



# NLP & APIs

Using Machine Learning  
to classify reddit posts

# Problem Statement:

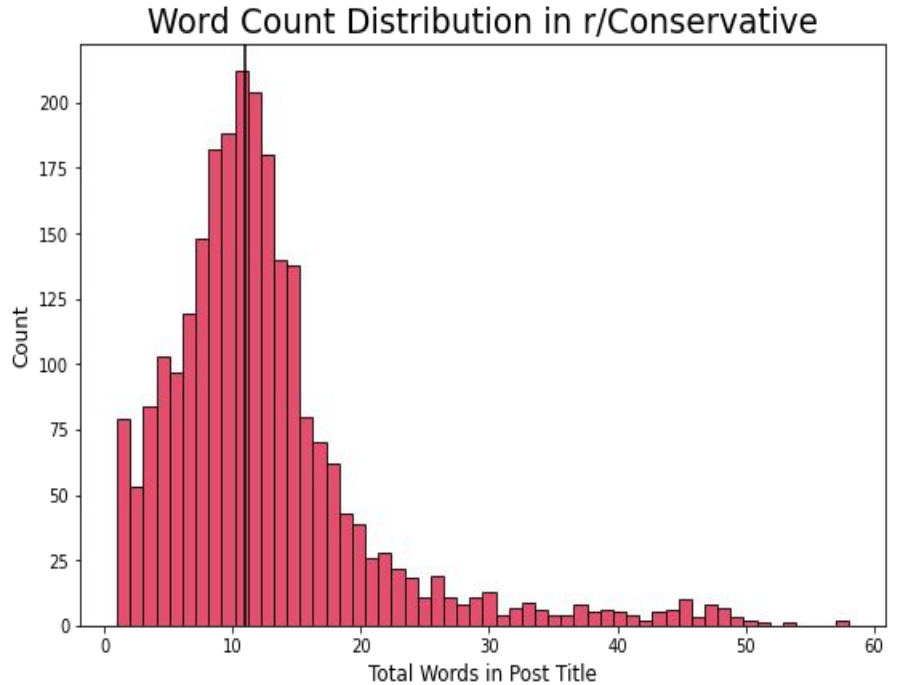
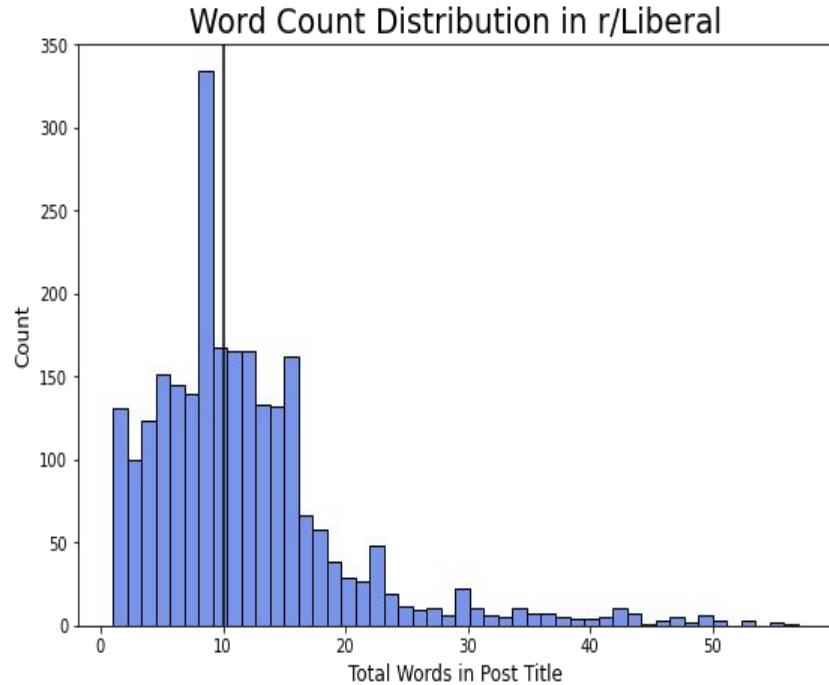
When given a dataframe of posts from 2 different subreddits, can Python use NLP to classify and differentiate between them?

Data Used:

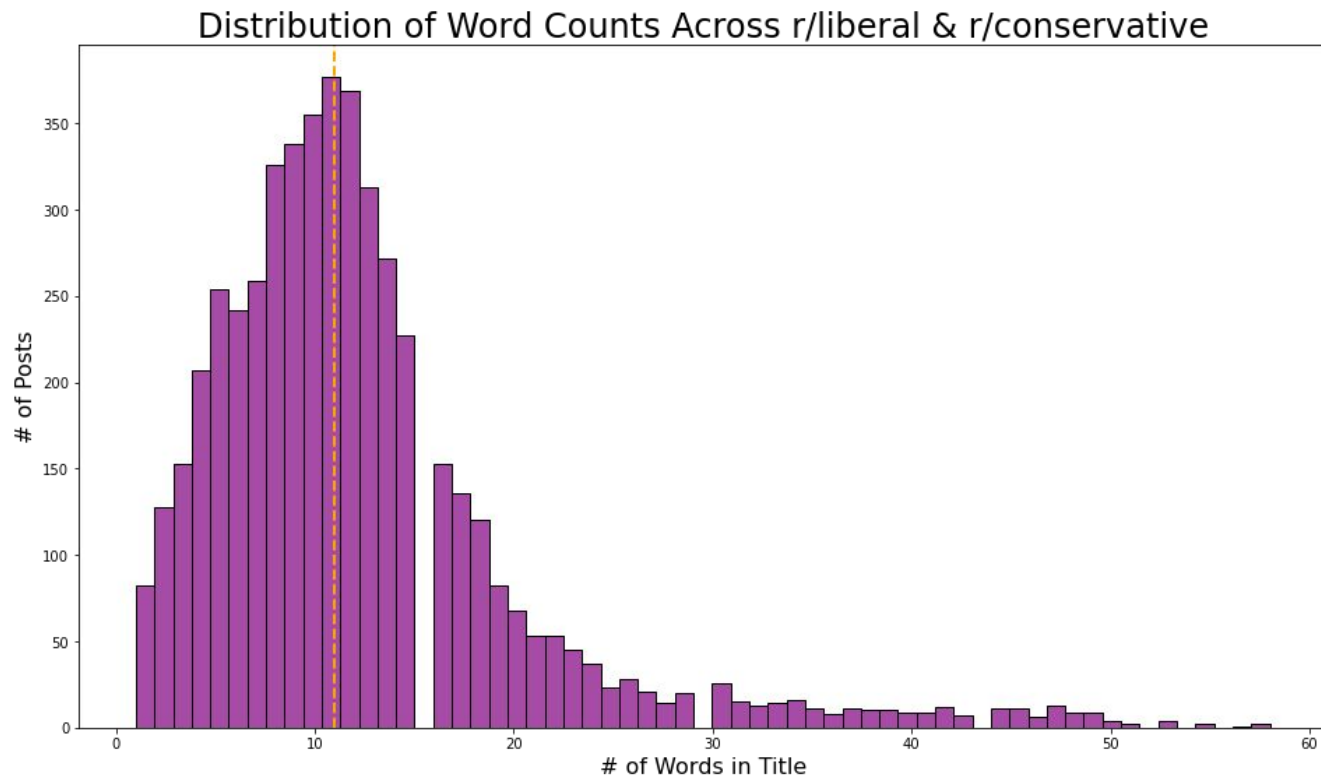
2500 posts taken from r/liberal

2500 posts taken from r/conservative

# Word Counts:



# Word Counts:



# Repost Alert!

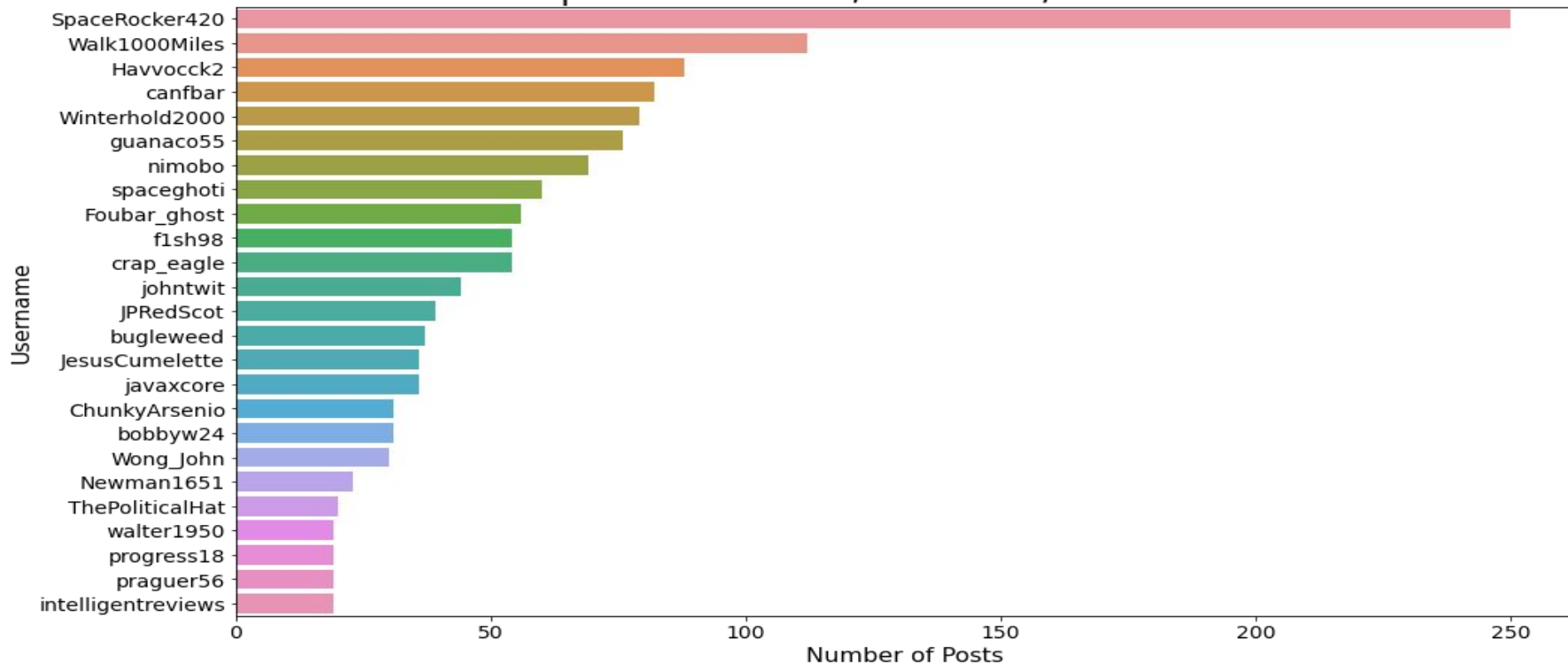
- Aggregated the number of reposts in both subreddits:
  - 94 total duplicated posts when extras were dropped
  - 61 came from r/conservative, 33 from r/liberal
  - Returns about 65% to 35% from one to another

# Unique Users:

- Total unique users from both subreddits:
  - 2,216 unique users
  - 44.32% of total user base from collected posts is unique
- 834 Unique Users in r/conservative
- 1,398 Unique Users in r/liberal
- 17 Users who post in both r/liberal and r/conservative

# User Posts:

Top 25 Posters in r/liberal & r/conservative



# Top Words:

r/liberal

<b>trump</b>	560
<b>biden</b>	258
<b>election</b>	145
<b>capitol</b>	110
<b>liberal</b>	88
<b>president</b>	75
<b>republicans</b>	71
<b>vote</b>	68
<b>people</b>	68
<b>joe</b>	67

r/conservative

<b>biden</b>	375
<b>trump</b>	182
<b>covid</b>	157
<b>cuomo</b>	113
<b>new</b>	107
<b>house</b>	107
<b>says</b>	102
<b>white</b>	100
<b>dr</b>	86
<b>texas</b>	86

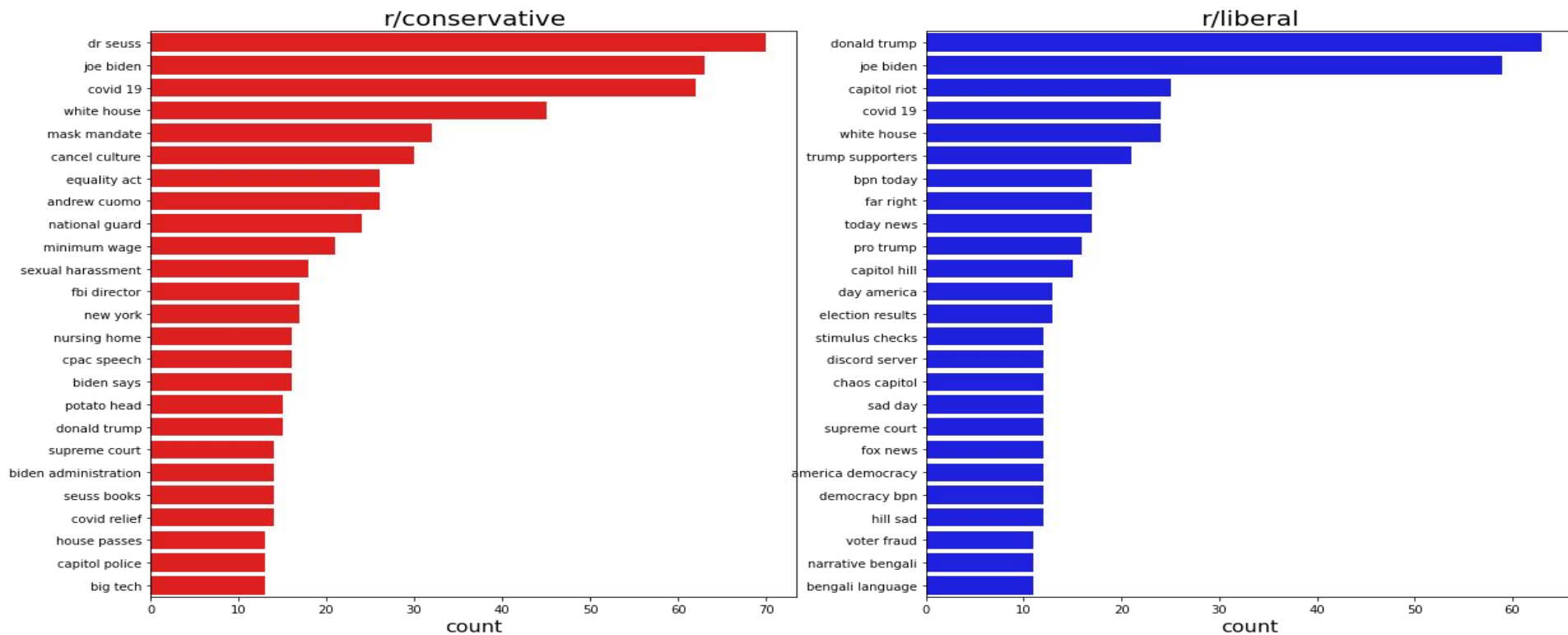
Total dataframe

<b>trump</b>	742
<b>biden</b>	633
<b>covid</b>	210
<b>election</b>	199
<b>capitol</b>	180
<b>new</b>	166
<b>says</b>	164
<b>white</b>	154
<b>house</b>	151
<b>joe</b>	146



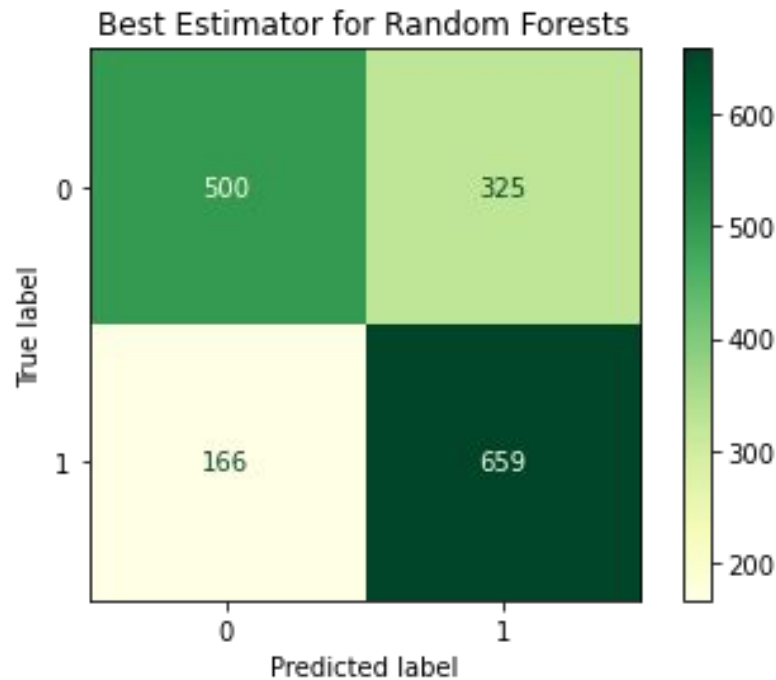
# Bipartisan Bigrams:

Top 25 Bigrams from Each Subreddit



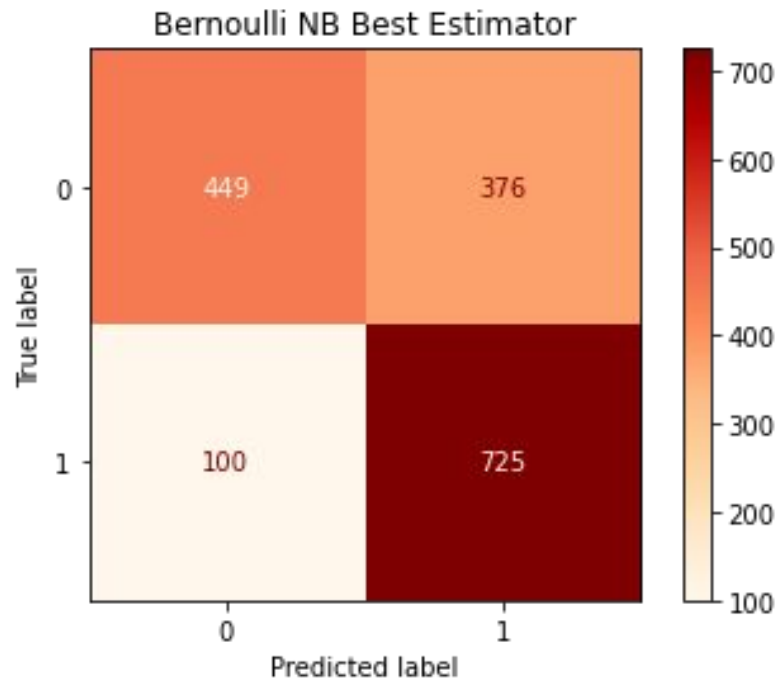
# Modeling: Random Forest

- Scoring:
  - Training: **99.76% accuracy**
  - Testing: **70.24% accuracy**
- Steps:
  - Count Vectorizer
  - Random Forest
- Best Parameters:
  - N-gram range : (1,1)
  - Stop words: English
  - # Estimators: 200



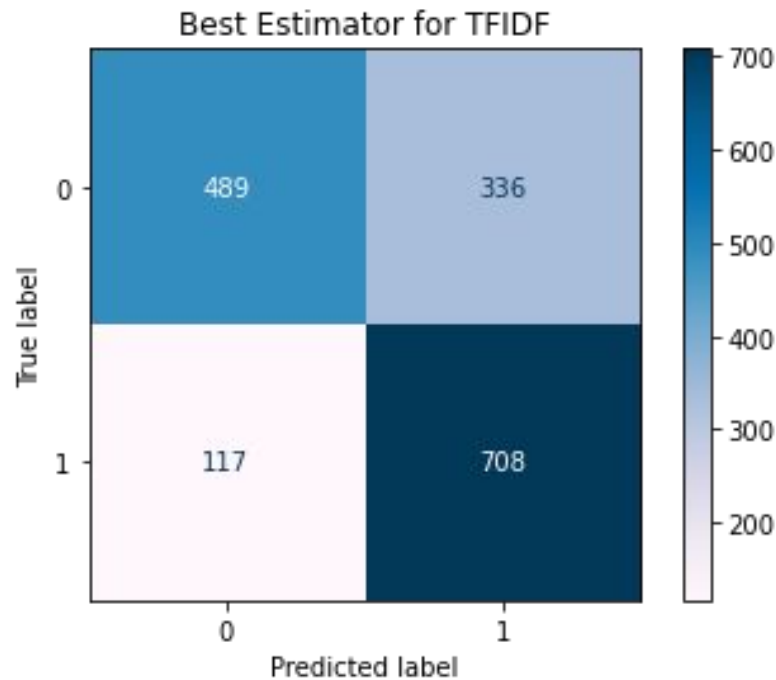
# Model: Bernoulli NB w/ Count Vectorizer

- Scoring:
  - Training: **83.40% accuracy**
  - Testing: **71.15% accuracy**
- Steps:
  - Count Vectorizer
  - Bernoulli Naive Bayes
- Best Parameters:
  - Min DF: 2
  - Max DF: 95%
  - Max Features: 4000
  - N-Gram Range: (1,2)



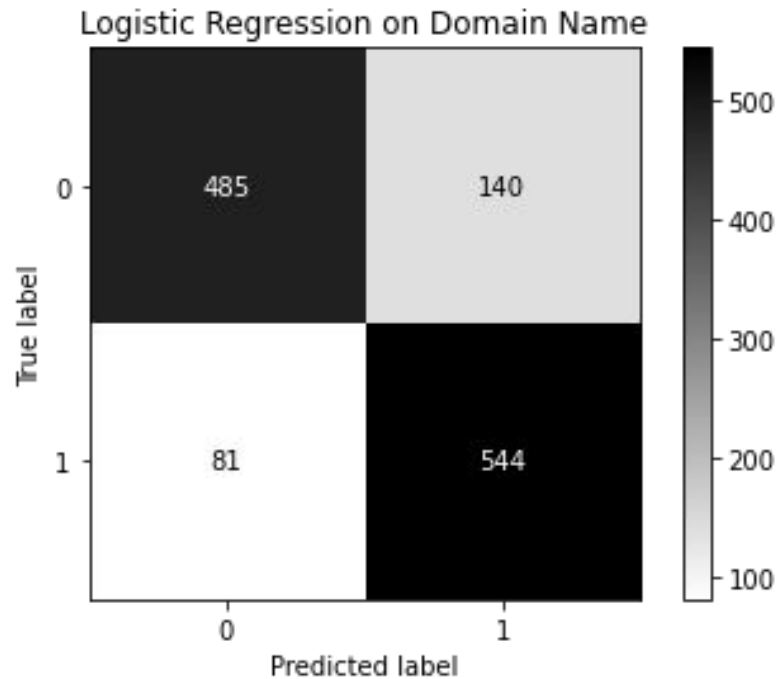
# Bernoulli NB w/ TFIDF

- Best Model
- Scoring:
  - Training: **88.14% accuracy**
  - Testing: **72.55% accuracy**
- Steps:
  - TFIDF
  - Bernoulli Naive Bayes
- Best Parameters:
  - Max Features: 6500
  - Stop words: None
  - N-Gram Range: (1,1)



# Model: Logistic Regression on Domain Name

- For Funsies
- Modeled on Domain name rather than text
- Scoring:
  - Training: **87.28% accuracy**
  - Testing: **82.48% accuracy**
- Best Testing Score of all Models



# Conclusions & Recommendations:

- More difficult to predict than anticipated
- Try more parameter tuning and other models
- Polynomial Features