

# Capstone 1: Milestone Report

Rika Yatchak

March 26, 2019

## 0 Introduction

In this project, I examine the Twitter US Airline Sentiment for my Springboard Capstone project. To view the code used to generate the results in this report, see the corresponding Git repository.

In this report, we will

1. Define a guiding question for the project.
2. Identify hypothetical clients.
3. Describe the dataset and the cleaning/wrangling steps needed to prepare for analysis.
4. List other potential datasets that could be used.
5. Explain the findings.

## 1 Project Question

Can we perform a sentiment analysis on customer tweets directed at selected U.S. Airlines?

## 2 Hypothetical Client

The hypothetical client is a an airline which would like to support its social media managers in quickly identifying tweets that should be prioritized for immediate response. They would like to identify upset customers who may need additional support using an automated system that flags tweets that are negative and directs customer service representatives to review them for possible issue resolution.

## 3 Dataset Description and Preparation

The Twitter US Airline Sentiment was uploaded to Kaggle in 2016 by Figure Eight. It consists of two files: There are two different files:

- `Tweets.csv` – Tweet text with tweet ID, date, timezone information, as well as labeled sentiments provided by Crowdfunder.
- `database.sqlite` – Same contents as `Tweets.csv` but in a sqlite database.

The dataset appears to have been collected sometime in 2015. Unfortunately, the methodology for data collection and postprocessing was not published with the dataset. From my EDA, it seems clear that the dataset was collected by looking for tweets that @ mention the official accounts of various U.S. airlines. There are some columns that were added after data collection: the columns `airline_sentiment` and `negative_reason` appear to be the output of Crowdfunder's own sentiment analysis and topic categorization.

As this dataset contains tweet text, significant postprocessing was needed: emojis were identified and removed, various regex expressions were utilized to remove undesired parts of speech or text snippets such as URLs and personal email addresses, and hashtags were also post-processed. The data cleaning notebook has more details about the steps that were undertaken to prepare text for sentiment analysis.

## 4 Other Datasets

Since the main bulk of the data is output from the Twitter API, there is always the possibility to obtain more relevant tweets from Twitter itself. To obtain more data, we would need to pull tweets that @ mention Virgin America, United, US Airways, JetBlue, and Southwest Air. Since the dataset was collected, Virgin America has ceased operation. We could also consider adding airlines such as Delta, Frontier, and Alaska Airlines, which are not included in the current dataset.

The main reason that more tweets were not pulled for this project is the dramatic increase in data wrangling and data cleaning that this would require: the tweets made available by Crowdfunder in this dataset appear to have been selected especially for sentiment analysis and have been partially cleaned.

## 5 Findings

For the final sentiment analysis, I tried three types of input features:

1. The first type of input features are TF-IDF features output from scikit-learn's `TfidfVectorizer`. This function runs a `CountVectorizer` to return a matrix of token counts, and then performs a TF-IDF transform on this matrix. This transform takes into account both the term frequency in the individual document, as well as the document frequency of the token. In this way, words that are found in most tweets such as "the" are penalized as they are common words that do not add much meaning.
2. The second type of input features are generated using Word2Vec. First, we use Word2Vec to represent our large corpus in a smaller-dimensional space. Each word in our corpus is a vector in this space, and words that are related to one another are closer to each other in the vector space. Once I generated the word vectors, I used two separate types of embeddings: the sum of words embedding, in which a tweet is represented by the sum of the word vectors corresponding to the words in the tweet, and the mean embedding, in which we take the mean of the word vectors corresponding to the words in the tweet.
3. The third type of input features are paragraph vectors generated by Doc2Vec, which is a small modification of Word2Vec in which we add a unique document ID to the Word2Vec model. Each vector output is a representation of a specific tweet in our dataset.

For all three types of features, I used SMOTE to oversample the underrepresented "neutral" and "positive" sentiment classes in the data as the sentiment of the tweets in the dataset was heavily skewed towards the negative.

Finally, I evaluated multinomial logistic regression and random forest classifier methods for predicting the sentiment of the tweets. I used a 70/30 train/test split and evaluated the models by using the test accuracy. Logistic regression performed rather well out of the box with a simple TF-IDF input. After tuning the random forest classifier with the best-performing sum of words embedding (SOWE) Word2Vec features, I obtained the best-case test accuracy of .754. Both classifiers performed similarly (after tuning) for all three classes of features. The code for the sentiment analysis can be viewed in this notebook.