N. F. Rosenthal, Sara and Preslav Nakov. 2017. Semeval-2017 task 4: sentiment analysis in twitter. In *11th International Workshop on semantic evaluation (SemEval '17)*, Vancouver, Canada.