

The background features two side-by-side movie posters. On the left is the poster for 'Star Trek: Discovery', showing several crew members in yellow uniforms against a purple and blue space background. On the right is the poster for 'Star Wars: The Force Awakens', featuring Rey, Finn, and BB-8 in a dynamic pose with lightsabers and a TIE fighter in the background.

Project 3 : Web APIs & NLP.

Star Wars vs. Star Trek

Author: Artem Lukinov



Problem Statement

A glitch in Reddit's system caused the Star Wars and Star Trek subreddits to drop their names and scrambled all posts. The project is tasked with analyzing how different the subreddits are and developing a model to classify posts into correct subreddits.

What was looked at?

A screenshot of the Star Wars subreddit homepage. The header features the "Star Wars" logo in blue and white. Below it is a banner with a dark background showing a Star Destroyer and several TIE fighters against a red nebula. The main title "Star Wars - A long time ago, in a galaxy far, far, away..." is displayed in large white text. Below the title is the subreddit name "r/StarWars". There are three navigation links: "Posts" (underlined in red), "Upcoming SW Media", and "The Mandalorian ▾". To the right is a red "Join" button.

STAR WARS

Star Wars - A long time ago, in a galaxy far, far, away...

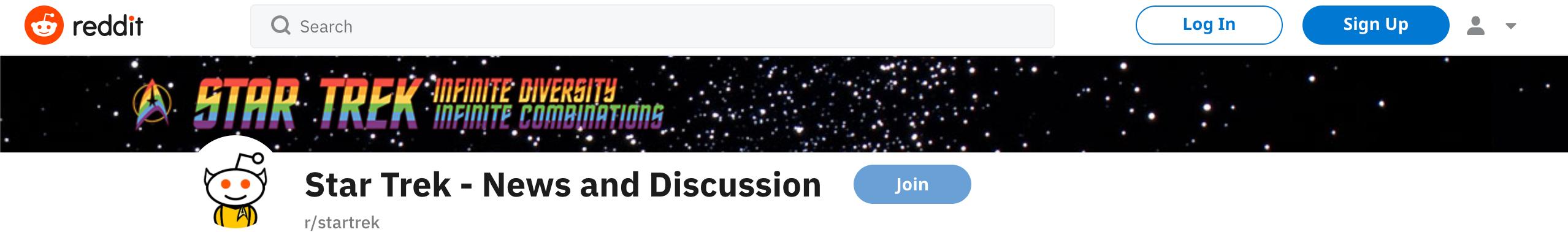
r/StarWars

Join

Posts Upcoming SW Media The Mandalorian ▾

- 1.8 million users
- Created May 27, 2008

What was looked at?



The screenshot shows the homepage of the Star Trek subreddit. At the top, there's a navigation bar with the Reddit logo, a search bar, and buttons for "Log In" and "Sign Up". Below the navigation is a banner featuring the Star Trek logo and the text "STAR TREK INFINITE DIVERSITY INFINITE COMBINATIONS". The main content area has a dark background with a starry pattern. On the left, there's a small image of the Reddit alien logo wearing a Starfleet uniform. To its right, the title "Star Trek - News and Discussion" is displayed in large, bold, black letters, followed by a "Join" button. Below the title, the subreddit name "r/startrek" is shown.

- 300,000+ users
- Created June 17, 2008



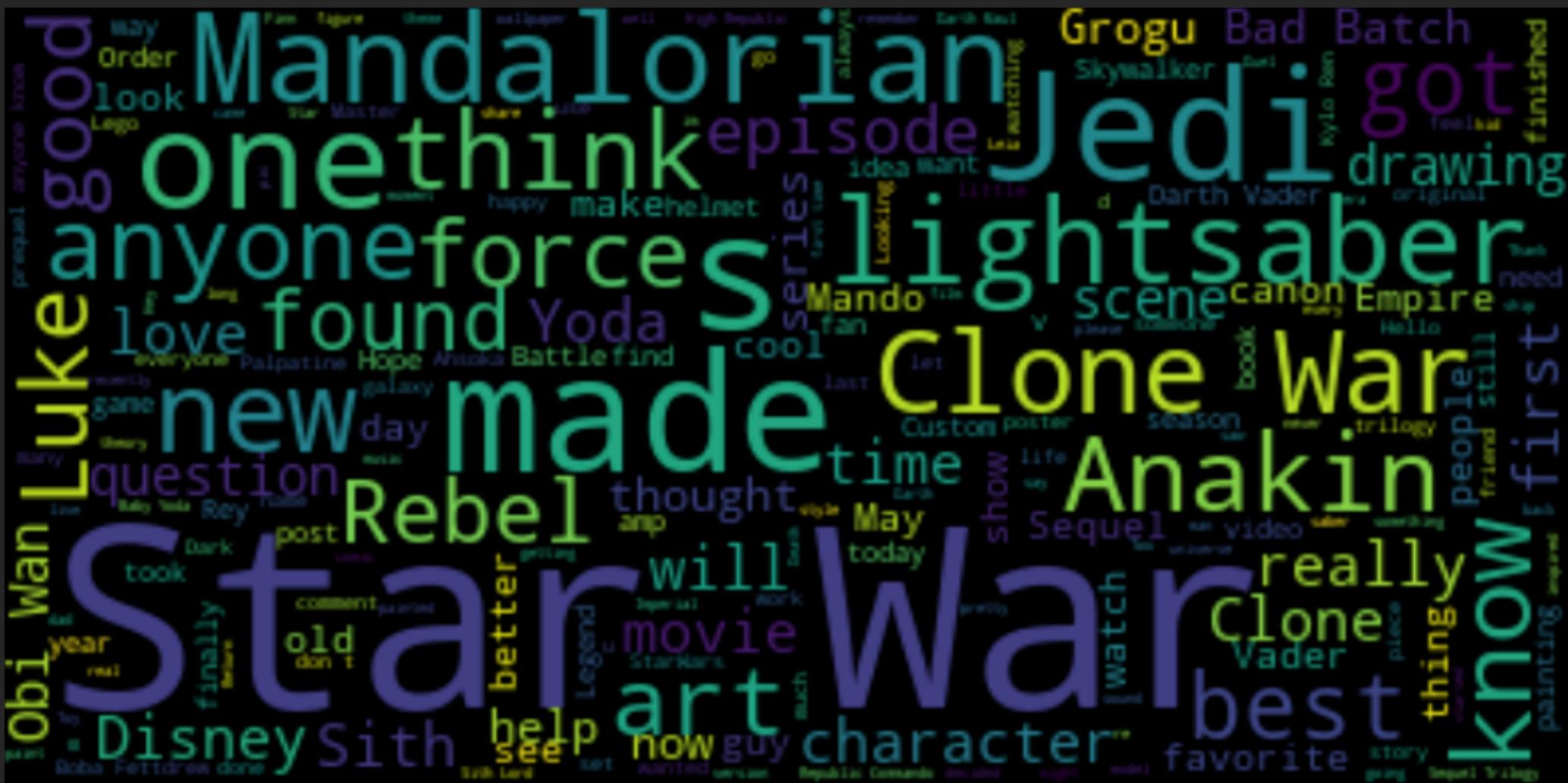
How was the data cleaned?

- two dataframes merged into one
- thorough analysis of 80+ columns
- three columns remained: title, content and subreddit's name
- rows with empty data were removed



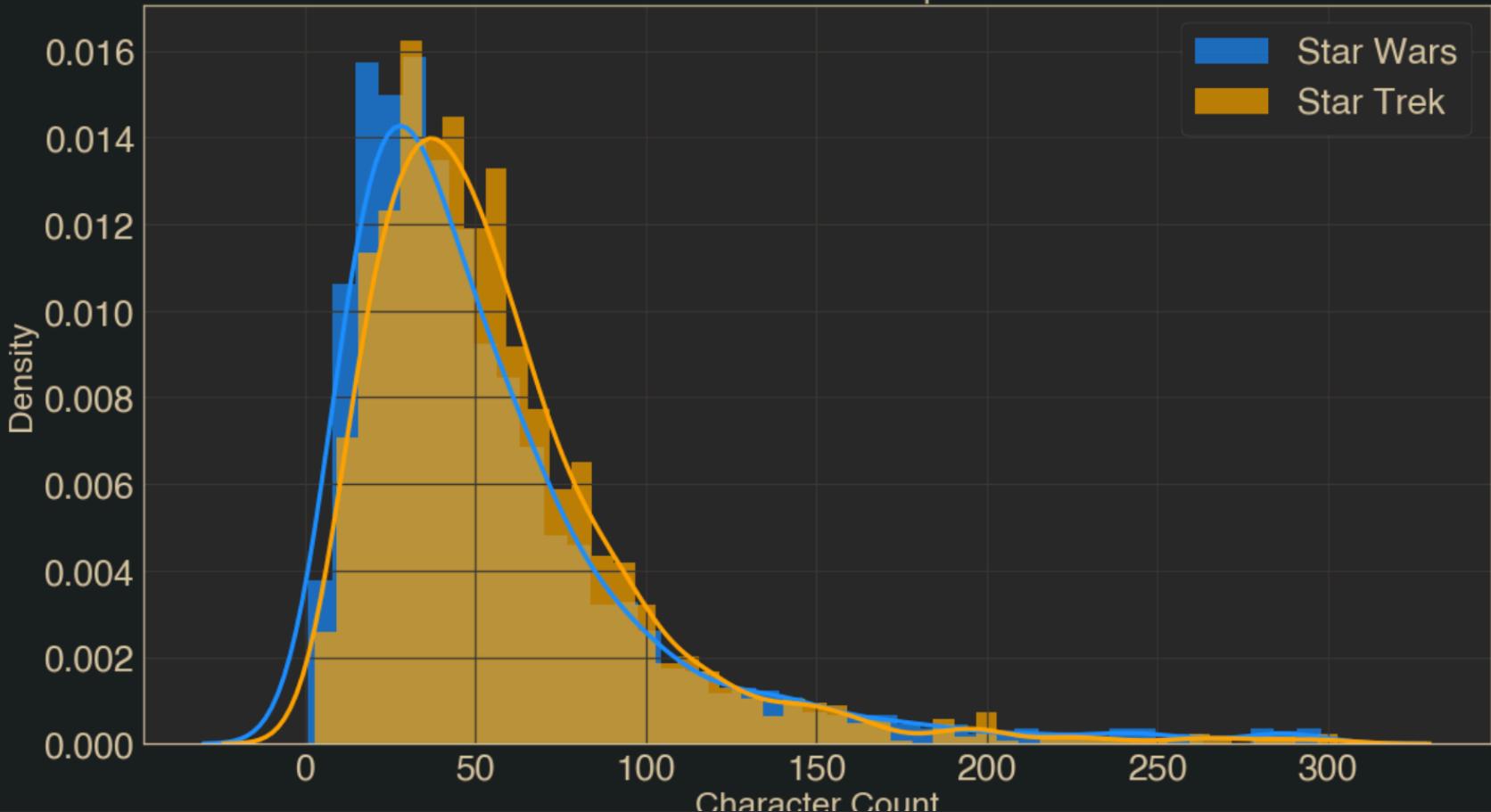
Pre-processing

- emojis removed
- “Star Trek” and “Star Wars” as phrases removed
- URLs removed
- Lemmatization applied
- Stopwords list was augmented with additional elements

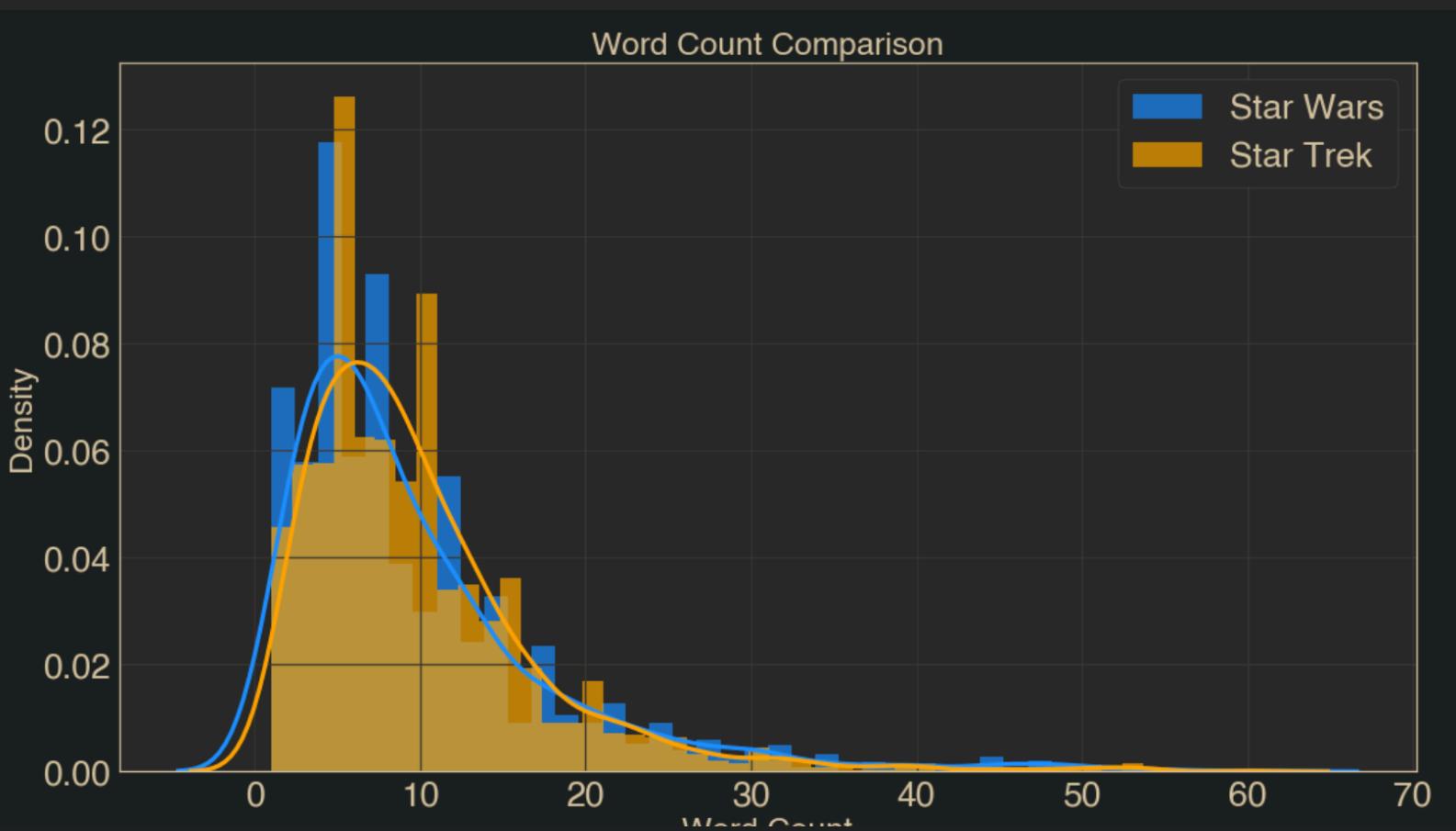




Character Count Comparison

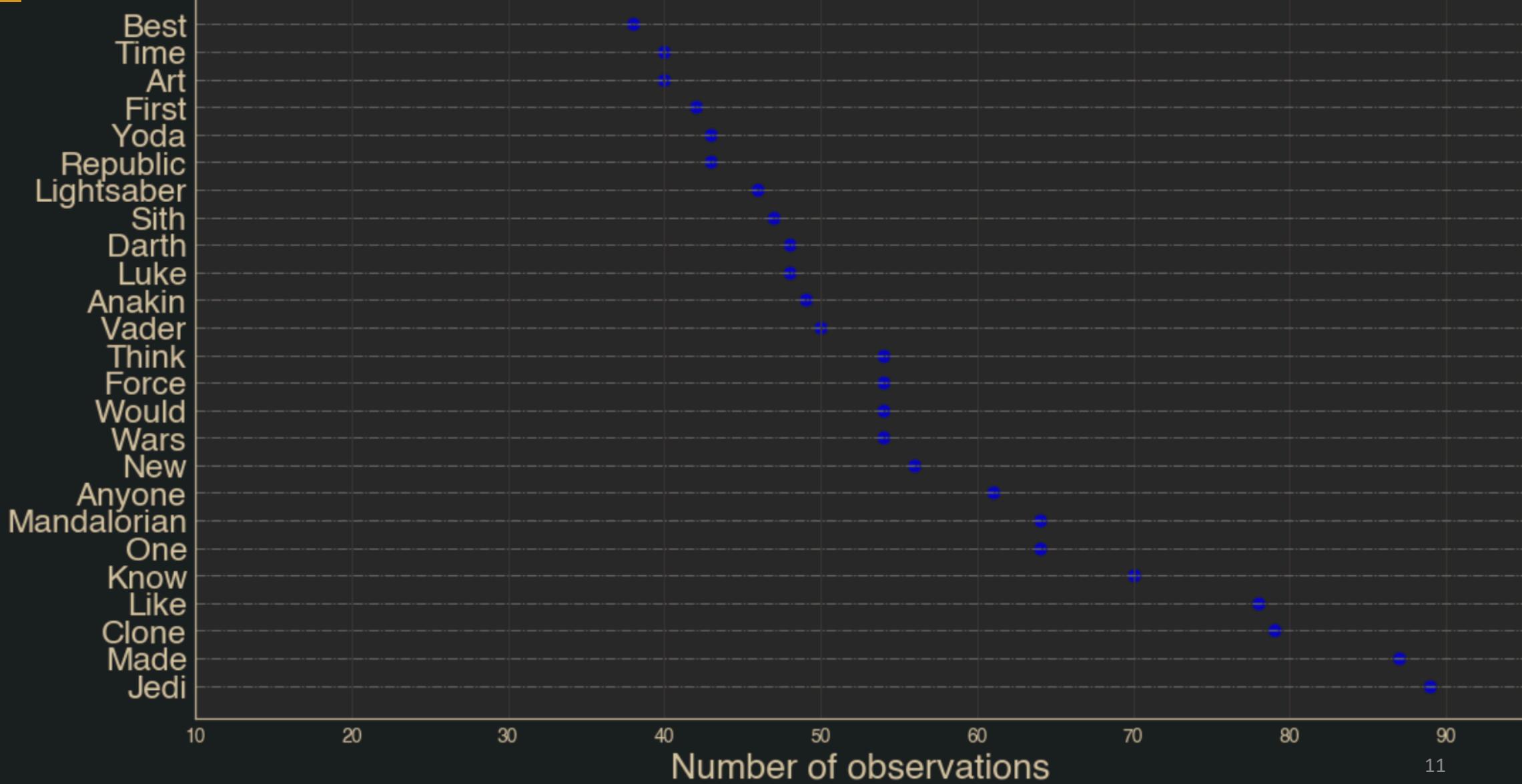


	count	mean	std	min	25%	50%	75%	max
 subreddit								
StarWars	1989.0	54.332831	46.435670	1.0	24.0	41.0	68.0	298.0
startrek	1989.0	57.182001	40.685735	3.0	31.0	48.0	71.0	303.0



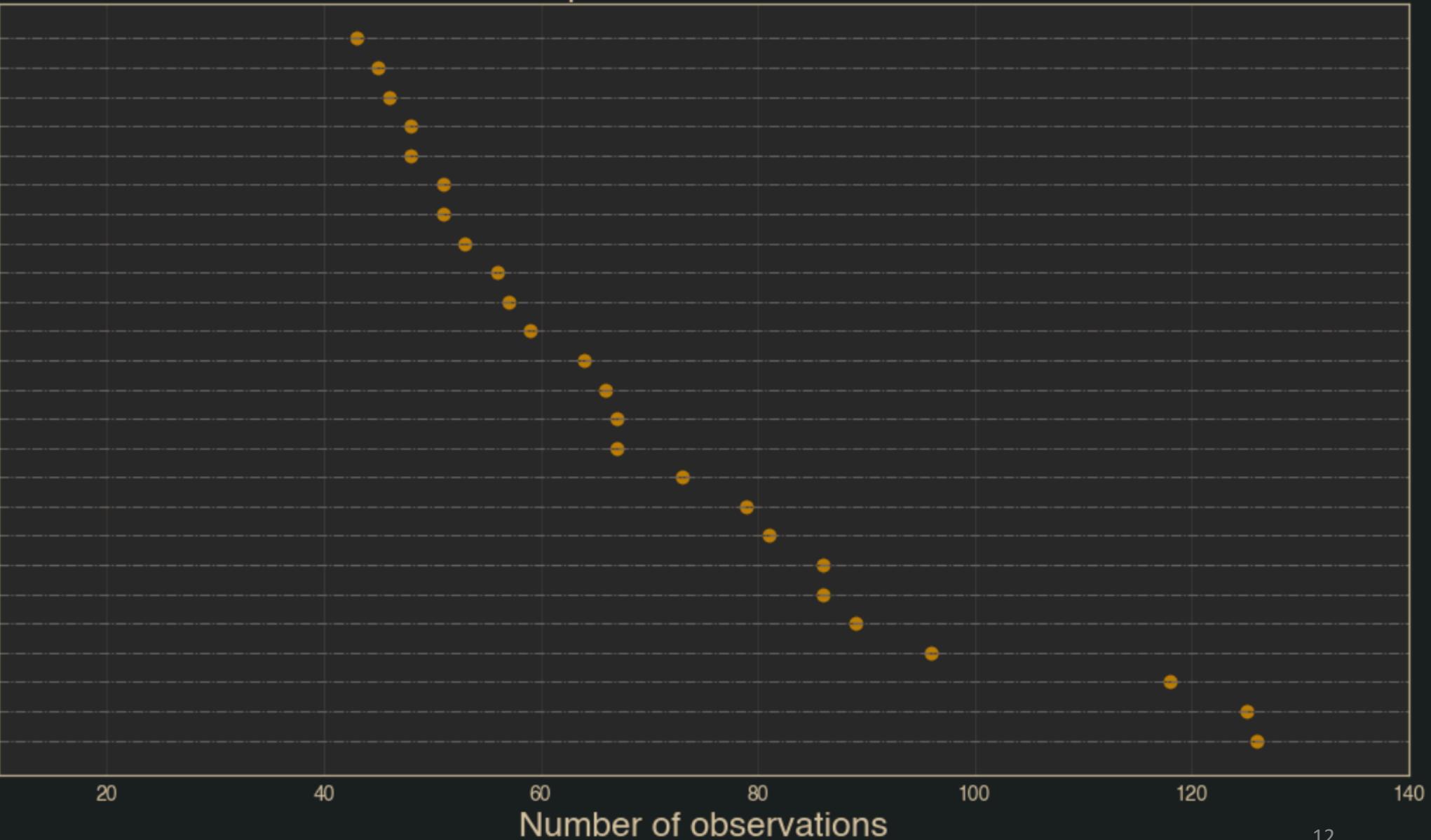
	count	mean	std	min	25%	50%	75%	max
 subreddit								
StarWars	1989.0	9.918552	8.569247	1.0	4.0	7.0	13.0	61.0
startrek	1989.0	10.024635	7.397312	1.0	5.0	8.0	13.0	60.0

Top words in Star Wars



Top words in Star Trek

Best
Decks
Lower
Watch
One
Anyone
Tos
Question
New
Think
Captain
Trek
First
Series
Time
Like
Would
Picard
Season
Voyager
Discovery
Enterprise
Ds9
Tng
Episode





Word overlap

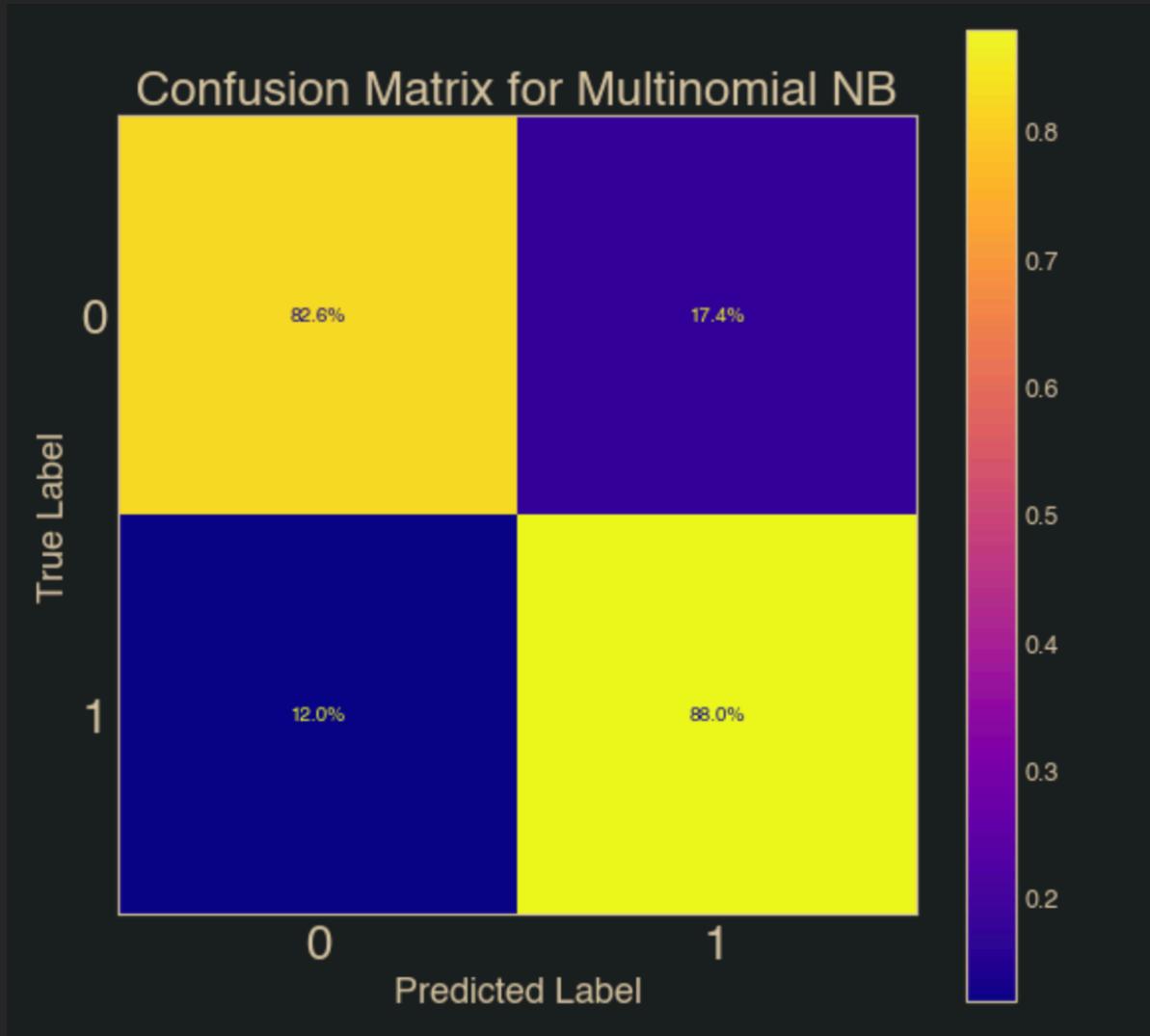
'like'
'one'
'anyone'
'new'
'would'

'think'
'first'
'time'
'best'

The Modeling process

Model	Training Score	Testing Score	Specificity	Sensitivity
Multinomial NB	93.6	85.3	82.6%	88.0%
Random Forest	98.5	81.4	75.6%	87.3%

The Modeling process



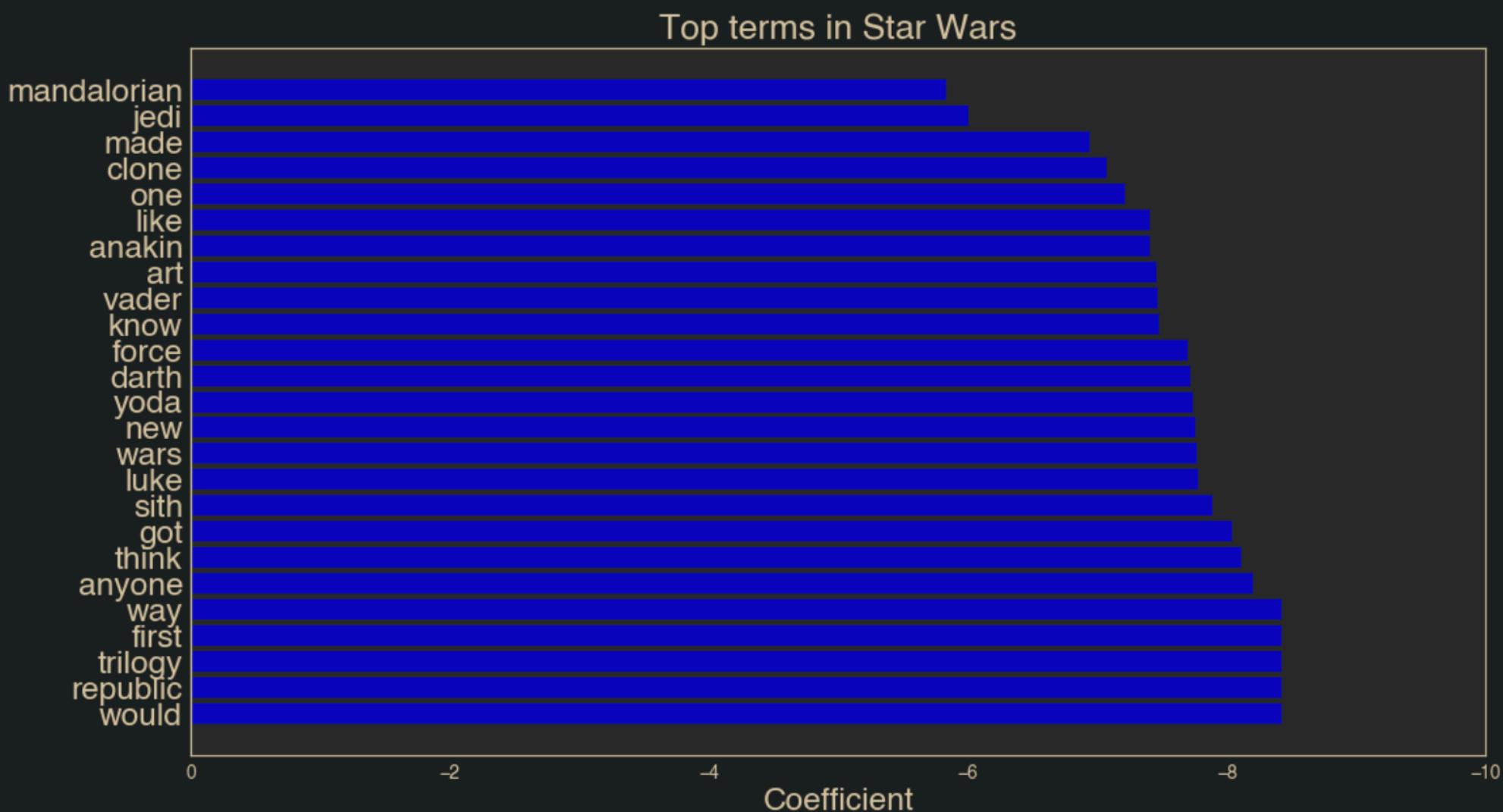
Specificity
(true negative rate):

82.6

Sensitivity
(true positive rate):

88.0

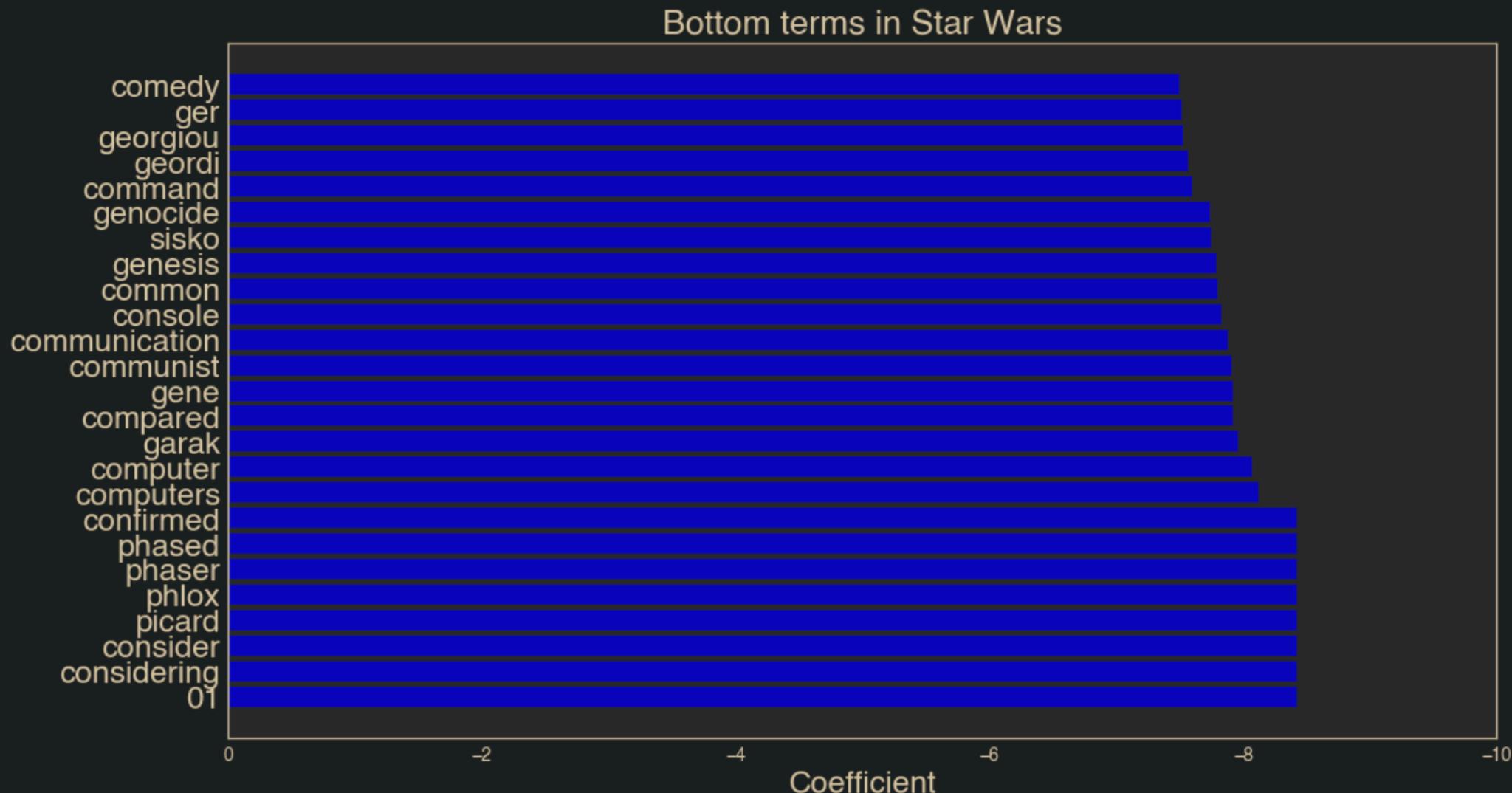
Conclusions and recommendations



Conclusions and recommendations



Conclusions and recommendations



Conclusions and recommendations





Conclusions and recommendations

- Unique terms help identify the subreddit
- There could be further tweaks to the modeling process
- Pipelines with grid search are recommended
- As of now, the project helped to learn difference between the two subreddits and help classify posts into correct categories.

Thank you



Q&A