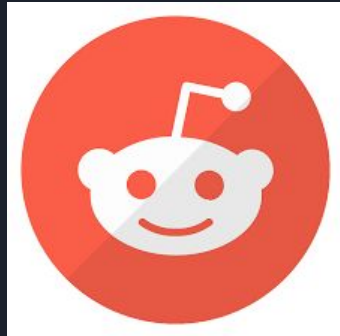# Social Media

Case Study: Subreddit Political Verbiage

# Background

- Political media has become essential to the exchange of political content on platforms the major players being facebook, twitter, and google

# Initial Strategic Plan for political candidate support

# Data Collection & Preliminary Cleaning

- Pulled about 3,000 posts from AskPolitics and Conspiracy Subreddits
- Removed duplicate titles, texts removed by moderators, blanks
- Cleaned titles and text with RegEx
- Feature Engineered word counts, punctuation count, and upper/lower case count

Final dataset:

- Ultimately left with approximately 1,000 subreddit posts per subreddit
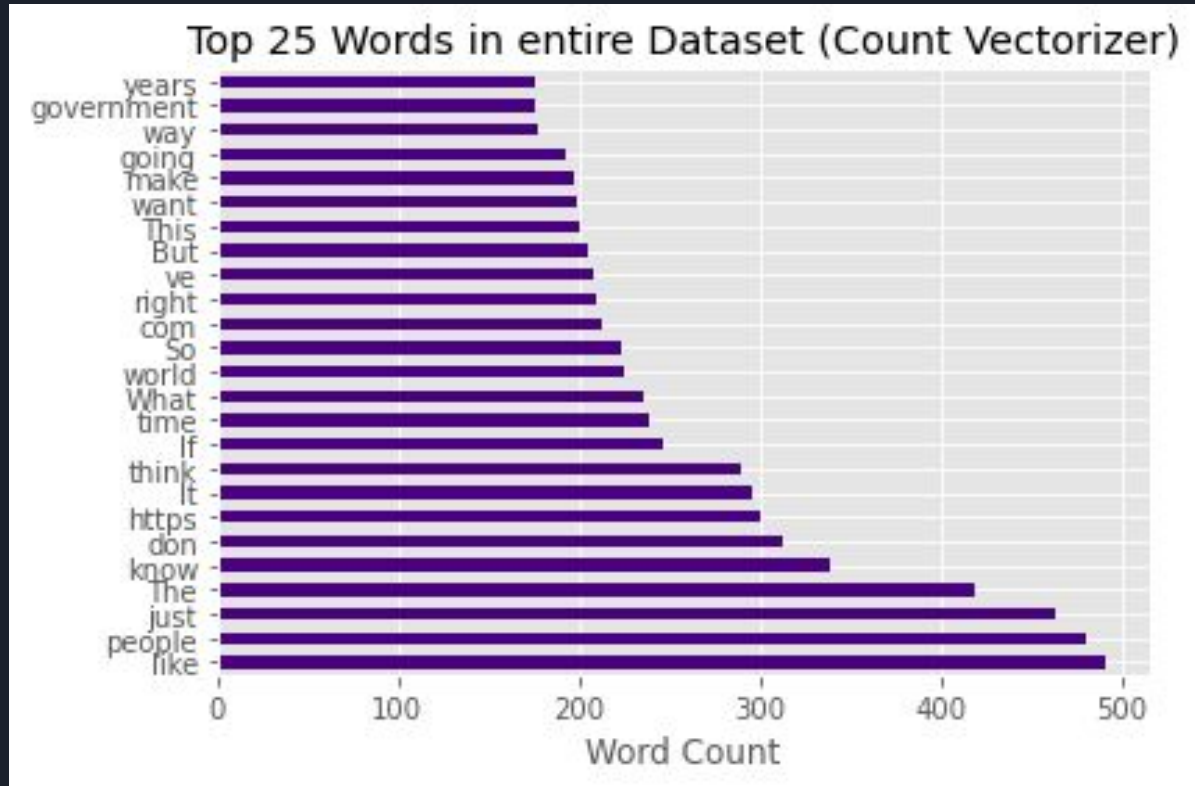
# Exploratory Data Analysis
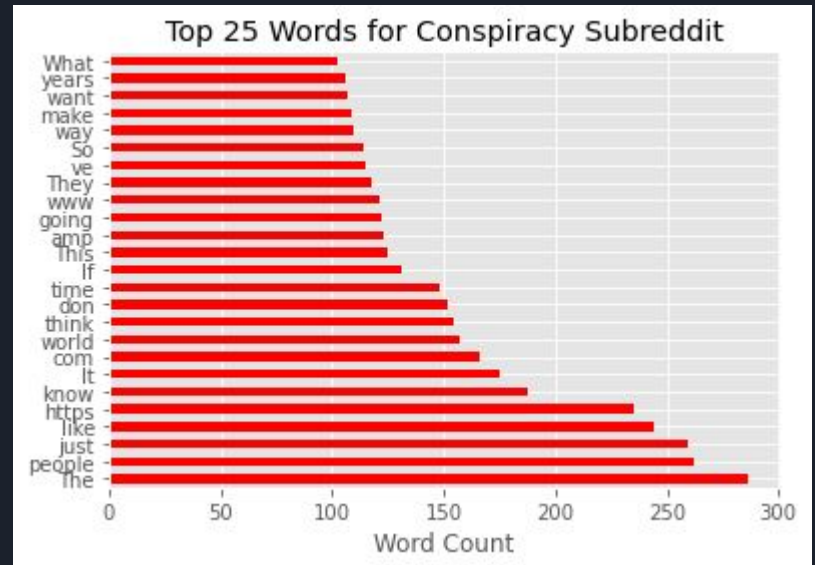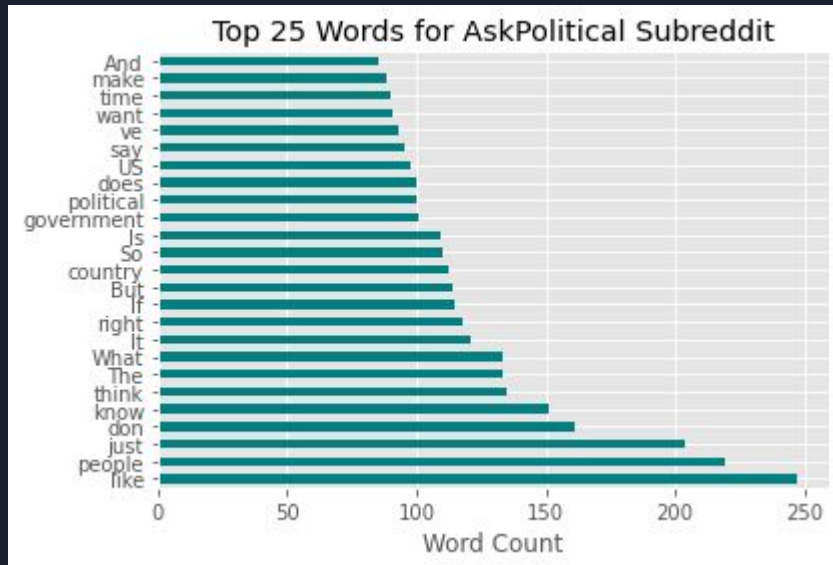
Overall and by subreddits

Characteristics explored:

- Top 25 Words
- Word Count
- Punctuation Count
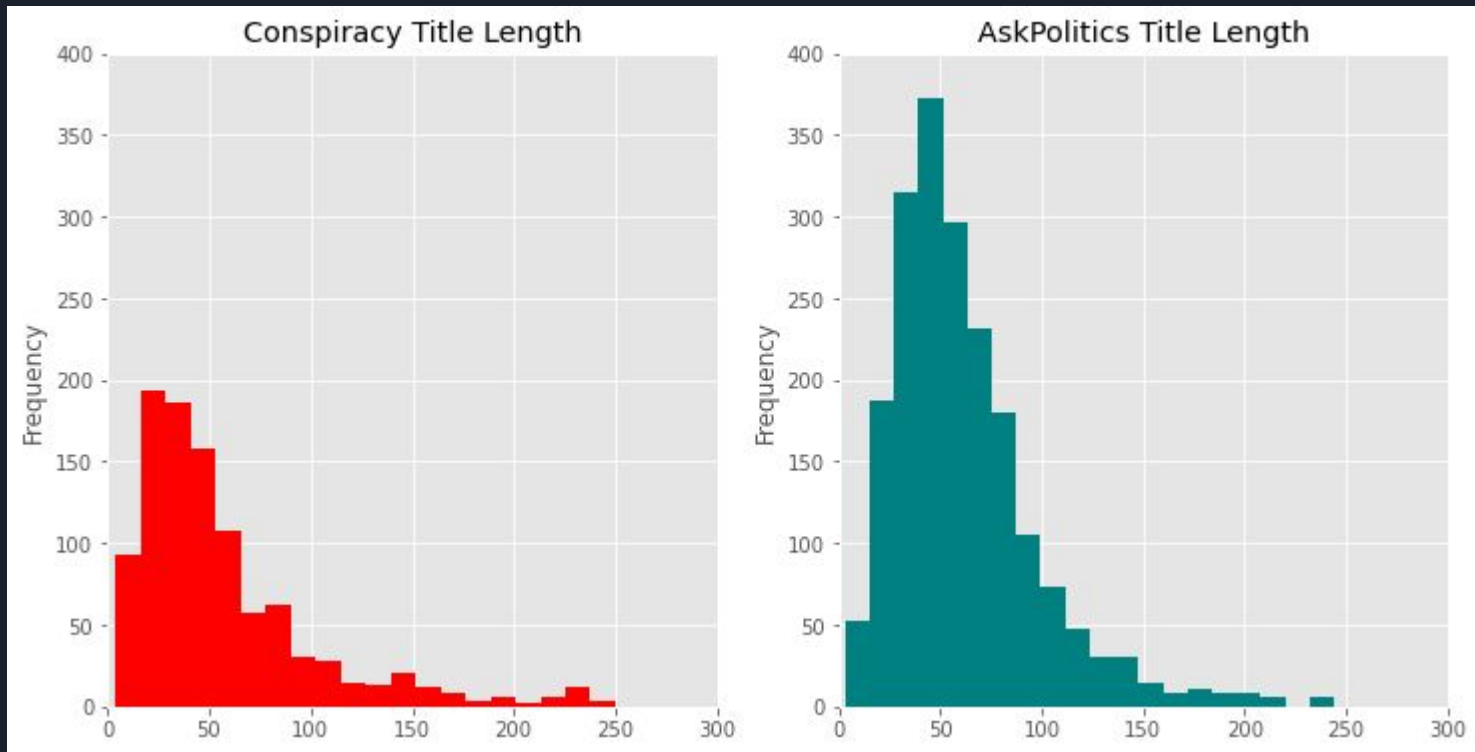- Correlations among features (including Sentiment Analysis Features)
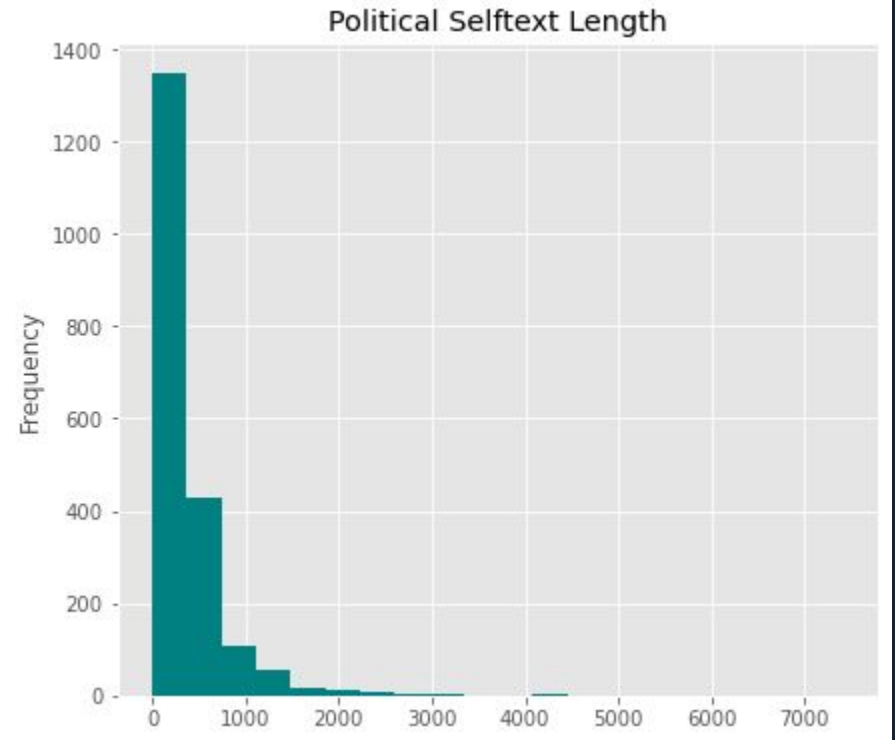
# Top 25 entire Dataset



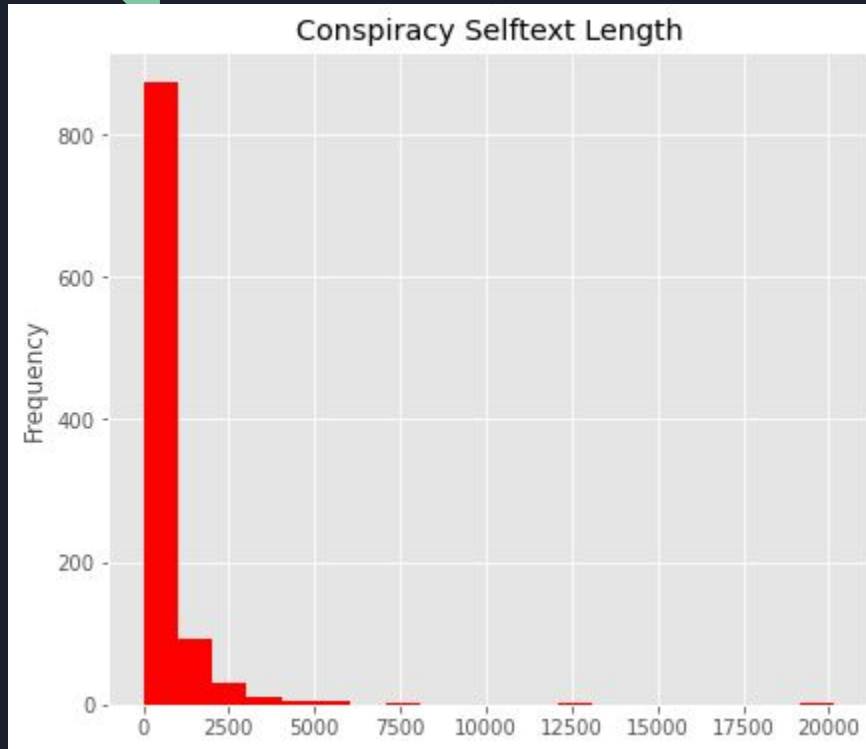Top 25 Words in entire Dataset (Count Vectorizer)

# Top 25 per Subreddit



Top 25 Words for AskPolitical Subreddit



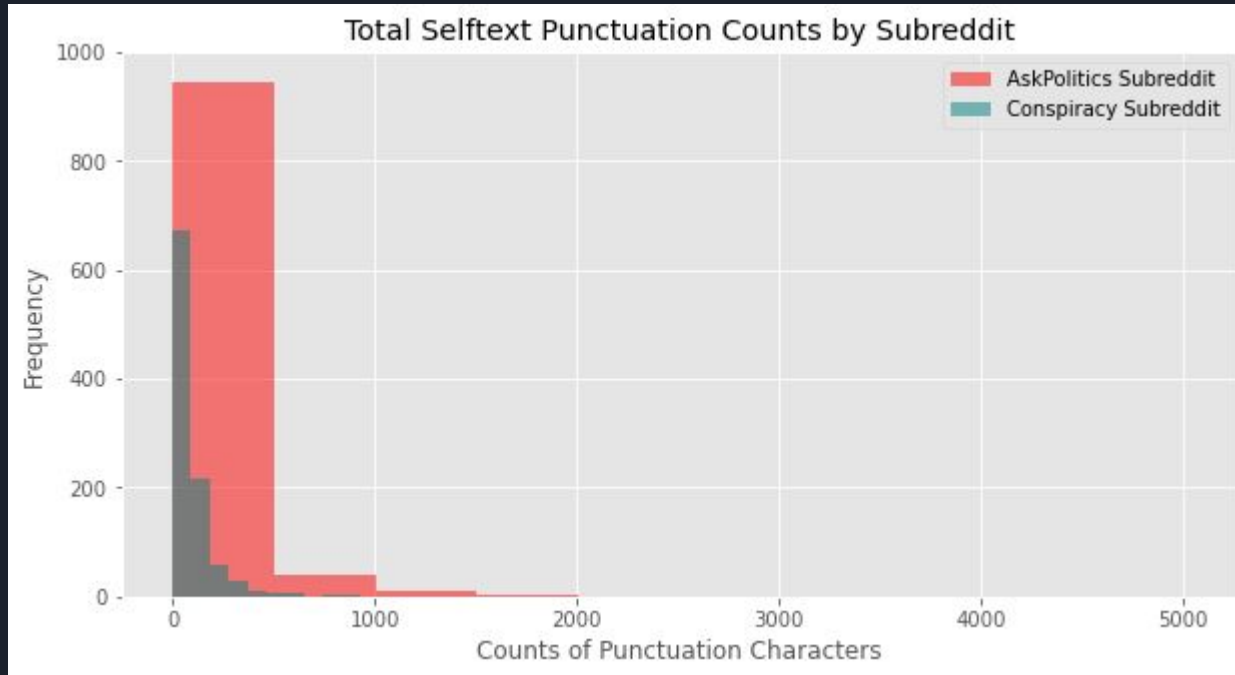Top 25 Words for Conspiracy Subreddit

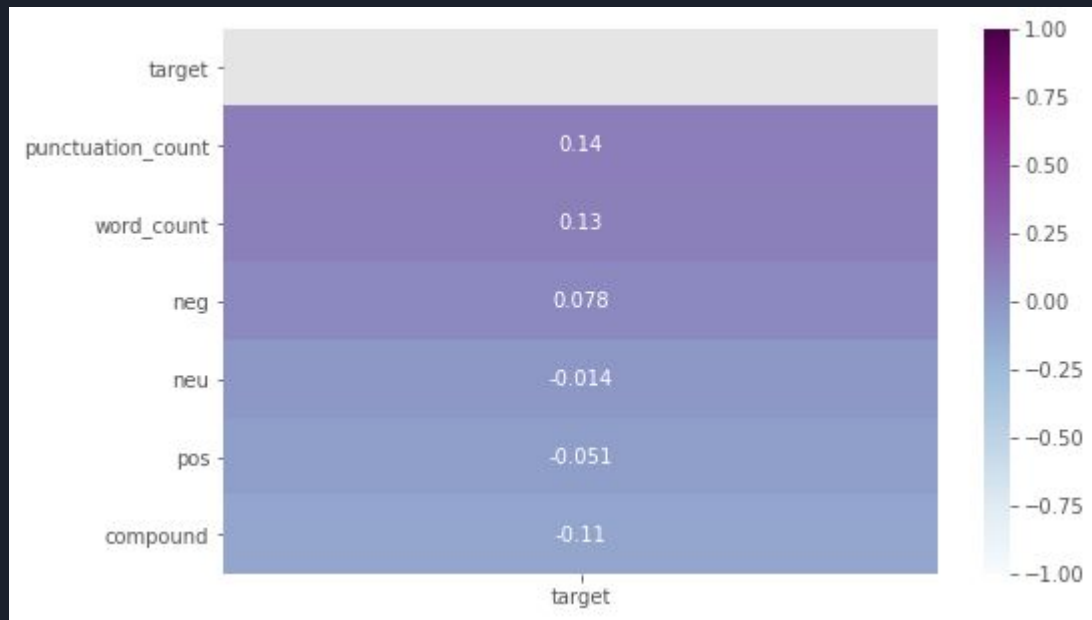# Title Length Distributions per Subreddit

# Self Text Length Distributions per Subreddit

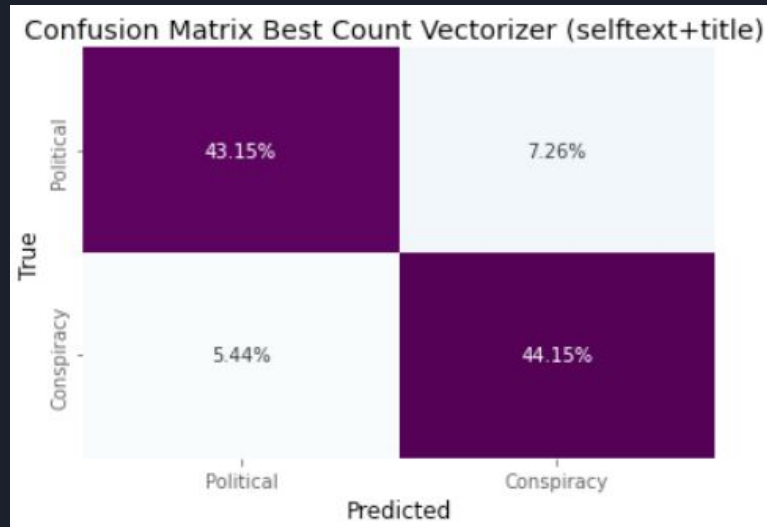# Self Text Punctuation Count Distributions per Subreddit
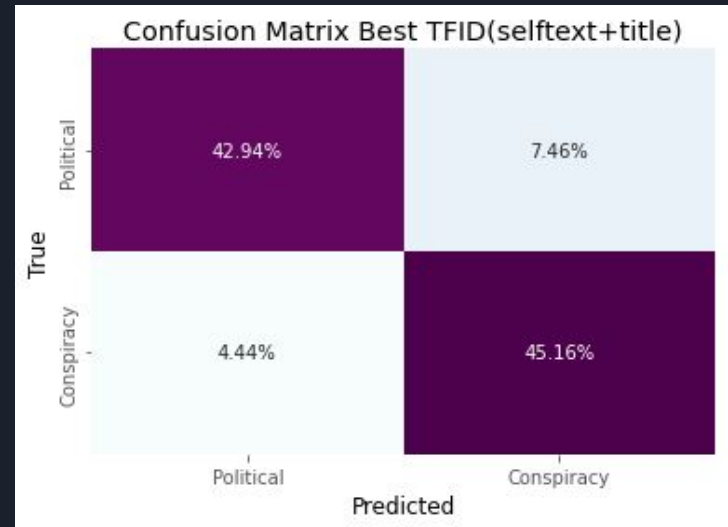
# Model Approach

- Logistic Regression prioritized
- Random Forest explored, ultimately not heavily optimized

# Misclassifications (TFID vs CV) before Word Cnt + Punc counter

Countvectorizer (First)

TFID (Second)



Confusion Matrix Best Count Vectorizer (selftext+title)

|  | Political | Conspiracy |
|---|---|---|
| **Political** | 43.15% | 7.26% |
| **Conspiracy** | 5.44% | 44.15% |

True / Predicted



Confusion Matrix Best TFID(selftext+title)

|  | Political | Conspiracy |
|---|---|---|
| **Political** | 42.94% | 7.46% |
| **Conspiracy** | 4.44% | 45.16% |

True / Predicted

# Best Log Ref(count vectorizer)

train 0.985858585858585859

test 0.8729838709677419



Entire Training Set (Count Vectorizer)
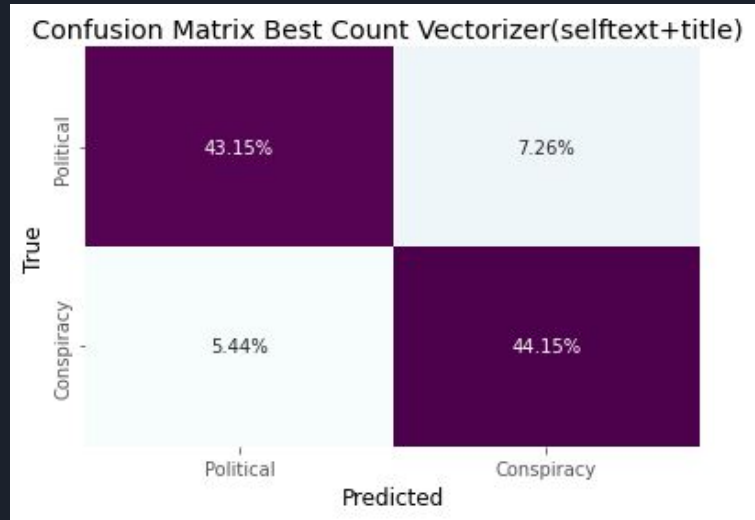
# TFID Vectorizer - 0.947(train), 0.88 (test)



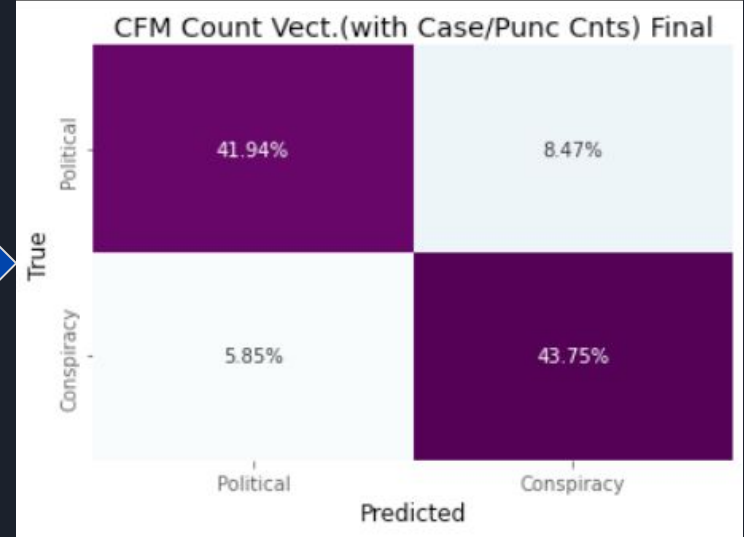Entire Training Set (TFID Vectorizer)

# Misclassifications (using best CV) & Word Cnt + Punc counts included

Countvectorizer (before)

Countvectorizer (after)



Confusion Matrix Best Count Vectorizer(selftext+title)

|  | Political | Conspiracy |
|---|---|---|
| **Political** | 43.15% | 7.26% |
| **Conspiracy** | 5.44% | 44.15% |

True / Predicted



CFM Count Vect.(with Case/Punc Cnts) Final

|  | Political | Conspiracy |
|---|---|---|
| **Political** | 41.94% | 8.47% |
| **Conspiracy** | 5.85% | 43.75% |

True / Predicted

# Random Forest
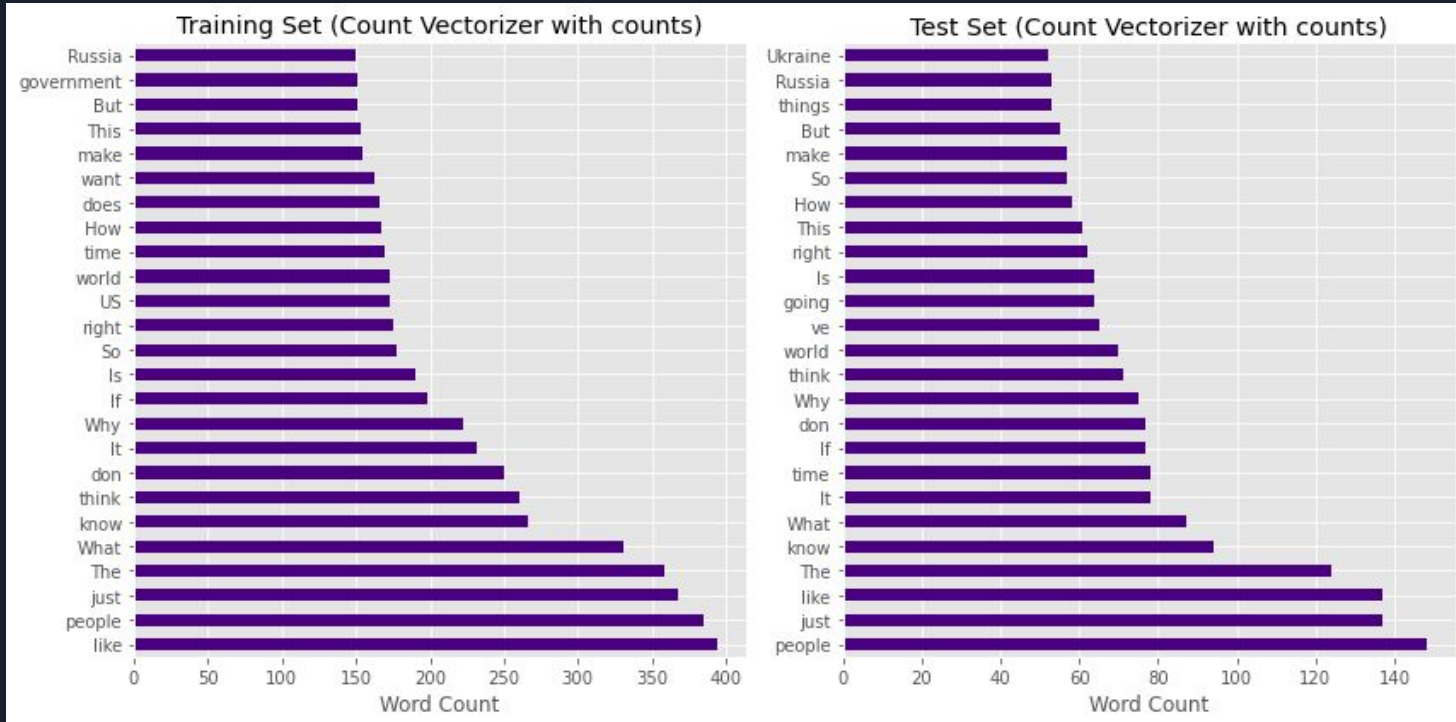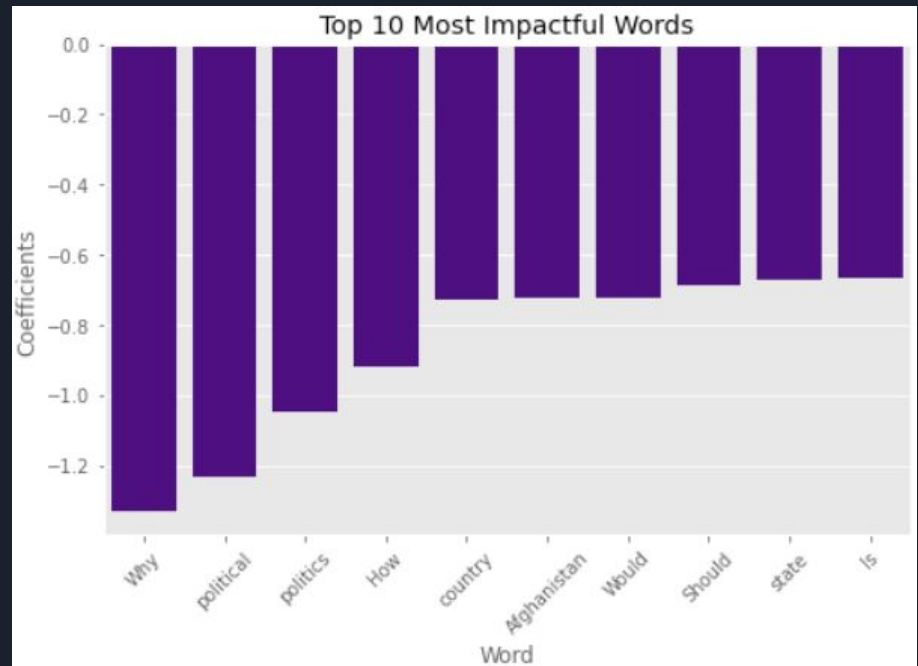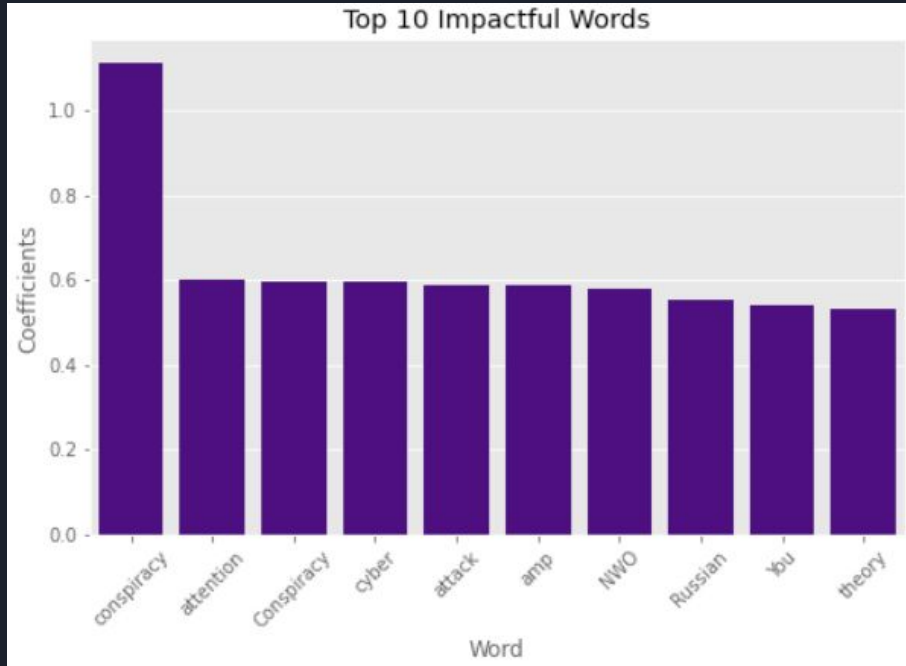


CFM Random Forest

# Best Count Vectorization Insights

# Best Logistic Regression Insights

# Conclusions & Recommendations

- Logistic Regression (without additional features) is best performer
- Coefficients provide insights
- Best text accuracy scores (slightly unbalanced)
- Communications group knows what words to avoid in messaging