



Letty Wu

Problem



- Reddit is one of the top 10 the most visited websites in the US
- 2019 Ad revenue – \$100 million
- Twitter generated more than 2.99 billion in advertising service revenues in 2019

Join Google AdSense!

- Big data with machine learning algorithm
- Using machine learning model to predict between subreddits, then use distinguished words to target advertising better.



Sample Data



Nutrition

r/nutrition

A subreddit for the discussion of nutrition science.



A place for nutrition professionals

r/dietetics

A place for current and future nutrition professionals to discuss all aspects of the profession.

Workflow

Data Collection

Used Pushshift API to pull posts from subreddits

1

Data Cleaning & EDA

Removed Outliers, used sklearn to preprocess text data

2

Modeling

Tried 13 different classifiers, and did Gridsearch on 3 classifiers

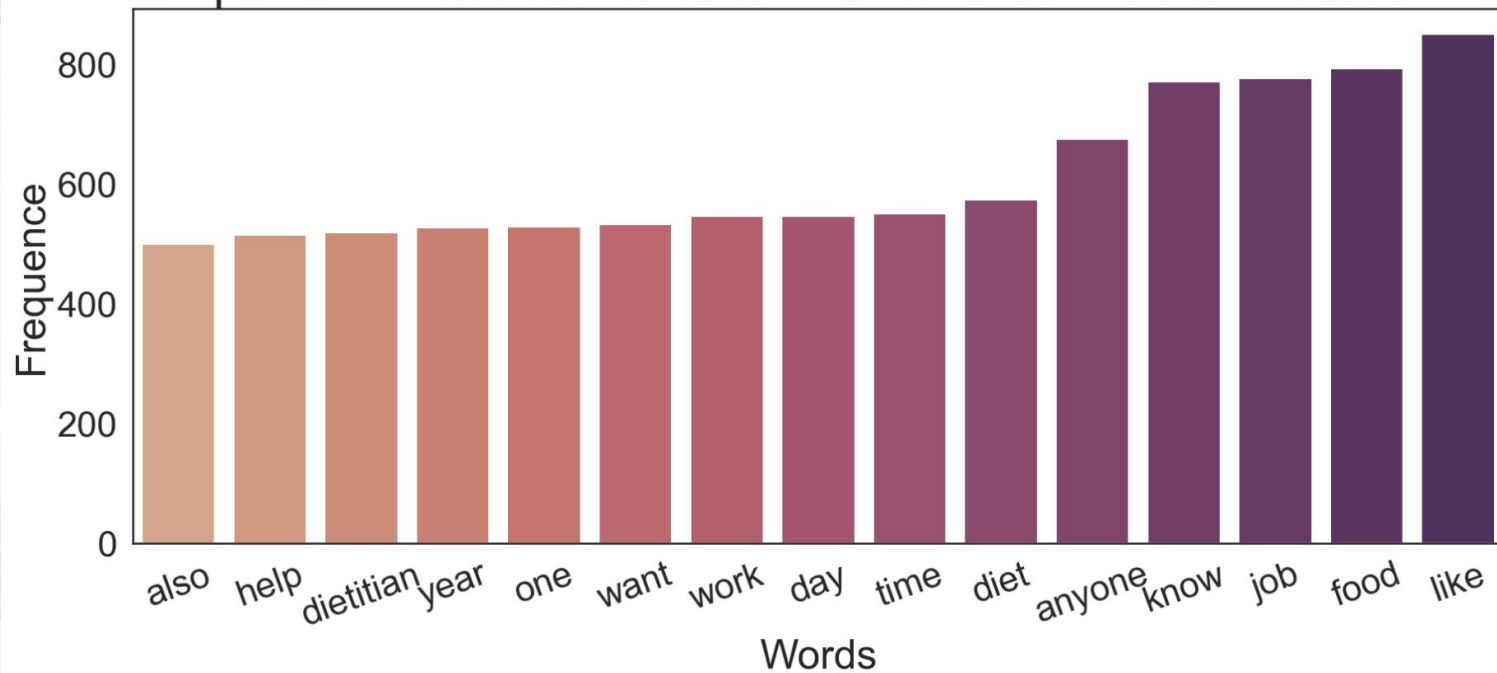
3

Evaluation

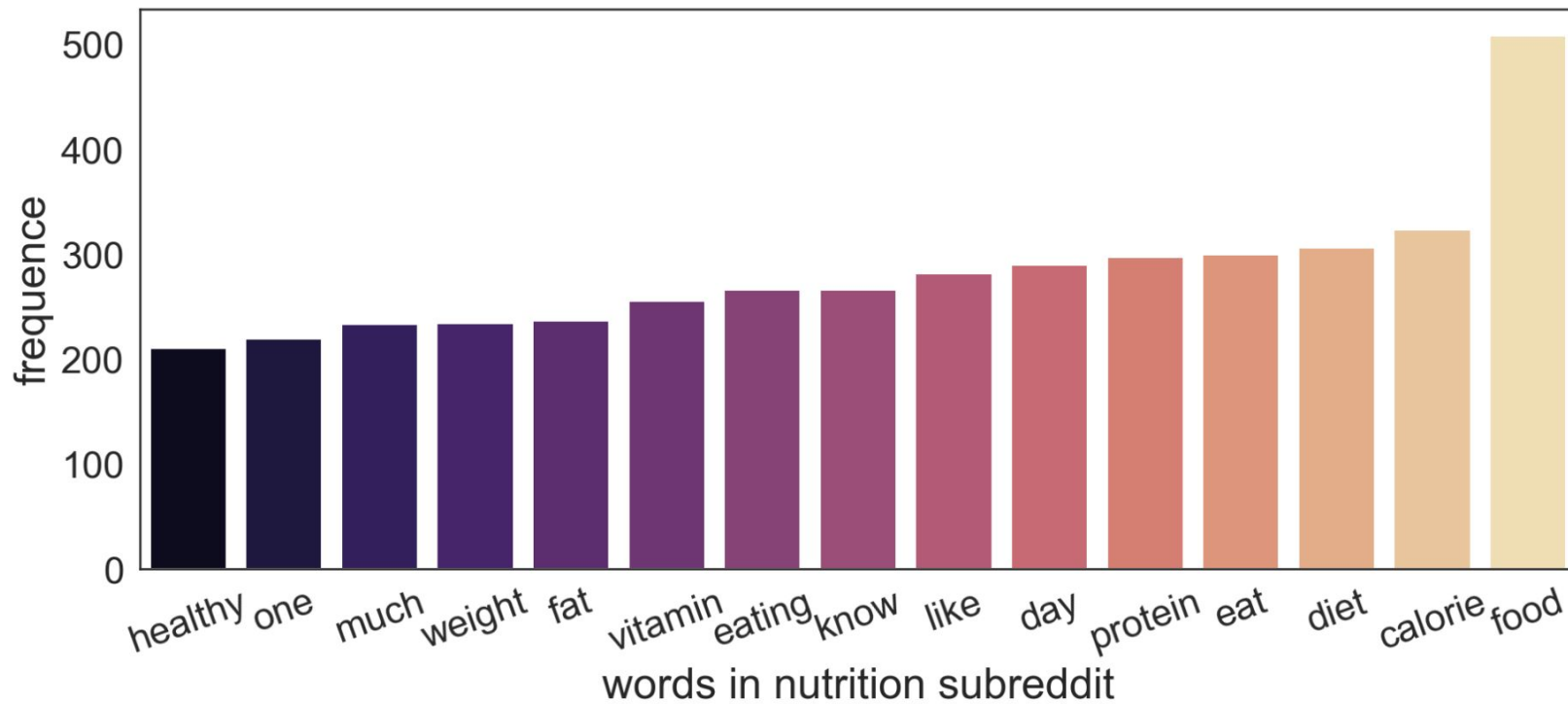
Used accuracy as the main metric

4

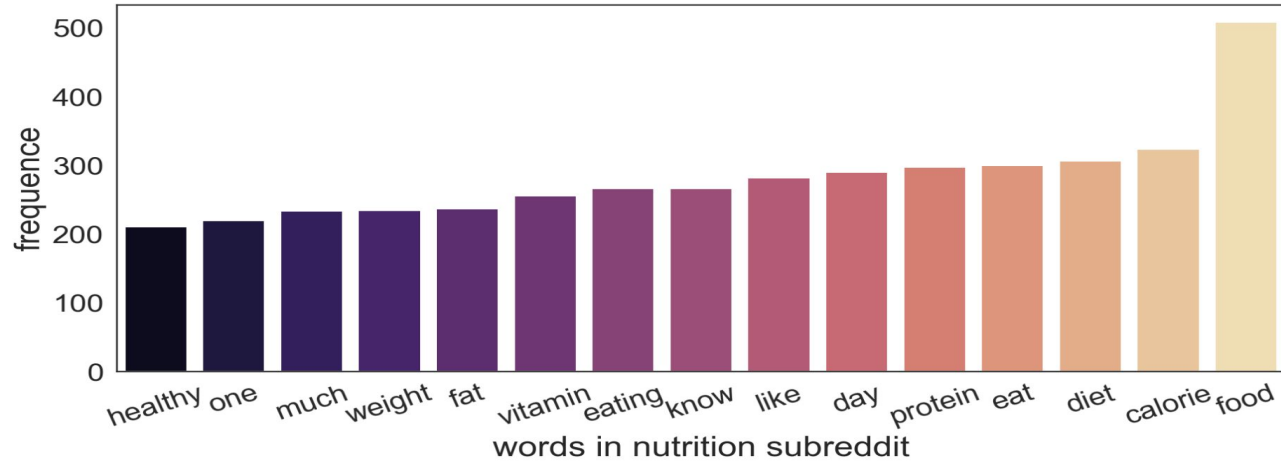
Top 15 Common Words In Nutrition and Dietetics Subreddits



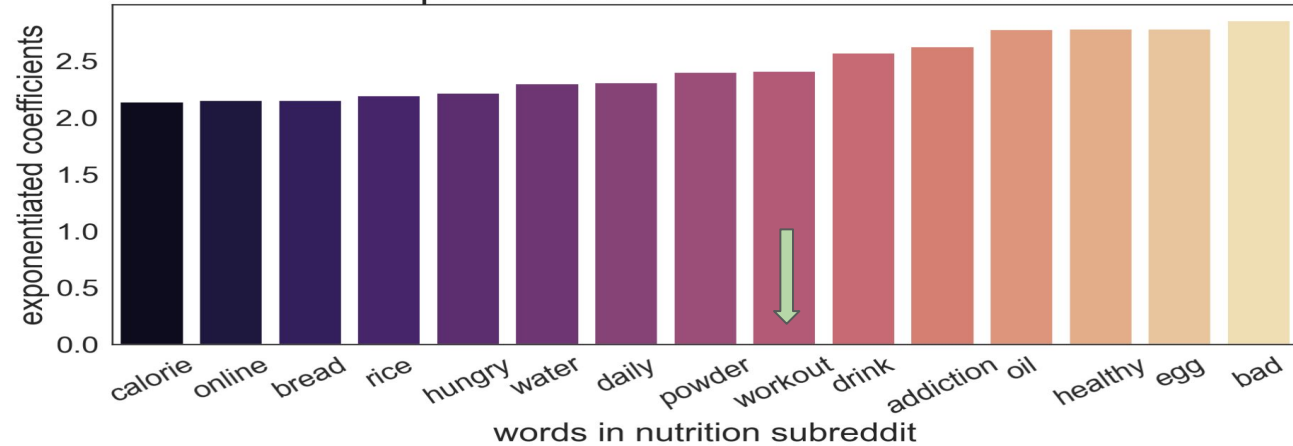
15 Common Words In Nutrition Subreddit



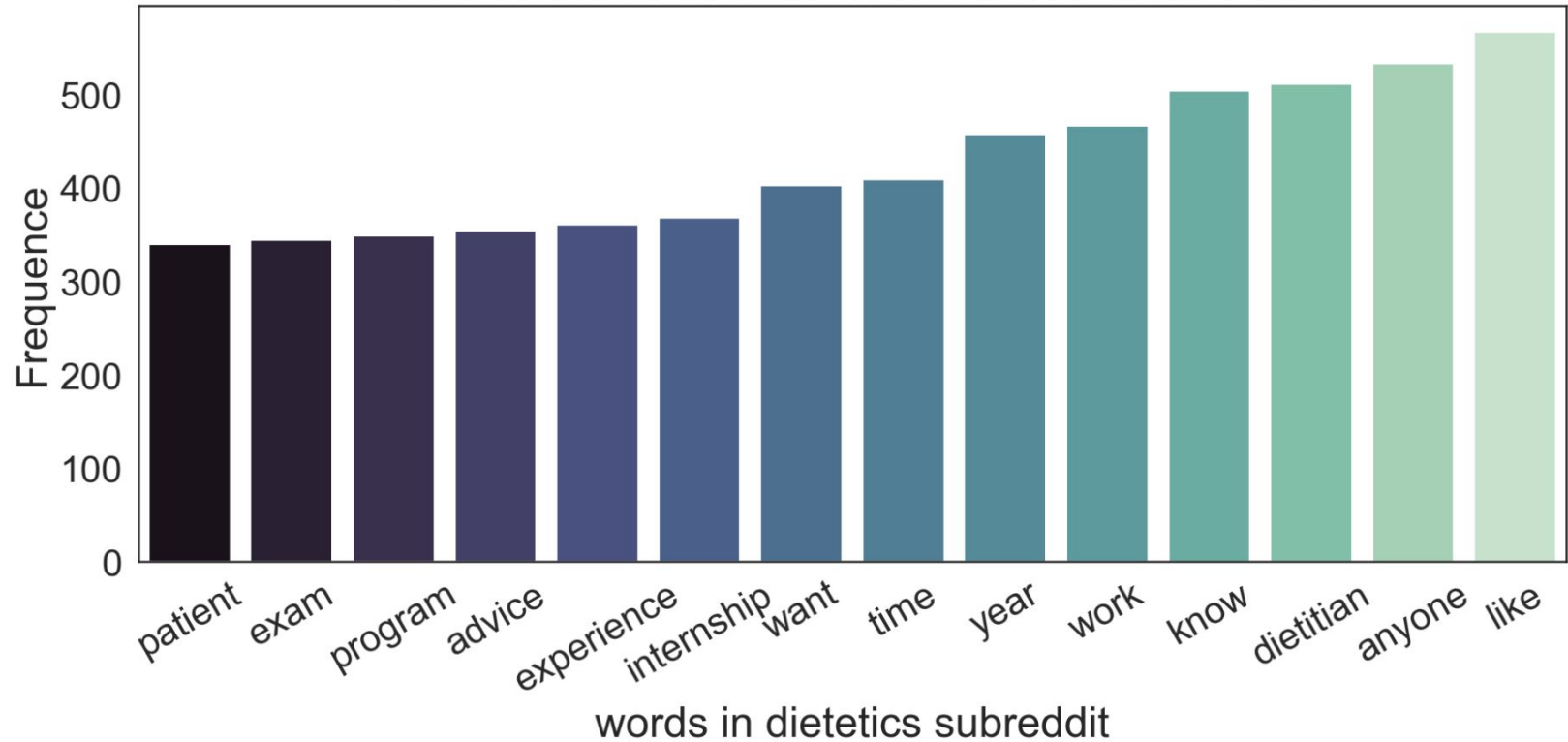
15 Common Words In Nutrition Subreddit



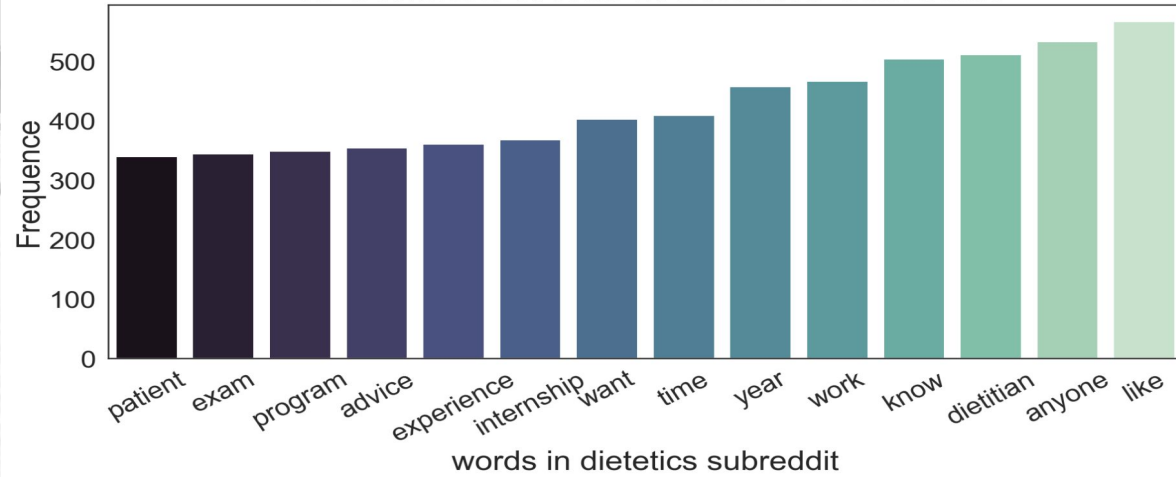
15 Important Words In Nutrition Subreddit



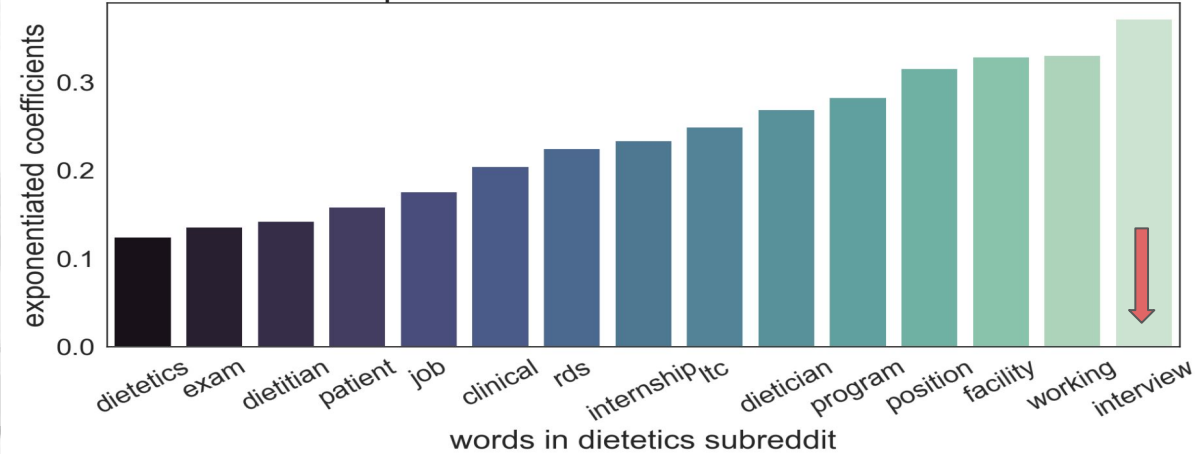
15 Common Words In Dietetics Subreddit



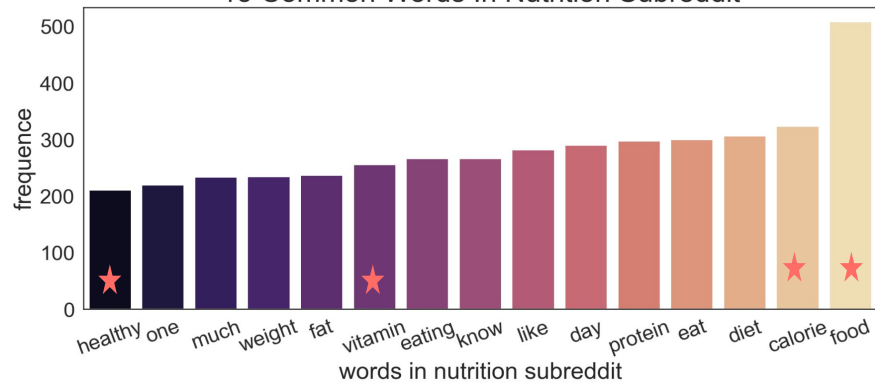
15 Common Words In Dietetics Subreddit



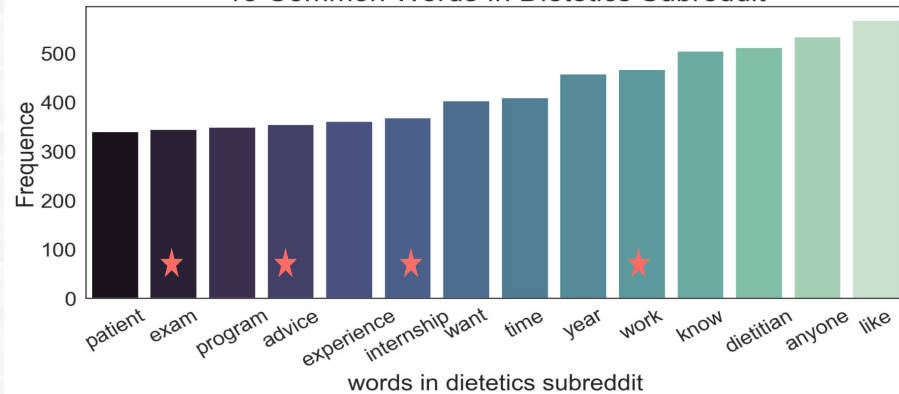
15 Important Words In Dietetics Subreddit



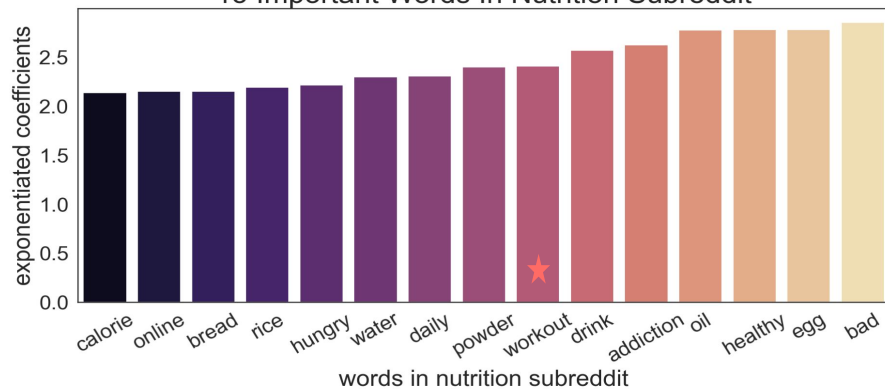
15 Common Words In Nutrition Subreddit



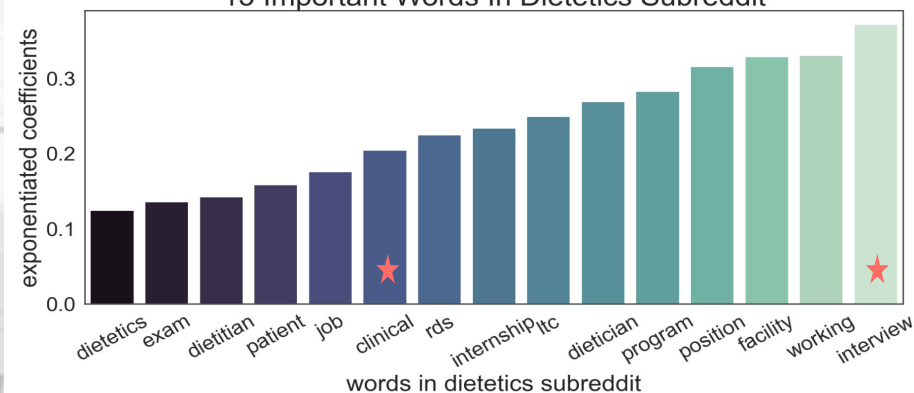
15 Common Words In Dietetics Subreddit



15 Important Words In Nutrition Subreddit



15 Important Words In Dietetics Subreddit



Model Accuracy

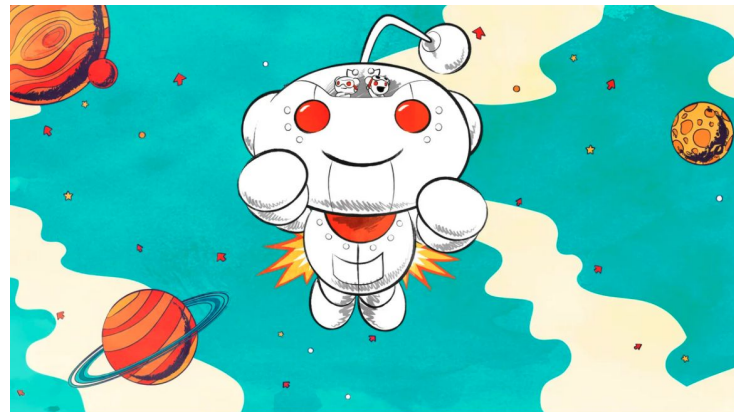
Multinomial Naive Bayes 90.08%

Voting Classifier 89.13%

LogisticRegression 88.22%

Conclusion

- Google AdSense has variety of advertisers
- Automate the process and customize the strategy for placing ads
- Join Google AdSense!





Thank you!
