

Identifying Fake News using NLP

Creating a model from r/TheOnion and
r/Politics



Outline

Executive Summary

Approach

Wireframes

Next Steps



Executive Summary

Fake News is a growing concern and there are significant efforts to identify fake articles and curtail them.

Reddit is a free public forum where anyone with an account can submit and link posts.

- r/TheOnion is a group dedicated to TheOnion articles.
- r/Politics is a group dedicated for discussing current world politics.



Approach

1. Webscraping

- Reddit API
- 2 sets of 500 posts

2. Preprocessing

- Removing http links
- Removing blank space
- Lemmatizing and Stemming

3. Modelling

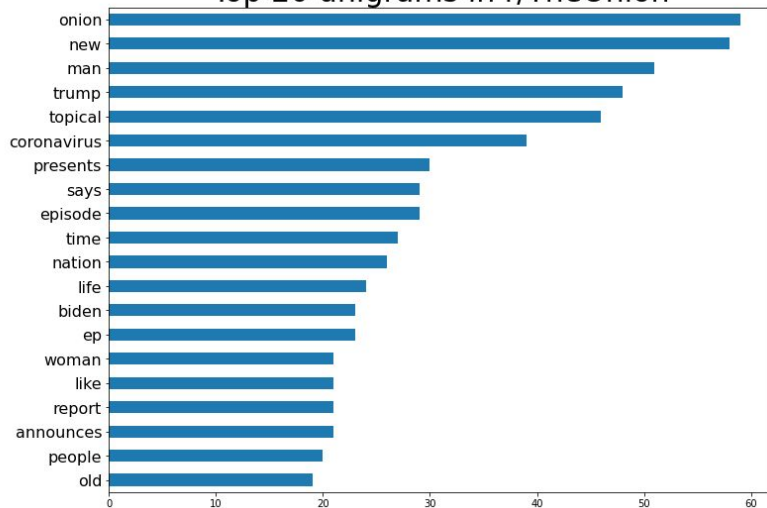
- Count Vectorizer / TFIDF
- Logistic Regression, Multinomial NB, Gaussian NB, Decision Tree, Bagging Classifier, K-Nearest Neighbours, Random Forest, Extra Trees, Support Vector Classifier
- Gridsearch



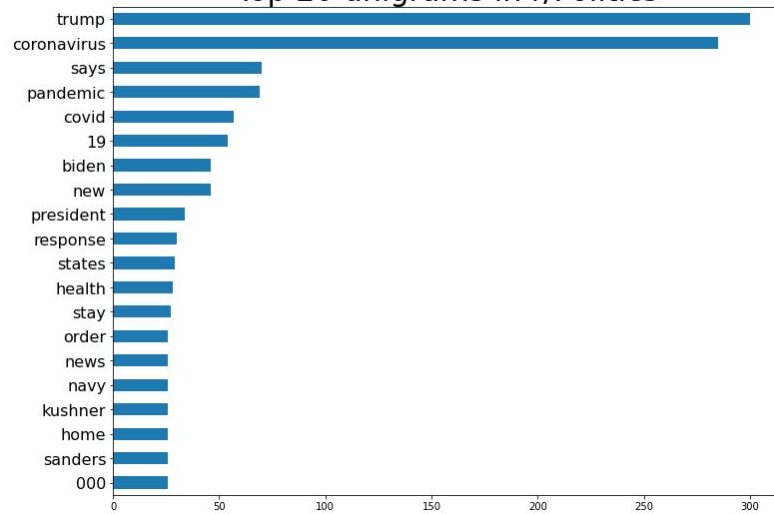
Observations

1. **Word Count**
 - Comparing 1&2 n-grams
2. **Accuracy**
 - Highest Score
 - Confusion Matrix
3. **ROC/ AUC curve**

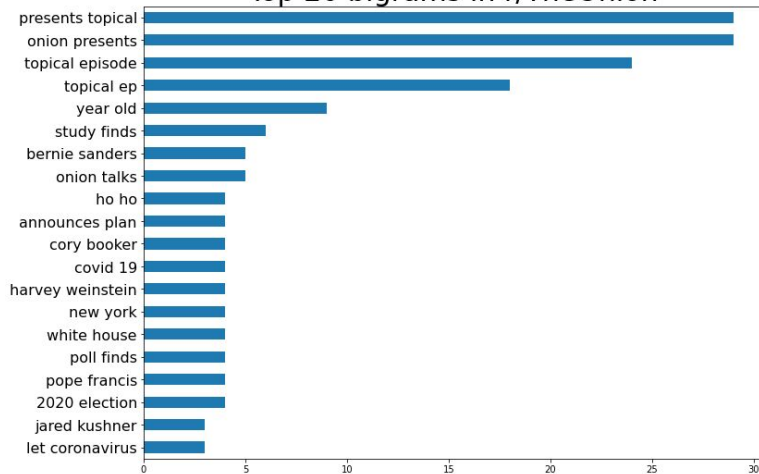
Top 20 unigrams in r/TheOnion



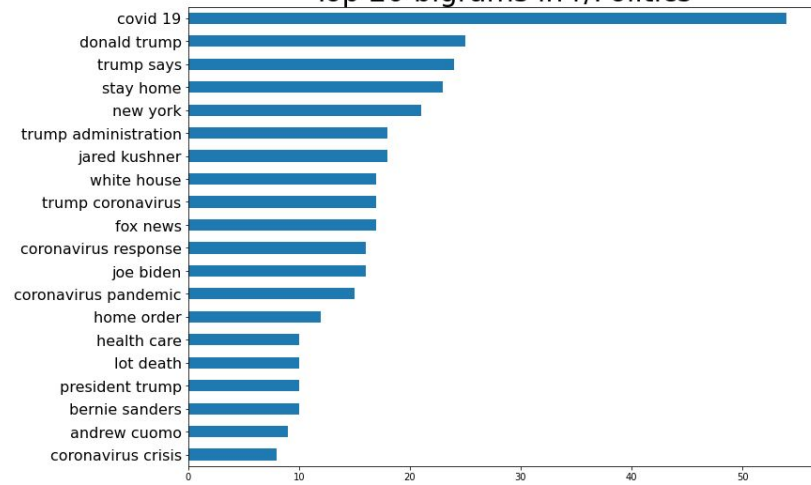
Top 20 unigrams in r/Politics



Top 20 bigrams in r/TheOnion

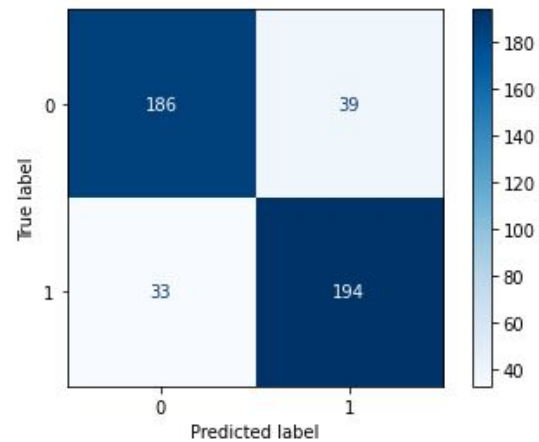


Top 20 bigrams in r/Politics

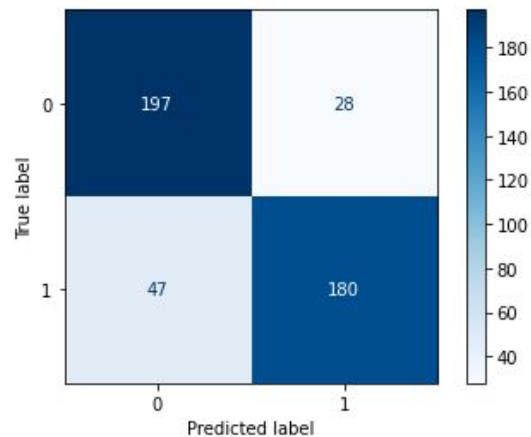


Model	Training Score	Test Score	ROC AUC Score
Logistic Regression	0.9845	0.8407	0.9119
SVC	0.9963	0.8407	0.9178

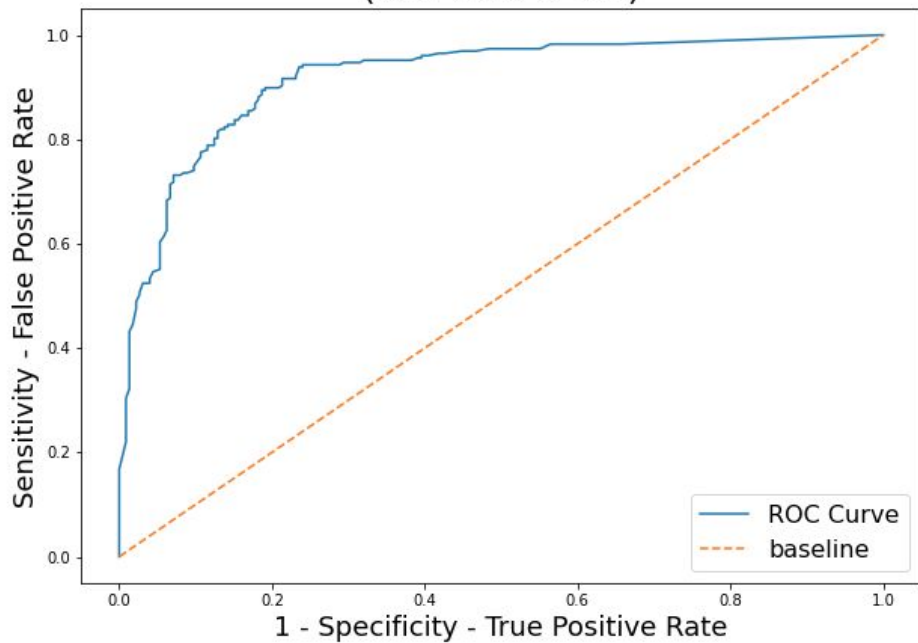
SVC:



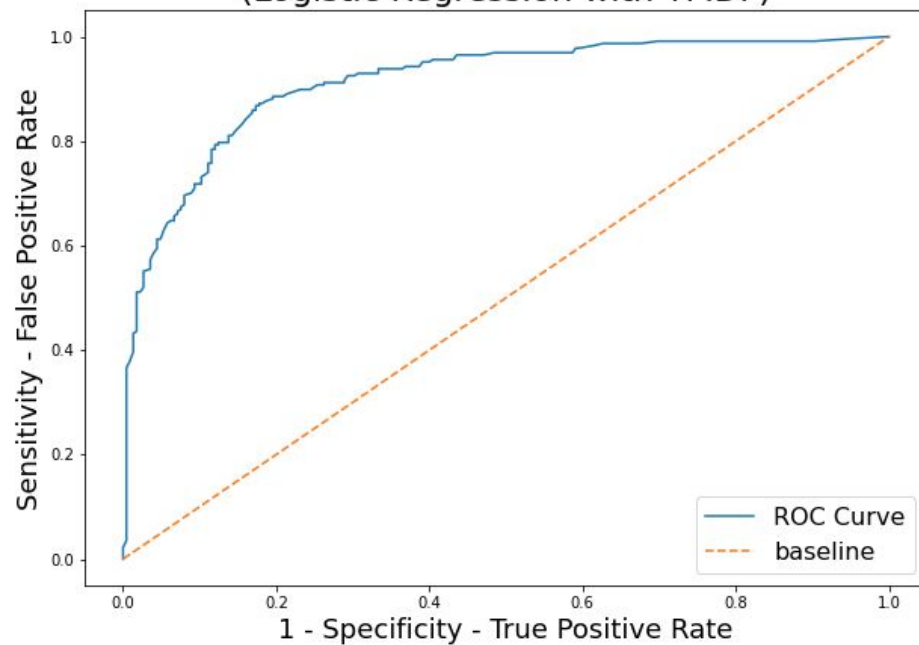
Logistic Regression:



ROC Curve with AUC = 0.918
(SCV with TFIDF)



ROC Curve with AUC = 0.913
(Logistic Regression with TFIDF)





Conclusions

- 17 models trained
- Highest Accuracy Score: 0.8407 | Highest AUC: 0.9178
- From these metrics, SVC can be considered as the best model to identify Fake News.

Improvements:

- Images, emojis, or any other non alphanumeric character.
- Context analysis
- Feature Engineering
- Stopwords

Questions?
