

PROJECT 3

QUANTIFYING REDDIT

By Evgeny Didenko

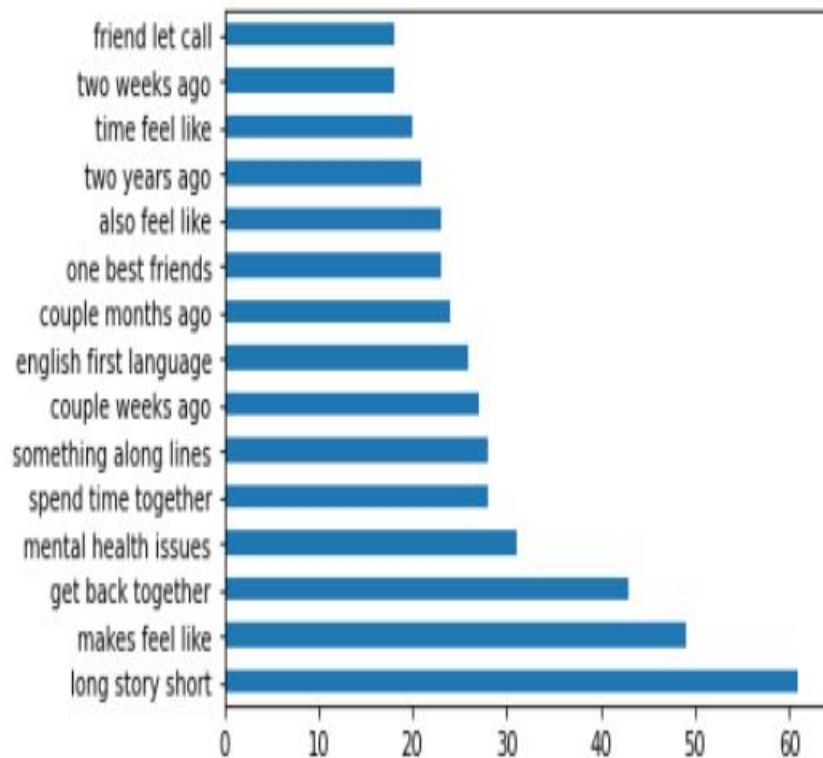
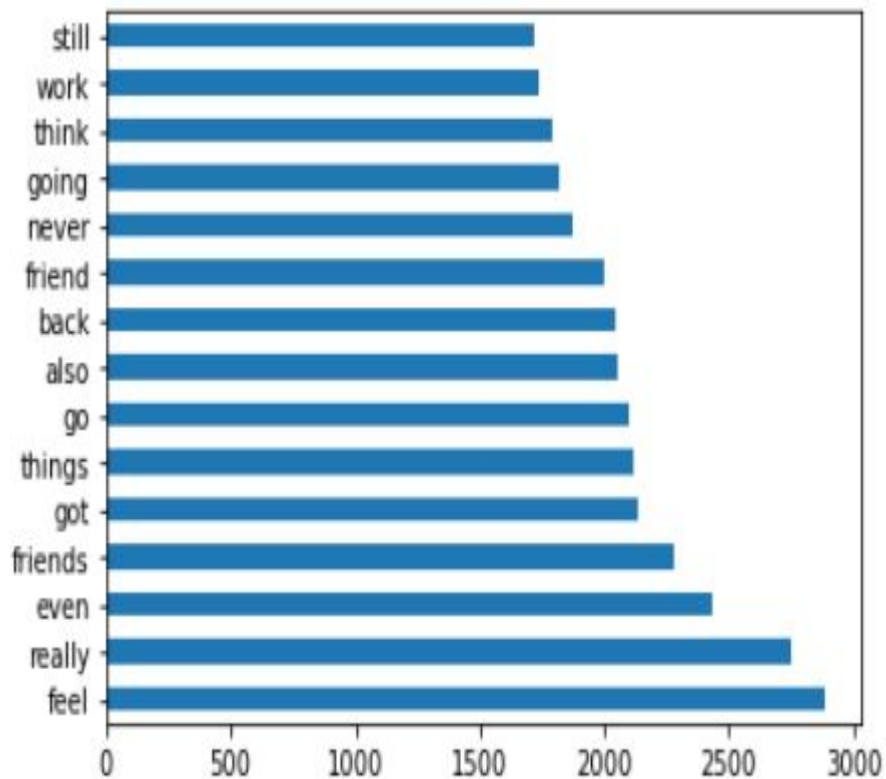
OVERVIEW

- ❑ The goal is to investigate how capable are the basic classification functions at distinguishing text data from two subreddits of similar general topic.
- ❑ The subreddits are "Relationship" and "Am I the Asshole", they both cover stories of conflict in interpersonal relationships and they are both very text heavy.
- ❑ Generally very large portion of Reddit is links, pictures or video clips, I don't think we have the tools to analyze these yet.

PROCURING THE DATA

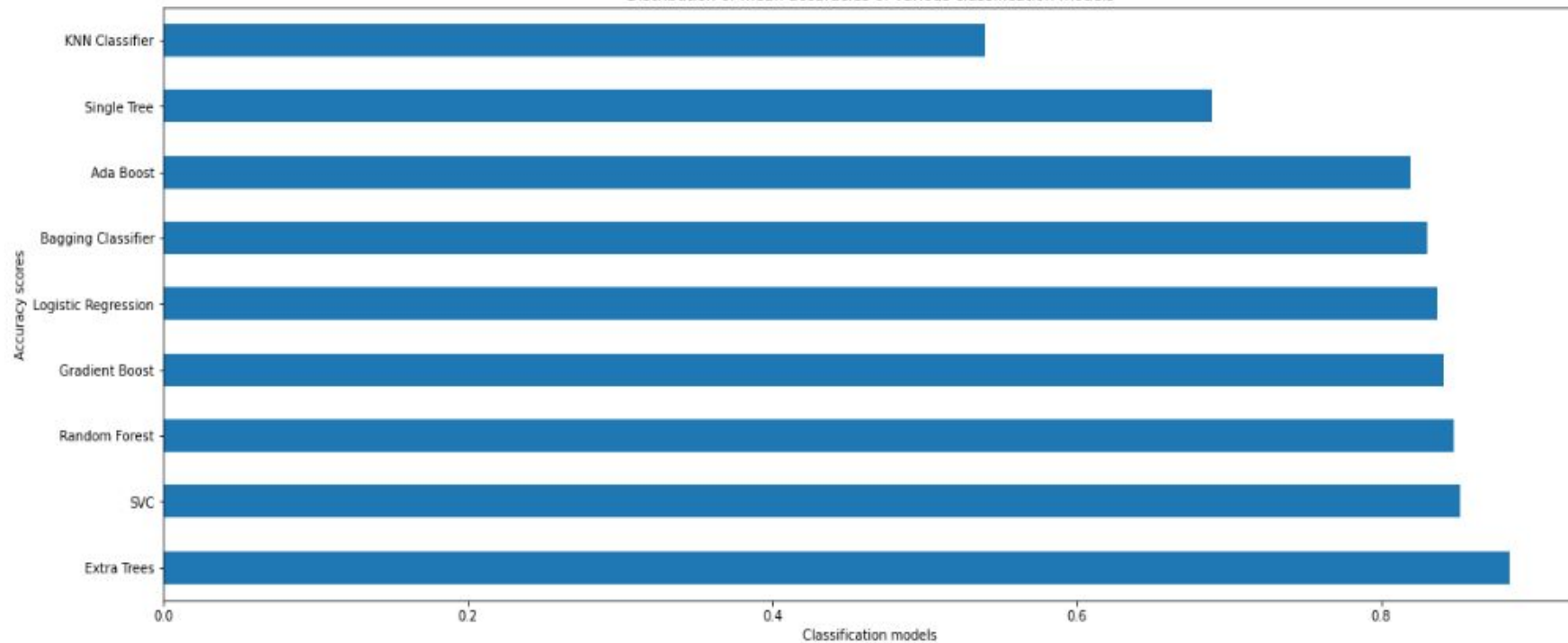
- ❑ We were given a link to API capable of downloading data from reddit as well as detailed instructions on how to use it.
- ❑ This API does not like handling too many posts or too many requests at once, so a loop with 2 second delay between cycles was implemented.
- ❑ Each loop after first used the date of the first (chronologically) post as the starting point, so although there were some duplicated, vast majority of posts was unique.
- ❑ Relevant (just the subreddit name, post title and post text) data from each subreddit was fed to a separate dataframes that were merged afterwards.

EDA



MODELS

Distribution of mean accuracies of various classification models



QUESTIONS