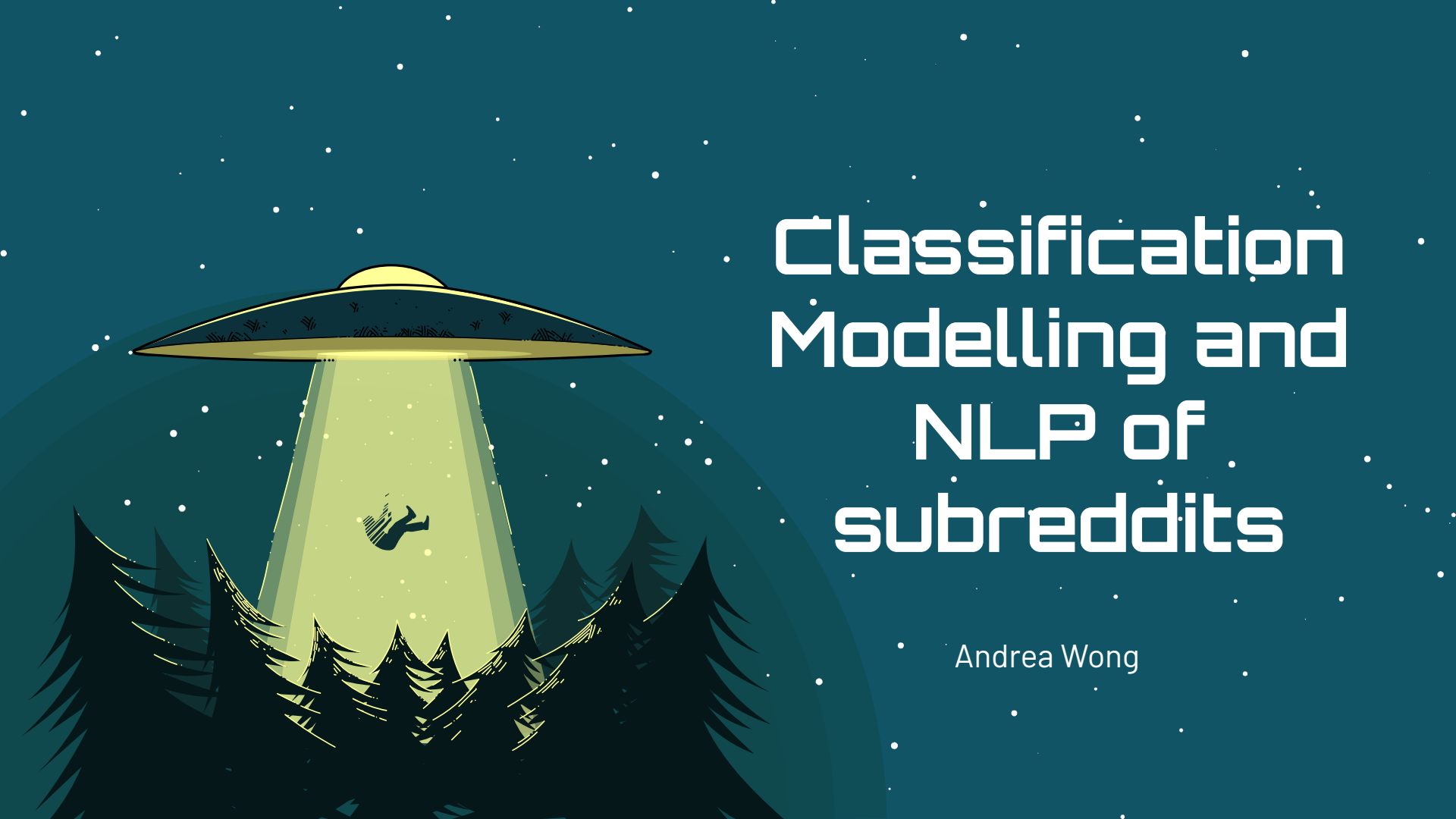


# Classification Modelling and NLP of subreddits

Andrea Wong



# TABLE OF CONTENTS

01

## INTRODUCTION

Description of chosen 2 subreddits, goals of project

02

## Pre-Process

EDA, cleaning data, transforming data



03

## Modelling

Build, train & evaluate model  
Tune Hyperparameters

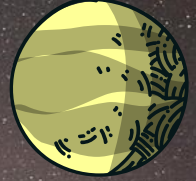
04

## CONCLUSION

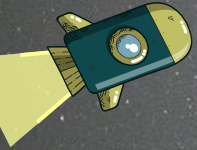
Performance measure



# Objectives



- 1) Collect posts from two subreddits (r/alien & r/conspiracytheories)
- 2) Use NLP to train a classification models to predict whether a post originated from a specific subreddit
- 3) The credibility of aliens' existence with conspiracy theories revolving around aliens



## r/alien



A community dedicated to discussion of the possibility of extraterrestrial life.

## r/conspiracytheories



A community to discuss every aspect of conspiracy theorism, from theories and current events to debunkings and popular culture.

# Data Aquisition

Aimed to collect 10,000 posts using Pushshift Reddit API



**r/alien**

9996 posts



**r/conspiracy**

9999 posts

# Cleaning and preprocessing data collected



## Stopwords

Remove stopwords,  
Added additional  
stopwords



## Links

Remove HTML links



## Emojis

Convert emojis to text



## Lowercase

Convert words to  
lowercase



## Stemming

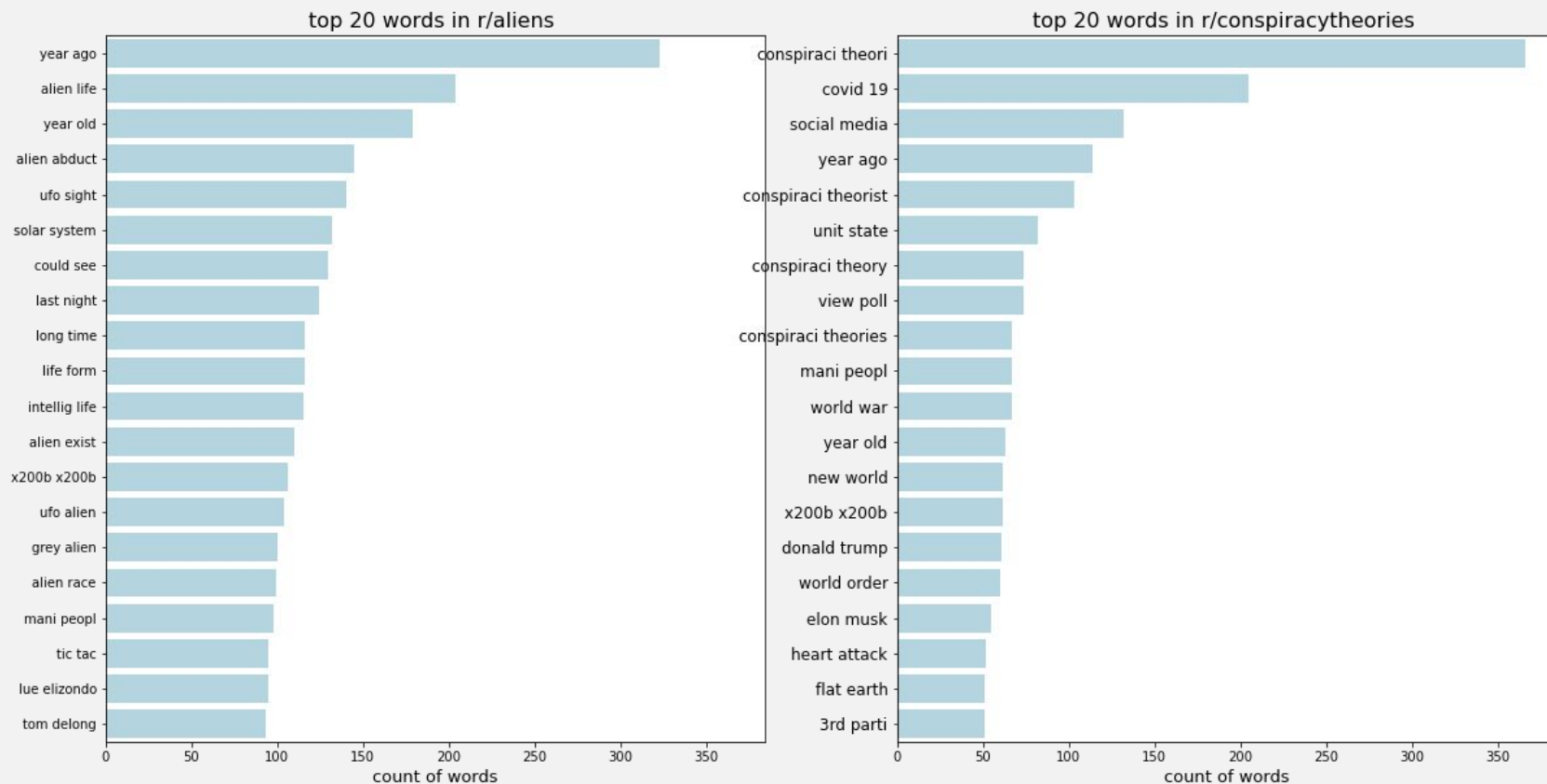
Stemming words to root  
form



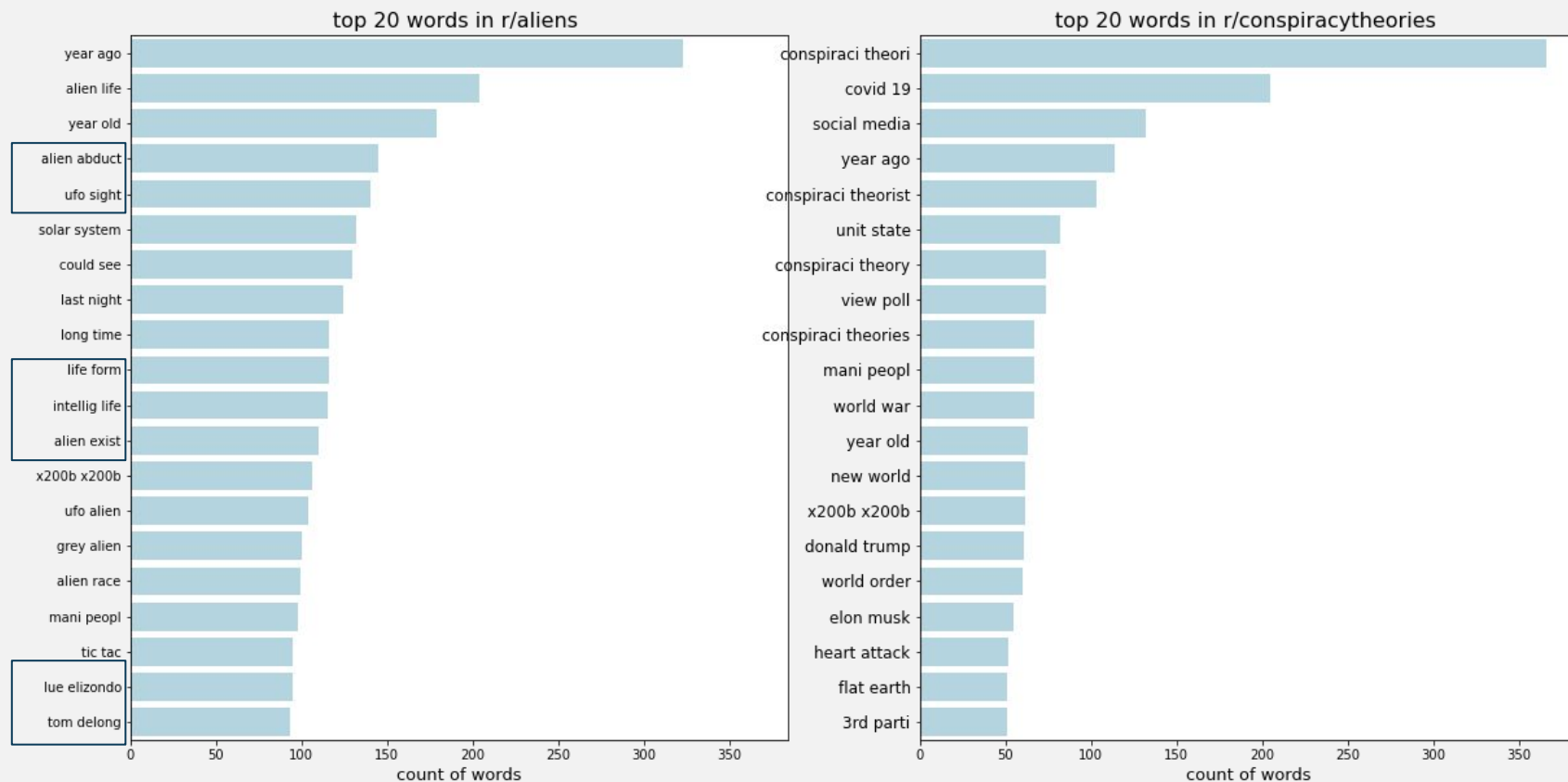
## Merged col

Merged title with selftext  
columns

# Top 20 bigrams in both subreddits

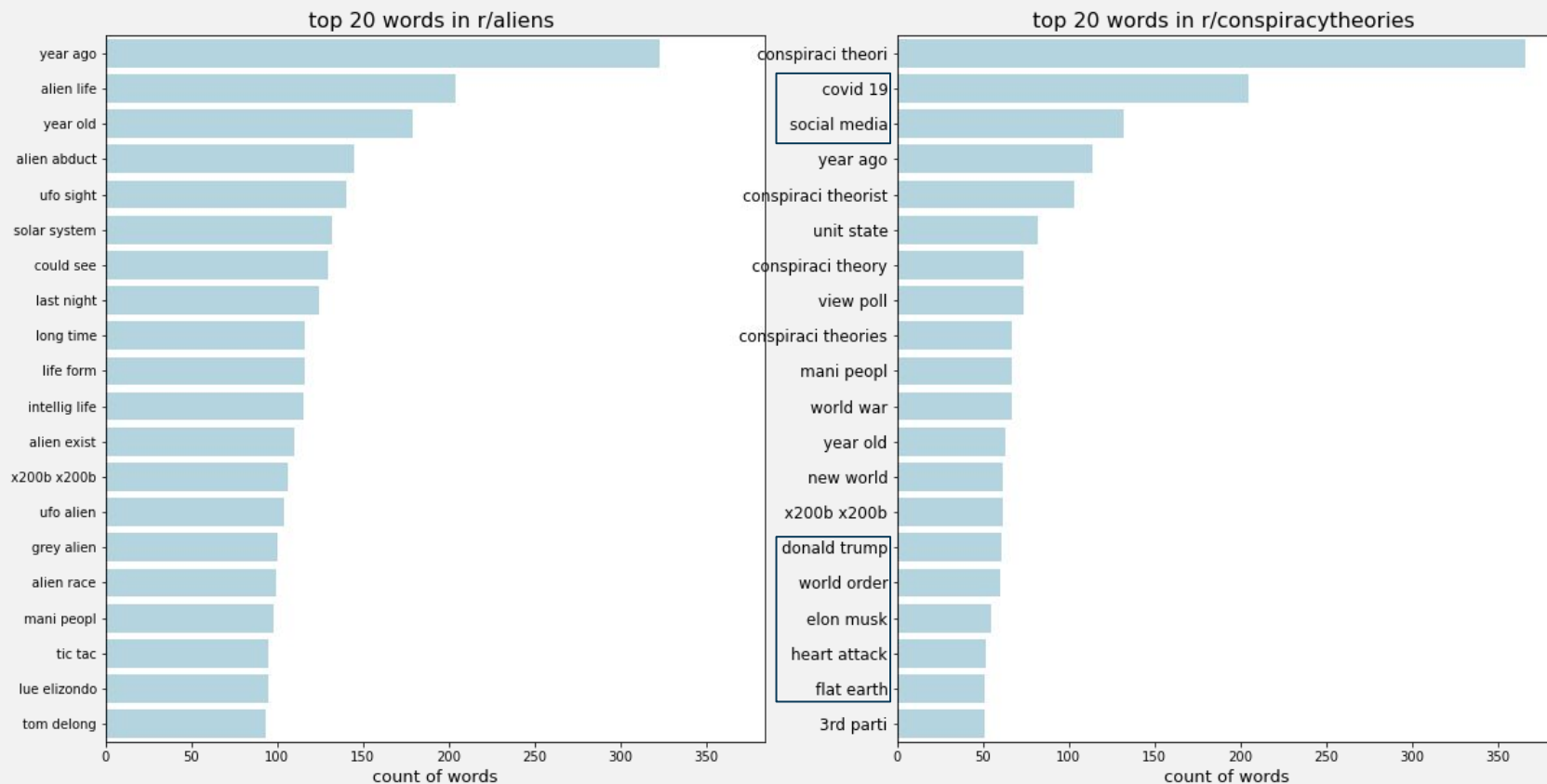


# Top 20 bigrams in both subreddits



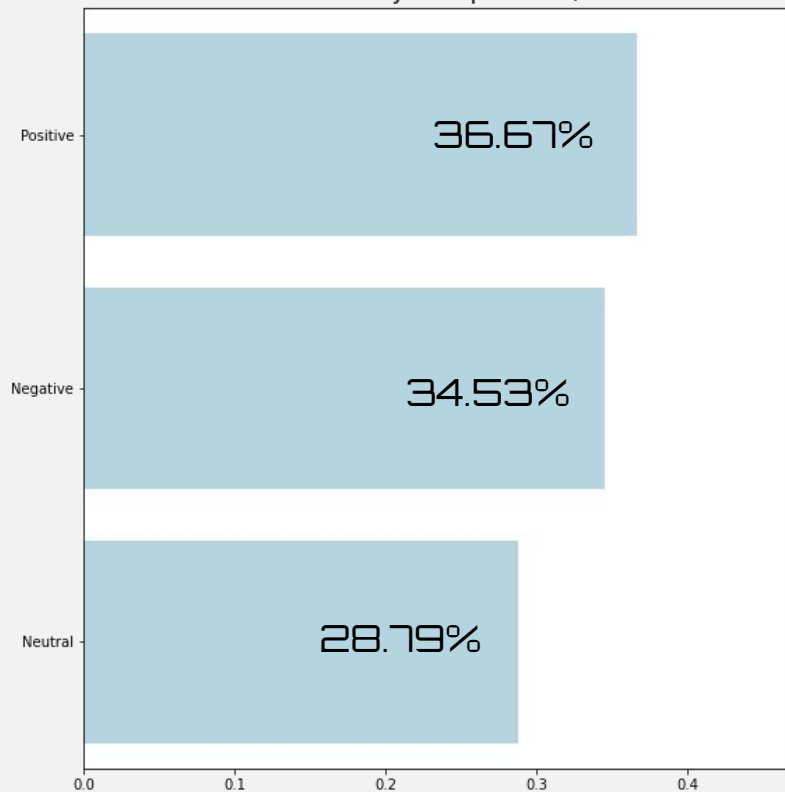


# Top 20 bigrams in both subreddits

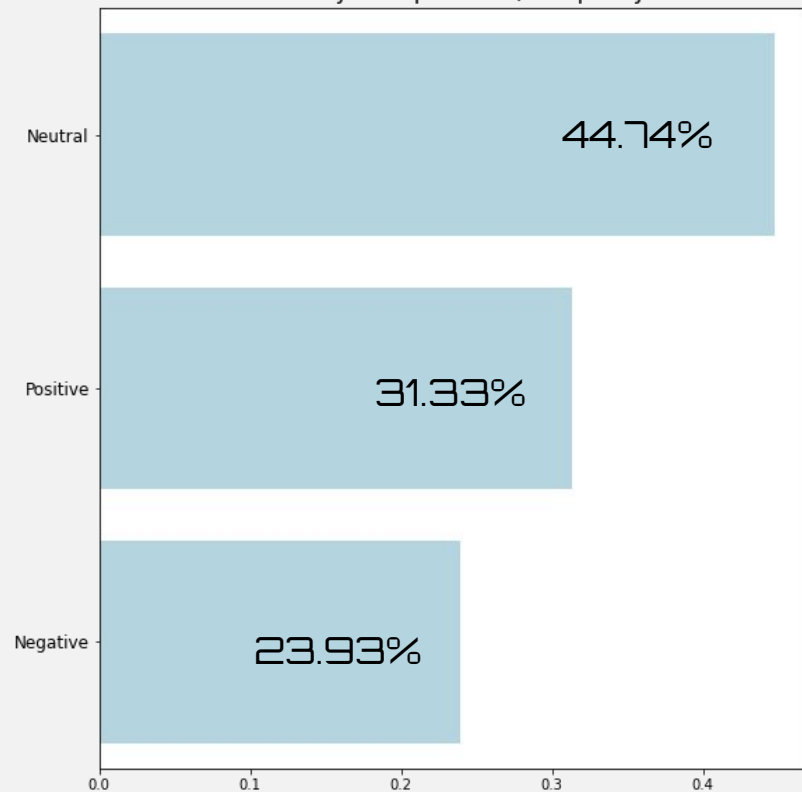


# Sentiment Analysis

Sentiment analysis of posts in r/aliens



Sentiment analysis of posts in r/conspiracytheories



# Modelling



## Train-Test-Split

Test-size = 0.2



## Classification Models

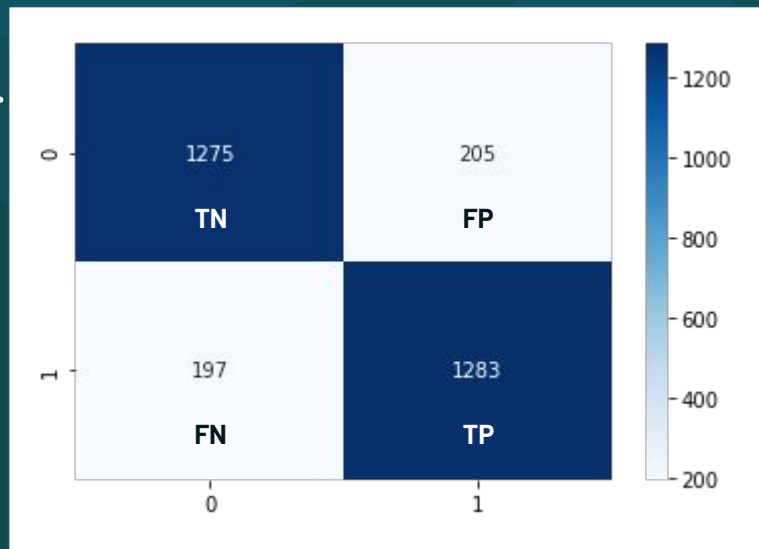
Logistic Regression  
Naive Bayes  
Random Forest Classifier



## Vectorizers

Count Vectorizer  
TF-IDF Vectorizer

# Best Performing Model



Logistic  
Regression



TF - IDF  
Vectorizer

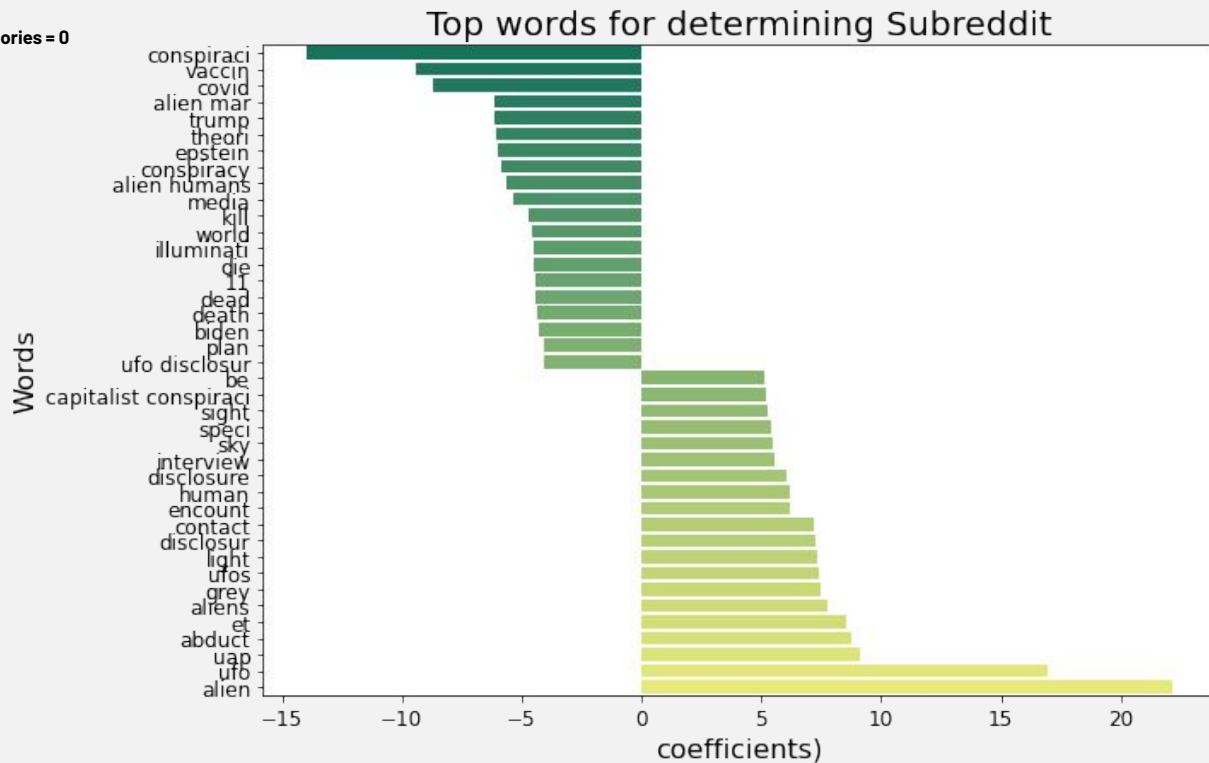


Accuracy - 0.864  
F1 Score - 0.865

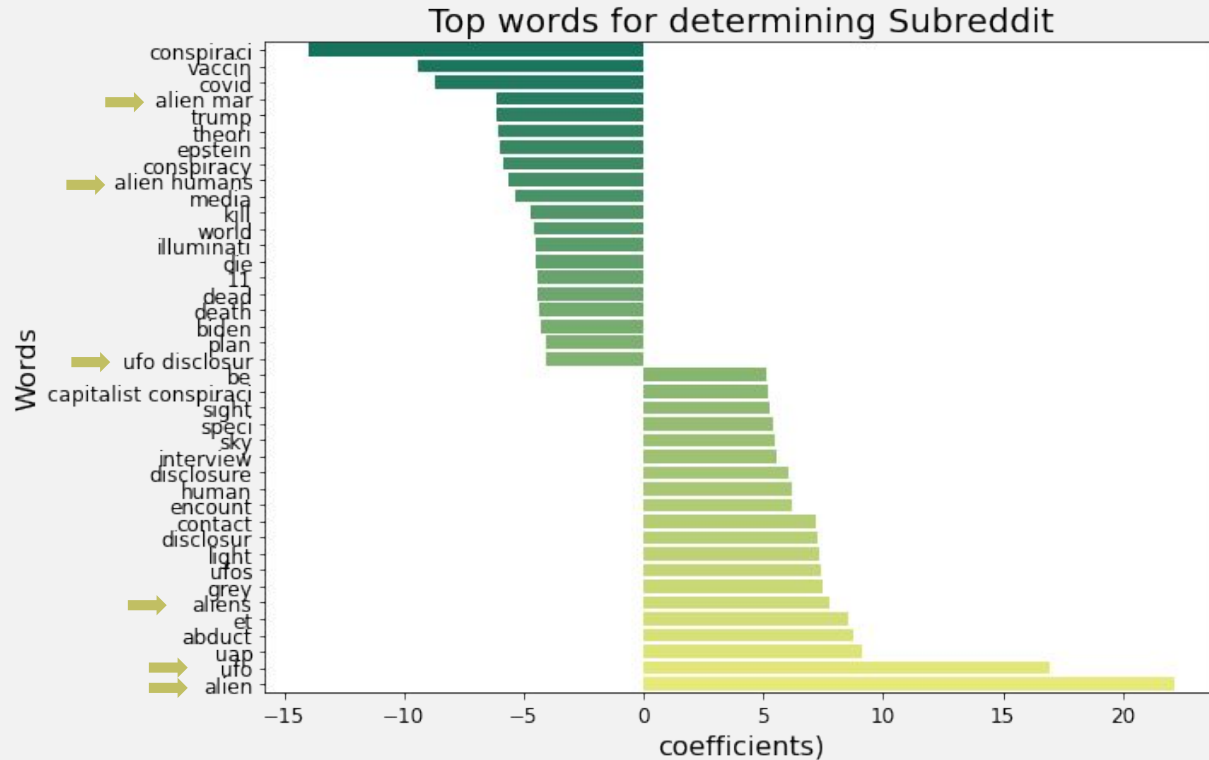
# Model Evaluation, Findings and Analysis

r/aliens = 1

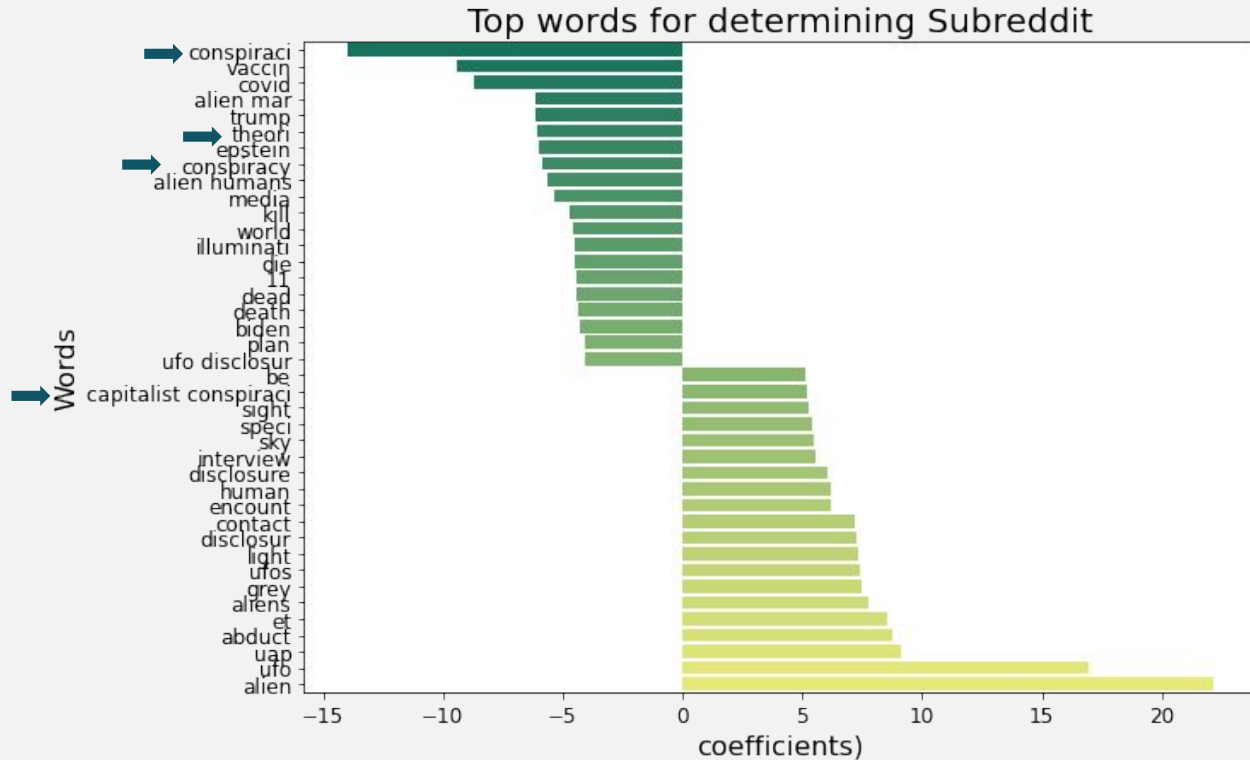
r/conspiracytheories = 0



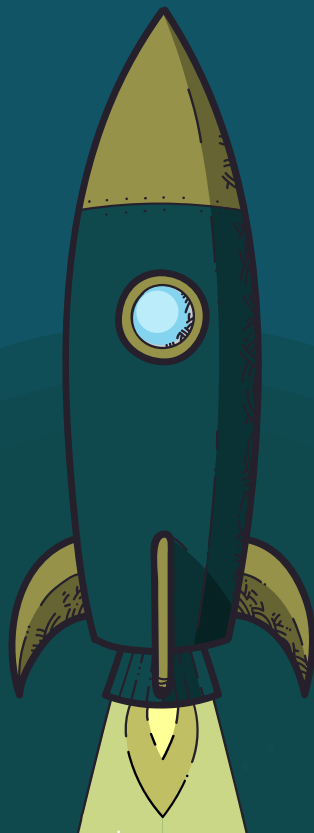
# Model Evaluation, Findings and Analysis



# Model Evaluation, Findings and Analysis



# Conclusion & Recommendations



- Relook at preprocessing methods (add more words in the list of stopwords) and other modelling methods to improve the score further