

# Subreddit Classification



r/harrypotter  
r/lotr

Boss Sarawongsuth

# PROBLEM STATEMENT

Reddit is a massive global online community forum. These communities are divided into subreddits with each being specific to a certain subject. Each subreddit usually has many moderators and admins moderating the posts to ensure that they meet the subreddit rules. With an ever-increasing userbase, the process of moderating these posts are becoming more cumbersome requiring greater amount of manpower to ensure the posts are appropriate for the subreddit.



This project aims to create a classification model that is able to classify whether a post belongs in either Harry Potter or Lord of the Rings subreddit by primarily looking at the title and post content.

- [r/harrypotter](#)
- [r/lotr](#)



# SUBREDDITS

**HARRY POTTER**

**r/harrypotter**

**993k subscribers**

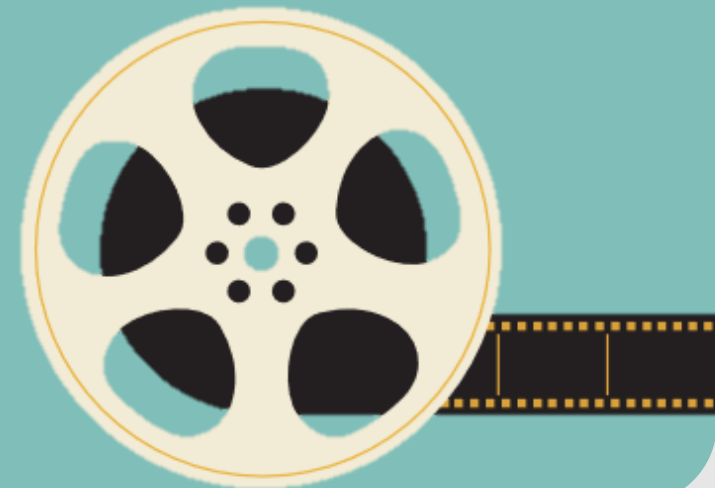
**57%**

**LORD OF THE RINGS**

**r/lotr**

**523k subscribers**

**43%**



# DATASETS

The datasets for r/harrypotter and r/lotr were obtained from  
Reddit API.

**1723 rows**  
**&**  
**19 columns**

# DATA TYPES

**6**

**STRING**

Features

---

title  
selftext  
subreddit

...

**3**

**BOOL**

Features

---

spoiler  
author\_premium  
is\_video

**3**

**FLOAT**

Features

---

upvote\_ratio  
num\_reports  
...

**7**

**INT**

Features

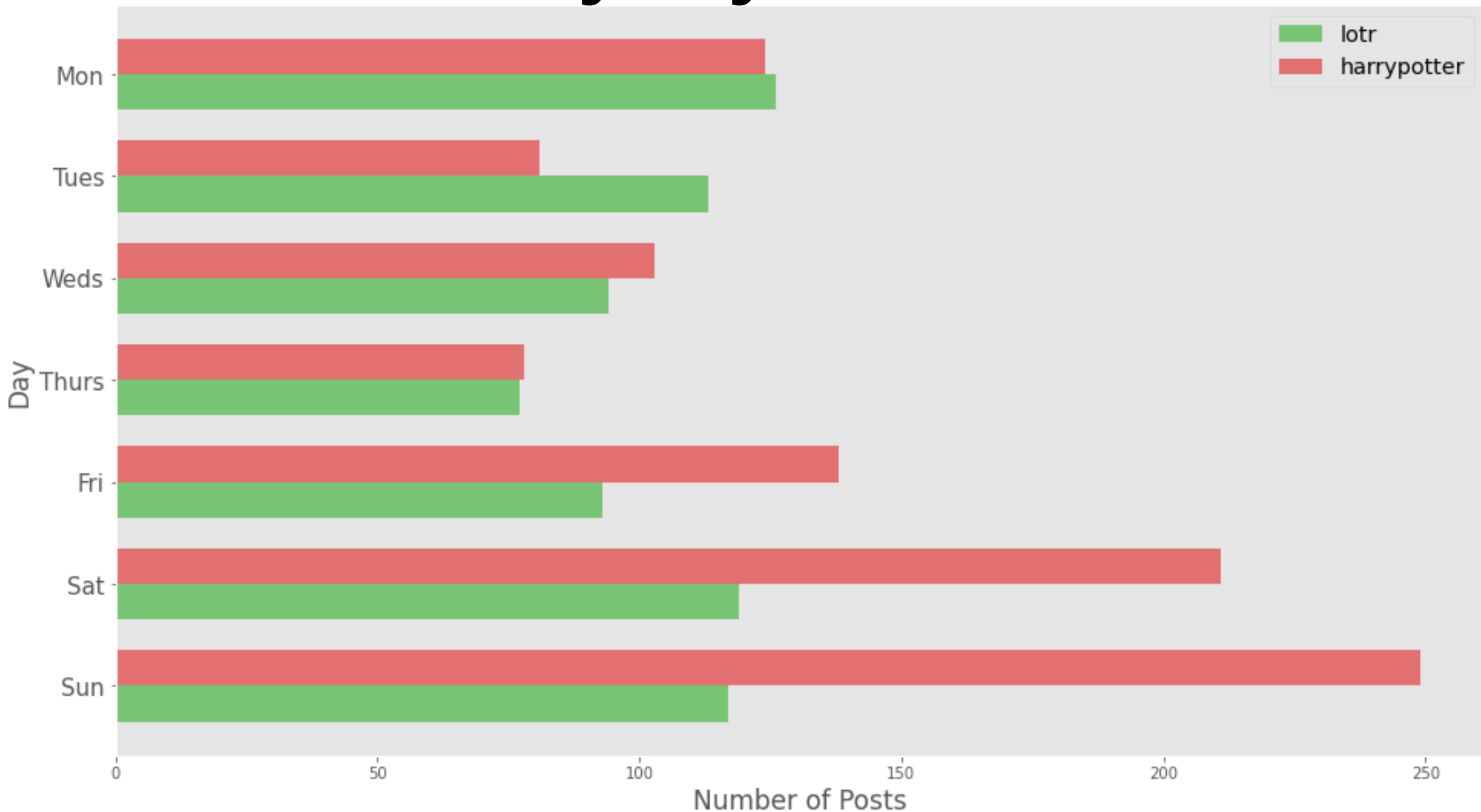
---

ups  
downs  
num\_comments  
...

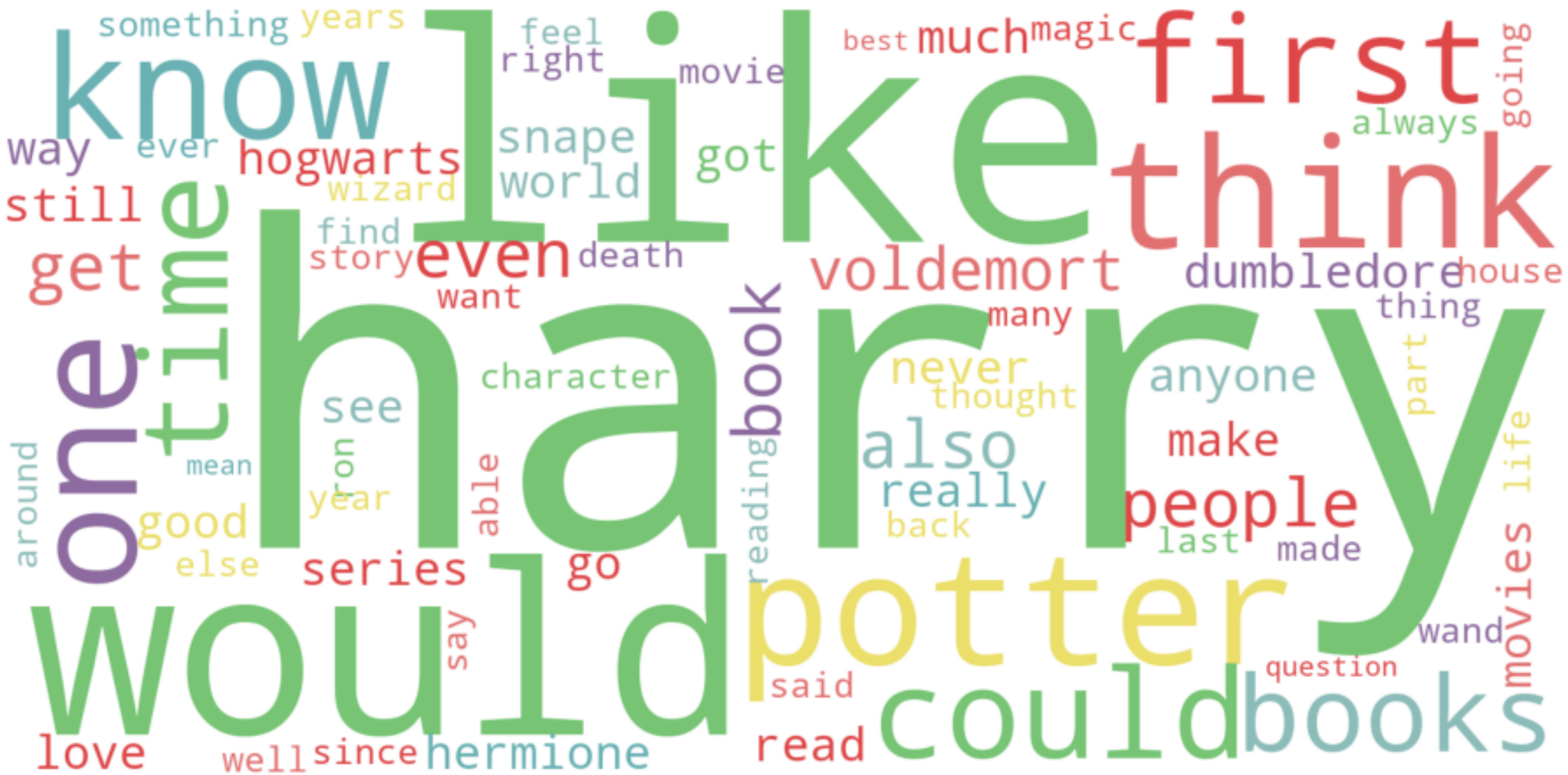


**EDA**

# Posts by Day & Subreddit

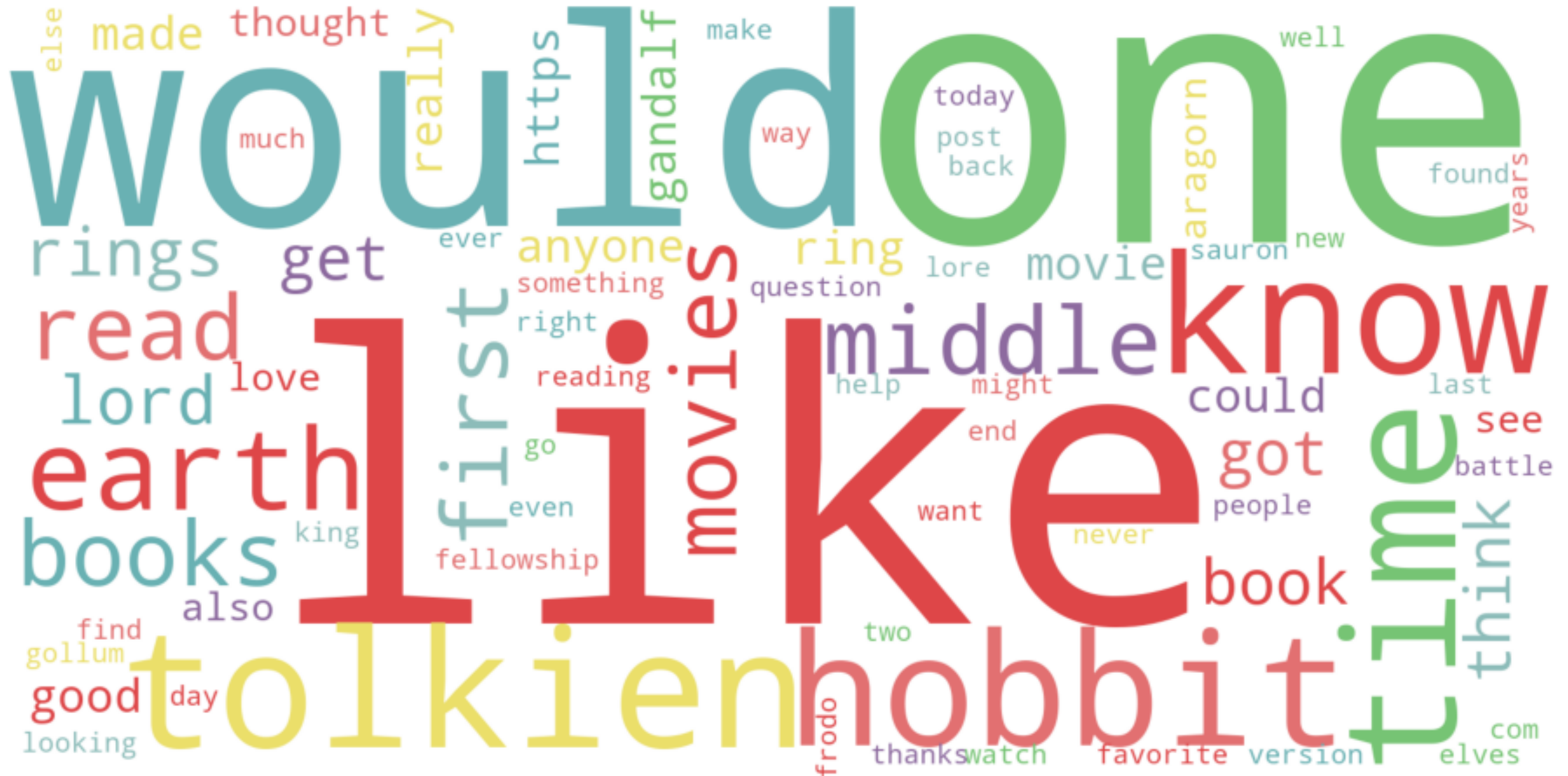


# HARRY POTTER

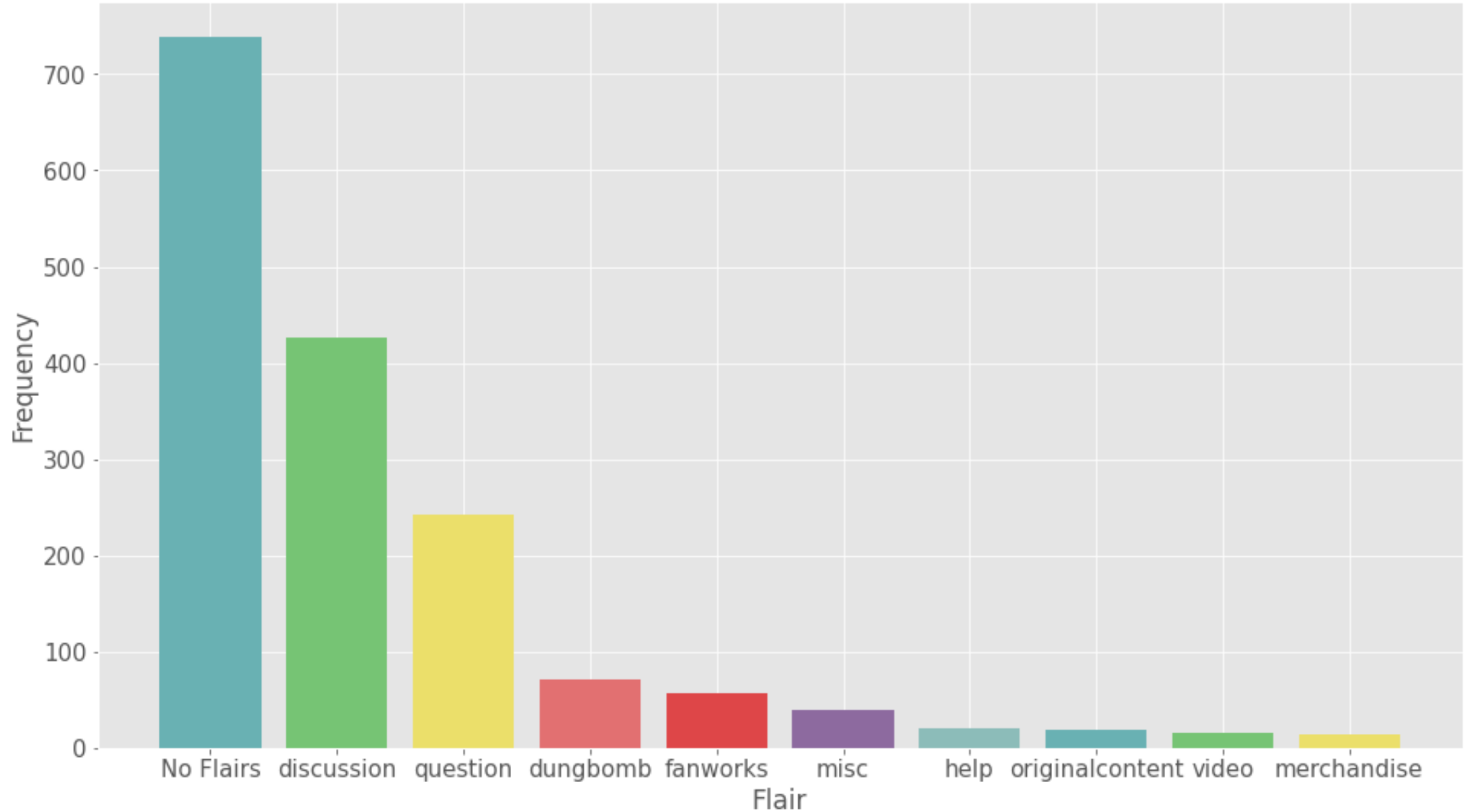




# LORD OF THE RINGS



# Reddit Flairs





# **CLASSIFICATION MODELS**

# CLASSIFICATION MODELS

## Logistic Regression

---

- 'logreg\_\_C': 50
- 'logreg\_\_penalty': 'l2'
- 'tvec\_\_max\_df': 0.9
- 'tvec\_\_max\_features': None
- 'tvec\_\_min\_df': 1
- 'tvec\_\_ngram\_range': (1, 1)

## kNN

---

- 'knn\_\_metric': 'minkowski'
- 'knn\_\_n\_neighbors': 15
- 'tvec\_\_max\_df': 0.9
- 'tvec\_\_max\_features': None
- 'tvec\_\_min\_df': 1
- 'tvec\_\_ngram\_range': (1, 2)

## Naïve Bayes

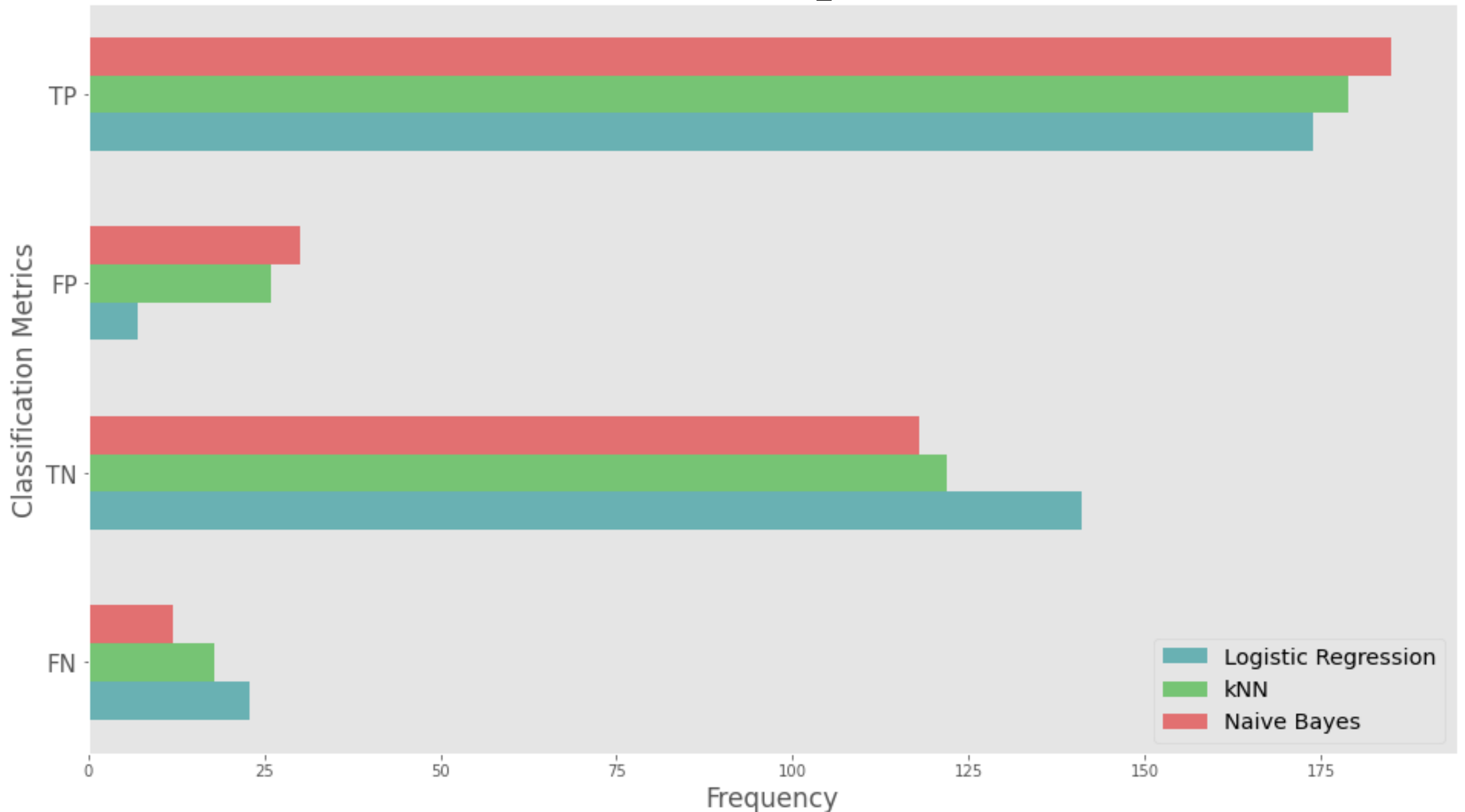
---

- 'tvec\_\_max\_df': 0.9,
- 'tvec\_\_max\_features': None,
- 'tvec\_\_min\_df': 3,
- 'tvec\_\_ngram\_range': (1, 2)

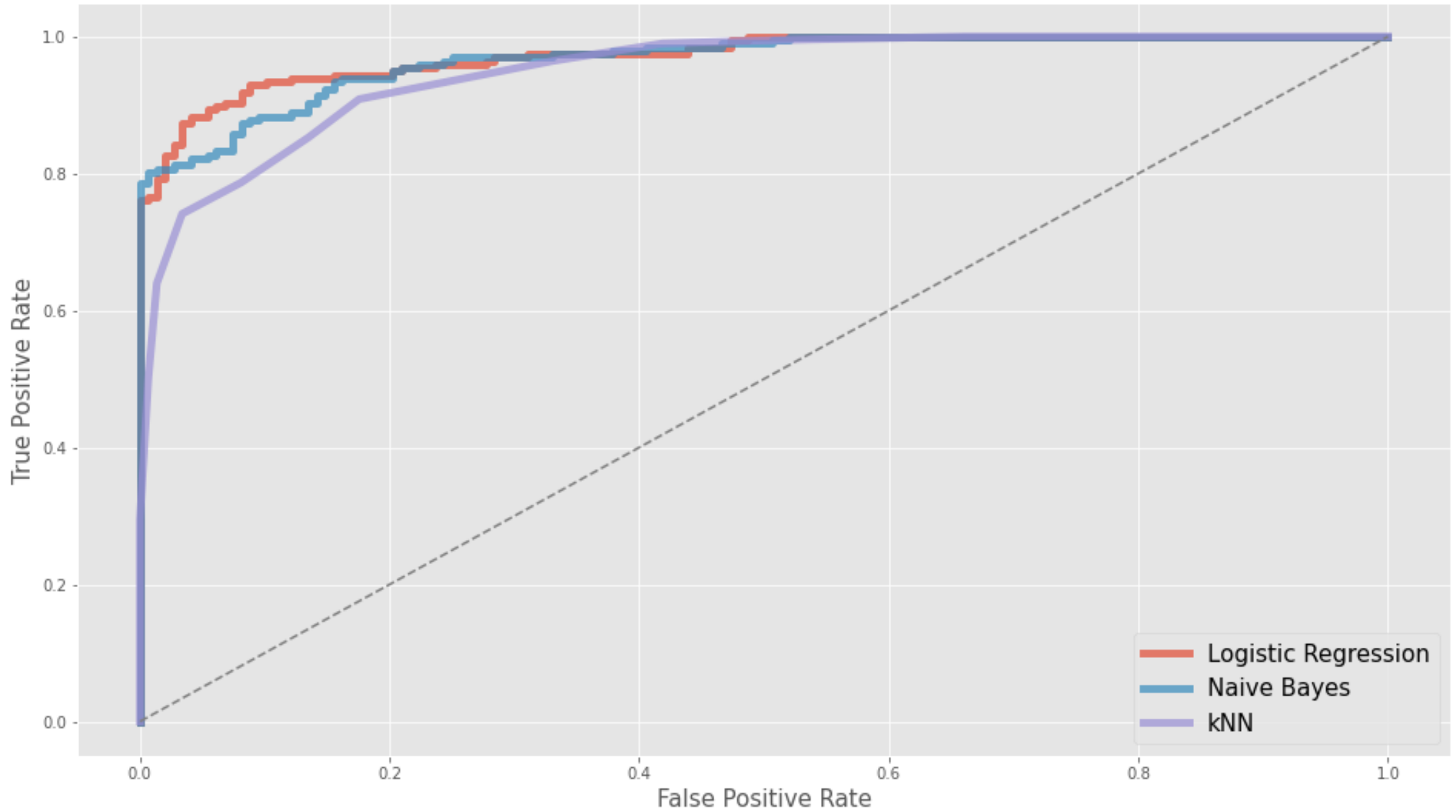


# **EVALUATION**

# Model Comparison



# ROC Curve



# Evaluation Metrics

Metrics	Logistic Regression	kNN	Naive Bayes
Accuracy (Train)	1	0.89	0.95
Accuracy (Test)	0.91	0.87	0.88
Misclassification Rate (Train)	0	0.11	0.05
Misclassification Rate (Test)	0.09	0.13	0.12
Sensitivity (Train)	1	0.95	0.98
Sensitivity (Test)	0.88	0.91	0.94
Specificity (Train)	1	0.80	0.91
Specificity (Test)	0.95	0.82	0.80
Precision (Train)	1	0.87	0.94
Precision (Test)	0.96	0.87	0.86
F1 (Train)	1	0.91	0.96
F1 (Test)	0.92	0.89	0.90
AUC	0.972	0.948	0.967





**INFERENCE**

# WORDS & SUBREDDIT LIKELIHOOD

---

## r/harrypotter

**‘harry’ : 4230000x**

**‘potter’ : 27000x**

**‘wizard’ : 5200x**

**‘hogwart’ : 2000x**

**‘weasley’ : 1100x**

## r/lotr

**‘today’ : 0.00271x**

**‘gandalf’ : 0.000546x**

**‘tolkien’ : 0.000360x**

**‘hobbit’ : 0.000153x**

**‘ring’ : 0.000122x**

**\*with regards to the positive class (r/harrypotter)**



# CONCLUSION



# CONCLUSION

The project has demonstrated that it is possible to classify whether or not a post belongs in certain subreddit like r/harrypotter and r/lotr. The best performing model (Logistic) uses the texts within a post to do the classification and was shown to be quite effective.

However, if it were to be implemented in production build, it would need to use the one-vs-all approach to classify if the post belongs in the subreddit or not instead of classifying whether or not a post belongs in r/harrypotter or r/lotr.

The main limitation of this model is that it only uses word frequency to determine the subreddits. Improvements could be achieved through the use of more advanced NTL techniques that can better understand the context of the texts within the post.

