# Reddit Content Analysis

## Group 22

| | |
|---|---|
| Markus Ritzer | Data Preprocessing, Submission Clustering |
| Thomas Aumüller | Community and Throw-away Analysis |
| Oana-Emilia Bodirnea | Sentiment Analysis |
| Gregor Yinan Potthast | Sentiment Analysis |

# Motivation

Analyze recent activities in r/conspiracy subreddit

- Topic Clustering
  - Comparison of TFIDF and Word Embeddings
  - Presentation of Clusters as Word Clouds
- Community Analysis
  - Track activity of users
  - detect Throw-Away accounts
- Sentiment Analysis
  - Detect positive/negative/neutral sentiment of
    - submissions
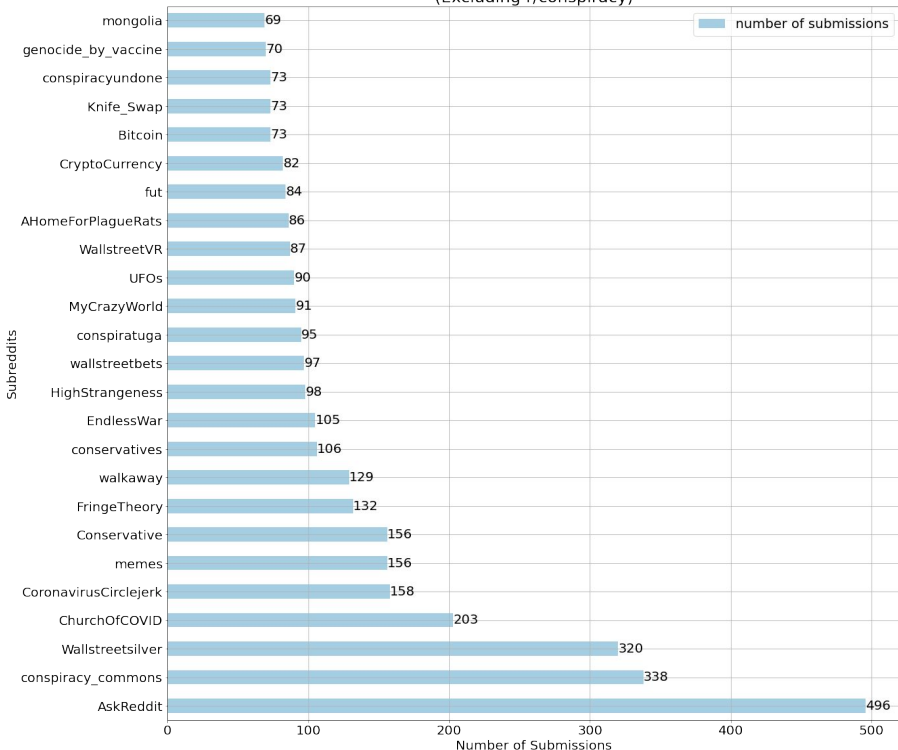    - comments

# Data and Methods

## Data Crawling

- PRAW Library
  - 1k latest/hottest submissions
  - unique authors
  - comments
- OpenAI
  - Word Embeddings
- Hugging Face
  - Sentiment analysis

## Methods

- TFIDF
- Word embeddings
- K-Means clustering
- Mannwitneju significance test
- General preprocessing
  - spell correction
  - stopword removal
  - punctuation removal
  - lemmatization

# Results: Topic Clustering

Submission to Vector via average pooling of

- TFIDFs (TfidfVectorizor)
- Word Embeddings (OpenAI text-embedding-ada-002 model)

Clustering with K-Means (k = 30)

- Averaging pairwise cosine similarities within clusters
- Averaging over cluster similarities

TFIDFs: **8.36 %**

Word Embeddings: **97.10 %**

# Results: Community Analysis / Throwaway Analysis

Top 25 Subreddits from the Community of r/conspiracy
(Excluding r/conspiracy)

| Subreddit | Number of Submissions |
|---|---|
| mongolia | 69 |
| genocide_by_vaccine | 70 |
| conspiracyundone | 73 |
| Knife_Swap | 73 |
| Bitcoin | 73 |
| CryptoCurrency | 82 |
| fut | 84 |
| AHomeForPlagueRats | 86 |
| WallstreetVR | 87 |
| UFOs | 90 |
| MyCrazyWorld | 91 |
| conspiratuga | 95 |
| wallstreetbets | 97 |
| HighStrangeness | 98 |
| EndlessWar | 105 |
| conservatives | 106 |
| walkaway | 129 |
| FringeTheory | 132 |
| Conservative | 156 |
| memes | 156 |
| CoronavirusCirclejerk | 158 |
| ChurchOfCOVID | 203 |
| Wallstreetsilver | 320 |
| conspiracy_commons | 338 |
| AskReddit | 496 |

number of submissions

Throwaway Account(TA) Def.:

*First Submission of the Redditor was in the observed Subreddit (r/conspiracy).*
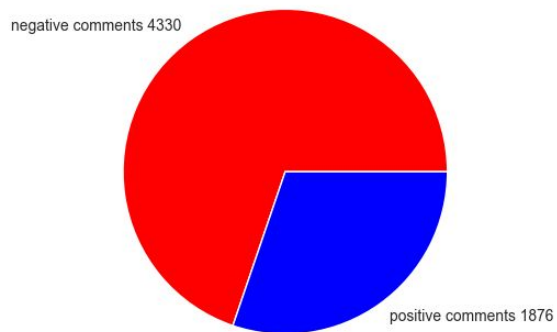
Results:

**38.57%** of our user base potentially used throwaway accounts.

There was <u>no significant difference in user-activity</u> when comparing normal users and TAs with MWU-Test.
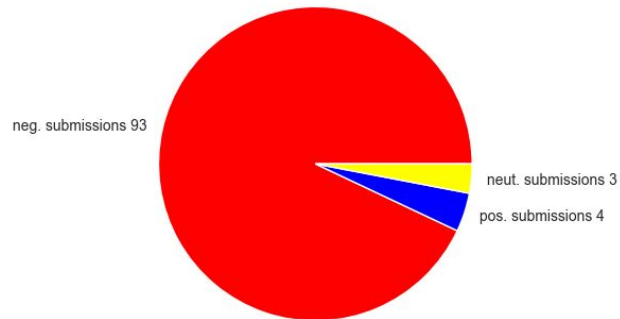
# Results: Sentiment Analysis

- Analysis of the comments from the top 100 submissions with the most comments

Sentiment analysis of all comments

negative comments 4330

positive comments 1876

Sentiment analysis of 100 submissions

neg. submissions 93

neut. submissions 3

pos. submissions 4

# Conclusion

## Conclusion

Data collection framework allows variety of analyses
Word Embeddings working better than TFIDF
Found activity behaviour and sentiment

## Future work

adapt data collection framework for more precise retrieval

## Limitations/Biases

- find best k for Clustering
- obtain adequate word embeddings
- more precise preprocessing (e.g. remove links)
- PRAW retrieval limits
- OpenAI retrieval limits