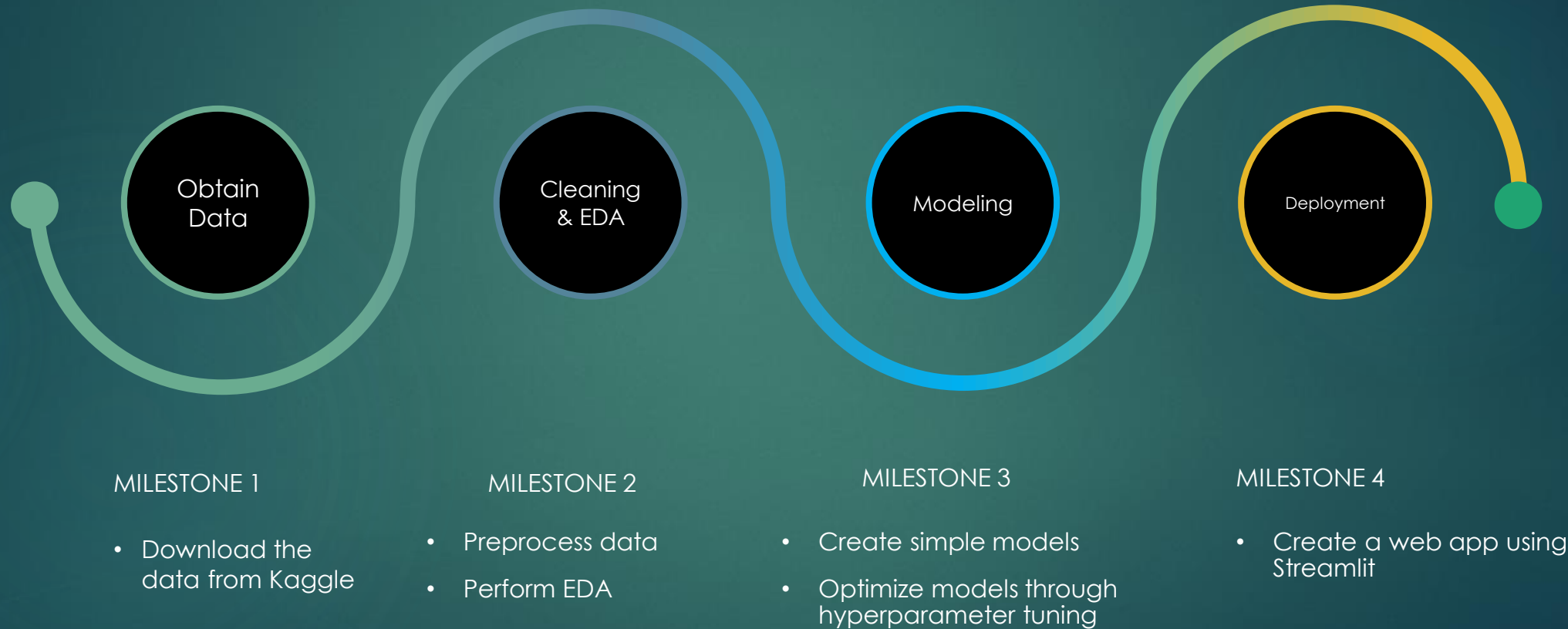# Fake News Detection

GETTING TO THE TRUTH IN AN ERA OF MISINFORMATION

*BY ADAM CUMURCU*

# Business Case

- Fake news has been on the rise this past decade
- Denying the 2020 election results
- Spreading misinformation about the COVID vaccine
- Breakdown of shared reality
- One of the main sources of fake news is social media, such as Facebook and Twitter
- 2019: 8 percent of engagement with the 100 top-performing news sources on social media was dubious
- 2020: this number more than doubled
- Most popular news platform on Facebook in 2021: The Daily Wire
- A fake news detection system like the one I aim to produce can be used by social media companies to filter out misinformation.
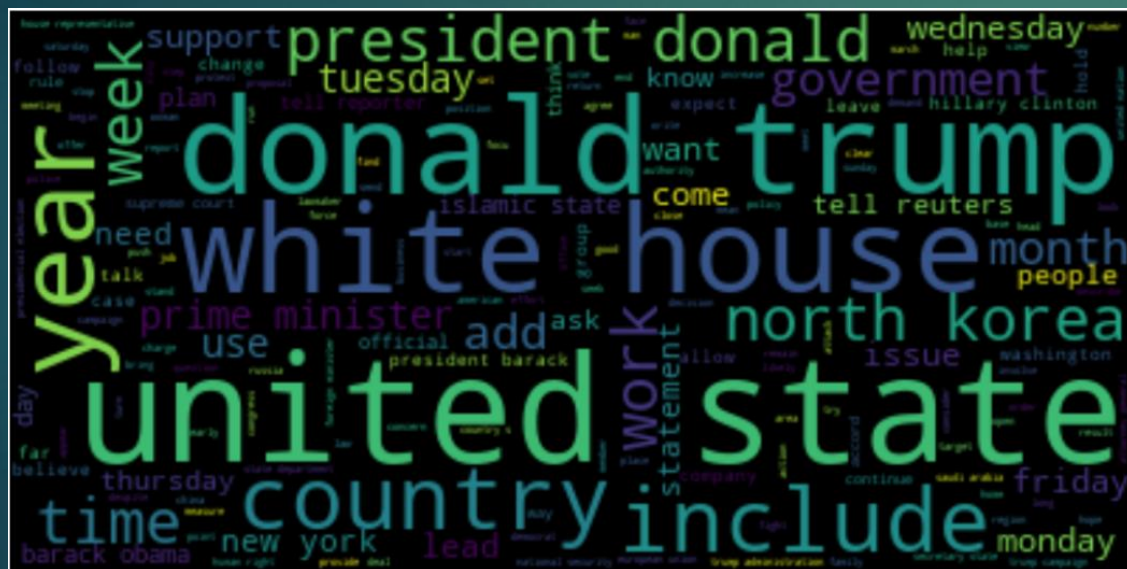
# Modeling Roadmap

Obtain Data

Cleaning & EDA

Modeling

Deployment

**MILESTONE 1**

- Download the data from Kaggle

**MILESTONE 2**

- Preprocess data
- Perform EDA

**MILESTONE 3**

- Create simple models
- Optimize models through hyperparameter tuning

**MILESTONE 4**

- Create a web app using Streamlit

# Data Understanding

- Date source: 'Fake and real news dataset' from Kaggle.
- 21,417 real articles from Reuters
- 23,481 fake articles from various sources.
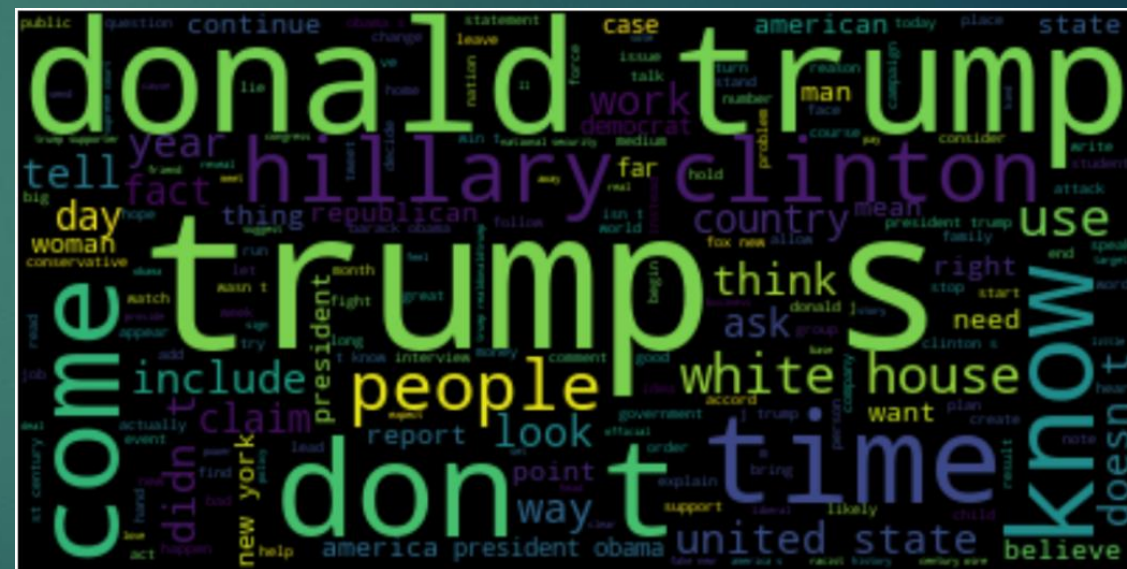- Period covered: January 2017 to December 2017.

# Data Preparation

▶ Only looking at the articles, not the titles, subject or date.

▶ Data preparation steps:

    1. Use regex to remove the names of the news outlet and city of origin from the   real articles.

    2. Use Spacy to tokenize the text of the articles and remove stop words.
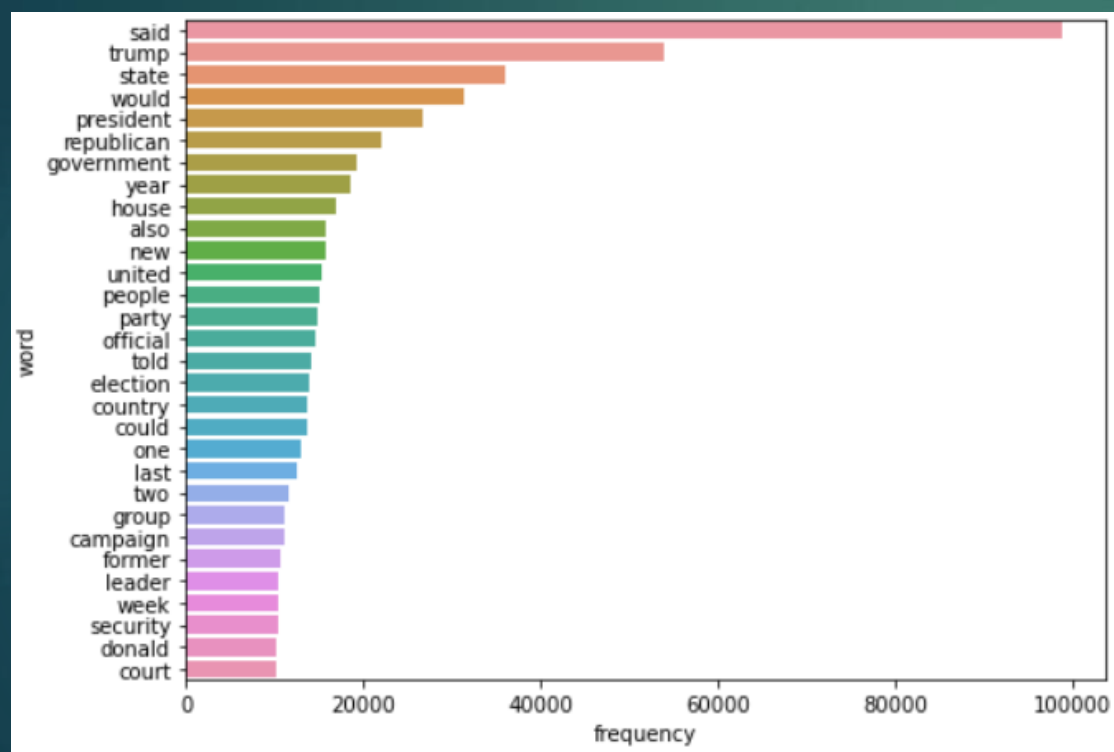
# World Cloud Visualization
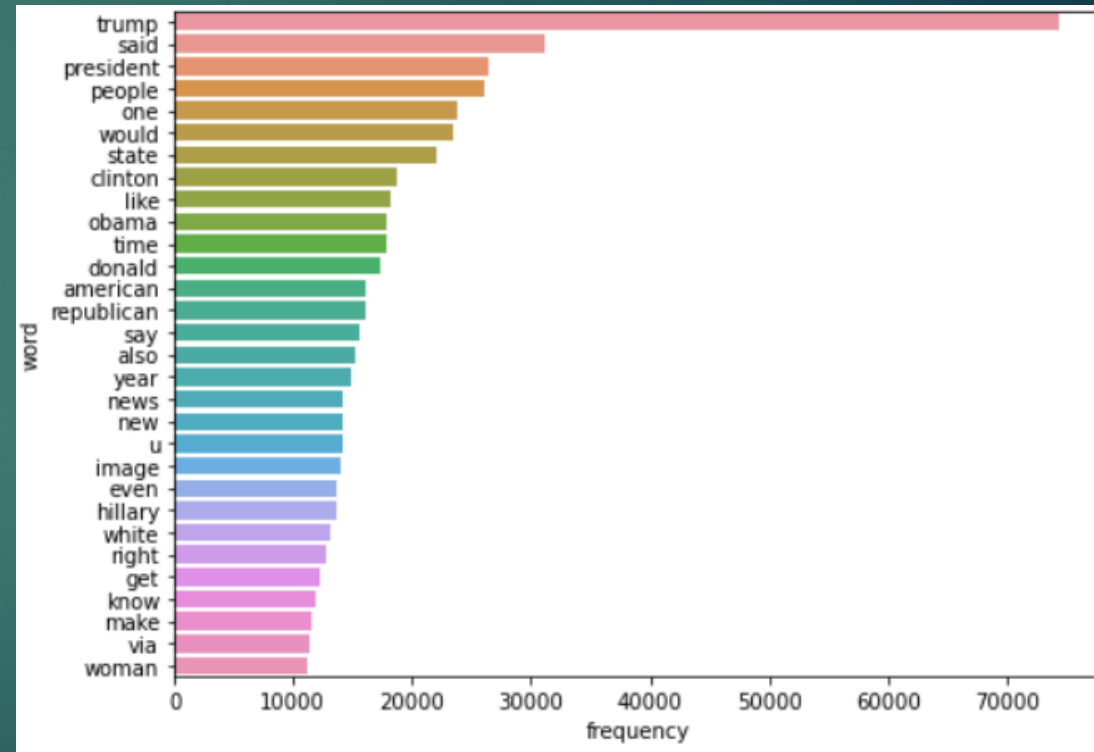
Real Articles

Fake articles
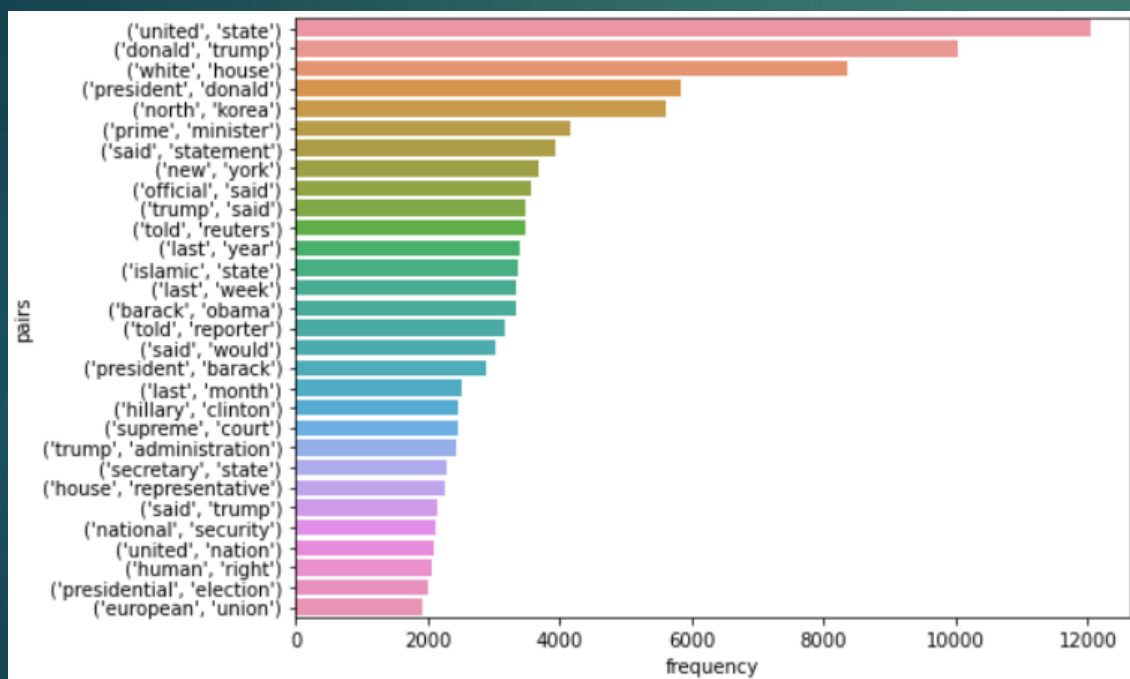
# Unigram Frequency Distribution
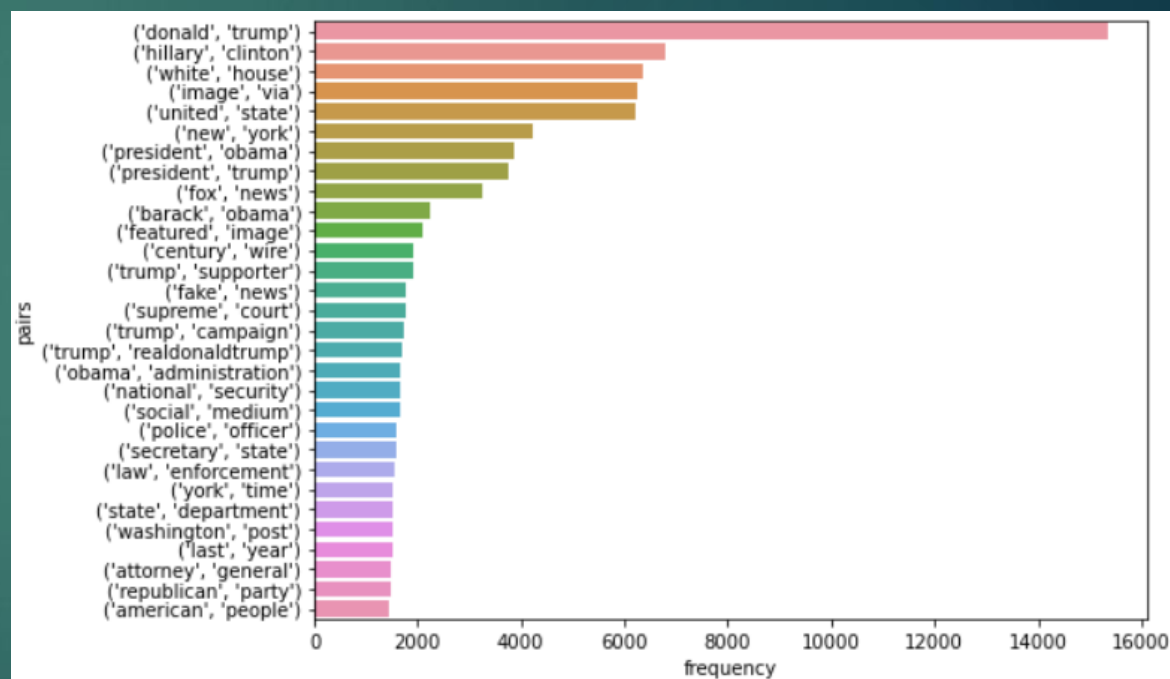
Real Articles

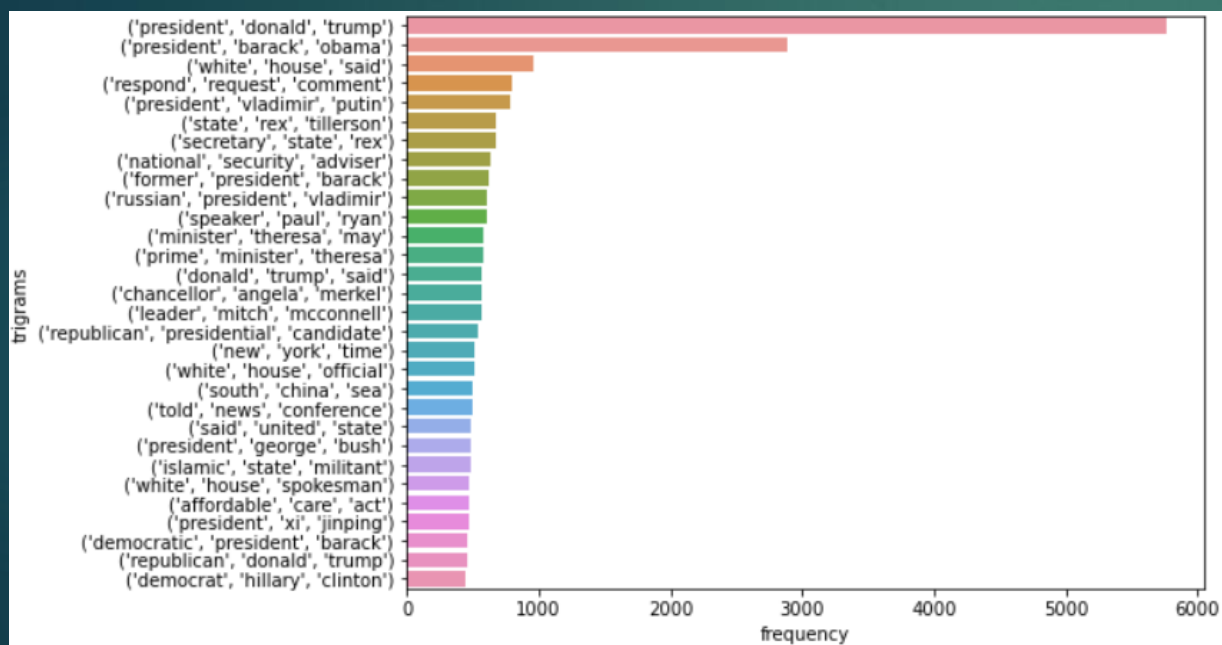Fake Articles
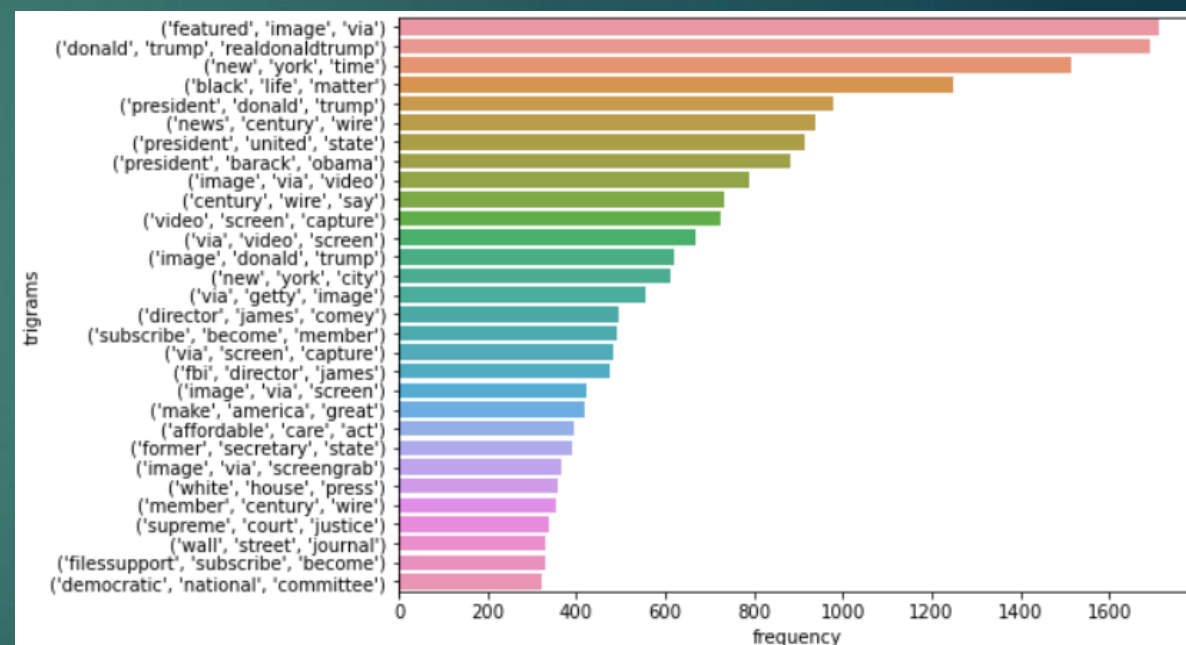
# Bi-gram Frequency Distribution

Real

Fake

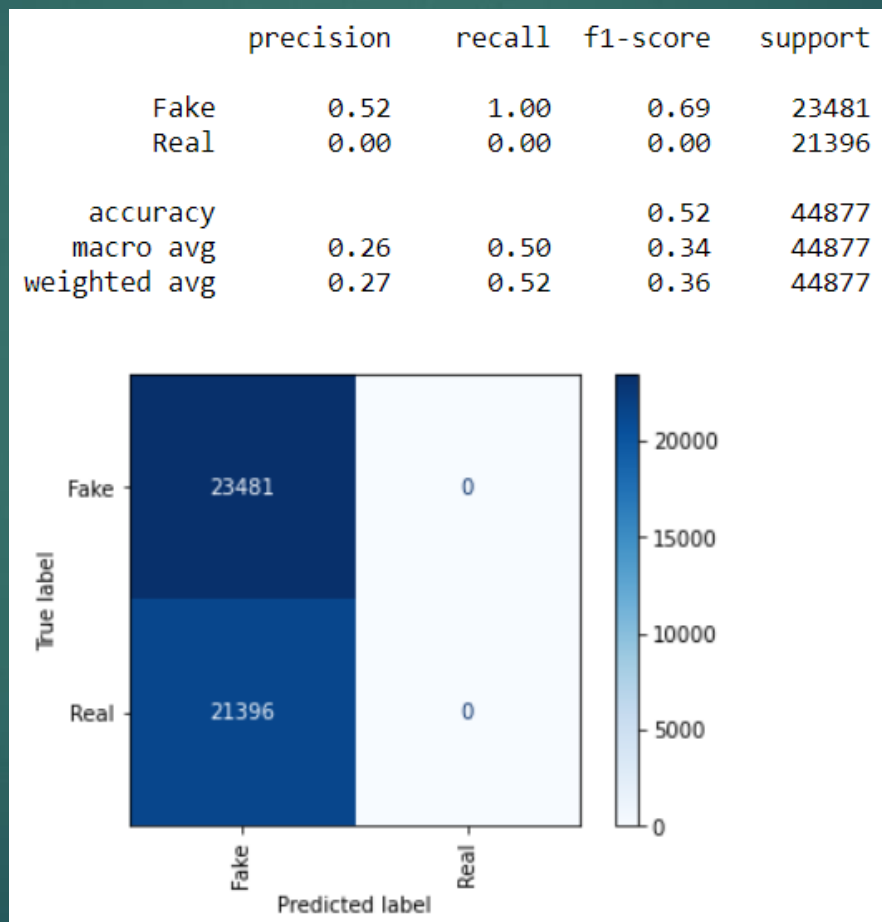# Tri-Gram Frequency Distribution

Real

Fake

# Modeling

- 1. dummy classifier to check the baseline score.
- 2. LogisticRegression
- 3. Multinomial Naïve Bayes
- 4. Random Forest
- 5. Voting classifier (LogisticRegression, Multinomial Naïve Bayes and Random Forest)
- Best Results: LogisticRegression
- Scoring: f1 (not overly concerned with false positives or negatives)
- Initial split: 0.25, which resulted in an f1 score of 0.99.
- Subsequent split (to address overfitting): 0.5, which resulted in an f1 score of 0.98

# Modeling Results

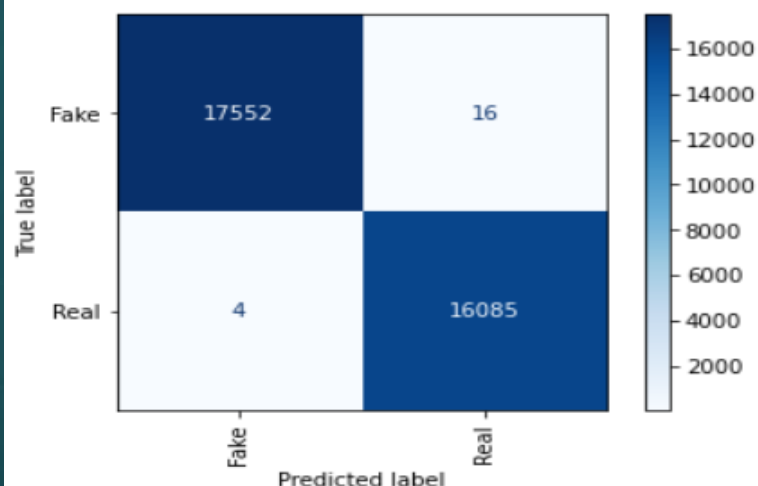| Model | Training F1-Score | Testing F1-Score |
| --- | --- | --- |
| Logistic Regression | 1 | 0.99 |
| Multinomial Naïve Bayes | 0.96 | 0.94 |
| Random Forest Classifier | 1 | 0.95 |
| Voting Classifier | 1 | 0.98 |

# Dummy Classifier

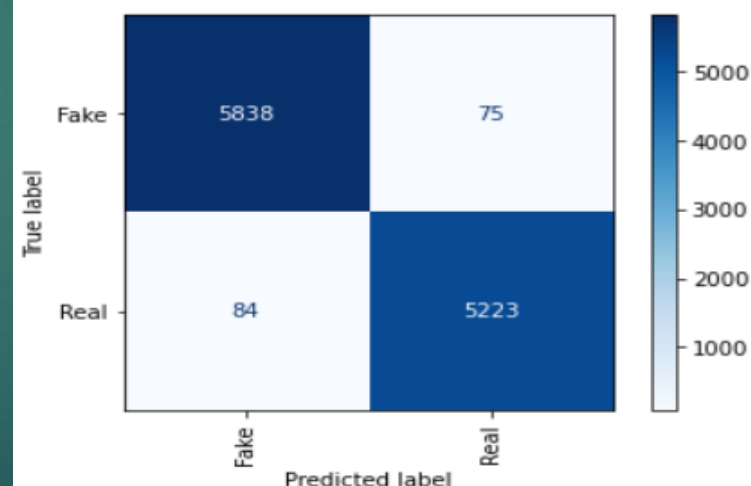# Logistic Regression Gridsearch Results with Split Set to 0.25

## Training Set

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Fake | 1.00 | 1.00 | 1.00 | 17568 |
| Real | 1.00 | 1.00 | 1.00 | 16089 |
| accuracy |  |  | 1.00 | 33657 |
| macro avg | 1.00 | 1.00 | 1.00 | 33657 |
| weighted avg | 1.00 | 1.00 | 1.00 | 33657 |



## Test Set

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Fake | 0.99 | 0.99 | 0.99 | 5913 |
| Real | 0.99 | 0.98 | 0.99 | 5307 |
| accuracy |  |  | 0.99 | 11220 |
| macro avg | 0.99 | 0.99 | 0.99 | 11220 |
| weighted avg | 0.99 | 0.99 | 0.99 | 11220 |

# Logistic Regression Gridsearch Results with Split Set to 0.5

## Training Set

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Fake | 1.00 | 1.00 | 1.00 | 11765 |
| Real | 1.00 | 1.00 | 1.00 | 10673 |
| accuracy |  |  | 1.00 | 22438 |
| macro avg | 1.00 | 1.00 | 1.00 | 22438 |
| weighted avg | 1.00 | 1.00 | 1.00 | 22438 |

## Test Set

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Fake | 0.99 | 0.98 | 0.98 | 11716 |
| Real | 0.98 | 0.98 | 0.98 | 10723 |
| accuracy |  |  | 0.98 | 22439 |
| macro avg | 0.98 | 0.98 | 0.98 | 22439 |
| weighted avg | 0.98 | 0.98 | 0.98 | 22439 |

# Next Steps

- Deploy web app so that anyone can verify if an article is real or fake.
- Train model on newer articles to keep it up to date.

# Thank you

ADAM CUMURCU

EMAIL: ACUMURCU@GMAIL.COM

GITHUB: HTTP://GITHUB.COM/ADAMCUMURCU

LINKEDIN: HTTP://LINKEDIN.COM/IN/ADAM-CUMURCU

BLOG: HTTP://ACUMURCU.MEDIUM.COM/