

# Reddit Classification

By Erik Mercado



reddit ads



# Problem Statement

We are running a ad compagine and want to learn what people are talking about in r/books and r/movies.

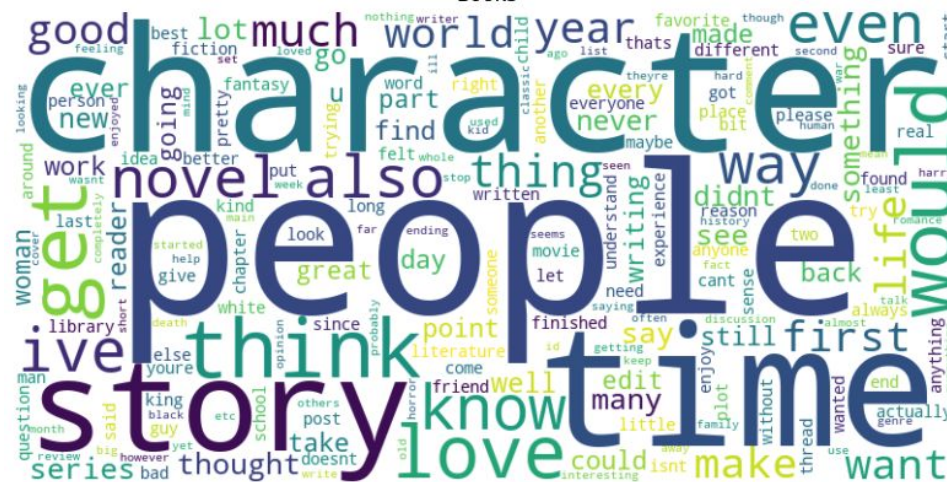
# Processes

- Scrape Reddit
  - Hot, new, top, rising, controversial
  - 1227 posts from r/Movies
  - 1179 posts from r/Books
- Clean Data
  - Remove URLs, Special Chars, and Nums, Stop
  - Remove obvious words from each subreddit
    - Book, Reading, Author, Page
    - Movie, Film, Director, Scene
- Modeling
  - Random Forest
  - SVC

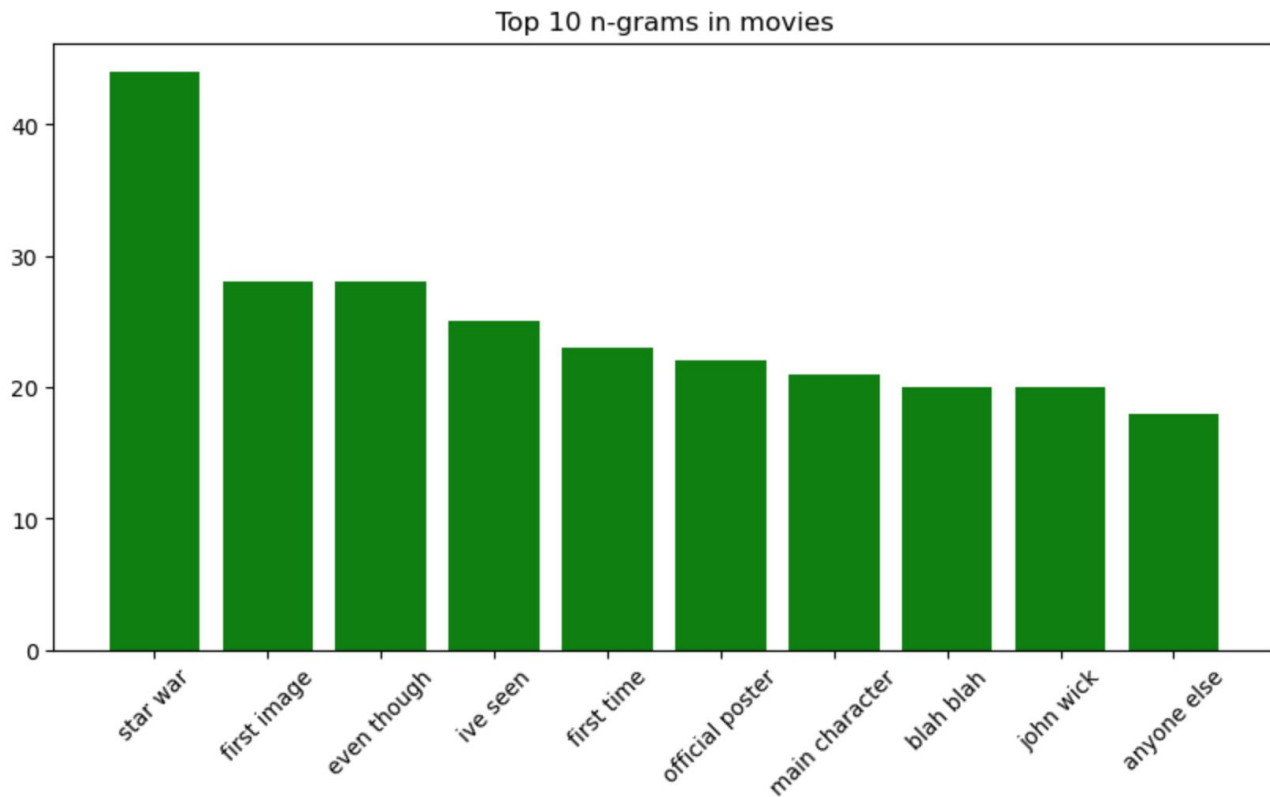
## Movies



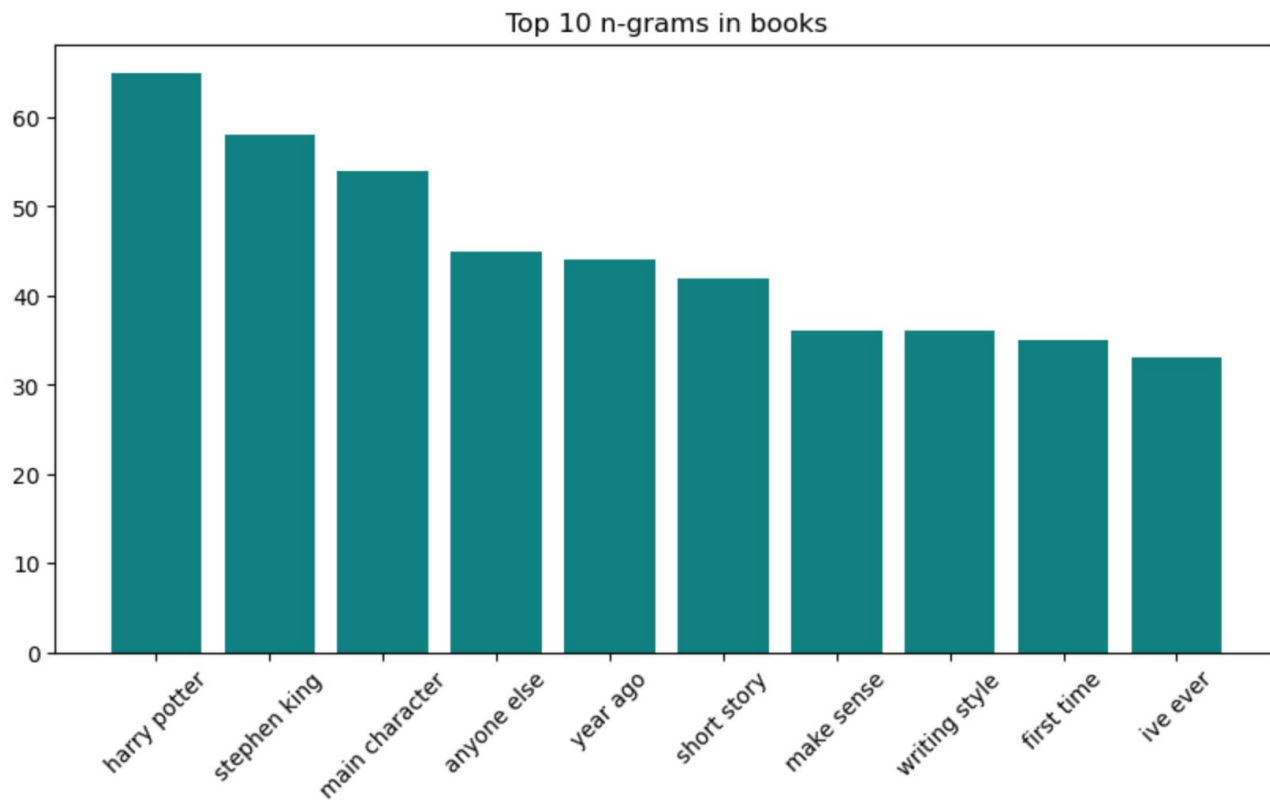
## Books



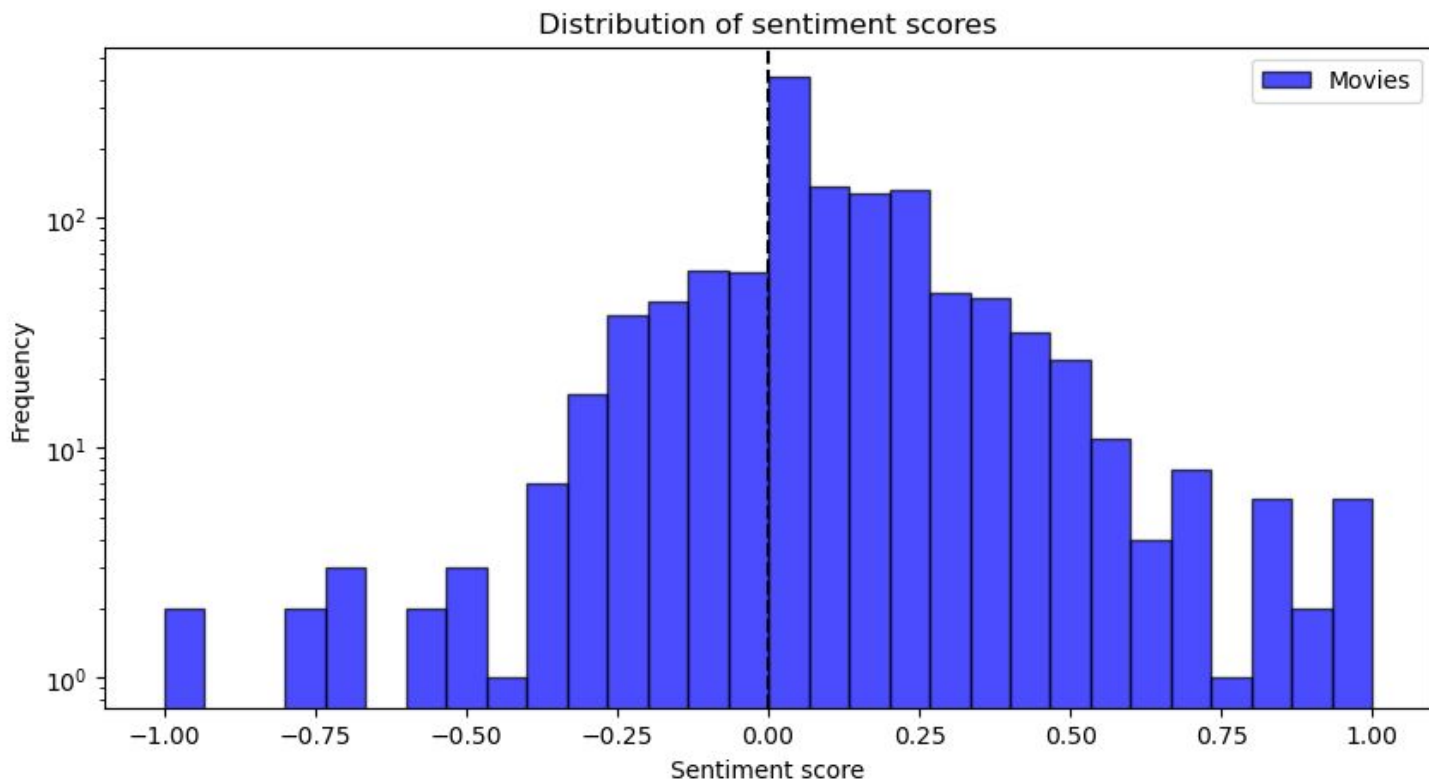
# EDA



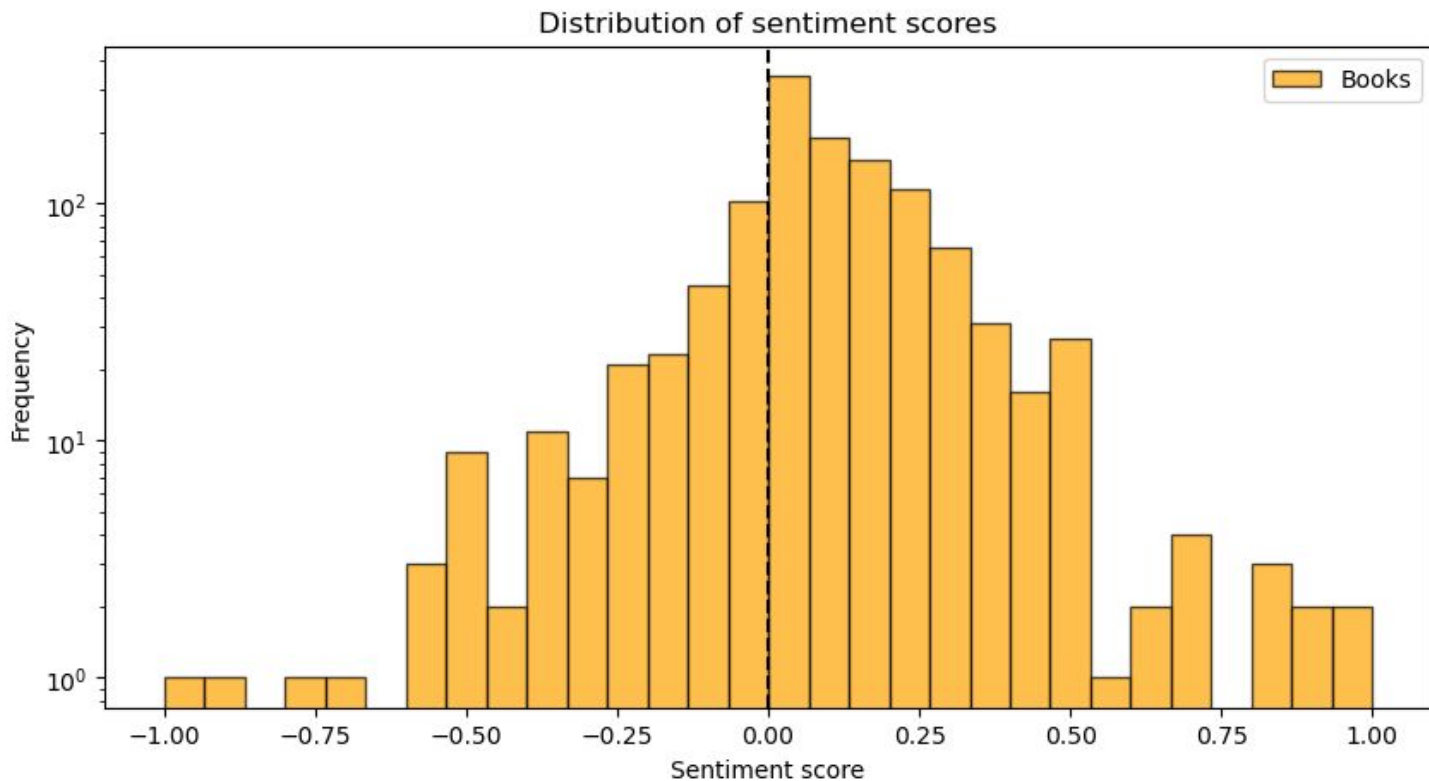
# EDA



# EDA



# EDA



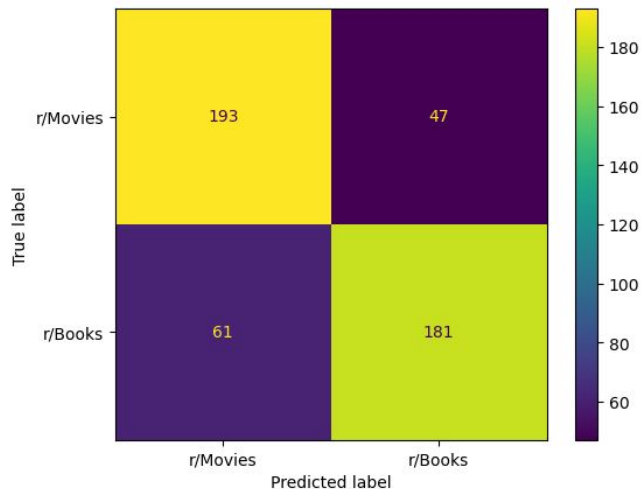


# Evaluations

## Random Forest

Best Accuracy: 0.8040449134199135

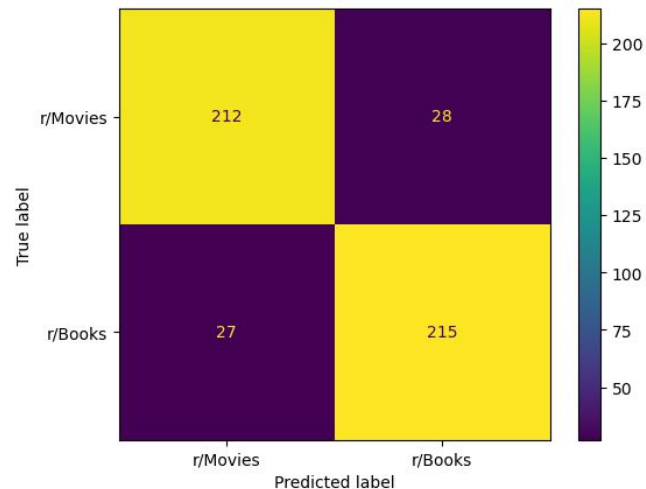
Test Accuracy: 0.7759336099585062



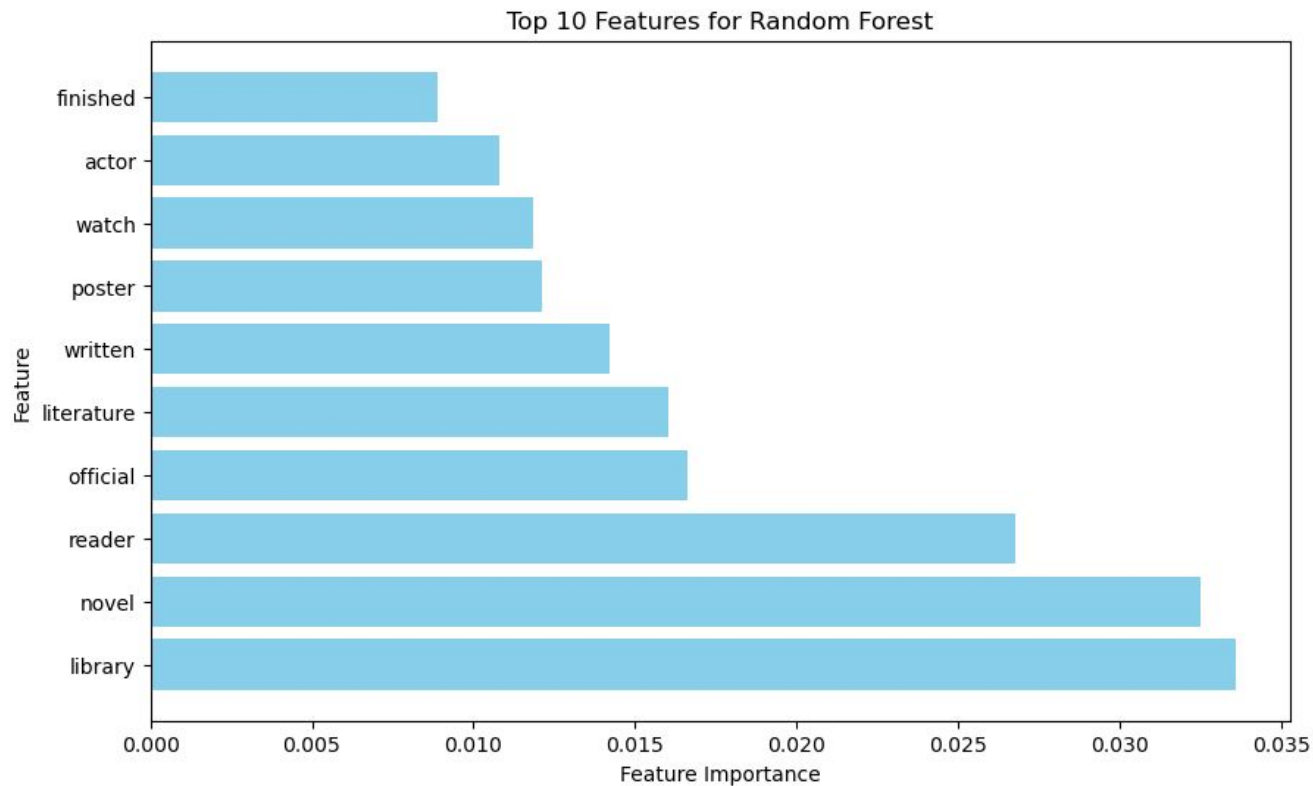
## SVC

Best Accuracy: 0.8762987012987011

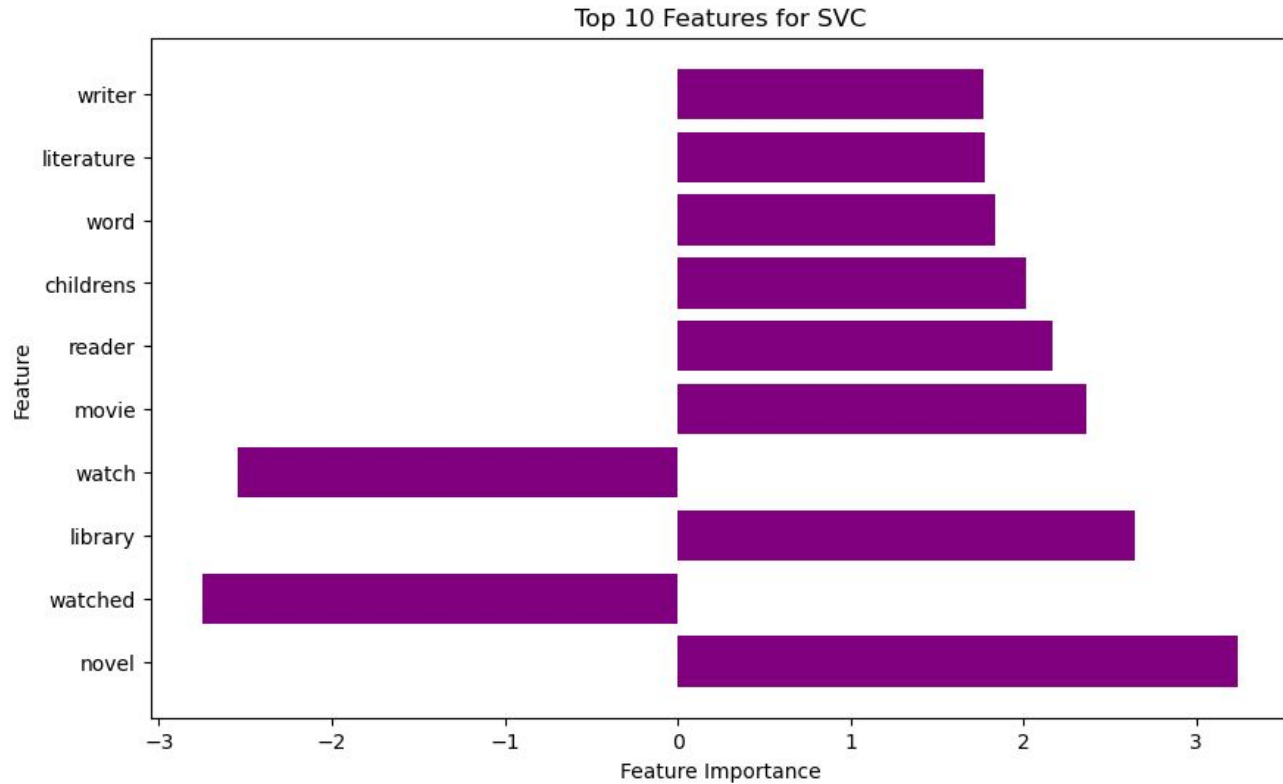
Test Accuracy: 0.8858921161825726



# Feature Importance



# Feature Importance



Thank You