# The Influence of Sentiment Polarity on News Popularity

## Kate Chen & Christopher Wolff
Duke University

ISS110
Patrick Herron
January 16, 2018

### Abstract

We conducted a quantitative study to illustrate the relationship between the sentiment polarity of a news article and its popularity. A collection of articles from the New York Times was aggregated and labeled to train a multinomial naive Bayes classifier and a support vector machine. The results were compared with an existing model and showed an improvement in classification accuracy from 58% to 71%. We also show that the output of a specific classifier sees an increase in standard deviation for sparse input vectors. Additionally, we created a Python framework that provides automated API querying, web scraping, data analysis, and visualization tools for future research. In contrast to our initial hypothesis, we found evidence that there may be a positive correlation between positive sentiment polarity and popularity, through a variety of third factors must also be taken into consideration.

## Introduction

In Orson Welles' classic film *Citizen Kane*, Charles Kane revives a dying newspaper and subsequently builds a news empire through yellow journalism often solely focused on scandal or other negative events. While Welles' depiction of corruption and fraudulent information in the news industry is exaggerated, his inspiration is not mislead - the saying "bad news sells" often propels news sources and journalists to focus disproportionately negatively on current events rather than maintain a neutral approach. Psychologists suspect that readers' attraction to negative news is the cause of an innate negativity bias, but past studies have primarily determined this phenomenon through psychological reasoning rather than quantitative evidence (Stafford, 2014). Therefore, the goal of our research was to create a software application that not only determines a sentiment score for an article based solely on its text, but also examine the general correlation of an article's popularity with its sentiment. Our hypothesis was that there will be a positive correlation between the negativity of an article and its popularity. In other words, we hypothesize to see an average trend of the most popular news being the most negative, and subsequent articles increasing in average positivity as the article popularity decreases. Using our computed sentiment scores and other data, we aimed to make a conclusion based on quantitative evidence.

## Background

### Previous Work

Sentiment analysis has been applied in a large range of domains. For example, it has been used to gauge the negativity or positivity of news coverage focused on a specific person or topic (Godbole, Srinivasaiah, & Skiena, 2007). The majority of this study focused on the effectiveness of the sentiment analysis model through analyzing coverage sentiment in news and blogs. In addition, news and blog sentiments over time about specific topics were compared, and top positive and top negative subjects were recorded. Studies have also been conducted to analyze the correlation between news coverage of businesses and their performances in financial stock markets (Atreya, Cohen, & Zhai, 2011; Kalyani, Bharathi, & Jyothi, 2016). This sentiment model was able to predict with 70 percent accuracy movements in the stock market based on preceding news and their sentiments. In addition, sentiment scores were confirmed through comparisons to manual human sentiment scoring. However, while many studies involving sentiment analysis exist, research regarding the validity of the claim that negativity in news from accredited news organizations results in higher view counts is not prominent, if existing. Nonetheless, many processes in sentiment analysis are similar regardless of the research objective, and therefore, these studies served as useful resources for our approach.

### Word Vector Representations of Documents

One of the key challenges in building a sentiment analyzer is generating a numerical representation of any given article, also referred to as a vector space representation (Salton, Wong, & Yang, 1975). The two most common representations are the bag of words model and the term frequency - inverse document frequency (tf-idf) model. Both represent the article as a fixed-size vector, which is necessary so that articles with different lengths can be used equivalently for computation. The bag of words model first tokenizes the document string, meaning it splits all of the words contained in it, and stores them in a list. It then counts how often each of the words occur and converts these frequencies into a vector $\vec{x}$. The vector has a size corresponding to the total number of words in the article collection. The words are always arranged in the same order, so that the $i^{th}$ element $x_i$ always corresponds to the frequency of the $i^{th}$ word. Since the capitalization of a word usually does not change its semantics, all words are often converted to lower case beforehand. Optionally, the resulting frequencies are normalized by the document length so that longer documents do not result in a word vector with a larger overall magnitude.

One problem with the bag of words model is that common words such as "the" or "a" tend to occur much more often in the English language, and will therefore carry a large weight in the word vectors. However, these words should have little to no effect on the overall sentiment score. The tf-idf model takes the overall occurrence of each term in the entire document collection into account as well and assigns a larger weight for rarer terms. The tf-idf score for some term $t$ is defined as the product of the term frequency TF(t) and the inverse document frequency IDF(t):

$$\text{TFIDF}(t) = \text{TF}(t) \cdot \text{IDF}(t) \qquad (1)$$

with

$$\text{IDF}(t) = \ln\left(\frac{n_d}{1 + \text{DF}(t)}\right) \qquad (2)$$

where $n_d$ denotes the total number of documents and DF(t) is the number of documents containing the term $t$. The 1 is added to the denominator so that it will never be zero. As a result of the inverse document frequency factor, terms will be weighted according to their "rarity" within the document collection in addition to their raw frequency. In some models the resulting tf-idf scores are normalized, typically by the Euclidean norm:

$$\vec{x}_{norm} = \frac{\vec{x}}{\|\vec{x}\|_2} = \frac{\vec{x}}{\sqrt{\sum x_i^2}} \tag{3}$$

**Naive Bayes Classification**

A naive Bayes sentiment classifier can be thought of as a function that takes a feature vector, such as one of the two described above, as an input, and assigns probabilities to each possible outcome. In our case, the outcomes were chosen to be binary labels that correspond to the articles sentiment – positive and negative. The model is based on Bayes' theorem (Berkson, 1930), which gives the relationship

$$P(y|x_1,...,x_n) = \frac{P(y)P(x_1,...,x_n|y)}{P(x_1,...,x_n)} \tag{4}$$

$P(y|x_1,...x_n)$ denotes the conditional probability of y given the set of features $x_1,...,x_n$; in other words, the likelihood of y occurring given that $x_1,...,x_n$ are true. In our case, y might be the probability that a given article has negative sentiment and the input probabilities $x_i$ are true if the $i^{th}$ word in the feature vector is present in the article. This term is also called the posterior probability, because it represents the desired output of the system. $P(y)$ is the prior probability that a given class occurs. For instance, if 60% of all articles have a negative sentiment, $P(y)$ is 0.6, regardless of the article. Similarly, $P(x)$ is the prior probability that a given term $x$ occurs in any given document. The term $P(x_1,...,x_n|y)$ is the likelihood that every term in the article occurs in the input vector given the output class $y$. At first, it might seem that this term is difficult to compute, and that converting the challenge of finding the posterior probability to finding its "inverse" has not made it significantly easier. However, that is where the main idea behind the naive Bayes model comes in:

**Assumption 1.** *The conditional probabilities $P(x_i|y)$ are independent of each other.*

That is also the origin of the word "naive" in naive Bayes, because the assumption is usually just false, especially in news articles. For example, the probabilities of the terms "white house" and "politics" occurring in an article are clearly dependent on each other. Nonetheless, the independence assumption proves useful for computational purposes, because it lets us compute $P(x_1,...,x_n|y)$ as

$$P(x_1,...,x_n|y) = \prod P(x_i|y) \tag{5}$$

It simply turns into the product of all individual conditional probabilities, which are easy to compute. Lastly, a multinomial naive Bayes classifier is a special version of naive Bayes designed to work well with tf-idf weighting.

**Support Vector Machines**

A support vector machine (SVM) is another model that can be used for classification. An SVM creates a plane of separation between the two output classes in the n-dimensional feature

space of the input vector. If we imagine that the input vector only has two features, which would correspond to having only two distinct words in the document collection, then a SVM would separate the positive and negative labels using a line.
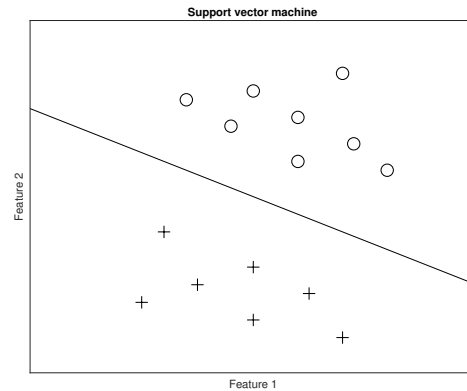


*Figure 1*. Two-dimensional SVM example (created with MATLAB$^{TM}$)

Here, each point represents an article with two features. The numerical value of each feature corresponds to each articles weight in its word vector, so either its frequency count or tf-idf weight. The plus signs and circles correspond to positive and negative sentiment labels, respectively. Once this line has been "learned," any new article can easily be classified as positive if it lies below the line and negative if it lies above the line. However, an input word vector usually has about 10,000 different features, so the data can not simply be separated by a one-dimensional object. Instead, an SVM creates an (n-1)-dimensional hyperplane through the n-dimensional space that the word vectors are in. For instance, if there were three total features, the separator would be a plane.

## Procedure

The procedure of our study included two main tasks. The first one was collecting an archive of news articles along with relevant information: title, abstract, full story text, and view count. We predicted that the title and abstract would most likely have the largest impact on an article's popularity, because they are directly exposed to the reader. The second component was determining a sentiment value for each of the articles. We experimented with the Python `TextBlob` library to assign sentiment values. However, for our final results, we trained our own sentiment analysis model using hand-labeled training data.

### Defining Sentiment and Popularity

Before gathering any quantifiable measures of sentiment and popularity, we needed to define these two attributes first. The term sentiment can span a multitude of dimensions, such as opinion, emotions, and affect (Munezero, Montero, Sutinen, & Pajunen, 2014). However, we were mostly interested in sentiment polarity - whether a given article is positive, negative, or neutral. We defined a sentiment score attribute of an article, which is a continuous number ranging from -1 to +1, where scores of -1, 0, and +1 correspond to extremely negative polarity, neutral polarity, and extremely positive polarity, respectively. In our own classifier, we used discrete binary labels because that

choice synergized best with our mathematical models. We defined popularity as the number of views a given article received over its lifespan. While other measures, such as the number of times an article was shared via online media or the total revenue it generated, were also viable measures, using the view count was convenient because it was directly available in the generated dataset. One downside to this form of measuring popularity is that older articles will likely have a higher view count. However, we reasoned that articles are generally featured for a relatively constant time period online, after which they are less likely to be read and view count stagnates. Hence, the extra views that arise from the additional time beyond the time period in which the article was featured would not significantly contribute to the total.

**Data Collection**

To approach collecting articles for the data, multiple factors were taken into consideration. In addition to the title and abstract of the article, popularity and time of publication needed to be taken into consideration. In the initial research phase, we investigated the ”trending” pages of news sources such as Facebook or Mashable. However, upon consideration, the New York Times offered an open source API that provided the necessary features we were looking for. In addition, the New York Times is an accredited news source with the second largest distribution in the United States and was therefore appropriate to be chosen as the data source for this project.

The New York Times Most Popular API offers three measures of popularity: most emailed, most shared and most viewed. Because the research question is focused on the popularity of an article in the context of how many people are attracted to read it, the view count was selected as the feature of choice. The API requested two inputs: a String representing the article section, and an integer denoting a time period during which the article was live. The section refers to the category within the newspaper, such as Education, Fashion, Food, Health, U.S, or World. The option to query all sections was also available. The time period values were limited to 1, 7, and 30, corresponding to the past day, week, and month of news. Surprisingly, when querying a time period value of 30, a few of the articles returned had a `published_date` older than one month. We assumed that these articles were originally published earlier, but were still available on the New York Times website. Additionally, the API allowed an offset parameter k, which specified that the k most viewed articles should be ignored and only articles starting from k+1 should be returned.

Using the inputs given by the user, the API returned a JSON file containing 20 articles. For each article, the API returns its URL, column, section, byline, title, abstract, published date, and source. The articles were listed in order of popularity. Therefore, the first article returned in the array would be the most popular article in the specified time period, defined by view count. Additional documentation about the New York Times Most Popular API can be found at New York Times Developer Website (Times, n.d.). Within the returned JSON file and features of each article, only title, and abstract were given, but the full document text was missing. We therefore used the web scraping tool `BeautifulSoup4` (Nair, 2014) to automatically extract the story content of each article from its given URL by parsing the websites HTML source code. In the end, a total of 1908 articles from all sections were collected from a period of 30 days, ranging from October 25, 2017 to November 24, 2017 and taking up 14.2 MB of memory storage space.

**Sentiment Analysis**

The sentiment analysis of the retrieved articles was initially done using the Python `TextBlob` library. For each article, the title, abstract, and full text were separately fed through the `TextBlob` sentiment analyzer. Then, each article was plotted against its relative popularity rank using `matplotlib` (Hunter, 2007). We used `numpy` (van der Walt, Colbert, & Varoquaux, 2011) to calculate the minimum, maximum, mean, and standard deviation for each of the three datasets.

Next, we trained our own sentiment classifier by hand-labeling over 600 news articles with binary sentiment labels: positive or negative, using a user interface (UI) we developed. The UI showed us each article's title and abstract and enabled us to enter a corresponding label, relabel previously labeled articles, or assign random labels to all articles for testing purposes. When labeling the articles, topic of the article in addition to any noticeable tone in the title or abstract were primarily taken into consideration; however, some cases were difficult to label as purely negative or purely positive. In these scenarios, guidelines for how these articles would be handled were established so that these articles would be judged equally and fairly. For example, an article about a positive aspect of a negative topic, such as a public figure releasing an apology after sexual harassment accusations, was judged as a negative article because of the core issue at hand, which would generally evoke negative emotion. In addition, because articles could not be labeled as neutral, or have a sentiment score of 0, seemingly neutral news was labeled as positive. An example of this would be an article including photographs and captions of Trump's activities in Asia, with no subjective information. Furthermore, political viewpoints were not taken into consideration - while the New York Times is generally viewed as more liberal, an article involving conversative topics did not automatically result in a negative sentiment score.

The articles were first preprocessed by filtering out stopwords and converting all terms to lowercase. We separately used both the bag of words representation as well as the tf-idf weights for the full story text of each article as input vectors. The support vector classifier and the multinomial naive Bayes classifier from the `scikit-learn` machine learning library (Pedregosa et al., 2011) were both utilized. Prior to training, the labeled articles were randomly split into 75% used for training and 25% used for testing. The testing data aided as a tool for validating our models performance to see how well they perform on unseen data. We recorded the resulting accuracy and stored the trained models using Python's object serialization tool `pickle`. Then, we compared the performance of our own models with the `TextBlob` baseline classifier. Lastly, we repeated the initial experiment of plotting sentiment against popularity but with our own classifiers.

### Results and Discussion

The plots for the title, abstract, and story sentiment obtained from the `TextBlob` classifier can be seen in Figure 2.
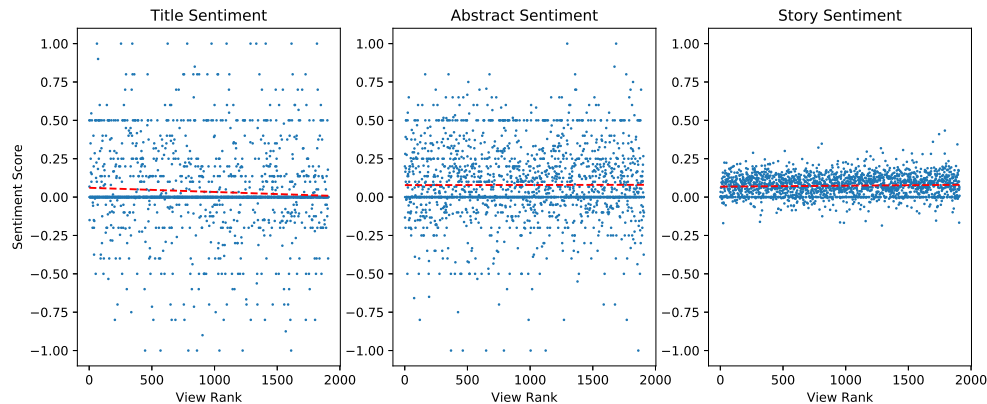


*Figure 2*. `TextBlob` sentiment score results (created with `matplotlib`)

The dotted red lines represent a linear fit for each of the datasets, created using `numpy`. From the flat slope of these trend lines, it is clear that there is no obvious pattern or correlation. Table 1 shows the general properties of each dataset.

|      | Title | Abstract | Story |
|------|-------|----------|-------|
| min  | -1    | -1       | -.19  |
| max  | 1     | 1        | .71   |
| mean | .14   | .22      | .23   |
| std  | .30   | .30      | .19   |

Table 1
*Minimum, Maximum, Mean, and Standard Deviation of TextBlob Sentiment Scores*

We noticed that both the range and standard deviation were much lower when the whole story was used as an input vector. This lead us to believe that the outputs for the `TextBlob` classifier tend to regress towards a mean score as the number of input features gets bigger. Our prediction is that `TextBlob` uses a bag of words model and only takes the existence of each term into account without regard for its frequency, so that its word vector representation is simply a multi-hot vector over all possible terms. This would imply that sparse input vectors have a much higher similarity in comparison.

Upon further investigation, we realized that `TextBlob` is based on the `PatternAnalyzer` from the `nltk` library (Loper & Bird, 2002), which was trained using movie reviews. The fact that the model was trained on a completely different type of text data lead us to believe that it might be mislabeling the news articles. An example of this is the phrase "Saudi Arabia Arrests 11 Princes, Including Billionaire Alwaleed bin Talal" with the abstract "The sweeping campaign of arrests appears to be the latest move to consolidate the power of Crown Prince Mohammed bin Salman, the favorite son and top adviser of King Salman," which received positive sentiment scores of 0.997 and 0.988, respectively. While this specific news is certainly negative, the sentiment analyzer return

an extremely high score. We believe that this occurs because words like favorite, prince, power, or billionaire would return high sentiment scores, potentially outweighing the more negative words, such as arrest.

Labels like the one described compromised the data, which warranted and necessitated a new way to obtain the sentiment data. Further research showed that many other sentiment analyzers were also trained on movie reviews or otherwise inadequate for the purposes of this study. Therefore, in order to achieve the most ideal results for the specific needs of the sentiment analysis in this study, we created our own model. The following snippet shows our final results for combinations of different word vector representations (bag of words and tf-idf) with different classifiers (naive Bayes and SVM).

```
>>> train_model()
Found 622 labeled articles
221 +, 401 -
Vectorizing using bag of words...
Naive Bayes: 67.31% accuracy
SVM: 64.10% accuracy
Vectorizing using tfidf...
Naive Bayes: 69.87% accuracy
SVM: 71.15% accuracy
```

The `TextBlob` classifier was used on the full set of labeled articles. Since the returned labels from `TextBlob` were non-binary, we rounded all positive scores to +1 and all negative scores to -1. For scores of 0, the training example was disregarded and not counted towards the overall accuracy. The `TextBlob` classifier achieved a total accuracy of **57.76%**, which is significantly less than our model.

The plots below show the relationship between view count and sentiment score based on the new sentiment scores computed with our models.
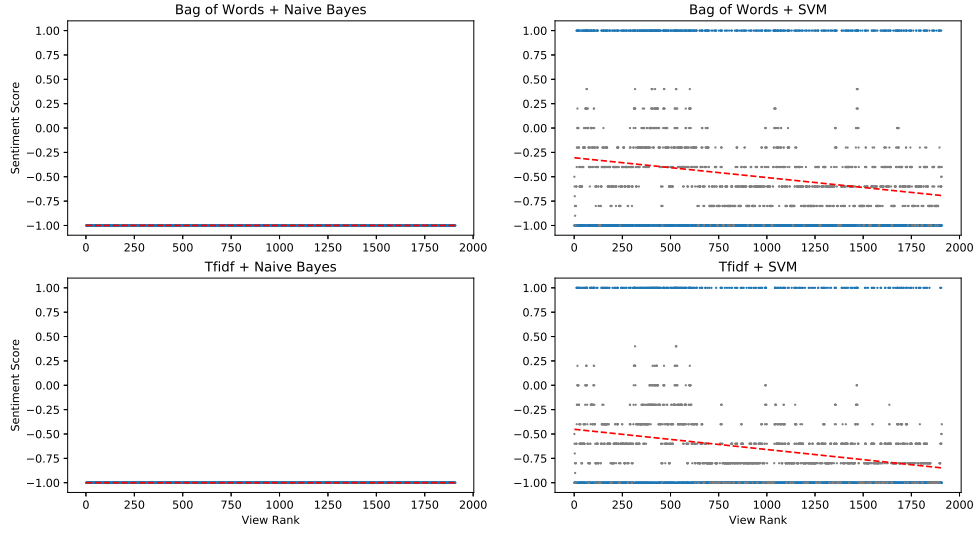
*Figure 3.* Our sentiment score results (created with `matplotlib`)

The blue data points are exact predictions made by the classifiers. The red line is a linear fit for each graph, representing the general trend of the data. The orange scatter plot is a ten-point moving-average of the data. This makes it easier to visualize the progression of the sentiment values over time. An unexpected observation is that the naive Bayes classifier guessed a negative sentiment for all data points. We predict that this relates to the fact that a majority of labeled articles (401) were negative, compared to the 221 positive articles. When computing the class probability $P(\text{negative}|t_1,...,t_n)$, the prior $P(y)$ dominates all other terms. When the final guess it computed using argmax{P(negative), P(positive)}, the positive probability will never be high enough, regardless of the word vectorizer chosen. The trend line for the SVM data shows a clear negative trend, contrary to our initial hypothesis. As the popularity of the articles decreases, the negativity increases. However, the hypothesis was supported in regards to the negativity of news articles, in that a clear majority of the articles returned and ranked were negative as opposed to positive. Therefore, the data supports that negative news is indeed popularly read and more so than positive news.

There are multiple third factors that may have influenced this data. For example, the conclusion may be influenced by a simple possibility that, in general, there is more negative news available to be read than positive news. It may also be possible that in general, negative news is more commonly relevant to readers' lives, as they may cause feelings of fear or urgency that may make them more important to read than positive news. Negative news may also lead to more complex, controversial, or nuanced issues that require the reader to click on the article and view it to further understand the article rather than being able to understand the jist from reading the title, and not opening the article to read it. Regardless, the data presents a conclusion that indicates a more complex relationship between sentiment and article popularity beyond the simple assumption that negative news is more popular.

## Conclusions

Though the results of this study did not strictly support nor refute the hypothesis and further research would be necessary to confirm our results or otherwise discover other correlations, the surprising outcome of the data raises an interesting question into the assumption that negative news is more popular than positive news. While the data does seem to support that humans naturally gravitate more towards negative news, the true correlation between sentiment and popularity does not seem to be as simple or clear. In addition, third variables such as the numerical prominence of negative news or cultural context in current events may heavily influence the data. However, the conclusion of this study does present a different perspective and possibly contradict a long-held common belief in the news industry. Therefore, further investigation may uncover facts and trends that could bring insight into the true impact of negativity in news popularty or otherwise call industry changes to action.

Due to the time and resource limitations of this project, there are a number of ways that this research can be expanded upon, confirmed, or disputed. Further investigation with larger data pools, more detailed sentiment analysis and potential introduction of other factors such as changes over time can help researchers form a more nuanced and accurate hypothesis. In addition, as discovered in the initial sentiment analysis model, the standard deviation of the `TextBlob` sentiment scores was lower when the entire story was analyzed for sentiment. Therefore, a potential extension to this study would be to conduct the sentiment analysis on the entire articles rather than just the titles or abstracts. However, the methods and models used in this study can provide a solid foundation upon which further research can be conducted to reach a stable and replicatable conclusion regarding the correlation between sentiment and news popularity.

References

Atreya, A., Cohen, N., & Zhai, J. (2011). Sentiment analysis of news articles for financial signal prediction. doi: https://nlp.stanford.edu/courses/cs224n/2011/reports/nccohen-aatreya-jameszjj.pdf

Berkson, J. (1930, 02). Bayes' theorem. *Ann. Math. Statist.*, *1*(1), 42–56. Retrieved from `https://doi.org/10.1214/aoms/1177733259` doi: 10.1214/aoms/1177733259

Godbole, N., Srinivasaiah, M., & Skiena, S. (2007). Large-scale sentiment analysis for news and blogs. *International Conference on Weblogs and Social Media*. doi: http://icwsm.org/papers/3–Godbole-Srinivasaiah-Skiena.pdf

Hunter, J. D. (2007). Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, *9*(3), 90-95. Retrieved from `http://aip.scitation.org/doi/abs/10.1109/MCSE.2007.55` doi: 10.1109/MCSE.2007.55

Kalyani, J., Bharathi, H., & Jyothi, R. (2016). Stock trend prediction using news sentiment analysis. doi: https://arxiv.org/abs/1607.01958

Loper, E., & Bird, S. (2002). Nltk: The natural language toolkit. In *In proceedings of the acl workshop on effective tools and methodologies for teaching natural language processing and computational linguistics. philadelphia: Association for computational linguistics.*

Munezero, M. D., Montero, C. S., Sutinen, E., & Pajunen, J. (2014, April). Are they different? affect, feeling, emotion, sentiment, and opinion detection in text. *IEEE Transactions on Affective Computing*, *5*(2), 101-111. doi: 10.1109/TAFFC.2014.2317187

Nair, V. G. (2014). *Getting started with beautiful soup*. Packt Publishing.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., . . . Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, *12*, 2825–2830.

Salton, G., Wong, A., & Yang, C. S. (1975, November). A vector space model for automatic indexing. *Commun. ACM*, *18*(11), 613–620. Retrieved from `http://doi.acm.org/10.1145/361219.361220` doi: 10.1145/361219.361220

Stafford, T. (2014). Psychology: Why bad news dominates the headlines. *BBC - Future*. doi: http://www.bbc.com/future/story/20140728-why-is-all-the-news-bad

Times, N. Y. (n.d.). The new york times developer network. doi: https://www.developer.nytimes.com

van der Walt, S., Colbert, S. C., & Varoquaux, G. (2011). The numpy array: A structure for efficient numerical computation. *Computing in Science & Engineering*, *13*(2), 22-30. Retrieved from `http://aip.scitation.org/doi/abs/10.1109/MCSE.2011.37` doi: 10.1109/MCSE.2011.37

Appendix
Source Code
This appendix contains the main source code used for this project. Additional files as well as a detailed README that explains how to run this project can be found at https://github.com/christopher-wolff/news-analysis.

**__main__.py**

```
1    """Sample use of analyzer module."""
2
3    from analyzer import query
4    from analyzer import scrape_stories
5    from analyzer import label_articles
6    from analyzer import train_model
7    from analyzer import analyze
8    from analyzer import visualize
9
10   if __name__ == '__main__':
11       print('Requesting articles from NYT API')
12       print('===============================')
13       query(num_queries=1)
14       print()
15
16       print('Scraping full article texts from NYT website')
17       print('============================================')
18       scrape_stories()
19       print()
20
21       print('Labeling articles')
22       print('=================')
23       label_articles(reset=True, rand_labels=True)
24       print()
25
26       print('Training classifiers')
27       print('====================')
28       train_model()
29       print()
30
31       print('Analyzing data')
32       print('====================')
33       analyze()
34       print()
35
36       print('Visualize results')
37       print('=================')
38       visualize()
```

**analyzer.py**

```
1    """Main module."""
2
3
4    __authors__ = 'Christopher Wolff, Kate Chen'
5    __version__ = '1.0'
6    __date__ = '9/10/2017'
7
8
9    import json
10   import os.path
11   import pickle
12   import random
13   import urllib
14
15   from bs4 import BeautifulSoup
16   from nltk.corpus import stopwords
17   from sklearn.feature_extraction.text import CountVectorizer
18   from sklearn.feature_extraction.text import TfidfVectorizer
19   from sklearn.model_selection import train_test_split
20   from sklearn import naive_bayes
21   from sklearn import svm
22   from sklearn.metrics import accuracy_score
23   from textblob import TextBlob
24   import matplotlib.pyplot as plt
25   import requests
26   import numpy as np
27
28
29   SETTINGS_PATH = 'settings.json'
30   RAW_PATH = 'data/raw.json'
31   STORIES_PATH = 'data/with_stories.json'
32   LABELS_PATH = 'data/with_labels.json'
33   SENTIMENTS_PATH = 'data/with_sentiments.json'
34   MNB_PATH = 'models/mnb.pkl'
35   SVM_PATH = 'models/svm.pkl'
36   COUNT_VECT_PATH = 'models/count_vect.pkl'
37   TFIDF_VECT_PATH = 'models/tfidf_vect.pkl'
38
39   BASE_URI = 'http://api.nytimes.com/svc/mostpopular/v2'
40   TYPE = 'mostviewed'
41   SECTION = 'all-sections'
42   TIME_PERIOD = '1'
43   RESPONSE_FORMAT = 'json'
```

```
44
45
46   def query(num_queries=1):
47       """Request data from NYT and store it as a json file.
48
49       Args:
50           num_queries (int): The number of queries
51       """
52       # Load API key
53       settings = json.load(open(SETTINGS_PATH))
54       API_KEY = settings['API_KEY']
55
56       # Send requests
57       URI = f'{BASE_URI}/{TYPE}/{SECTION}/{TIME_PERIOD}.{RESPONSE_FORMAT}'
58       articles = []
59       for k in range(num_queries):
60           print(f'Running query {k+1}...')
61           offset = k * 20
62           payload = {'api_key': API_KEY, 'offset': offset}
63           response = requests.get(URI, params=payload)
64           articles += response.json()['results']
65
66       # Save to file
67       with open(RAW_PATH, 'w') as output_file:
68           json.dump(articles, output_file)
69
70
71   def scrape_stories():
72       """Get full document texts from urls."""
73       # Load articles
74       articles = json.load(open(RAW_PATH))
75
76       # Submit GET request and parse response content
77       for k, article in enumerate(articles):
78           print(f'Scraping article {k+1}...')
79           url = article['url']
80           f = urllib.request.urlopen(url)
81           soup = BeautifulSoup(f, 'html5lib')
82           story = ''
83           for par in soup.find_all('p', class_='story-body-text \
84                                               story-content'):
85               if par.string:
86                   story += ' ' + par.string
87           article.update({'story': story})
88
```

```
89          # Save articles
90          with open(STORIES_PATH, 'w') as output_file:
91              json.dump(articles, output_file)
92
93
94      def label_articles(reset=False, relabel=False, start=0, rand_labels=False):
95          """Run UI for sentiment labeling.
96
97          Loads all articles and presents those without a label.
98
99          Args:
100             reset (boolean): Delete all labels
101             relabel (boolean): Allow option to override existing labels
102             start (int): Article number to start from
103             rand_labels (boolean): Assign all random labels
104         """
105         # Load articles
106         if reset or not os.path.isfile(LABELS_PATH):
107             articles = json.load(open(STORIES_PATH))
108         else:
109             articles = json.load(open(LABELS_PATH))
110         if start >= len(articles):
111             raise ValueError(f'Invalid starting point: {start}')
112
113         # Label articles
114         sentiments = [-1, 1]
115         print(f'Available sentiments: {sentiments}')
116         for k, article in enumerate(articles[start:]):
117             if not relabel and 'sentiment' in article:
118                 continue
119             print(f'Article: {k+start+1}')
120             print(f"Title: {article['title']}")
121             print(f"Abstract: {article['abstract']}")
122             if rand_labels:
123                 sent = random.choice(sentiments)
124             else:
125                 try:
126                     sent = int(input('Label: '))
127                 except ValueError:
128                     break
129                 if sent not in sentiments:
130                     break
131             article.update({'sentiment': sent})
132             print('---------------------------')
133
```

```
134        # Save articles
135        with open(LABELS_PATH, 'w') as output_file:
136            json.dump(articles, output_file)
137
138
139    def train_model(random_state=None):
140        """Train a sentiment analyzer model.
141
142        Args:
143            random_state (int): Random seed for train_test_split used by numpy
144        """
145        # Load articles
146        articles = json.load(open(LABELS_PATH))
147        # Extract data
148        articles = [article for article in articles if 'sentiment' in article]
149        stopset = set(stopwords.words('english'))
150        titles = [article['title'] for article in articles]
151        labels = [article['sentiment'] for article in articles]
152
153        # Vectorize data
154        count_vect = CountVectorizer(lowercase=True,
155                                     strip_accents='ascii',
156                                     stop_words=stopset,
157                                     decode_error='replace')
158        tfidf_vect = TfidfVectorizer(use_idf=True,
159                                     lowercase=True,
160                                     strip_accents='ascii',
161                                     stop_words=stopset,
162                                     decode_error='replace')
163
164        # Analyze and display relevant information
165        num_total = len(articles)
166        num_pos = sum(article['sentiment'] == 1 for article in articles)
167        num_neg = sum(article['sentiment'] == -1 for article in articles)
168        print(f'Found {num_total} labeled articles')
169        print(f'{num_pos} +, {num_neg} -')
170
171        # Train using count vectorizer
172        print('Vectorizing using bag of words...')
173        x = count_vect.fit_transform(titles)
174        y = labels
175        if random_state is not None:
176            x_train, x_test, y_train, y_test = train_test_split(
177                    x, y, random_state=random_state)
178        else:
```

```
179             x_train, x_test, y_train, y_test = train_test_split(x, y)
180
181         mnb_clf = naive_bayes.MultinomialNB()
182         mnb_clf.fit(x_train, y_train)
183         y_pred = mnb_clf.predict(x_test)
184         mnb_acc = accuracy_score(y_test, y_pred) * 100
185         print('Naive Bayes: %.2f%% accuracy' % mnb_acc)
186
187         svm_clf = svm.SVC(probability=True)
188         svm_clf.fit(x_train, y_train)
189         y_pred = svm_clf.predict(x_test)
190         svm_acc = accuracy_score(y_test, y_pred) * 100
191         print('SVM: %.2f%% accuracy' % svm_acc)
192
193         # Train using tfidf vectorizer
194         print('Vectorizing using tfidf...')
195         x = tfidf_vect.fit_transform(titles)
196         y = labels
197         if random_state is not None:
198             x_train, x_test, y_train, y_test = train_test_split(
199                     x, y, random_state=random_state)
200         else:
201             x_train, x_test, y_train, y_test = train_test_split(x, y)
202
203         mnb_clf = naive_bayes.MultinomialNB()
204         mnb_clf.fit(x_train, y_train)
205         y_pred = mnb_clf.predict(x_test)
206         mnb_acc = accuracy_score(y_test, y_pred) * 100
207         print('Naive Bayes: %.2f%% accuracy' % mnb_acc)
208
209         svm_clf = svm.SVC(probability=True)
210         svm_clf.fit(x_train, y_train)
211         y_pred = svm_clf.predict(x_test)
212         svm_acc = accuracy_score(y_test, y_pred) * 100
213         print('SVM: %.2f%% accuracy' % svm_acc)
214
215         # Store vectorizers and trained classifiers
216         with open(SVM_PATH, 'wb') as output_file:
217             pickle.dump(mnb_clf, output_file)
218         with open(MNB_PATH, 'wb') as output_file:
219             pickle.dump(svm_clf, output_file)
220         with open(COUNT_VECT_PATH, 'wb') as output_file:
221             pickle.dump(count_vect.vocabulary_, output_file)
222         with open(TFIDF_VECT_PATH, 'wb') as output_file:
223             pickle.dump(tfidf_vect.vocabulary_, output_file)
```

```
224
225
226   def analyze():
227       """Analyze article data."""
228       # Calculate sentiment scores
229       articles = json.load(open(LABELS_PATH))
230       mnb_clf = pickle.load(open(MNB_PATH, 'rb'))
231       svm_clf = pickle.load(open(SVM_PATH, 'rb'))
232       count_vocabulary = pickle.load(open(COUNT_VECT_PATH, 'rb'))
233       tfidf_vocabulary = pickle.load(open(TFIDF_VECT_PATH, 'rb'))
234       stopset = set(stopwords.words('english'))
235       count_vect = CountVectorizer(lowercase=True,
236                                    strip_accents='ascii',
237                                    stop_words=stopset,
238                                    decode_error='replace',
239                                    vocabulary=count_vocabulary)
240       tfidf_vect = TfidfVectorizer(use_idf=True,
241                                    lowercase=True,
242                                    strip_accents='ascii',
243                                    stop_words=stopset,
244                                    decode_error='replace',
245                                    vocabulary=tfidf_vocabulary)
246       for k, article in enumerate(articles):
247           title = article['title']
248           abstract = article['abstract']
249           story = article['story']
250           print(f'{k+1}: {title}')
251           title_sent = TextBlob(title).sentiment
252           abstract_sent = TextBlob(abstract).sentiment
253           story_sent = TextBlob(story).sentiment
254           article.update({'title_sent': title_sent,
255                           'abstract_sent': abstract_sent,
256                           'story_sent': story_sent})
257           print(f'{title_sent} {abstract_sent} {story_sent}')
258
259           count = count_vect.fit_transform([title])
260           tfidf = tfidf_vect.fit_transform([title])
261           article.update({'count_mnb_sent': mnb_clf.predict(count).item(0),
262                           'count_svm_sent': svm_clf.predict(count).item(0),
263                           'tfidf_mnb_sent': mnb_clf.predict(tfidf).item(0),
264                           'tfidf_svm_sent': svm_clf.predict(tfidf).item(0)})
265
266       # Test TextBlob performance
267       num_total = 0
268       num_correct = 0
```

```
269        for article in articles:
270            if 'sentiment' not in article:
271                continue
272            title_sent = article['title_sent'].polarity
273            true_sent = article['sentiment']
274            if title_sent == 0:
275                continue
276            if _sign(title_sent) == true_sent:
277                num_correct += 1
278            num_total += 1
279        acc = num_correct / num_total * 100
280        print('=========================')
281        print('TextBlob accuracy: %.2f' % acc)
282        print('=========================')
283
284        # Determine min, max, mean, and std
285        title_sents = np.array([a['title_sent'] for a in articles])
286        abstract_sents = np.array([a['abstract_sent'] for a in articles])
287        story_sents = np.array([a['story_sent'] for a in articles])
288
289        print('Title Sentiments')
290        print('----------------')
291        print(f'min: {np.min(title_sents)}')
292        print(f'max: {np.max(title_sents)}')
293        print(f'mean: {np.mean(title_sents)}')
294        print(f'std: {np.std(title_sents)}')
295        print()
296
297        print('Abstract Sentiments')
298        print('-------------------')
299        print(f'min: {np.min(abstract_sents)}')
300        print(f'max: {np.max(abstract_sents)}')
301        print(f'mean: {np.mean(abstract_sents)}')
302        print(f'std: {np.std(abstract_sents)}')
303        print()
304
305        print('Story Sentiments')
306        print('----------------')
307        print(f'min: {np.min(story_sents)}')
308        print(f'max: {np.max(story_sents)}')
309        print(f'mean: {np.mean(story_sents)}')
310        print(f'std: {np.std(story_sents)}')
311        print()
312
313        # Save to file
```

```
314        with open(SENTIMENTS_PATH, 'w') as output_file:
315            json.dump(articles, output_file)
316
317
318    def visualize():
319        """Visualize the data."""
320        # Load data
321        articles = json.load(open(SENTIMENTS_PATH))
322        title_sents = [article['title_sent'][0] for article in articles]
323        abstract_sents = [article['abstract_sent'][0] for article in articles]
324        story_sents = [article['story_sent'][0] for article in articles]
325        count_mnb_sents = [article['count_mnb_sent'] for article in articles]
326        count_svm_sents = [article['count_svm_sent'] for article in articles]
327        tfidf_mnb_sents = [article['tfidf_mnb_sent'] for article in articles]
328        tfidf_svm_sents = [article['tfidf_svm_sent'] for article in articles]
329
330        view_rank = range(1, len(articles) + 1)
331
332        # Calculate trendlines
333        z1 = np.polyfit(view_rank, title_sents, 1)
334        p1 = np.poly1d(z1)
335        z2 = np.polyfit(view_rank, abstract_sents, 1)
336        p2 = np.poly1d(z2)
337        z3 = np.polyfit(view_rank, story_sents, 1)
338        p3 = np.poly1d(z3)
339
340        z4 = np.polyfit(view_rank, count_mnb_sents, 1)
341        p4 = np.poly1d(z4)
342        z5 = np.polyfit(view_rank, count_svm_sents, 1)
343        p5 = np.poly1d(z5)
344        z6 = np.polyfit(view_rank, tfidf_mnb_sents, 1)
345        p6 = np.poly1d(z6)
346        z7 = np.polyfit(view_rank, tfidf_svm_sents, 1)
347        p7 = np.poly1d(z7)
348
349        # Compute moving average
350        window_size = 10
351        window = np.ones(int(window_size))/float(window_size)
352        count_svm_sents_ma = np.convolve(count_svm_sents, window, 'same')
353        tfidf_svm_sents_ma = np.convolve(tfidf_svm_sents, window, 'same')
354
355        # Plot sentiment versus view rank
356        # TextBlob
357        plt.figure(1)
358        plt.subplot(1, 3, 1)
```

```
359        plt.scatter(view_rank, title_sents, s=5)
360        plt.plot(view_rank, p1(view_rank), 'r--')
361        plt.title('Title Sentiment')
362        plt.xlabel('View Rank')
363        plt.ylabel('Sentiment Score')
364        plt.ylim(-1.1, 1.1)
365
366        plt.subplot(1, 3, 2)
367        plt.scatter(view_rank, abstract_sents, s=5)
368        plt.plot(view_rank, p2(view_rank), 'r--')
369        plt.title('Abstract Sentiment')
370        plt.xlabel('View Rank')
371        plt.ylim(-1.1, 1.1)
372
373        plt.subplot(1, 3, 3)
374        plt.scatter(view_rank, story_sents, s=5)
375        plt.plot(view_rank, p3(view_rank), 'r--')
376        plt.title('Story Sentiment')
377        plt.xlabel('View Rank')
378        plt.ylim(-1.1, 1.1)
379
380        # sklearn classifiers
381        plt.figure(2)
382        plt.subplot(2, 2, 1)
383        plt.scatter(view_rank, count_mnb_sents, s=5)
384        plt.plot(view_rank, p4(view_rank), 'r--')
385        plt.title('Bag of Words + Naive Bayes')
386        plt.ylabel('Sentiment Score')
387        plt.ylim(-1.1, 1.1)
388
389        plt.subplot(2, 2, 2)
390        plt.scatter(view_rank, count_svm_sents, s=5)
391        plt.scatter(view_rank, count_svm_sents_ma, s=5, facecolor='0.5')
392        plt.plot(view_rank, p5(view_rank), 'r--')
393        plt.title('Bag of Words + SVM')
394        plt.ylim(-1.1, 1.1)
395
396        plt.subplot(2, 2, 3)
397        plt.scatter(view_rank, tfidf_mnb_sents, s=5)
398        plt.plot(view_rank, p6(view_rank), 'r--')
399        plt.title('Tfidf + Naive Bayes')
400        plt.xlabel('View Rank')
401        plt.ylabel('Sentiment Score')
402        plt.ylim(-1.1, 1.1)
403
```

```
404        plt.subplot(2, 2, 4)
405        plt.scatter(view_rank, tfidf_svm_sents, s=5)
406        plt.scatter(view_rank, tfidf_svm_sents_ma, s=5, facecolor='0.5')
407        plt.plot(view_rank, p7(view_rank), 'r--')
408        plt.title('Tfidf + SVM')
409        plt.xlabel('View Rank')
410        plt.ylim(-1.1, 1.1)
411
412        plt.show()
413
414
415    def _sign(x):
416        if x < 0:
417            return -1
418        elif x > 0:
419            return 1
420        else:
421            return 0
```