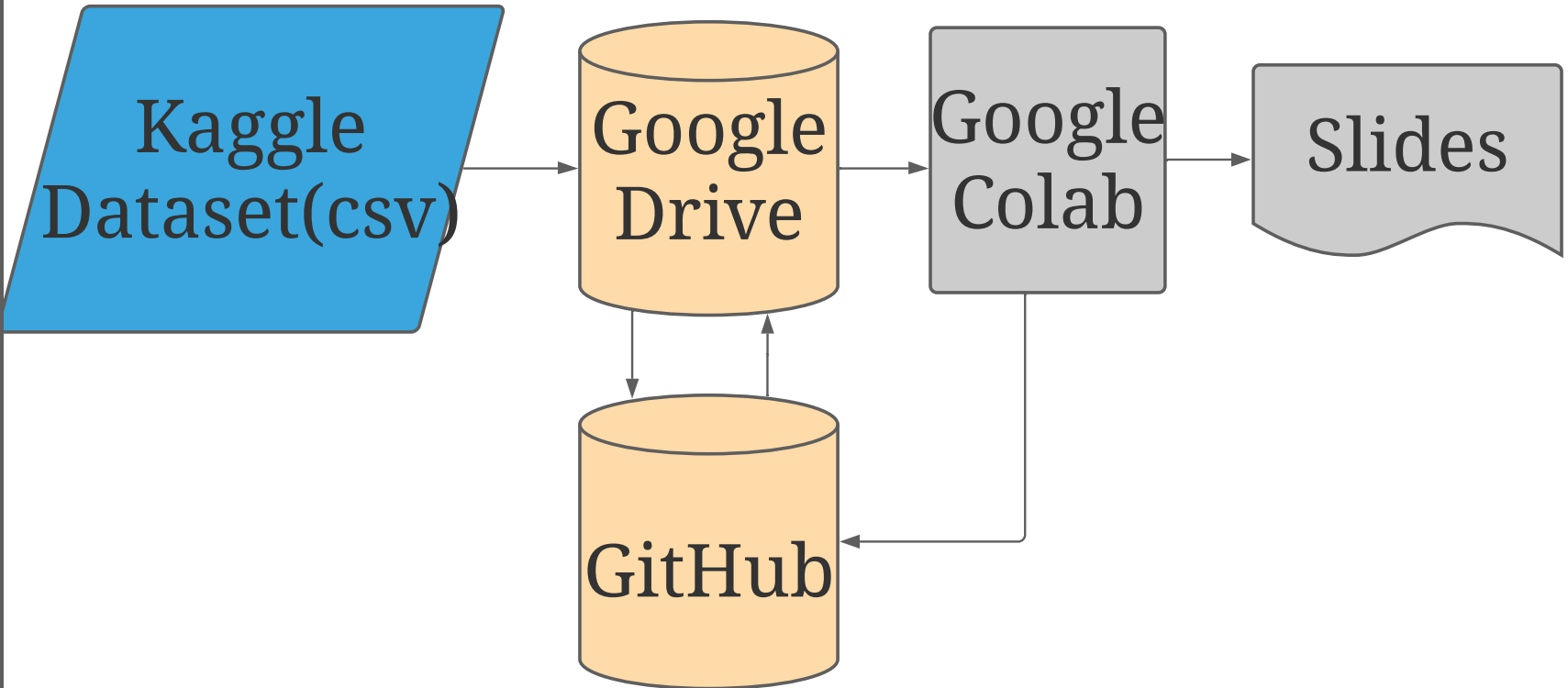


Research Questions:

- A) How has fake news changed?  
B) Can we predict whether something is fake news or not using NLP?

Workflow:



Introduction:

Fake news is not a recent or unique issue, although it has been exacerbated by the ease and speed in which information could be spread. Since the 1890s, fake news has plagued the profession of journalism. Dubbed "yellow journalism", newspaper publishers would run sensationistic news articles. However, fake news was eventually replaced by objective journalism, until recently due to the rise of the internet. (Source:UCSB)

- In terms of impacts fake news has on society, University of Michigan lists the following as possibilities:
- Anti-intellectualism
  - Antiscience
  - Widespread Mistrust

Works Cited:

- WordCloud Tutorial
- TextBlob Documentation
- Confusion Matrix tutorial
- Wikipedia for SVM
- Wikipedia for Logistic Regression

Real News WordCloud:



Fake News WordCloud:

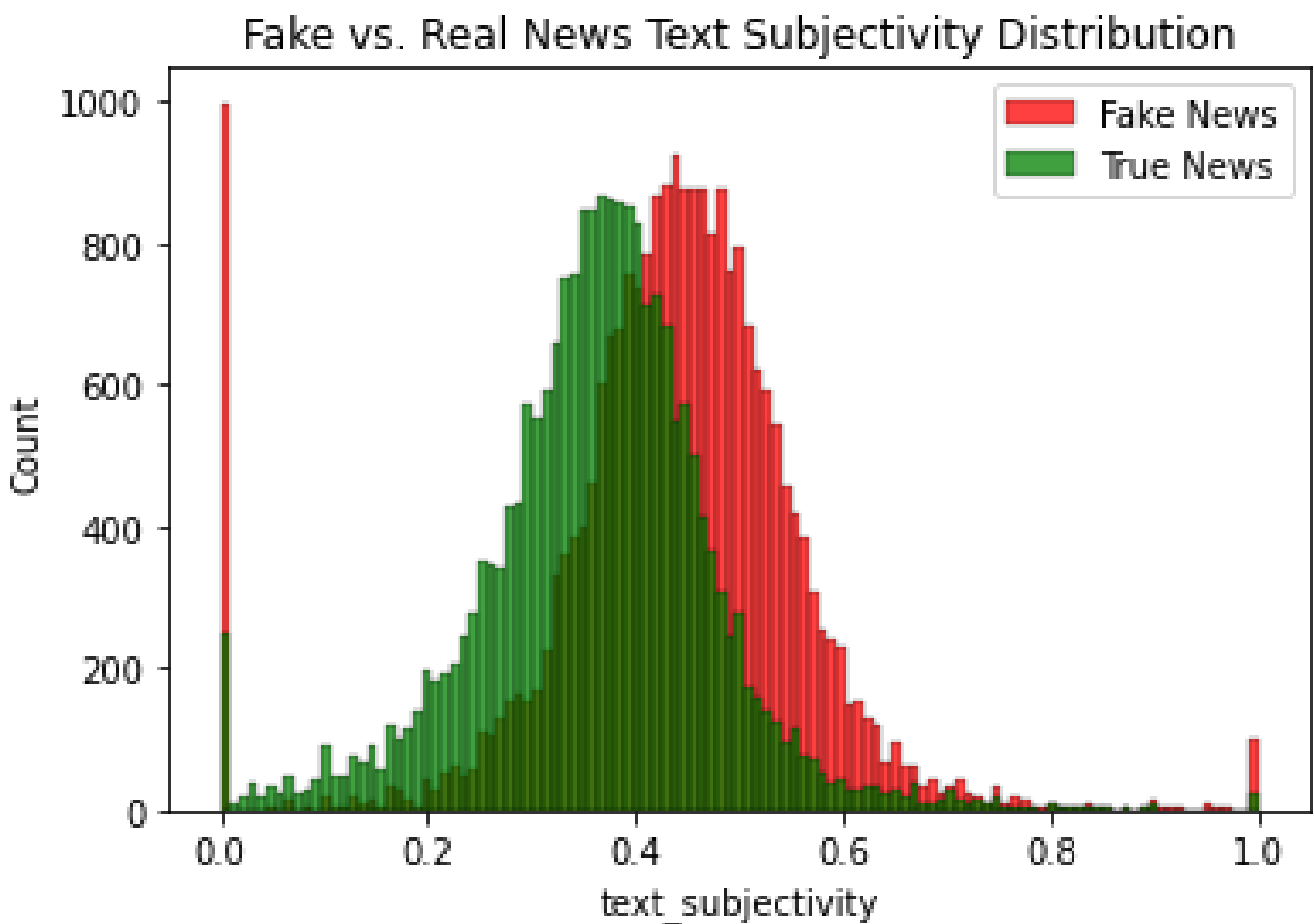
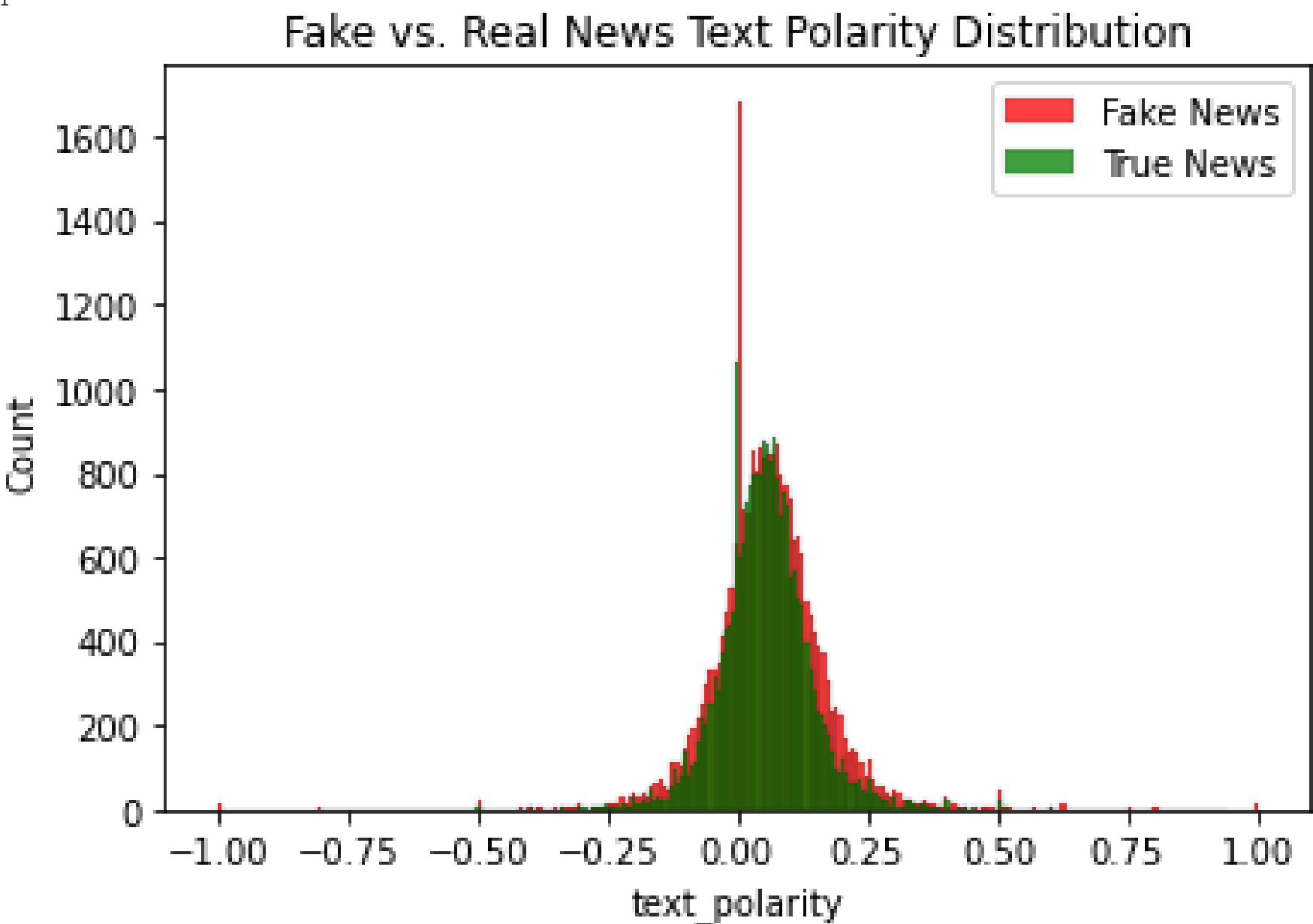


EDA: Continued on Page 2

- From wordcloud, can see differences in words used
- Fake news dataset primarily about Donald Trump
- Sources matter!
  - Reuters makes up a large part of the real news data
  - 21st Century Wire, a right wing conservative web blog makes the a large part of fake news data
- Difference in topics discussed:
  - Fake news tends to be nationalistic (america) and more focused on traditional issues (christmas, sheriff)
  - Real news tends to be more about international issues (nato, brussels, london) and more on progressive ideals (transgender)
- Interestingly, it seems real news is more interested in conservatives than fake news is.

Sentiment Analysis: Continued on Page 3

- Tracked polarity and subjectivity scores
- Distribution of polarity scores were not markedly different in true and fake news
- Fake news contains more subjective vocabulary than true news

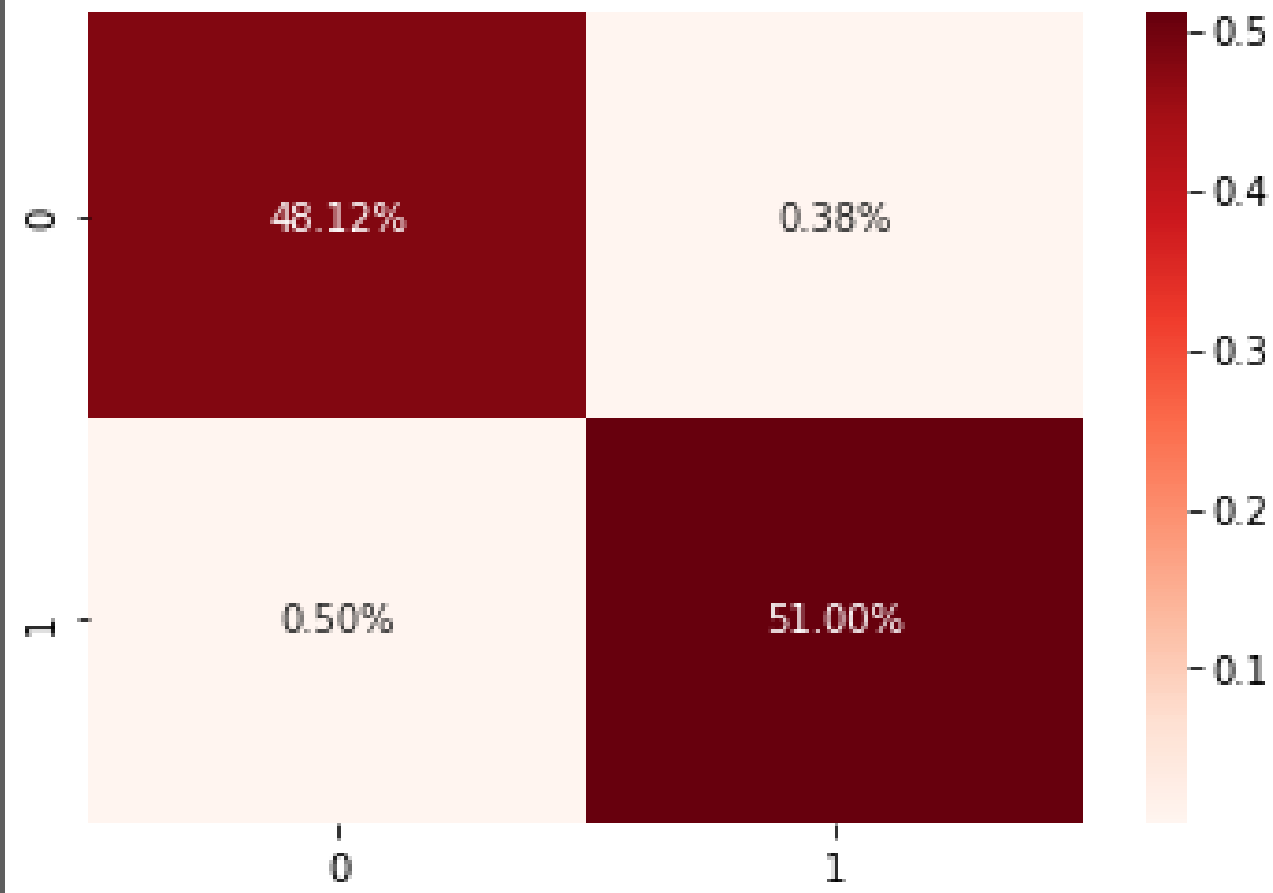


Predicting Fake News vs.

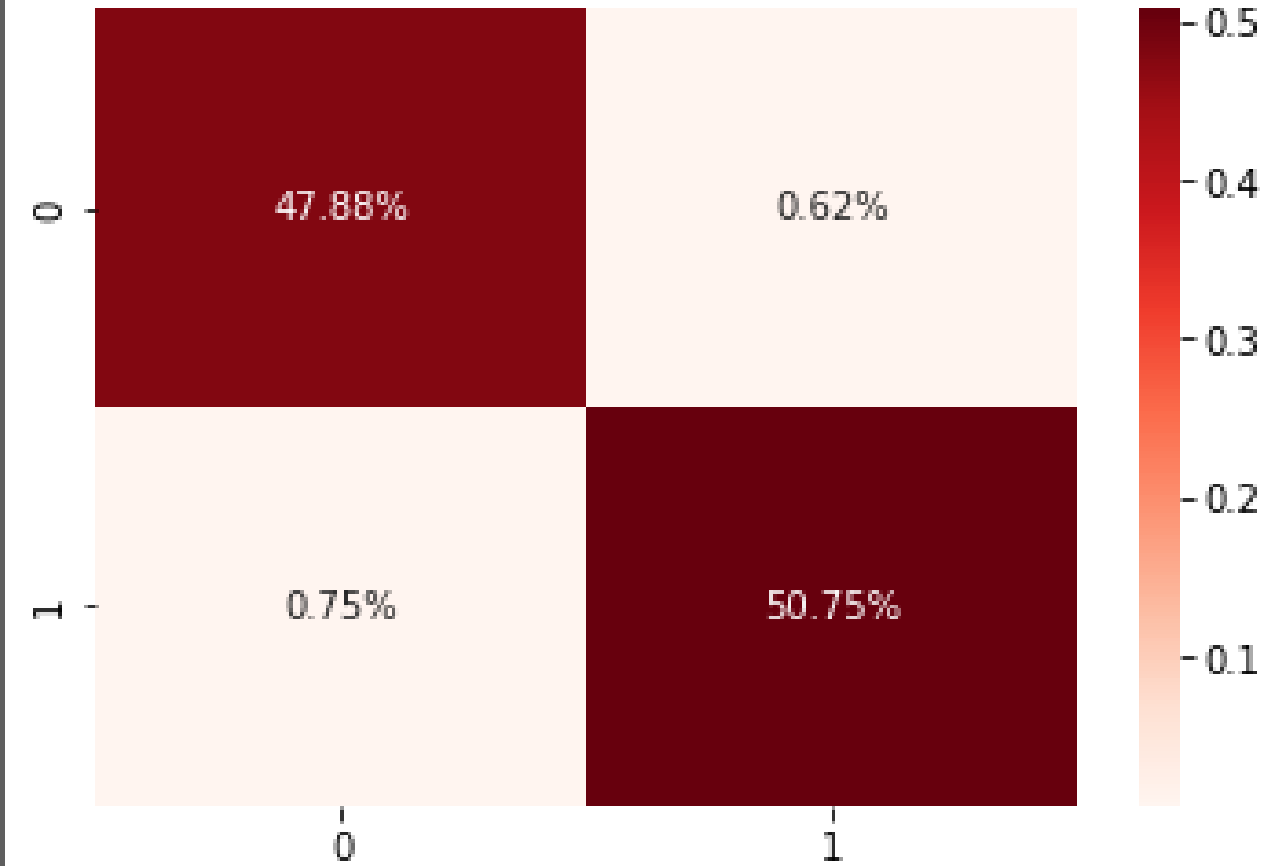
Real News: Continued on Page 4

- Fitted an SVM and Logistic Regression Model based on TF-IDF word frequencies
- Achieved very high accuracy:
  - Logistic Regression: **99.12%**
  - SVM: **98.63%**
- However, may be due to quirks in the data
  - Sources for fake and real news are very uniform
- 1 for Fake, 0 For Real

Logistic Regression Confusion Matrix:



SVM Confusion Matrix:









### What is Sentiment Analysis?

- Sentiment Analysis attempts to use computational methods to identify the tone of a text
- In our project, we are using the TextBlob package to determine the polarity and subjectivity of a sentence

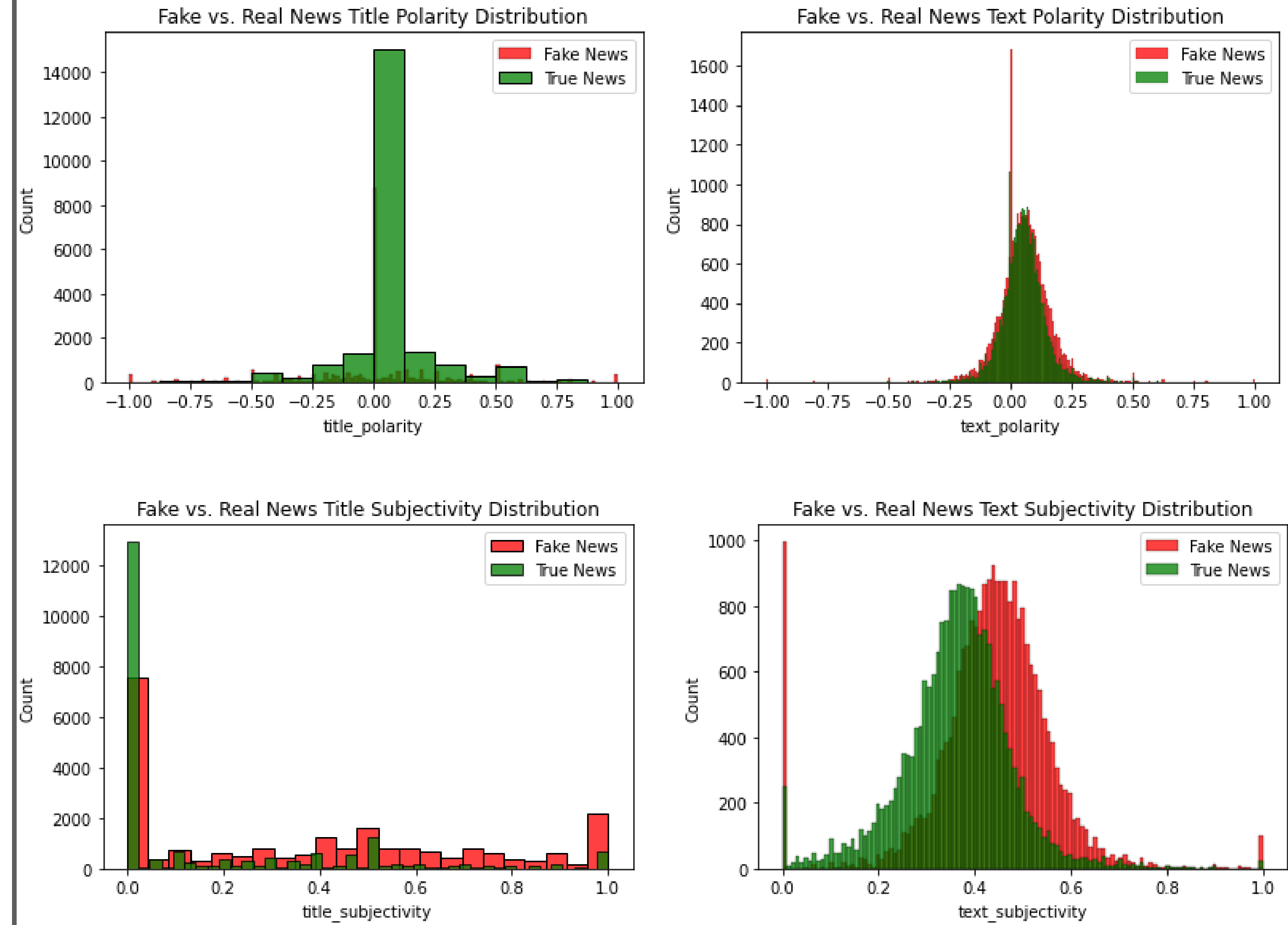
### What is Polarity and Subjectivity?

- **Polarity:** Polarity indicates how positive or negative a text is. In TextBlob, we are given a score from -1 to 1. A score of -1 indicates a sentence is very negative, while a score of +1 indicates a sentence is very positive
- **Subjectivity:** Subjectivity indicates how subjective a text is. Similarly, we are given a score from -1 to 1. A score of -1 indiciates a less subjective text, while a score of +1 indicates a very subjective text.
- TextBlob contains a dictionary of words and a score associated with each word. These words are matched to the sentence and then averaged to return a score for the entire sentence

### Sentiment Analysis Process:

- I wanted to look at the distribution of polarity and subjectivity of the following:
  - Polarity of the title for fake and real news
  - Subjectivity of the title for fake and real news
  - Polarity of the text for fake and real news
  - Subjectivity of the text for fake and real news
- **Why does this matter?**
  - Doing so allows individuals to potentially use sentiment analysis as a way to augment our classification process
  - Reveals insights into the nature of fake and real news
  - Evaluate the suitability of TextBlob for sentiment analysis of fake and real news

### Results of Sentiment Analysis



### Conclusion of Sentiment Analysis:

- The distribution of the polarity and subjectivity of the titles are not significantly different
- This may be due to the shortness of the titles, which may not allow many matches between the words in the title and the words in the TextBlob dictionary
- Interestingly, the polarity of fake and real news are very similar.
- However, significant difference between the subjectivity of fake and real news
- On average, fake news is more subjective than real news.
- Standard Deviations of both are similar as well



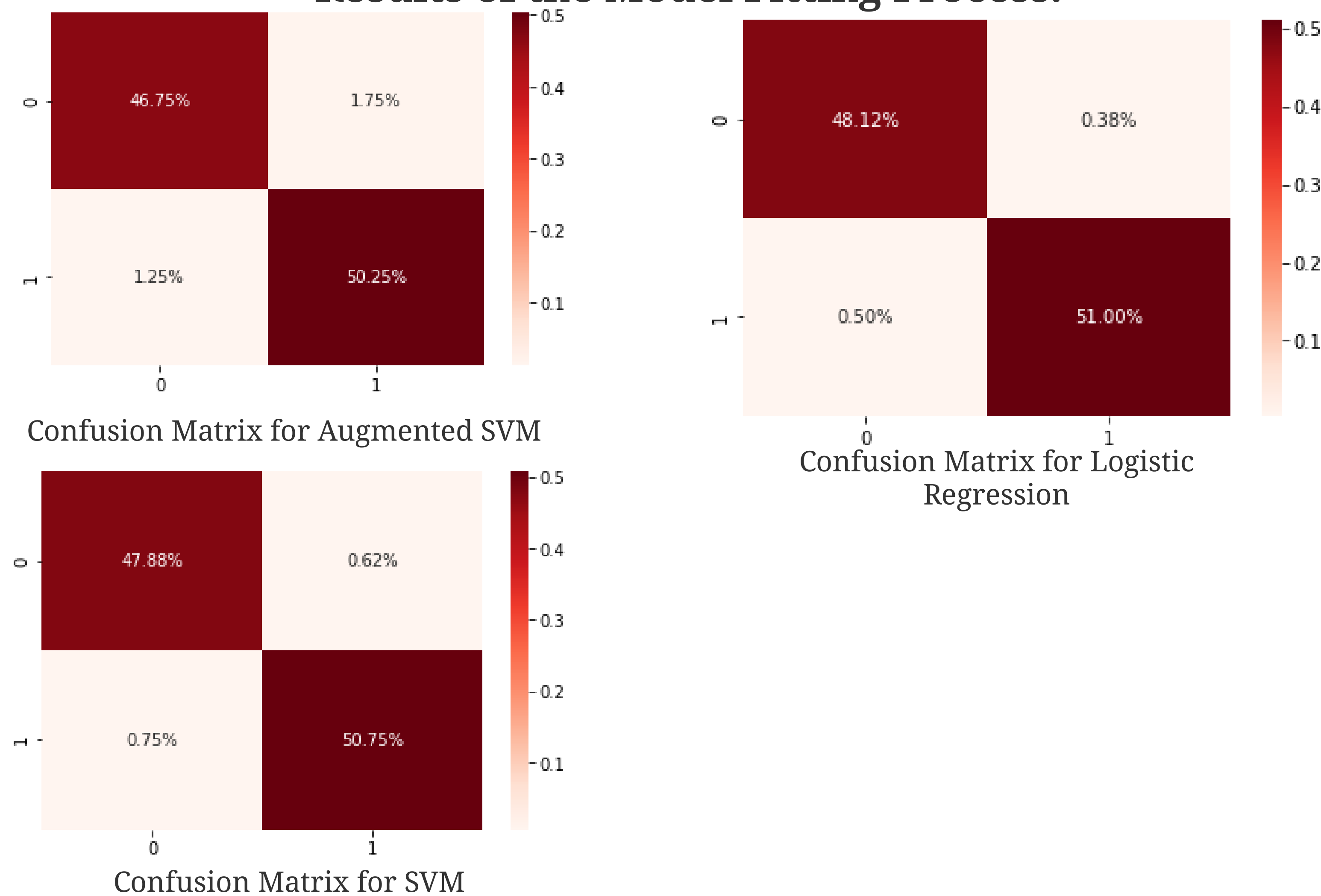
## Model Fitting Process:

1. Vectorize our text into a TF-IDF matrix
2. Create a training, validation and test set of our vectorized data
3. Fit two models on training set: Linear SVM and Logistic Regression Model
4. Evaluate which performs the best on validation set
5. Take best model, so far, augment the TF-IDF matrix with subjectivity score
6. Fit using training set to create a new model
7. Evaluate all models on the test set

## What is TF-IDF?

- TF-IDF stands for Term Frequency-Inverse Document Frequency
- Is a method to account for the frequency of words in a collection of documents
- Aims to weight words by importance
- TF-IDF consists of two parts: TF (Term Frequency and IDF (Inverse Document Frequency)
- $TF\text{-}IDF = TF \times IDF$ , where:
  - TF stands for the term frequency, which could be a simple count of each word in each document
  - IDF stands for Inverse Document Frequency, which is the inverse of the fraction of number of documents where term  $t$  appears among the total documents
  - Intuitively, a high TF-IDF score

## Results of the Model Fitting Process:

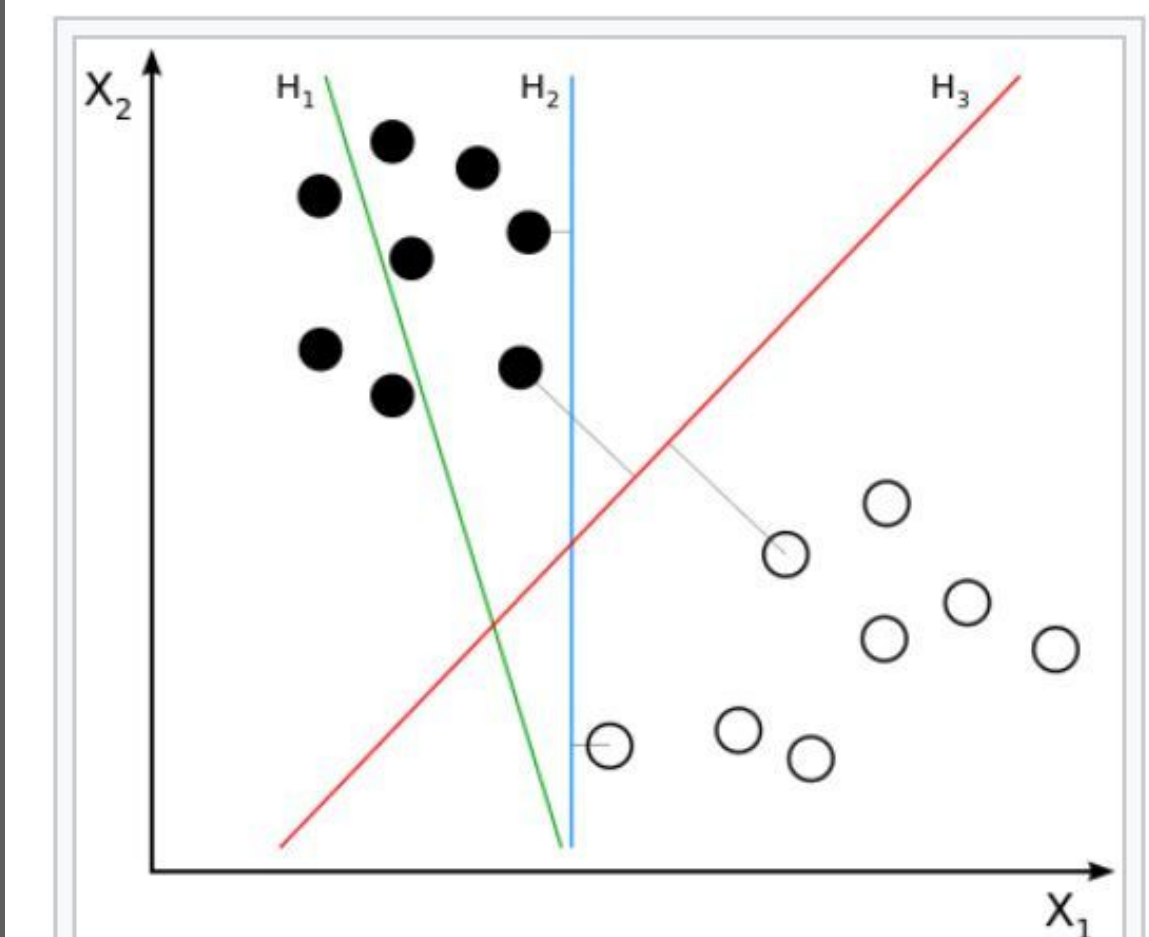


## Conclusions from the Model Fitting Process:

- Accuracy of the TF-IDF only Logistic Regression Model: 99.12%
- Accuracy of the SVM of the TF-IDF only SVM Model: 98.63%
- Accuracy of the Augmented SVM: 97.00%
- Interestingly, augmenting the TF-IDF matrix with subjectivity scores lowered the accuracy
- Unusually high accuracy of the TF-IDF only Logistic Regression
  - I believe this is due to the fact the sources were very homogenous
  - As a result, a very high weight was placed on the specific source
- For future work, two approaches are possible
  - Keep the current data, find all the sources, then remove all the text which mention the source in the article
  - Conduct own web scraping and follow the approach above. With own web scraping, can know the exact source and remove mentions of the source in the text

## What is SVM?

- Support Vector Machines
- Classifies data into binary categories (0 or 1)
- Aims to find a hyperplane (in 2D this would be a line) where the average distance between the hyperplane and each point is the largest



$H_1$  does not separate the classes.  $H_2$  does, but only with a small margin.  $H_3$  separates them with the maximal margin.

Source: Wikipedia

## What is Logistic Regression?

- Classifies data into binary categories
- Aims to find a linear model such that the log odds of an event category occurring (1) is maximized
- Where odds are defined as  $p/(1-p)$  where  $p$  represents the probability an event occurs
- Then, a threshold is defined (default = 0.5)
- If  $p \geq 0.5$ , then classify as 1
- if  $p < 0.5$ , then classify as 0
- Source: Wikipedia