

Reddit Subreddit analysis

By Derik Vo

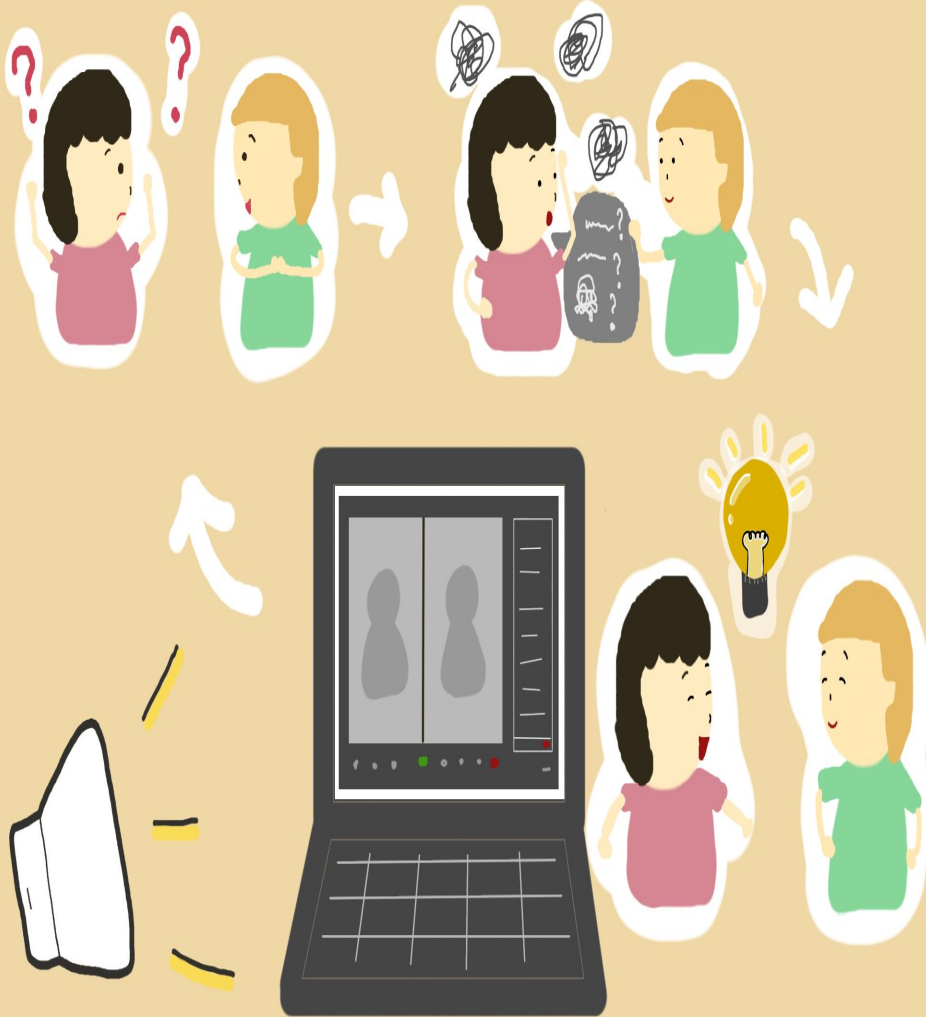
Background

1. Our team is studying social interaction by gender
2. We will incorporate this by using Reddit
3. We want to build a classification model

AskWomen vs. AskMen

Agenda

1. Data Collection
2. Exploratory Data analysis
3. Modeling
4. Conclusions
5. Recommendations
6. Q & A



Can we
categorize a
reddit post?

Data Collection

Record Count

Ask Men: 3995 records

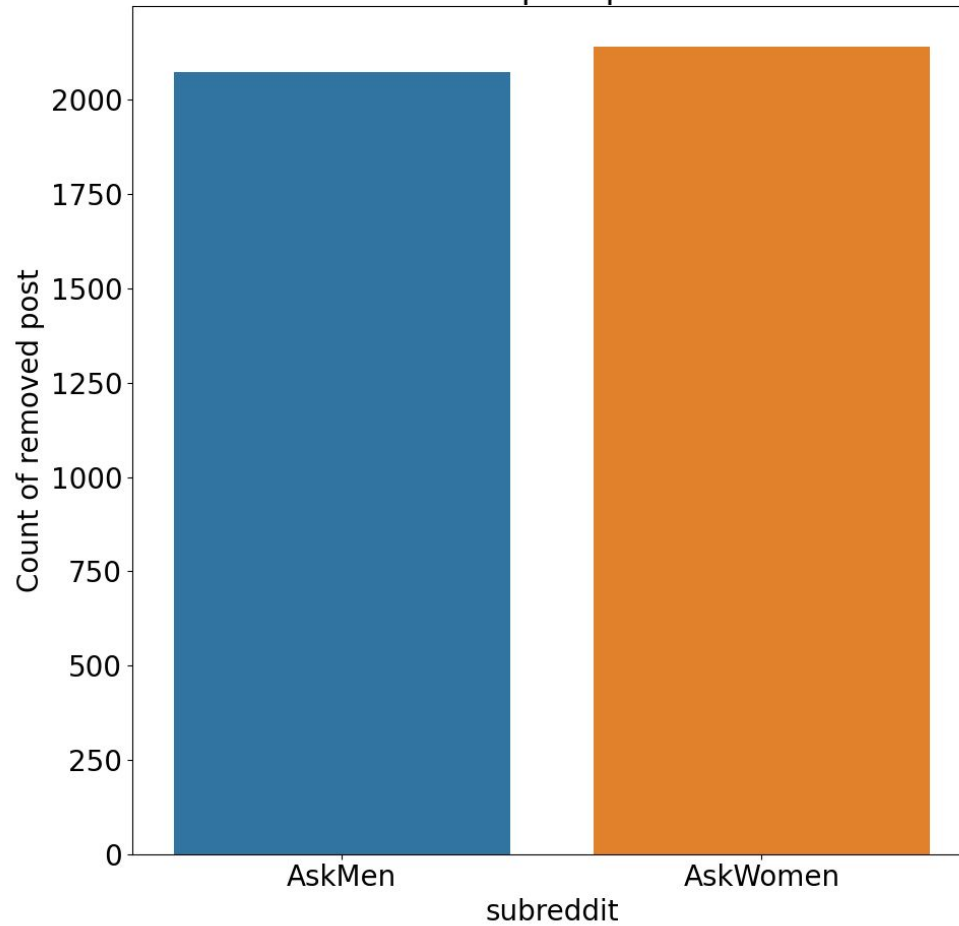
Ask Women: 3991 records

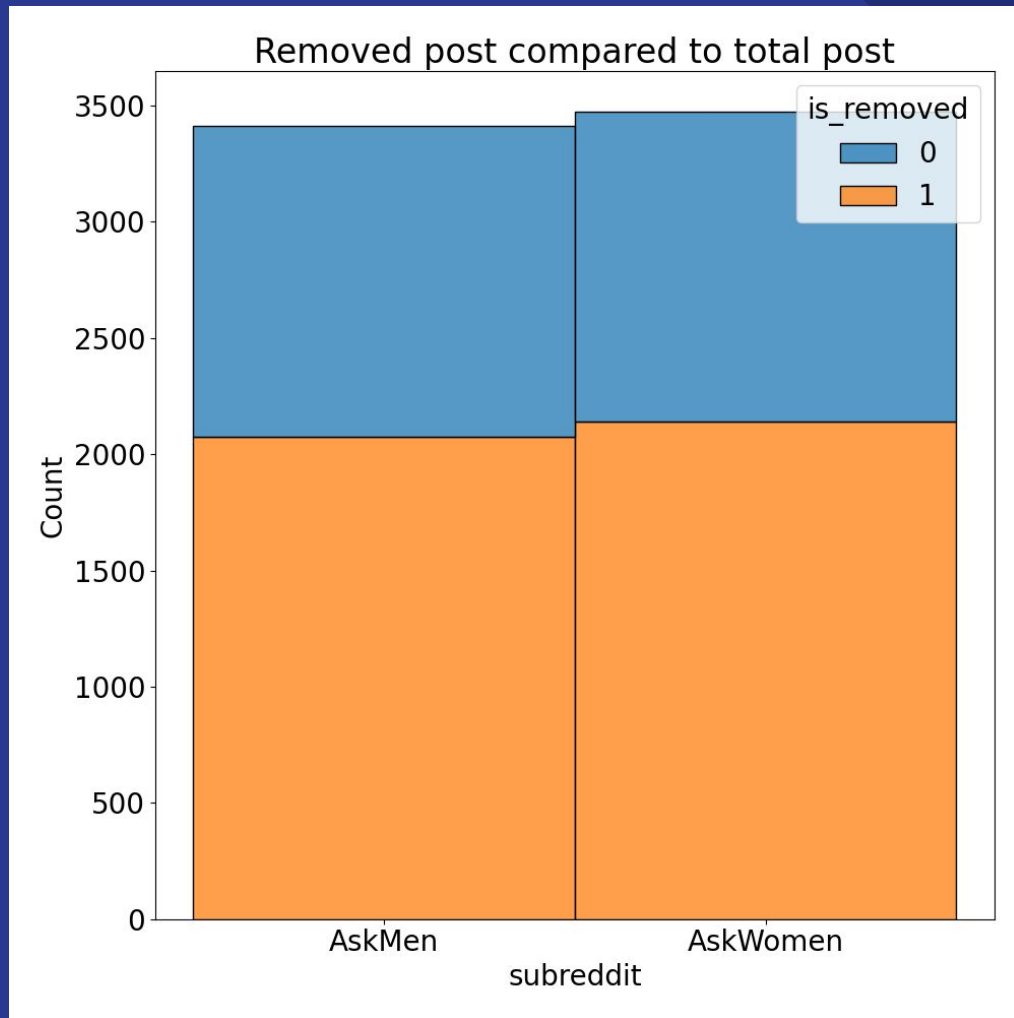
Filters used

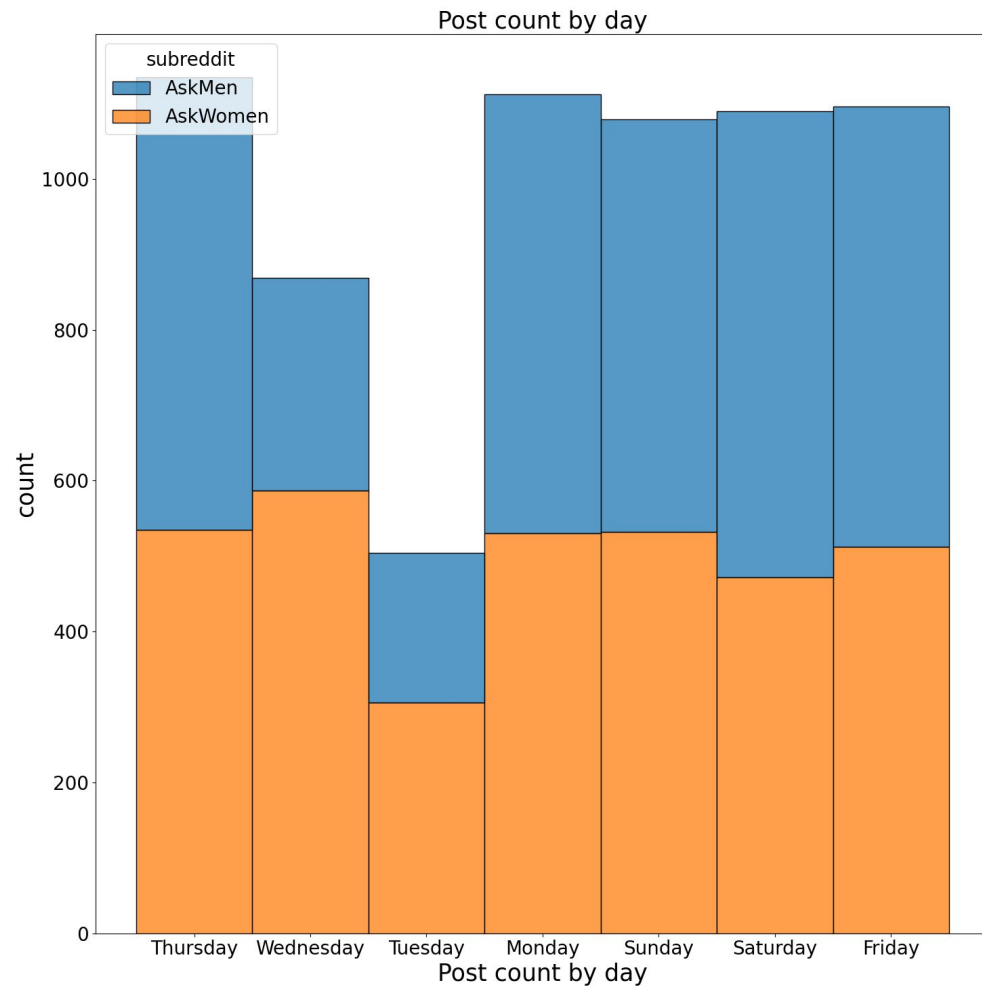
1. Since
2. min_comments

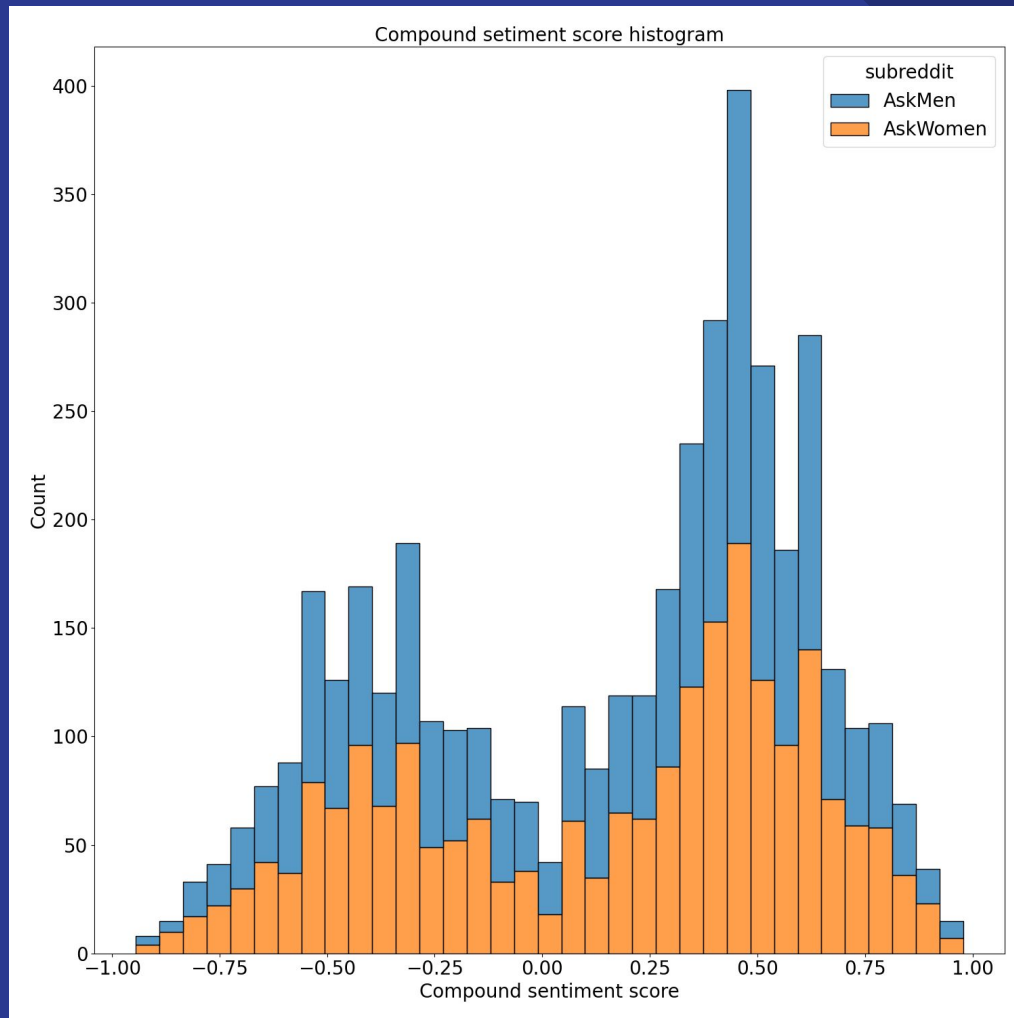
Exploratory Data analysis

Total removed post per subreddit

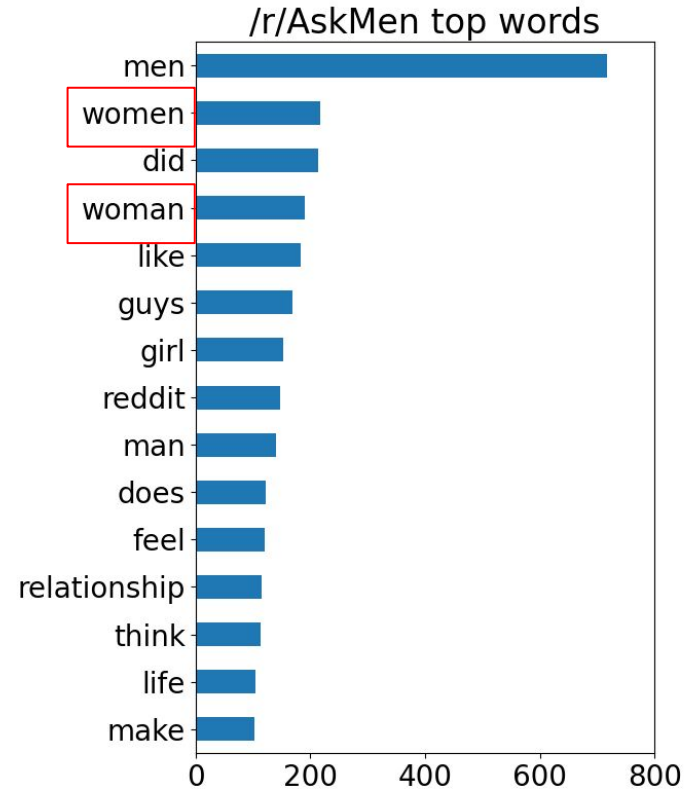
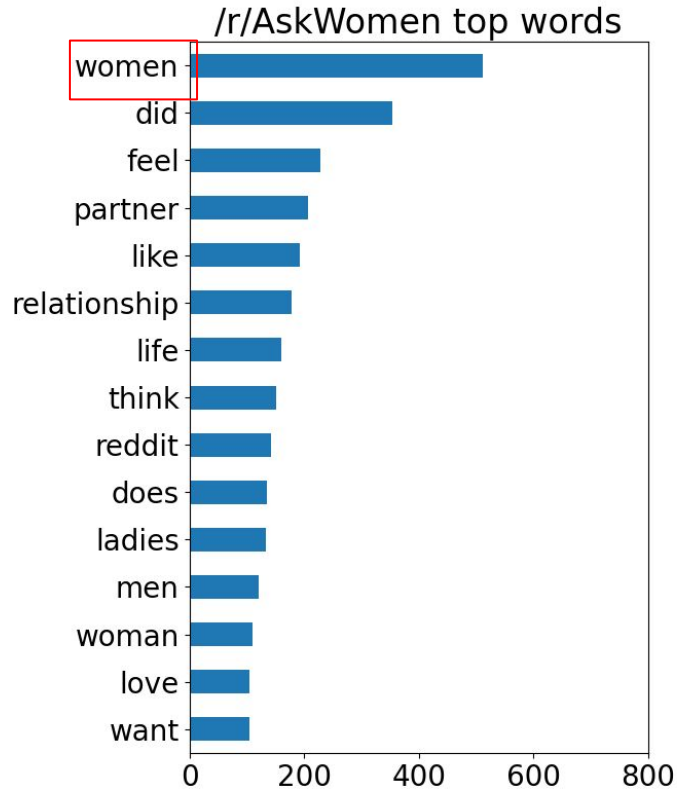




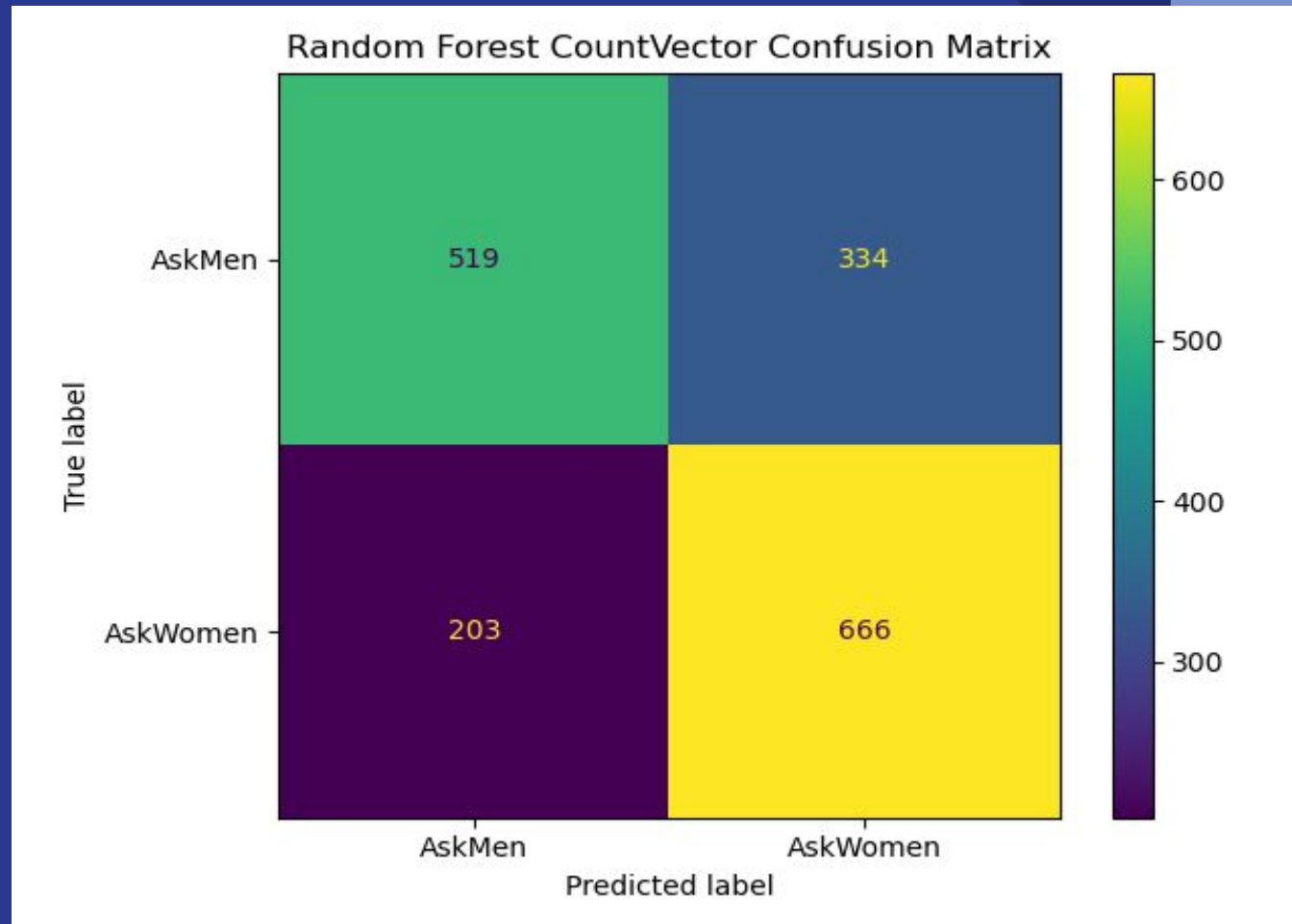


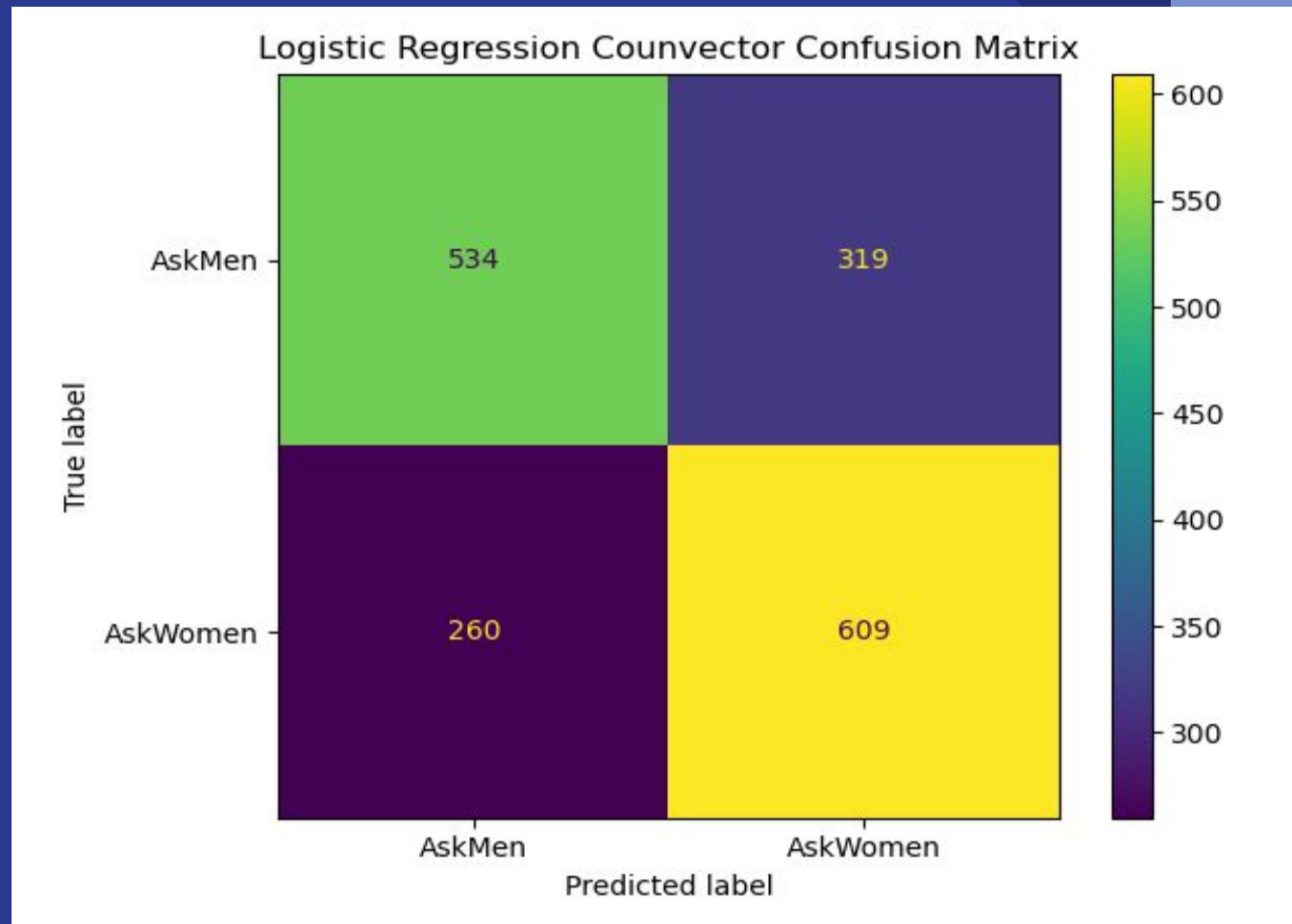


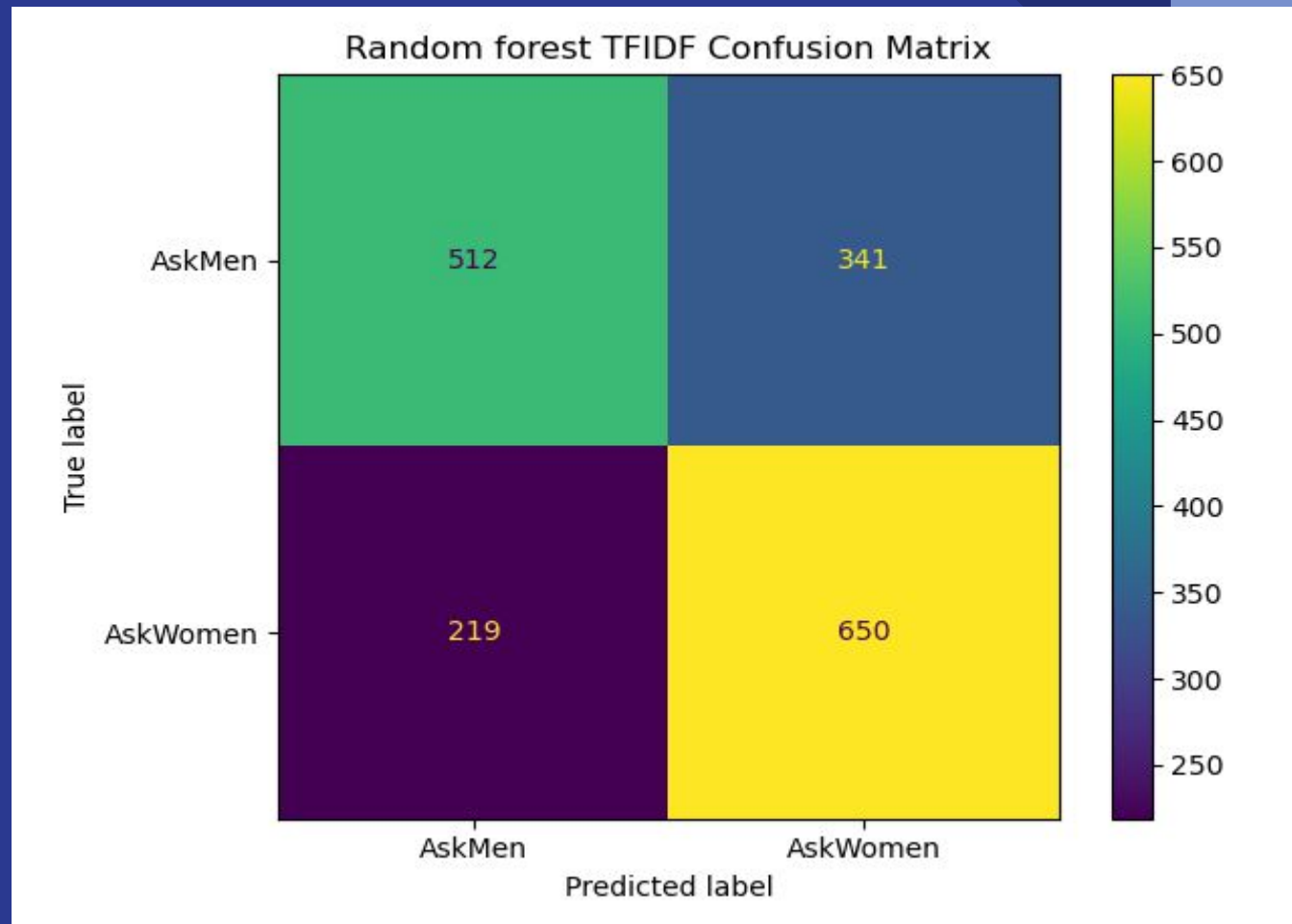
Top 15 words in each subreddit

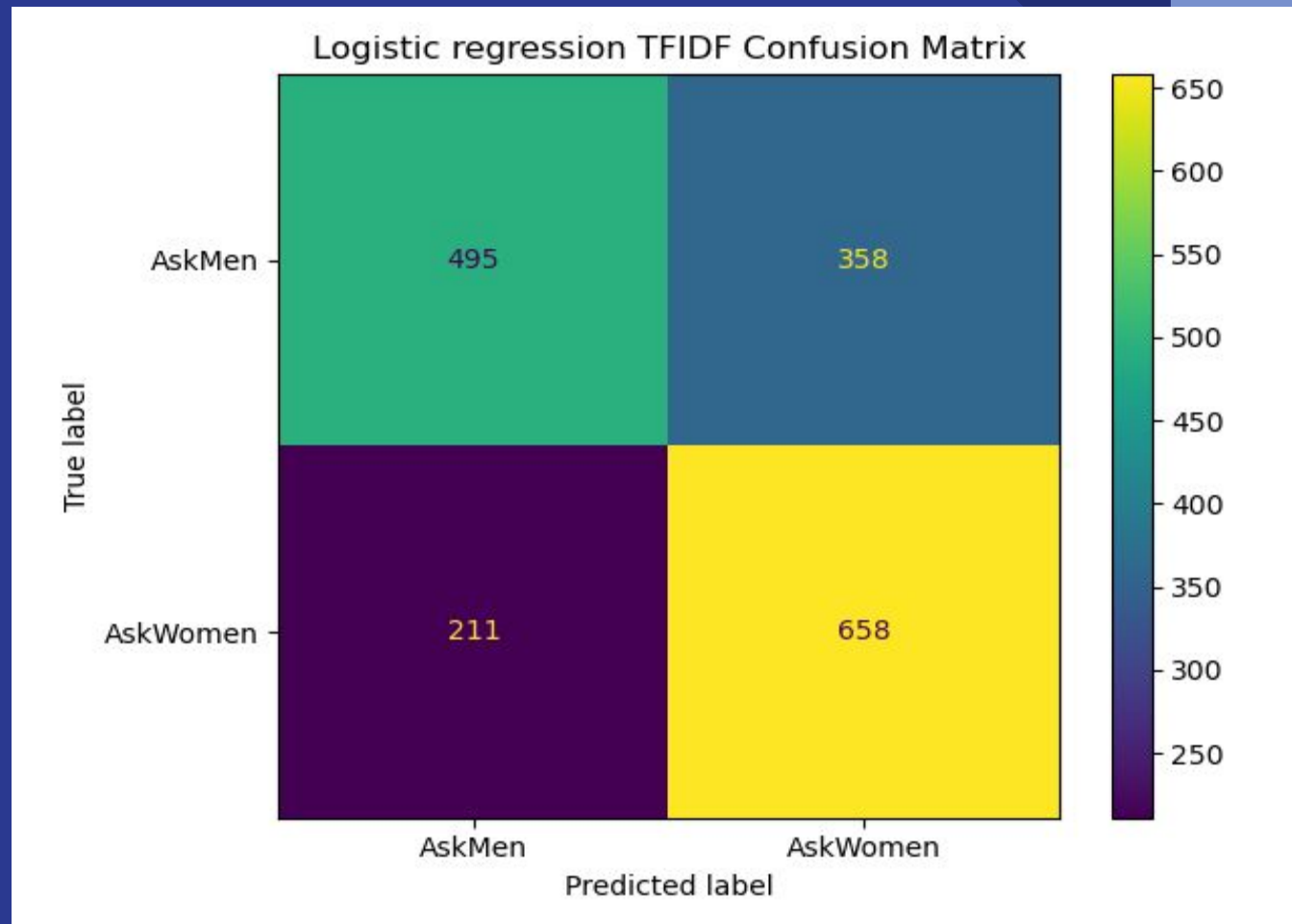


Modeling









Conclusion

Models could not categorize

1. Models were relatively consistent in the predictions
2. Models and transformer usage had marginal effect
3. The behavior of each subreddit had too much similarities

Recommendations

1. Filter out common words shared between subreddits
2. Incorporate the comments within the subreddits
3. Share our dataset to build a bot to identify post that should be removed
