



Project 3

WebAPI & NLP

by Radchenko Gleb

DSI-822

Project statement:

Our marketing company is looking for a way to improve effectiveness of advertising campaigns running on Reddit. For that purpose stakeholders would like to know how accurately we can predict, if a user a woman or a man based on the subreddits they are publishing submissions to.

Project success will be evaluated using resulting f1 score and accuracy score.



Subreddits

/r AskWoman



/r AskMan

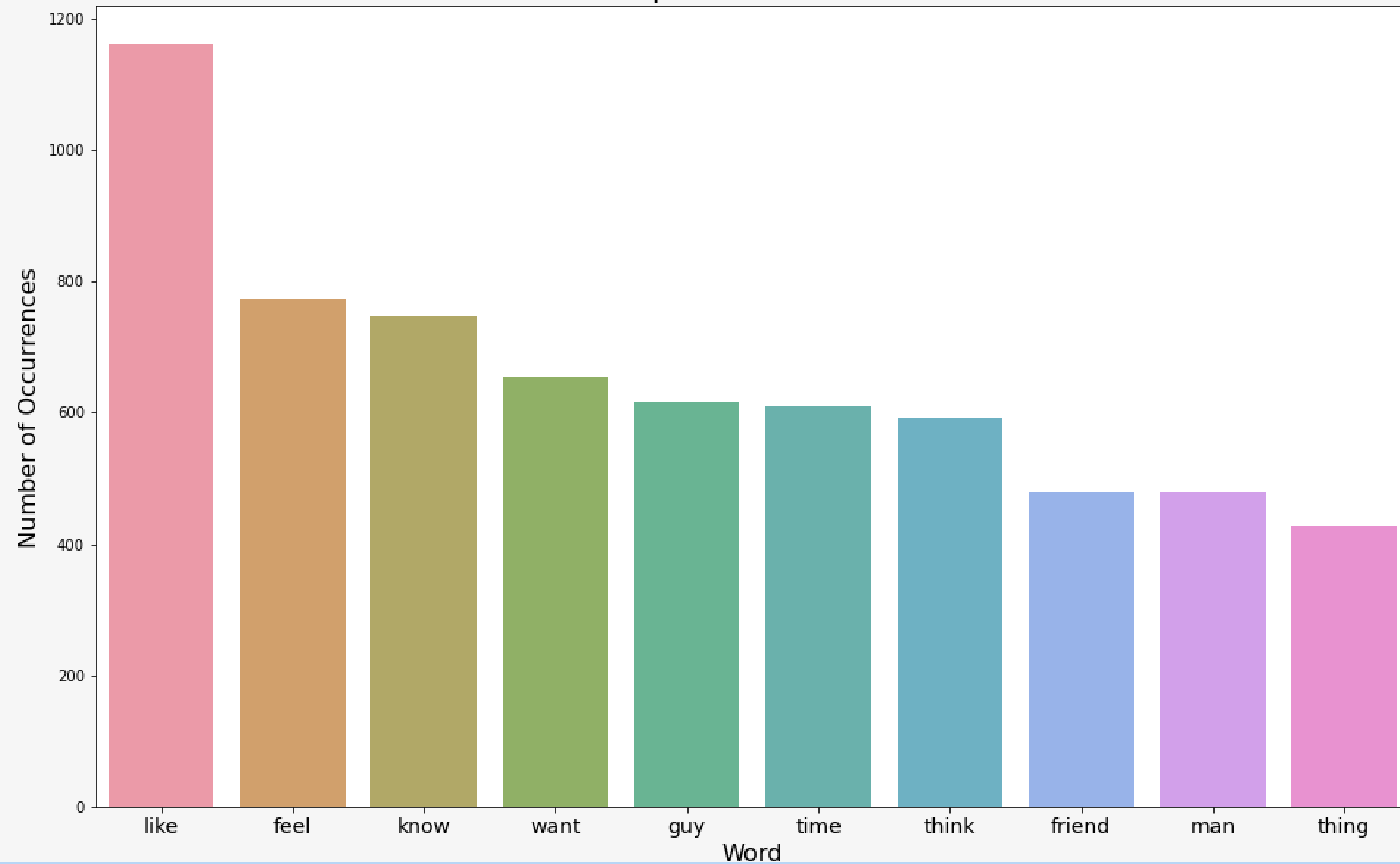
Reasons:

A lot of text data

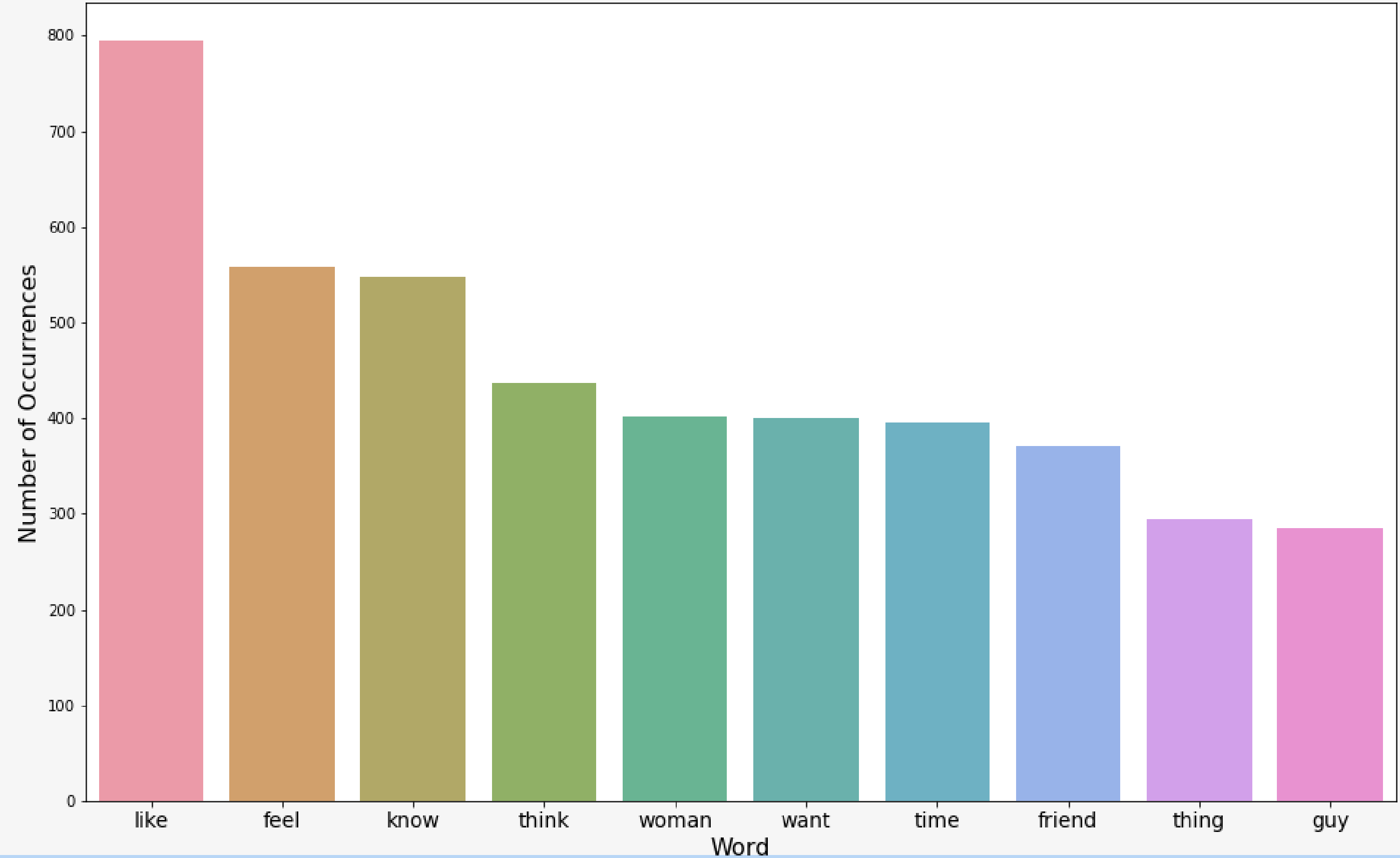
Significant amount of submission and commentaries

Interesting problem

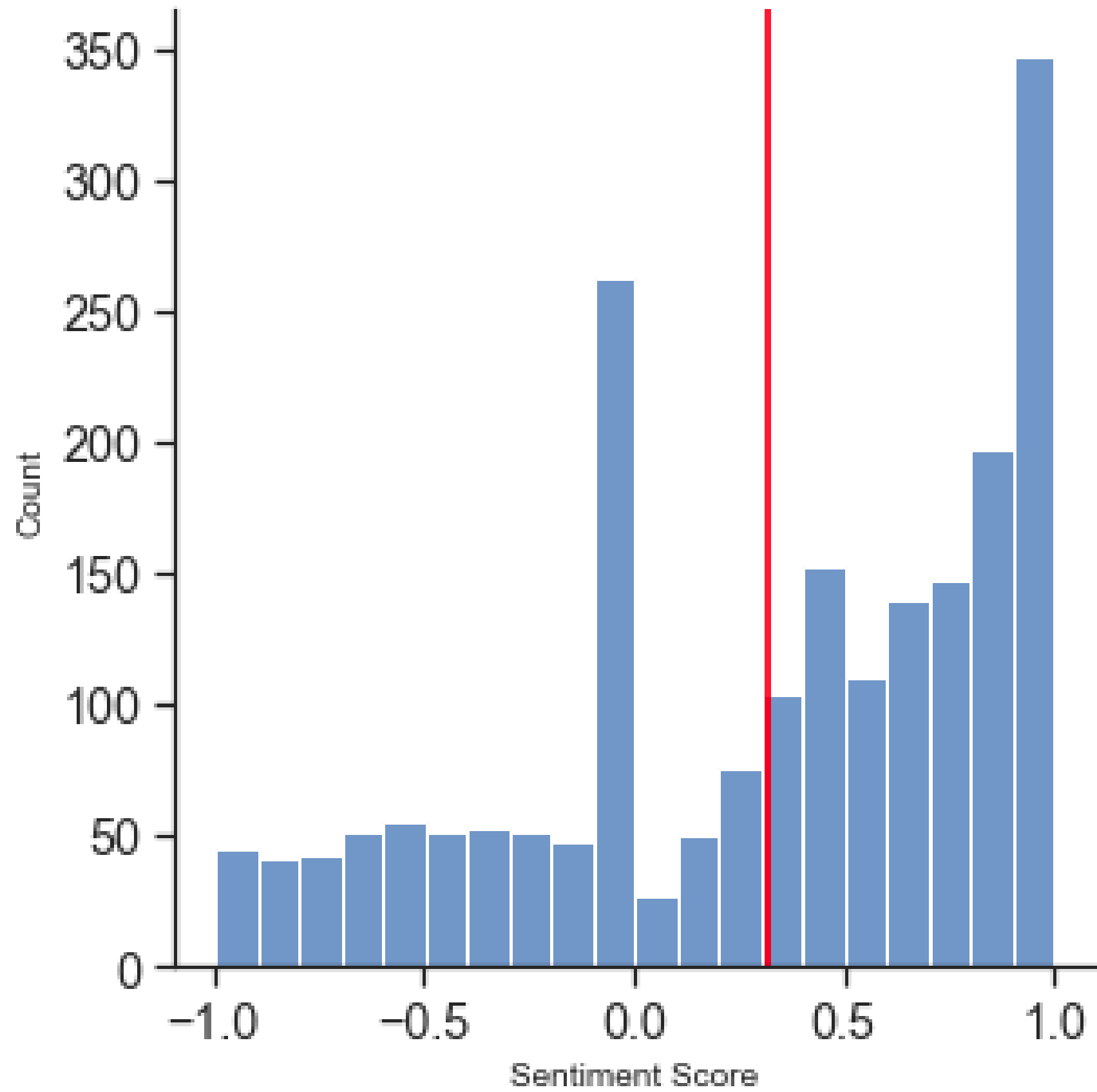
Most frequent words in AskMan



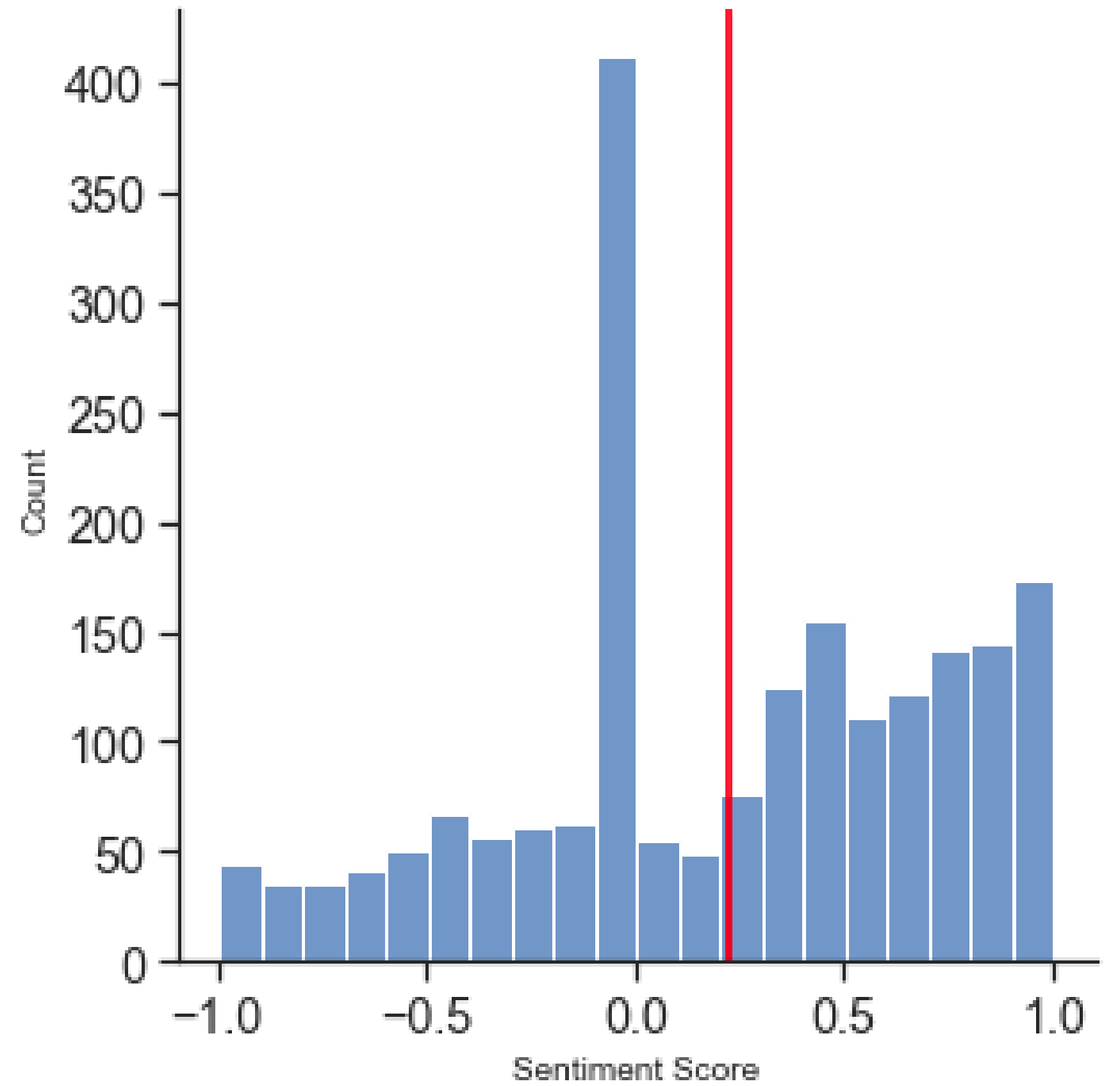
Most frequent words in AskWoman

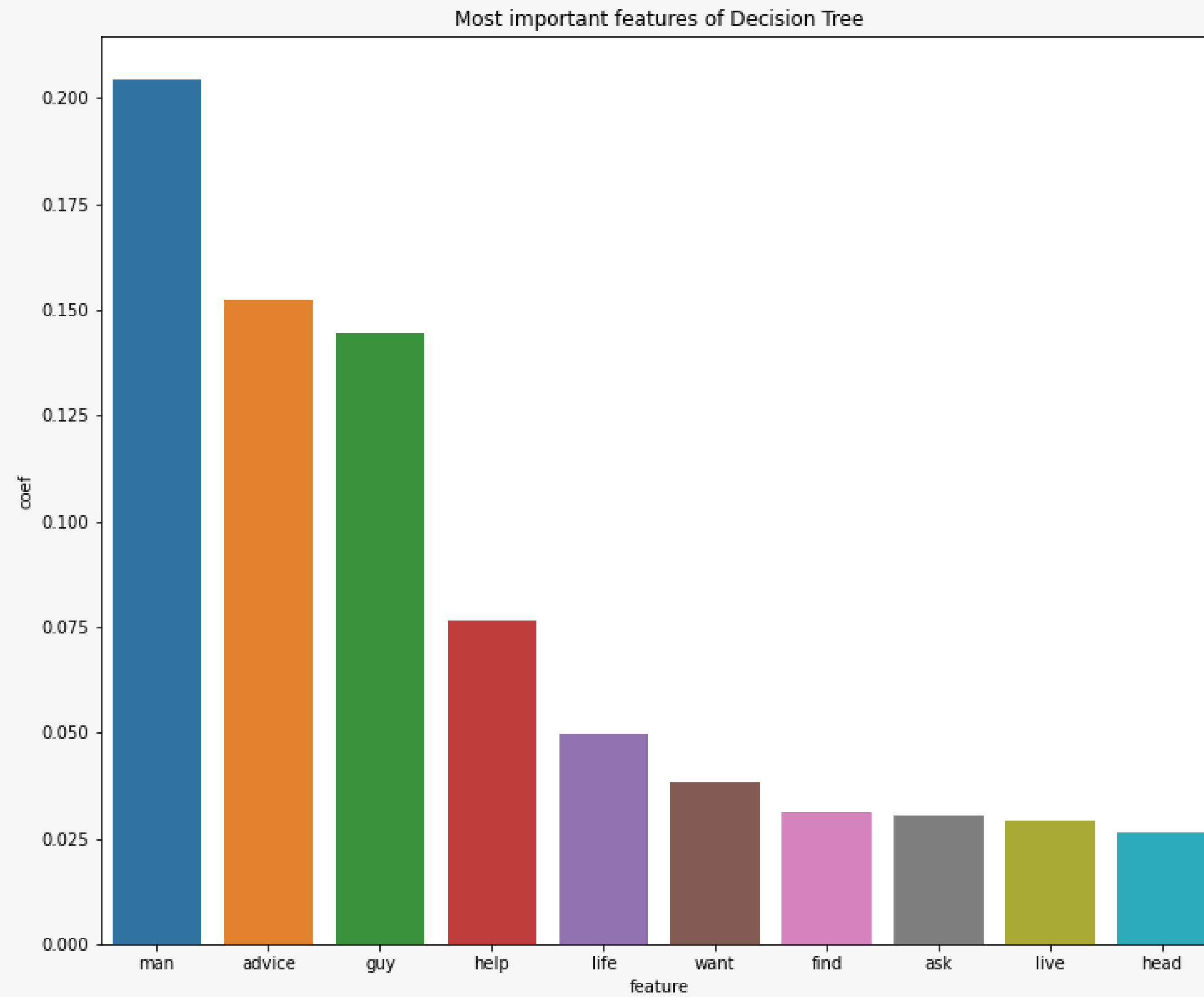


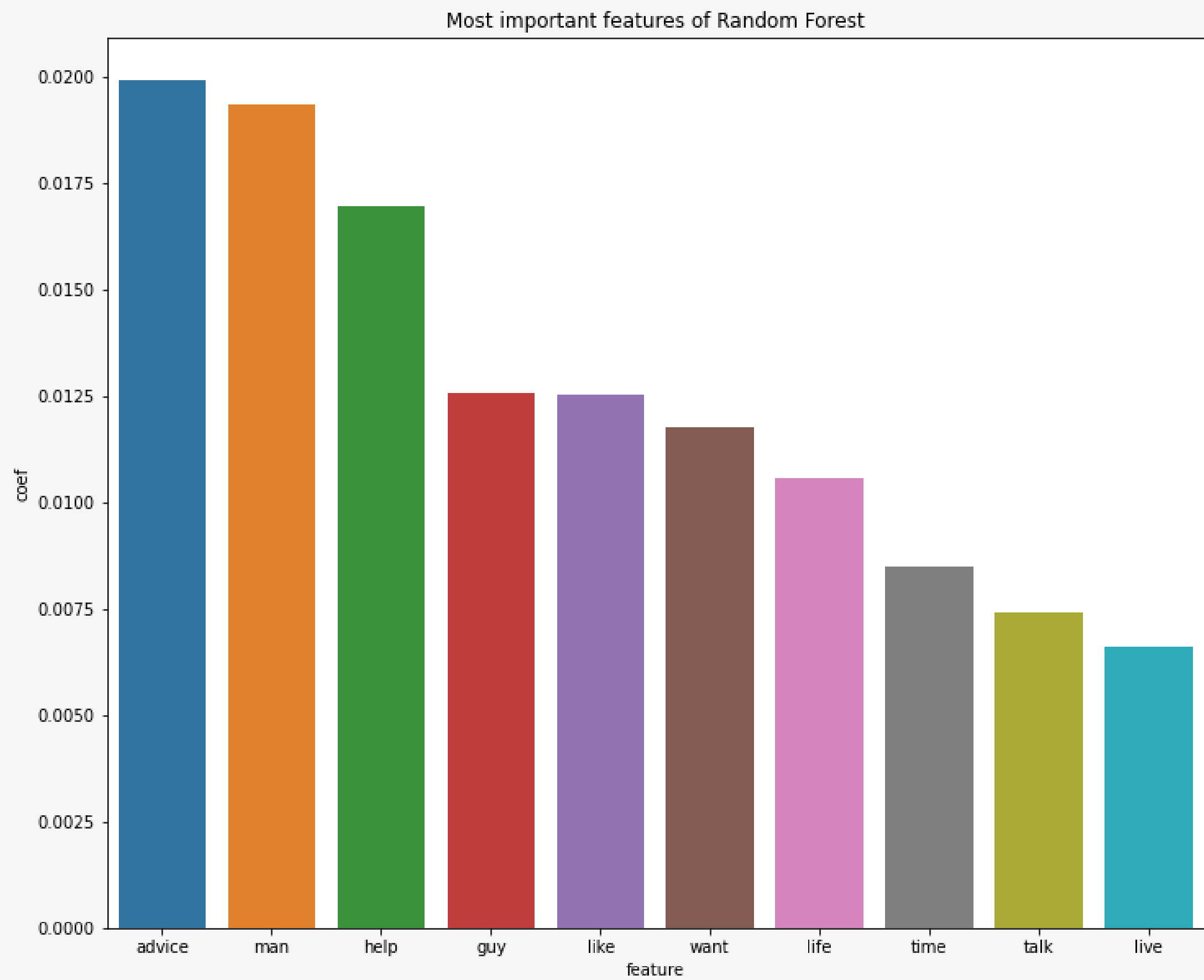
'AskMan' Sentiment Score Distribution

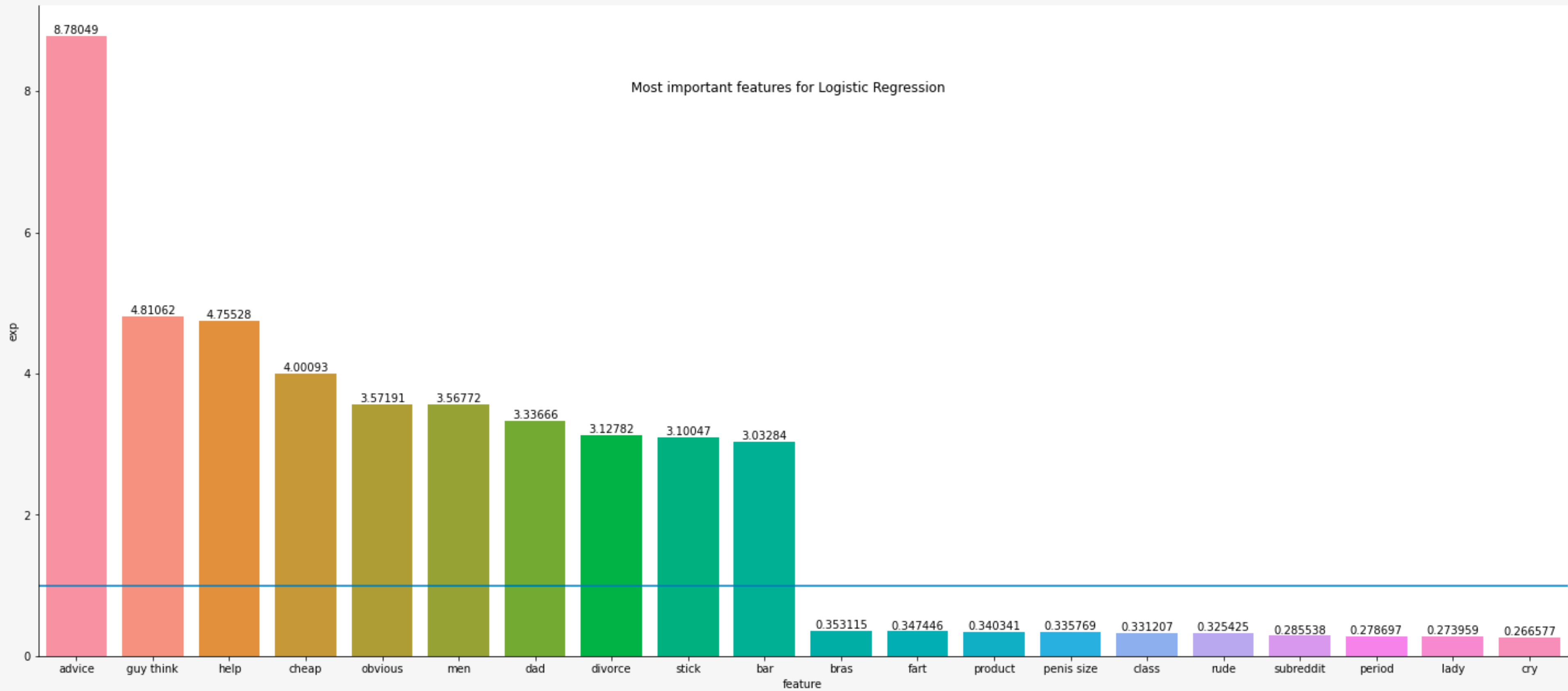


'AskWoman' Sentiment Score Distribution









AskWoman

Mean amount of words **21**

The longest sub **529** words

Mean amount of letters **122**

Max amount of letters **2702**

AskMan

Mean amount of words **30**

The longest sub **369** words

Mean amount of letters **167**

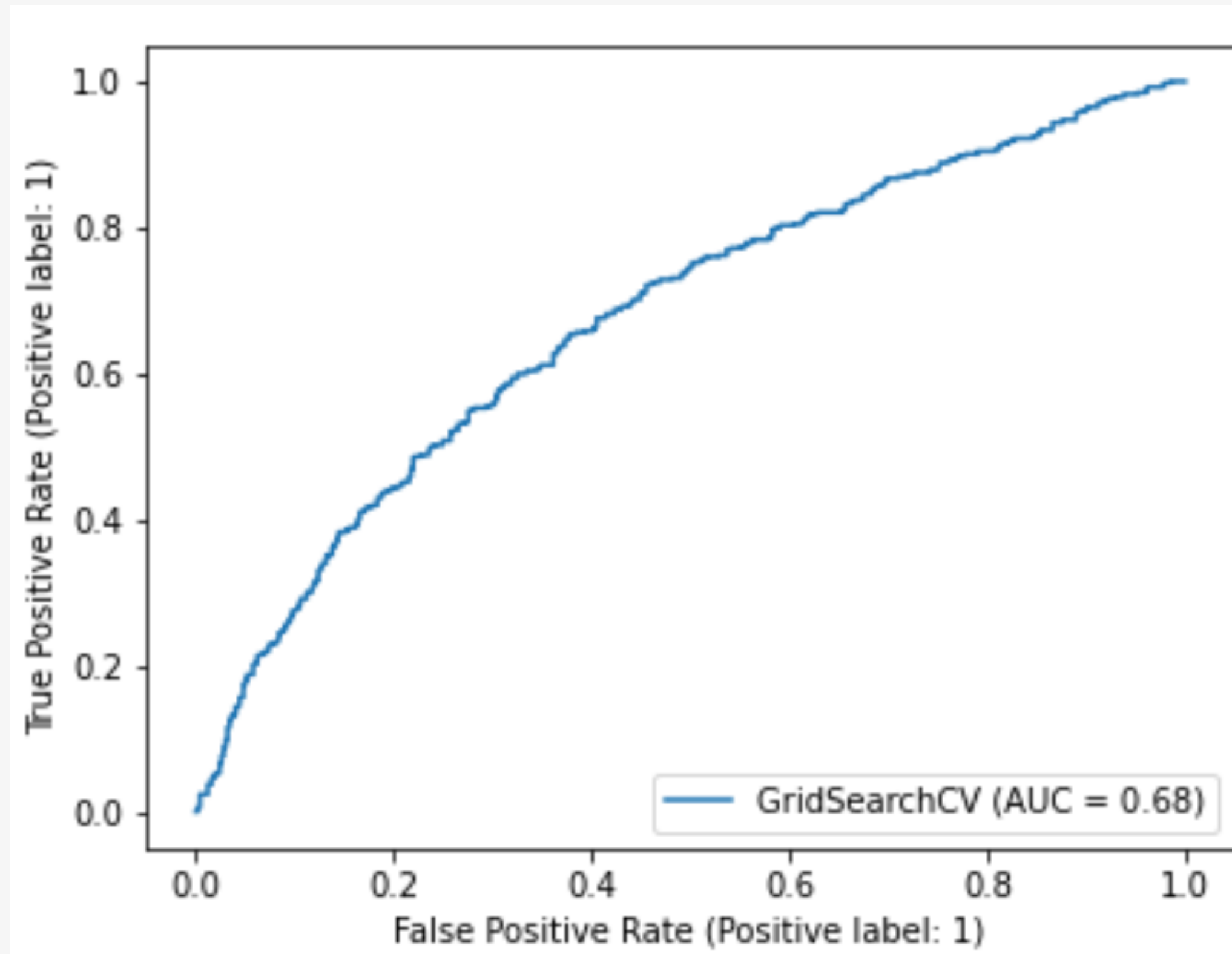
Max amount of letters **2062**

Top 3 models

Model	Accuracy	F1 score
BASELINE	0.5	0.504
AdaBoost	0.63	0.59
RandomForest	0.62	0.59
Naive Bayes	0.59	0.59



Logistic Regression ROC plot



Model	Submissions	Commentaries
AdaBoost	0.63	0.605
DecisionTree	0.62	0.579
Naive Bayes	0.59	0.635

Recommendations

**Model shows noticeable
better performance over
baseline model**



**Let's give it
a try!**



2 ways to improve results

**Dive deeper into the
dataset**

**Find new data which
would show more
differences between
man and woman**



Thanks!

Thanks!

Thanks!

Thanks!

Thanks!

Thanks!