

Table of Contents



01 Introduction

02 Workflow

03 Dataset

04 Tools

05 EDA

06 Final results

Introduction

We all see a lot of fake news almost on a daily. Often its purpose is to spread fear and corruption, especially in societies where the level of digital awareness is low.



Workflow

Data cleaning EDA Classification Topic Modeling

NLP

Remove Remove Remove **Emails** punctuation Stop words Lemmatize stemming Remove numbers

&

Dataset







News dataset

44k rows

4 columns



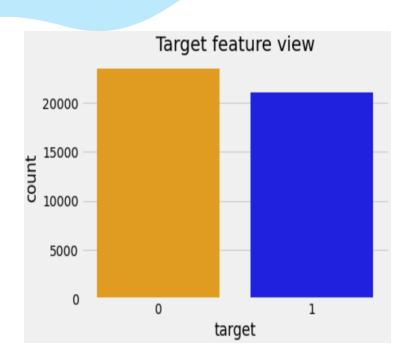


- Python, Jupyter notebook
- NumPy, Pandas
- Matplotlib, Seaborn and wordcloud
- Sklearn
- **NLTK** and gensim
- **Regular Expressions**
- **Ngrams**

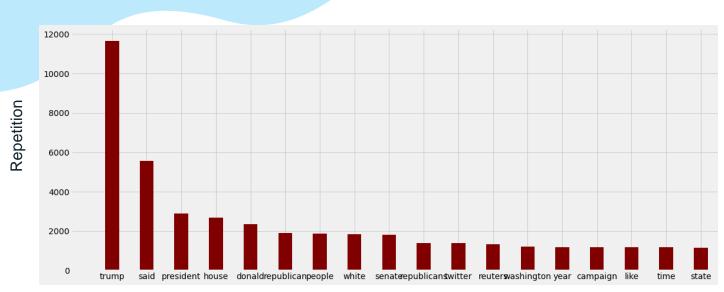
EDA



Fake and real counts



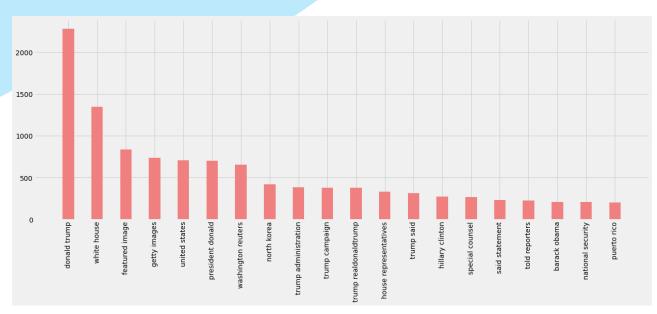
Unigrams Counts



Words in unigrams form

Bigrams counts





Words in Bigrams form

Trigrams Counts



Words in Trigrams form

Most common words in real news

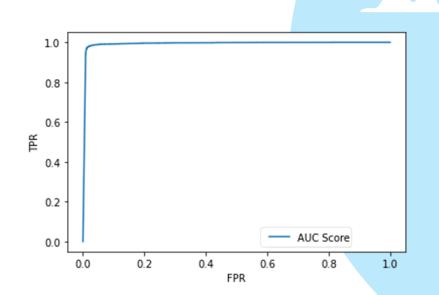


Most common words in fake news





Classification using naïve bayes models

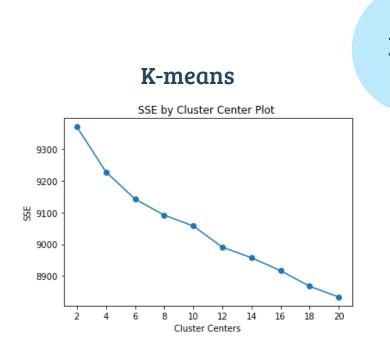


Accuracy Score

0.95	MultinomialNB Baseline model
0.965	BernoulliNB
0.98	BernoulliNB After Tuning
0.978	BernoulliNB Testing



What is the optimal number of topics?



LDA

Hyperparameter Tuning

10 Topics

LDA using different libraries

Bag of word

TF-IDF model

LDA

Count Vectorizer

LDA

Conclusion

- In this project we built a model that predict fake and real news by using classification technic
- The best score was 98%
- Our best method in topic modeling was LDA

Thank You