

Web APIs & NLP

Sara Zhou

Problem Statement :

Given a post from either
r/OutOfTheLoop and
r/explainlikeimfive, how well can we
predict which subreddit it came
from?

Background on the Subreddits

r/OutOfTheLoop

r/explainlikeimfive

“A subreddit to help you keep up to date with what’s going on with reddit and other stuff.”

People ask mostly about news or trends they’re missing out on. Ex. 2017 What’s going on with Net Neutrality? At the time the FCC was voting on whether or not NN should exist.

Posts asking about concepts that only people with extensive knowledge on the subject could answer.

People ask any type of concept, simple to complex. Ex. Why is the sky blue? Vs How after 5000 years of humanity surviving off of bread do we have so many people who are allergic to gluten?

“A subreddit to help you keep up to date with what’s going on with reddit and other stuff.”



How well could the model perform?

The Process

Obtaining Posts

Grabbed 100,000 Post

After cleaning out all the null and deleted posts I ended up with around 25,000 - 30,000 per subreddit.

Analysis & Preprocessing

Tokenizing

I tokenized all the self text and removed common words and decided to not look at titles because in ELIF users must have ELIF in the title name.

Modeling

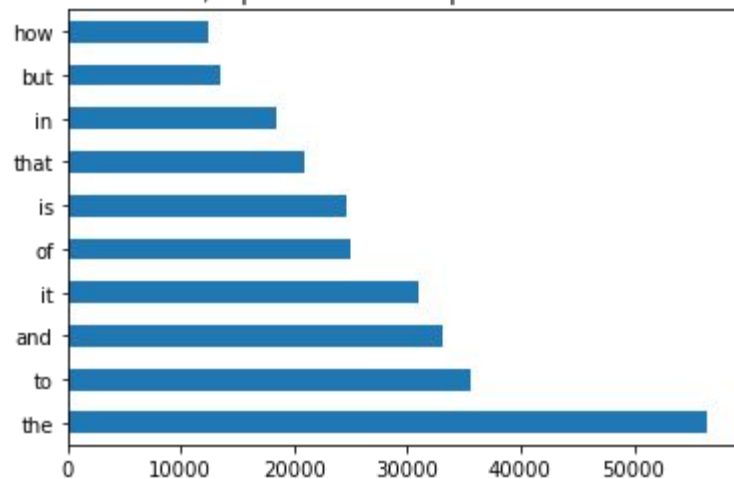
I used 3 different models

1. GS Over CountVectorizer BernoulliNB
2. GS Over TF-IDF Vectorizer and Multinomial NB
3. Random Forests

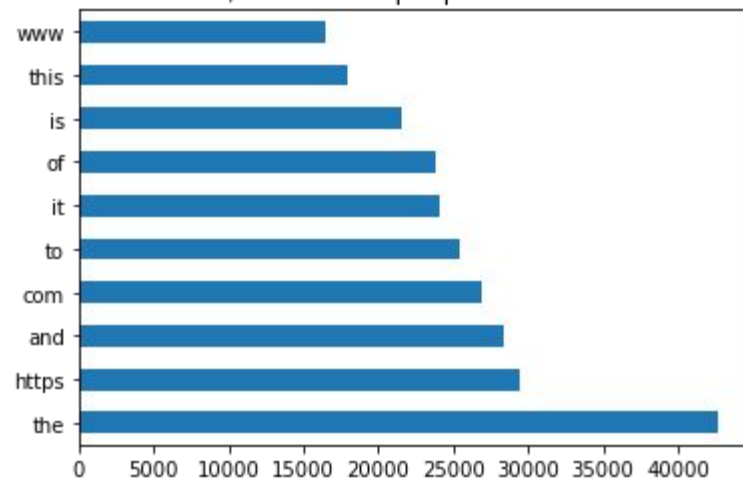
Model Performance

GS Over CountVectorizer BernoulliNB	GS Over TF-IDF Vectorizer and Multinomial Naive Bayes	Random Forest
Train: 94.16%, Test: 93.81%	Train: 94.53%, Test: 94.26%	Train: 99.93%, Test: 93.81%

r/explainlikeimfive Top Word Count



r/OutOfTheLoop Top Word Count



Conclusion:

My training and testing score stayed in the 90s. When choosing these subreddits I expected a lot of overlap because they both seemed very similar but after running these models it's clear that they had no problem differentiating the two. The list of top words for either were different enough for the most part, it seems that the users of these subreddits were both looking for answers for their questions but the types of questions asked differed enough to produce a fairly accurate model.
