

# FAKE NEWS

A Study of the Differences between Real and Fake News

Claudia Hidrogo  
Madeline Odom  
Jordan Sims  
Xiomara Winkler

# Questions Sought to Answer

1. Is there an association between length of article and the realness of news?
2. What words are most commonly used in real vs fake news?
3. Is there a time when more fake news started showing up (in the Harvard set)?
4. What categories (sports, politics, entertainment) are more prone to fake news?
5. What publications are more prone to publish fake news?

# Datasets

**Source:** Harvard Dataverse **URL:** <https://dataverse.harvard.edu>

*A Large Multi-Labelled News Dataset for the Study of Misinformation in News Articles*

**Source:** Kaggle: Clément Bisailon **URL:** <https://www.kaggle.com/clmentbisailon/fake-and-real-news-dataset?select=True.csv>

*Dataset consisting of fake and real news articles.*

**Source:** Kaggle: Samrat Sinha **URL:** <https://www.kaggle.com/samrat96/fake-news-detection?select=test.csv>

*Fake News Detection - Train dataset.*

**Source:** Kaggle: Paul Larmuseau and Deechain **URL:** <https://www.kaggle.com/c/classifying-the-fake-news/data?select=training.csv>

*Classifying the Fake News – Training dataset.*

**Source:** Data Flair **URL:** <https://data-flair.training/blogs/advanced-python-project-detecting-fake-news/>

*Detecting Fake News with Python and Machine Learning*

**Source:** Zenodo: Radu Prodan, Prateek Agrawal, Pawan Verma **URL:** <https://zenodo.org/record/4561253#.YWNCKnRMKUI>

Generic word embedding dataset for fake news detection consisting of real and fake news from popular news sources.

# Data Pre-processing

<b>Data Cleaning</b>	<p>Prepare and organize the raw datasets before doing the data integration.</p> <p>This step includes:</p> <ul style="list-style-type: none"><li>Delete any empty rows</li><li>Drop NaN values</li><li>Correct inconsistent data</li><li>Make sure that the label attribute only has a binary representation</li></ul> <pre>0      title      ... label 0  As U.S. budget fight looms, Republicans flip t...  ...  0</pre>
<b>Data Integration</b>	<p>Concatenate the datasets containing the fake news articles with the real news articles to form a large dataset.</p>



# Data Pre-processing

<b>Data Transformation</b>	<p>Stop Words: Identify and remove/drop all the words in each news article that do not have any significant meaning</p> <p>Tokenization: Convert the list of words into tokens (numbers) for easy manipulation of the data.</p>	<pre>['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', 'theirs', 'the', 'mselves', 'what', 'which', 'who', 'whom', 'this',</pre> <pre>[[1103, 93, 6335, 259, 2282, 1937, 5, 69, 26, 9, 18, 13, 1103, 106, 6335, 2399, 868, 2282, 1937, 500, 3201, 1233, 980, 753, 988, 52, 5, 69, 2, 87, 18, 13, 93, 1103, 106, 6335, 1351, 5, 69, 228, 5152, 6412, 17, 1246, 11300, 26, 2569, 93, 6412, 2, 18, 13, 433, 1880, 8453, 2473, 292, 4044, 9853, 167, 5, 110, 3631, 12022, 174, 9289], [6659, 4635</pre>
<b>Data Reduction</b>	<p>Use pad sequences to reduce the news articles to an equal length.</p>	<pre>[[ 868 2282 1937 ... 12022 174 9289]  [ 6774 857 165 ... 3120 1576 194]  [11055 17 22 ... 174 3925 2484]  ...  [ 2488 1318 502 ... 112 81 1835]  [ 73 103 58 ... 81 1047 1138]  [ 2 3824 286 ... 3528 98866 16604]]</pre>

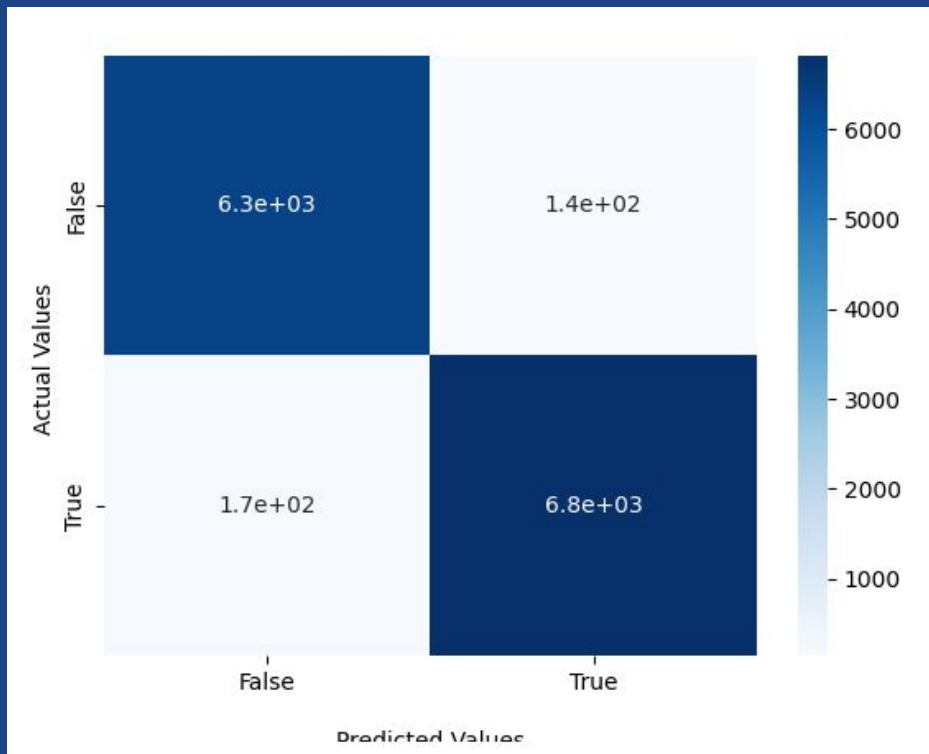
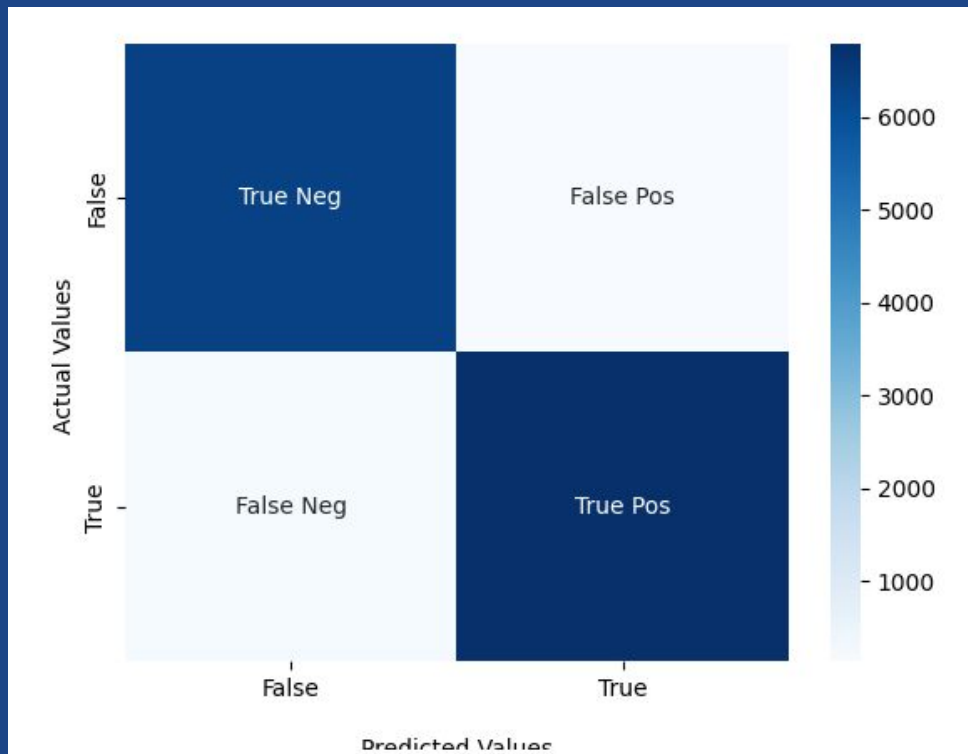
# Tools



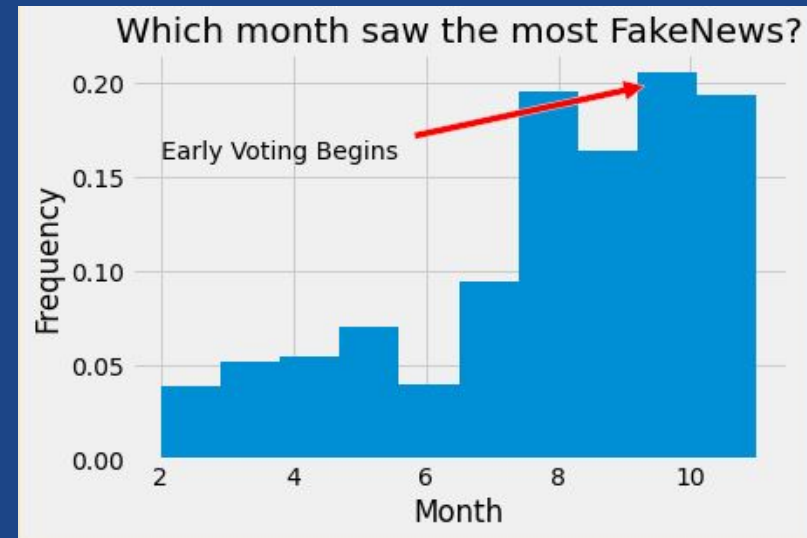
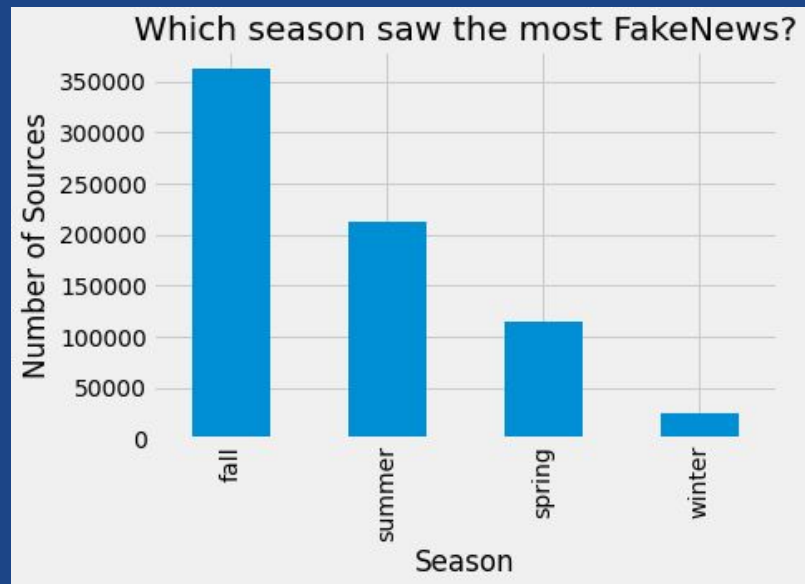
<i>Pandas</i>	For data frame manipulation
<i>Numpy</i>	For numeric analysis
<i>Matplotlib</i> <i>Seaborn</i>	For data visualization
<i>Plotly</i>	To create interactive plots – dynamic visualization
<i>WordCloud</i>	To represent frequency of words
<i>Keras</i> <i>Tensorflow</i>	For deep learning
<i>NLTK</i>	For language processing

# Evaluation

## *Confusion Matrix*



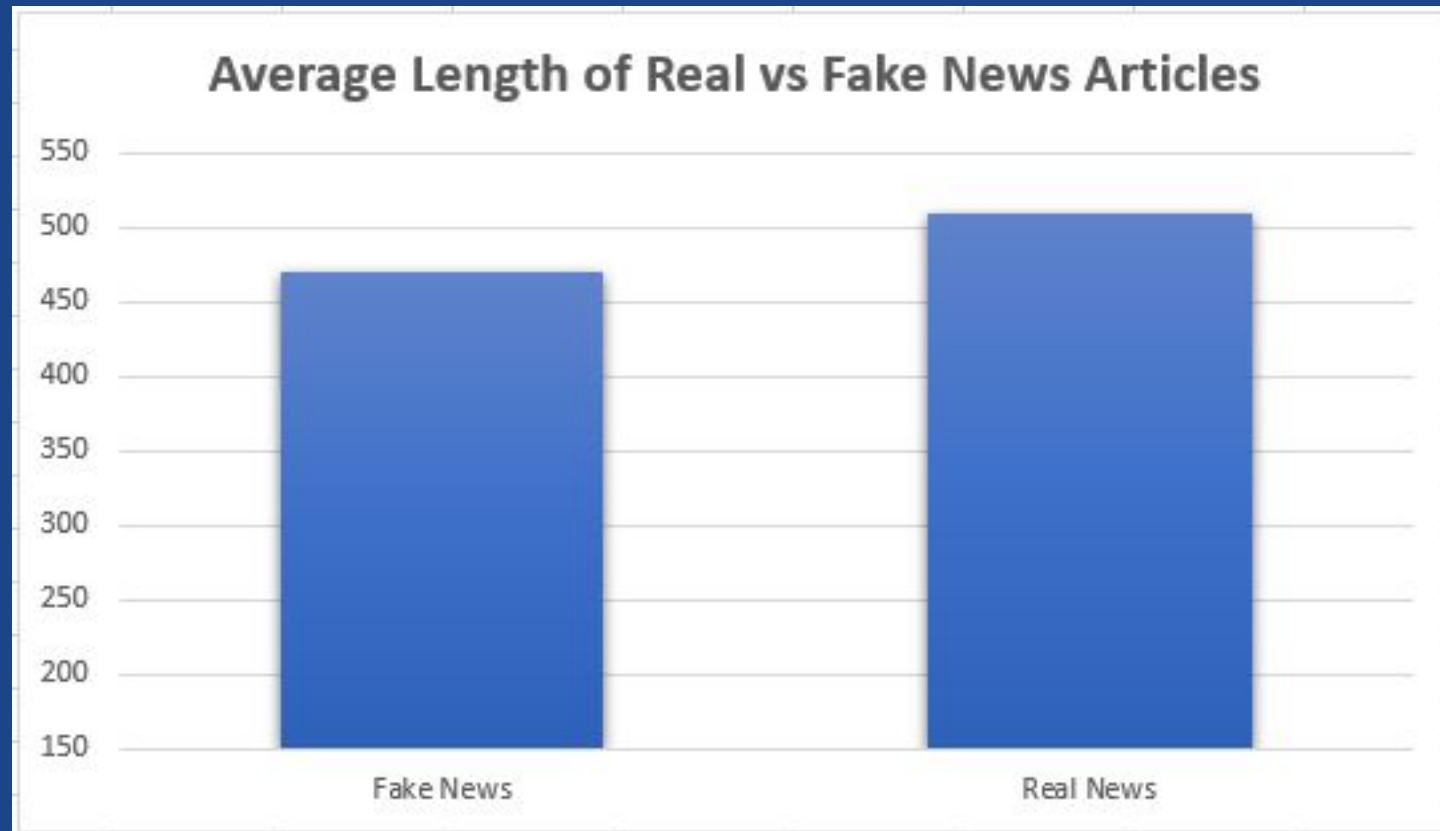
# Knowledge Gained





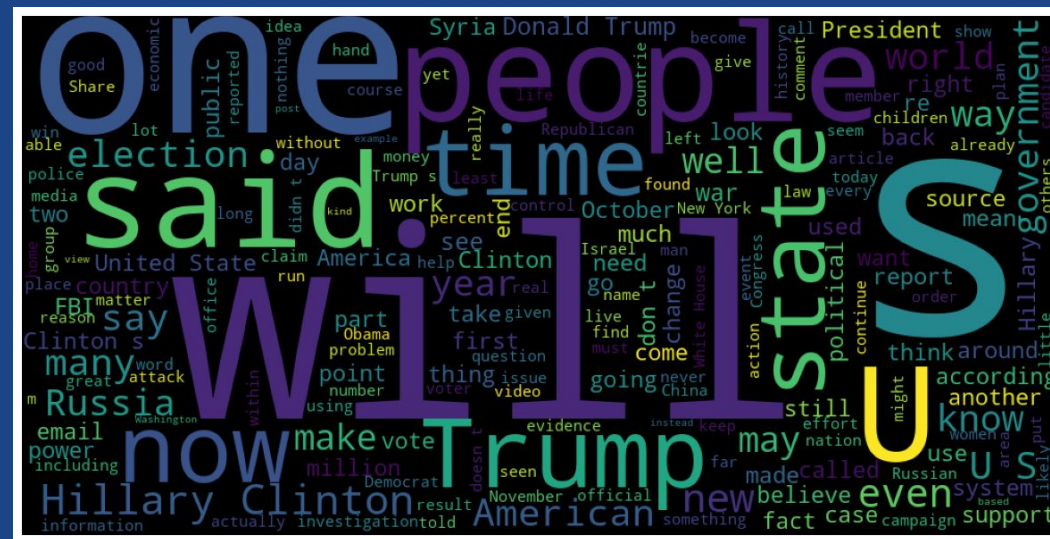
# Knowledge Gained

*Average Article Length*

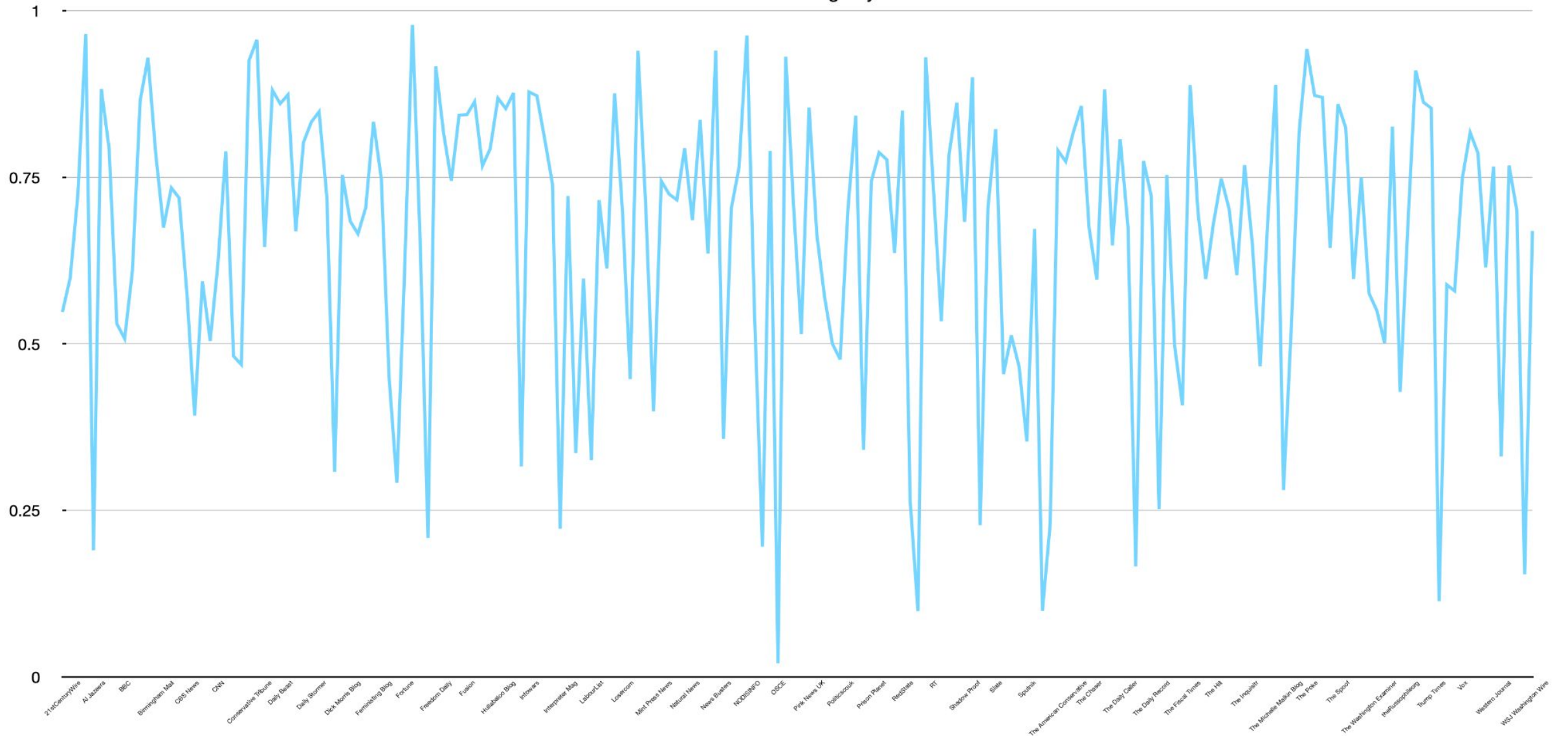


## Word Usage Per Article Type

# Fake



Grand Total Predicted (Average)  
Fake News on Average by Source



# Applications of Knowledge

## *Findings*

- The average reader will not be able to deduce if an article is fake based off article length alone
- The average reader will not be able to deduce if an article is fake based off word choice alone
- Fake news articles tend to be posted more frequently during election cycles in the United States
- Political news is by far the category most prone to fake news
- It is clear that certain publishers of news generate much more fake news than others, so it is important when getting news to verify the integrity of the publisher before trusting what you read