# Classification models on subreddits [r/Capitalism and r/communism]

By Genaro Garcia Pereyra

# Problem Statement

This projects seeks to demonstrate a classification model that is capable of classifying the subreddits of **r/Capitalism** and **r/communism** better than the baseline model given by the majority class in the dataset collected. In addition, the goal for this data analysis is to provide more information of the sentiment that users display on each of the subreddits and how sentiment and influence in some degree to the accuracy of the classifier.

The classification model uses ensemble methods in order to reduce the risk of overfitting and have a good perform in unseen data by combining predictions from several base models.

Accuracy is used as the metric for evaluation for all the models, and depending the model's performance on this metric the final model is determined.

# Data Collection

Data collection was performed with the help of the **pushshift API** which allows us to search all publicly available comments and submissions on Reddit. For this particular case the text to be classified  only considers the **comments section** of each post, this decision was made with the idea of having more available text that could allow a better prediction of the models implemented. In addition, the interactions between users in the comments section is a critical component on how to determine the general sentiment of the users that use that particular subreddit.
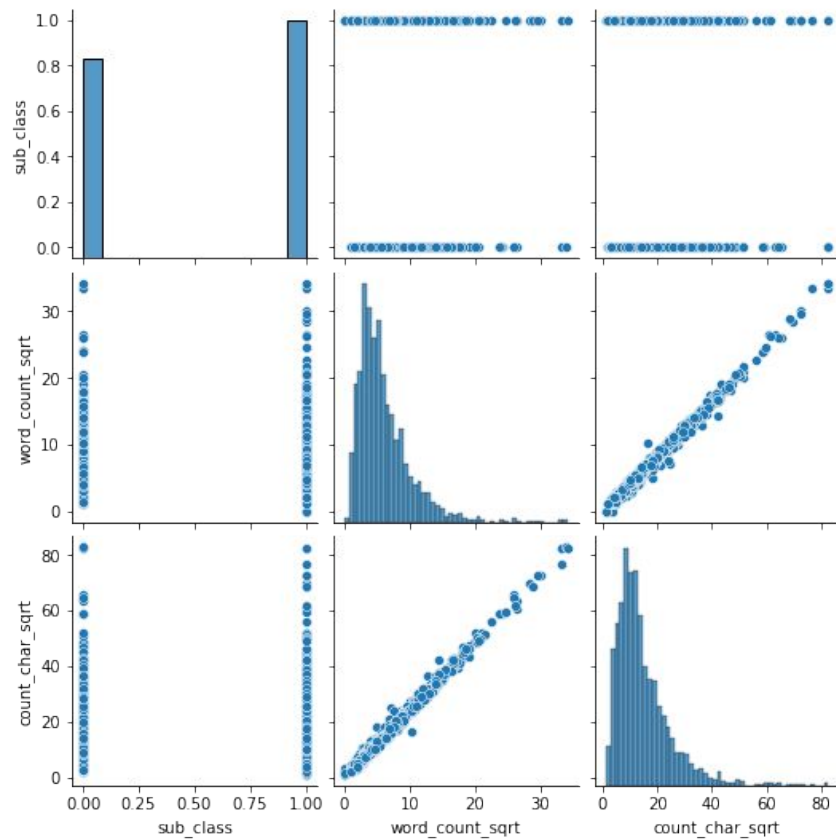
A size of 100 most recent comments was extracted for each of the last 11 months starting form (10/19/2020) for both the r/Capitalism and r/communism subreddits. This data was then used to created a single balanced dataset containing 2200 observations.
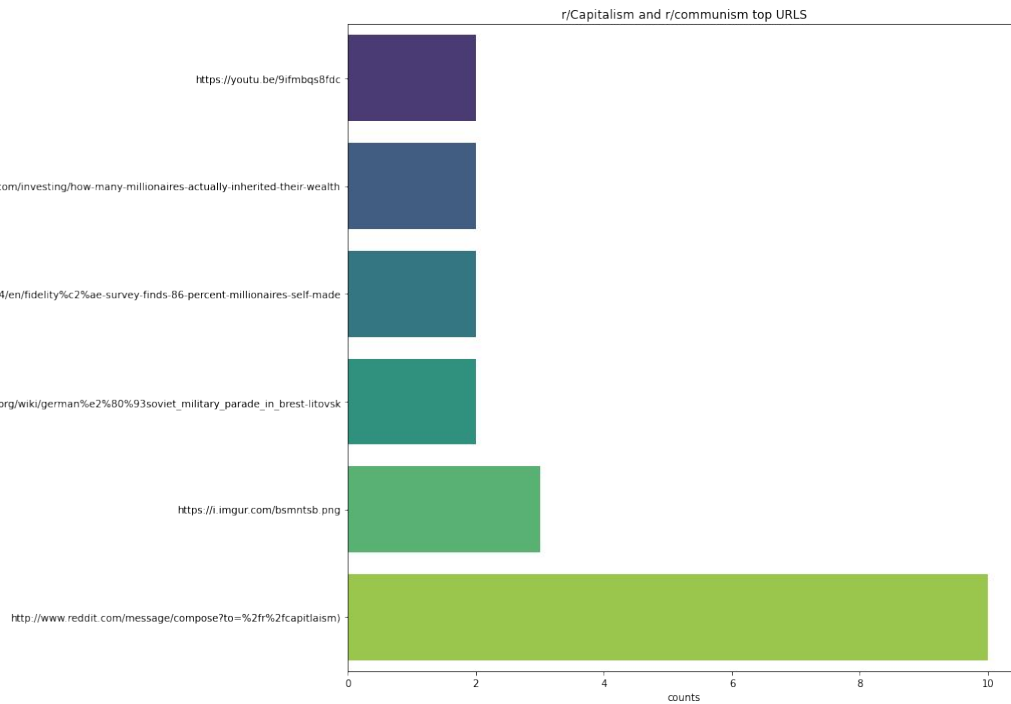
# Data Cleaning + EDA

Goals after performing data cleaning on 2200 comments extracted:

1. Remove presence of URLs
2. Remove presence of Russian words
3. Remove presence of Chinese words
4. Remove english stop words
5. Remove numerical values
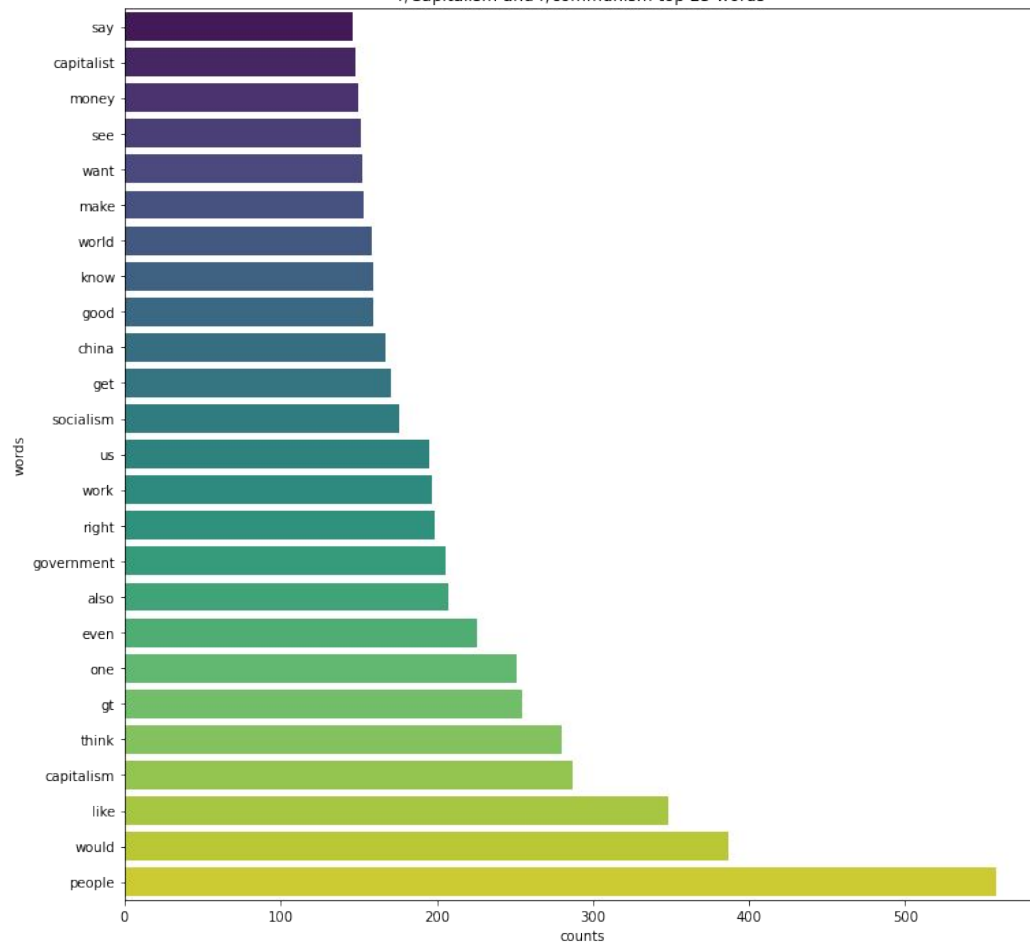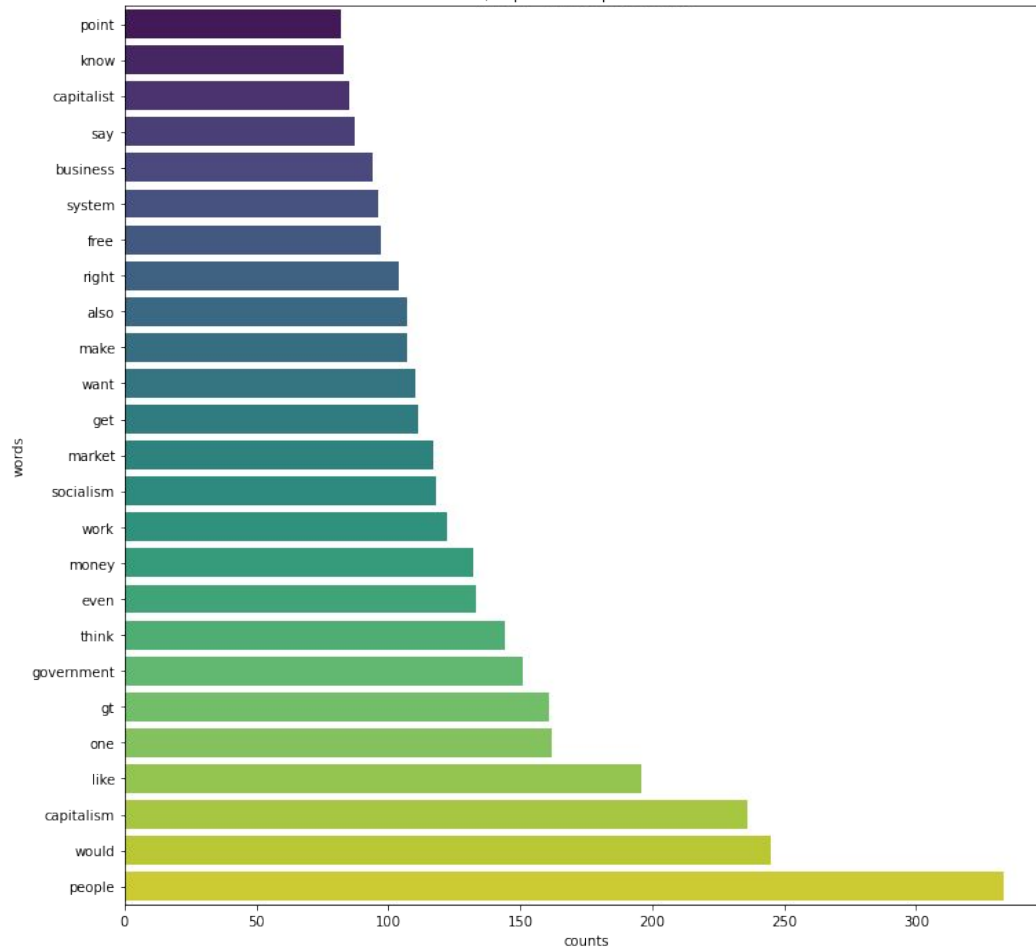6. Take care of all removed comments by users
7. Binarize class

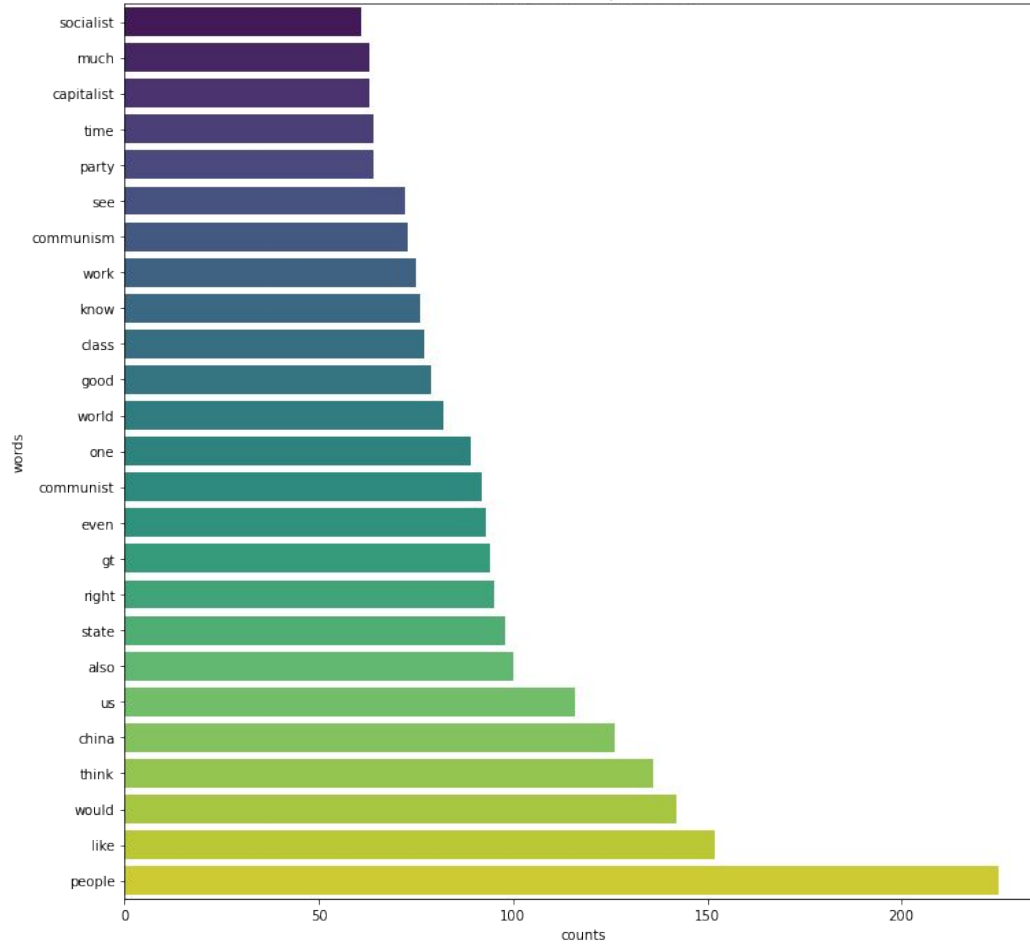r/Capitalism and r/communism top URLS

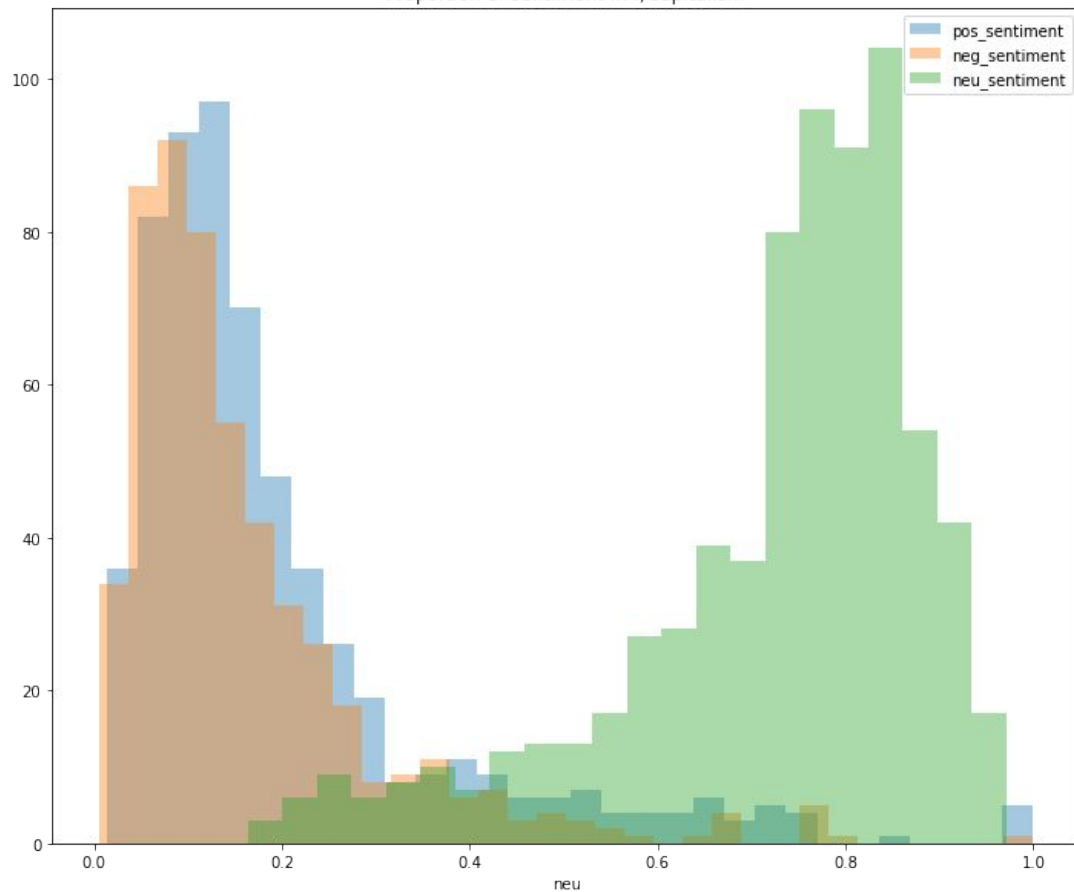r/Capitalism and r/communism top 25 words

r/Capitalism top 25 words

r/communism  top 25 words

r/Capitalism  top 25 words

r/communism  top 25 words
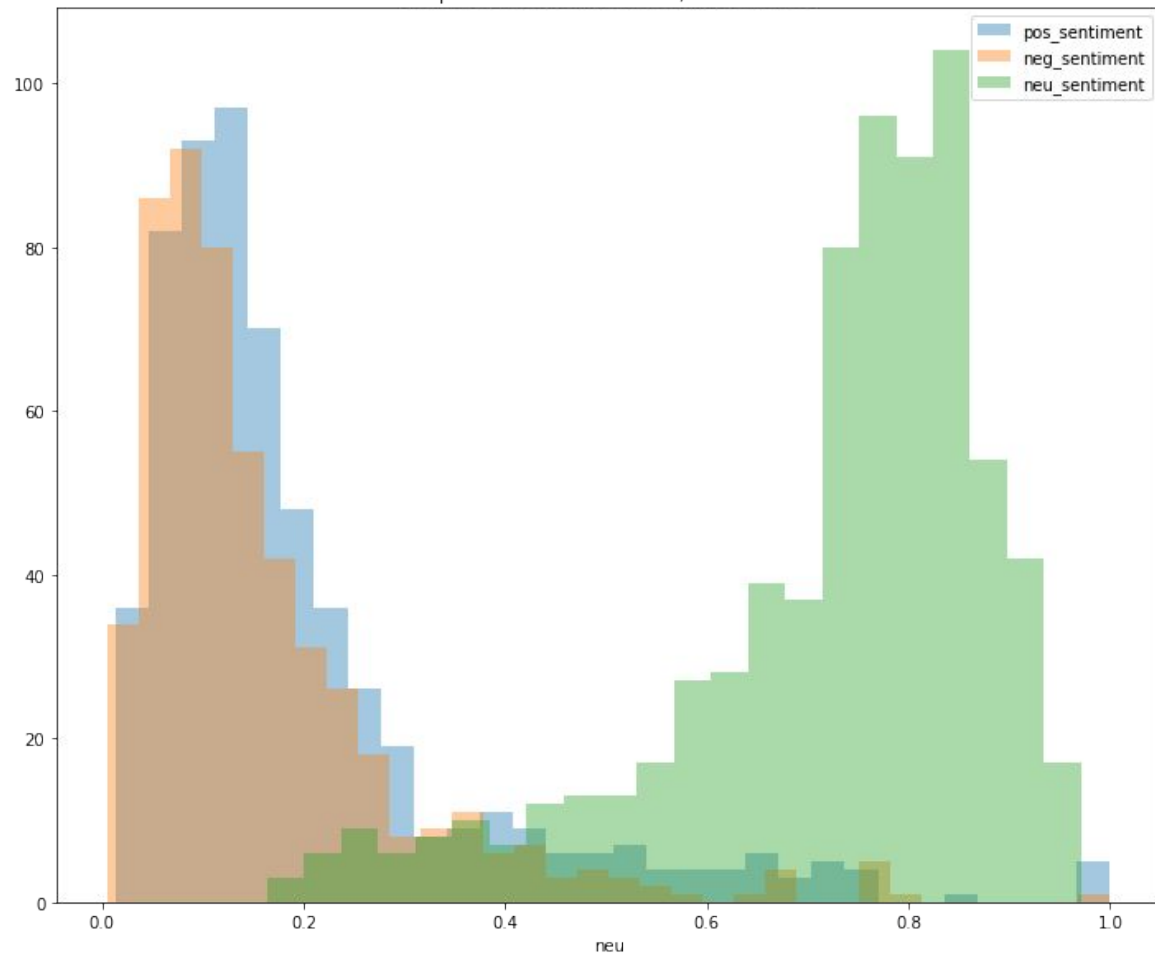
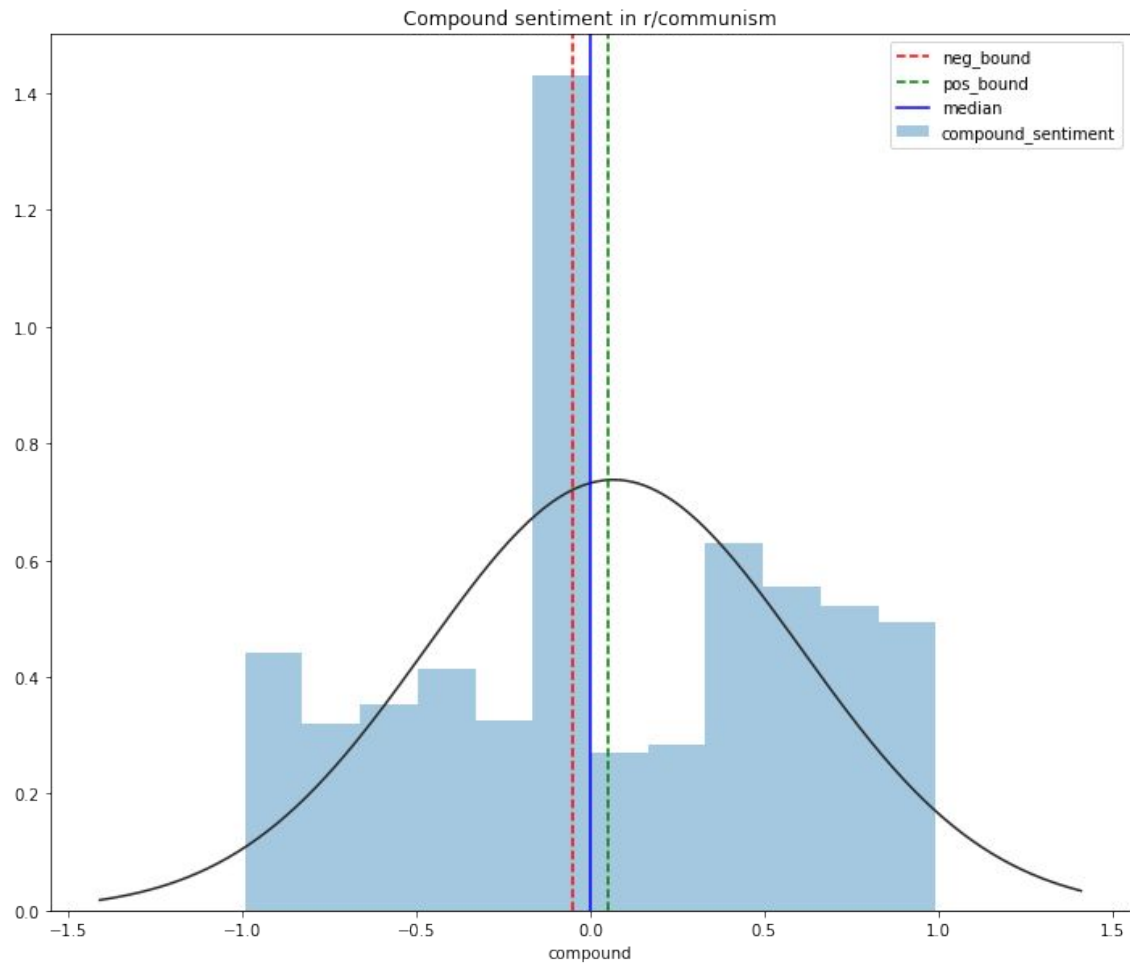Proportion of sentiment in r/Capitalism

Proportion of sentiment in r/communism

Compound sentiment in r/Capitalism

Compound sentiment in r/communism

# Modeling

| Model | Score on Training set | Score on Testing set | Vectorized text | VADER sentiment analysis |
|-------|----------------------|---------------------|-----------------|--------------------------|
| RFC + GS | ~ 0.9849 | ~ 0.7116 | Yes | No |
| Logistic Reg + GS | ~ 0.9071 | ~ 0.7382 | Yes | No |
| RFC + GS | ~ 0.9979 | ~ 0.7668 | Yes | Yes |
| Logistic Reg + GS | ~ 0.9098 | ~ 0.7525 | Yes | Yes |
| Ensemble [LogReg + RFC] | ~ 0.9938 | ~ 0.7689 | Yes | Yes |

# Conclusion and Recommendations

The use of sentiment analysis metrics as part of the features used in the model to classify the holdout test set of subreddits of c/Capitalism and c/communism was essential for the creation of a  model that could generalize better on unseen data. As seen in the word clouds for both subreddits there is high similarity of words used in both comment sections, and because of this it makes it harder for the model without VADER to obtain an acceptable accuracy score in the test set.

It is recommended to use tools such as VADER to have a better classification model. Furthemore, because of the presence of other languages it might be useful to incorporate corpus from Russian and Chinese languages.