

Project 3: Web APIs & NLP

Fallout vs. Star Trek

Regene M. DePiero

General Assembly

regenedepiero@gmail.com

October 9, 2020

Overview

1 Some EDA

2 Models

3 Conclusions

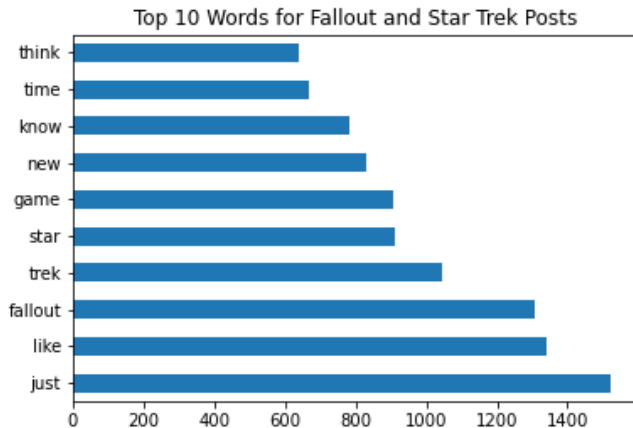
The Problem

Train a classifier that can determine which subreddit a post came from.

Subreddits of Interest

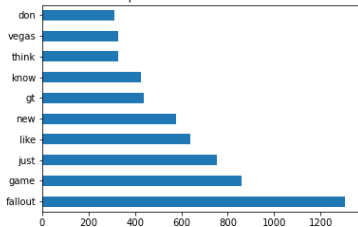


Top Words

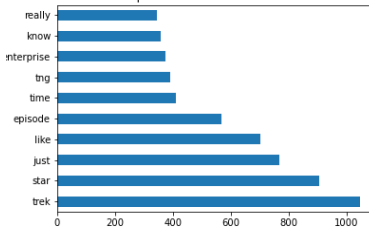


Top Words by Subreddit

Top 10 Words for Fallout Posts



Top 10 Words for Star Trek Posts



- Best model used $\alpha = 0.0001$
- $\text{max_df} = 0.9$

Scoring:

- F1 Score Training data: 0.9577
- F1 Score Testing data: 0.9437
- F1 Score on complete data: 0.9542

- No special parameter selection was done on this model.

Scoring:

- F1 Score Training data: 0.9872
- F1 Score Testing data: 0.9211
- F1 Score on complete data: 0.9704

Conclusions

- Increasing n-grams did not improve model performance.
- Use of TfidfVectorizer did not improve model performance.
- Further processing of data could be done to clean text.

The End