



An Analysis of Domestic Violence Using Natural Language Processing

Elie Park

June 22, 2022

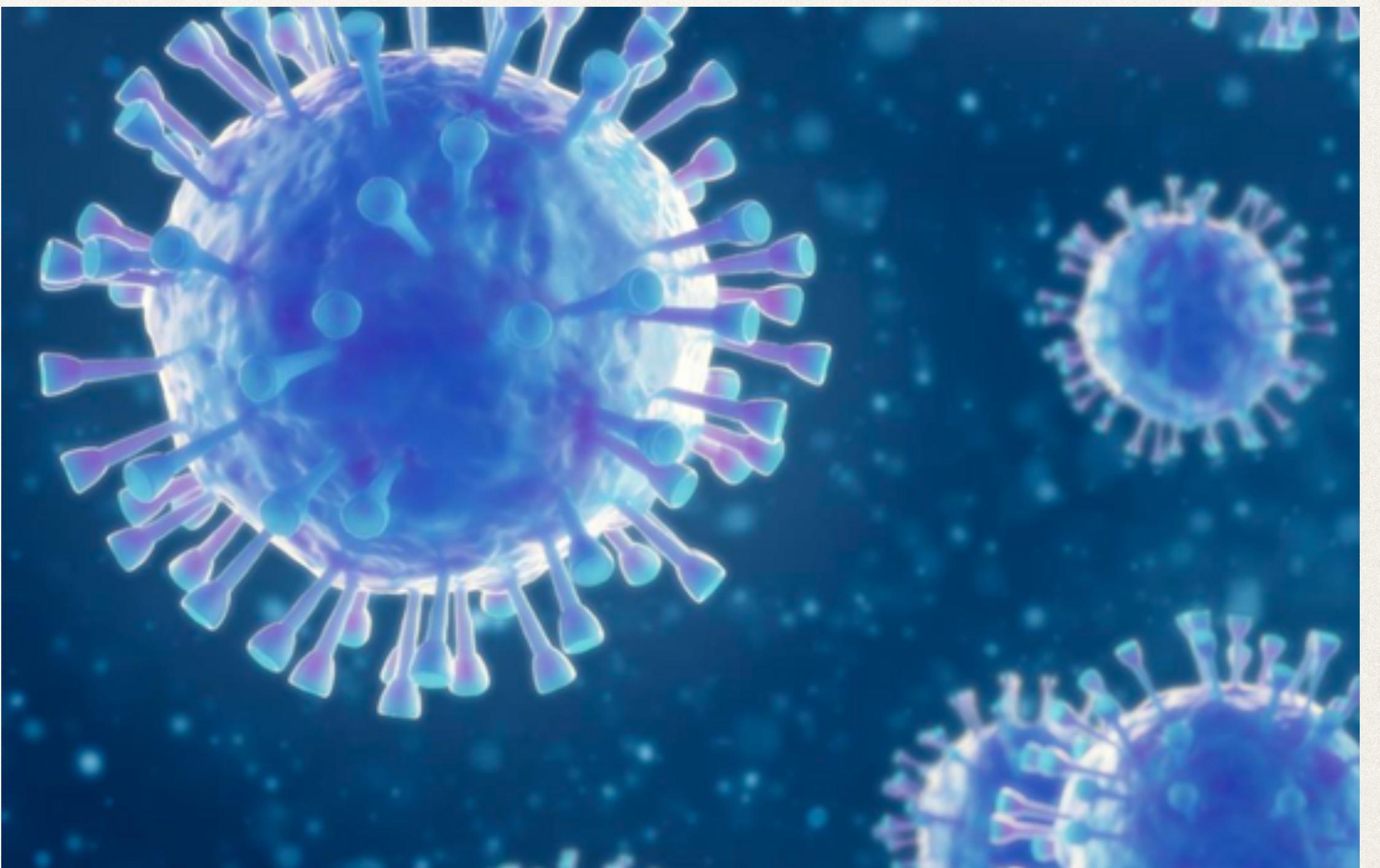
47% of women and
6% of men

who were victims of homicides were killed by an intimate partner in 2019.



71,563 calls

were received by Canada's Assaulted Women's Helpline from April to September 2020. (36,362 calls within the same months in 2019)



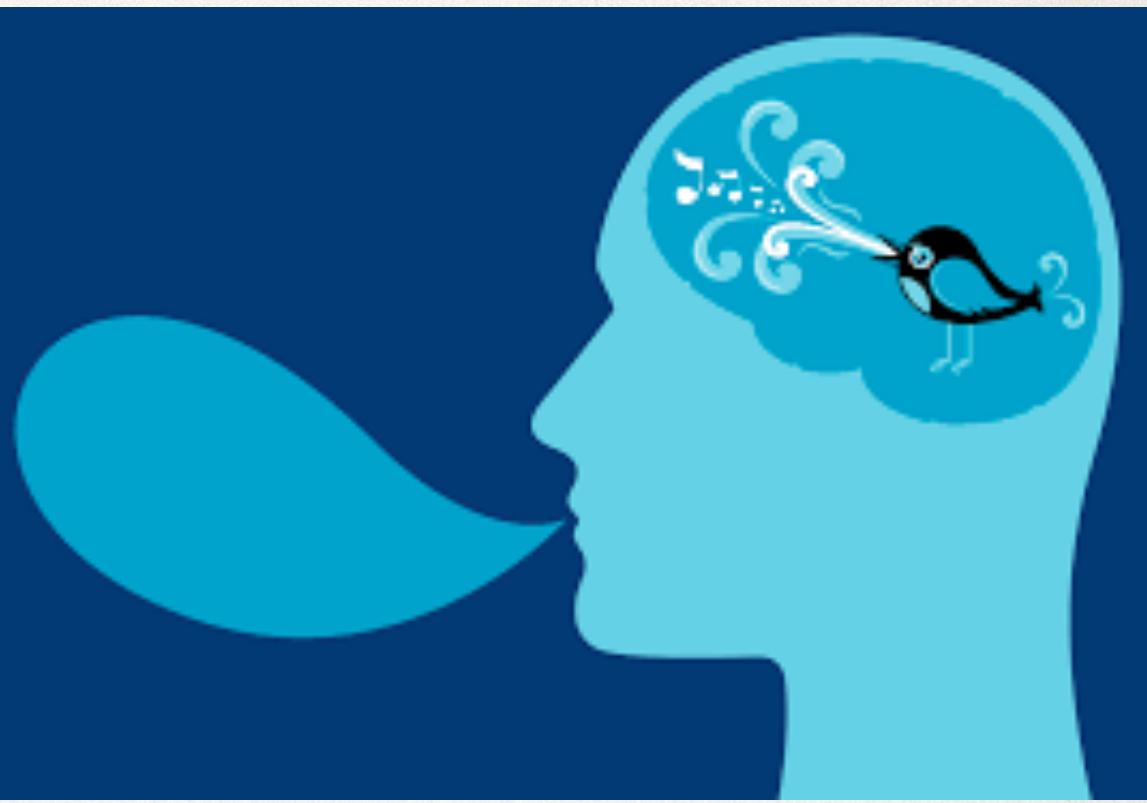
12%

is the increase in calls to police related to domestic violence between March and June of 2020, compared to the same time period in 2019.



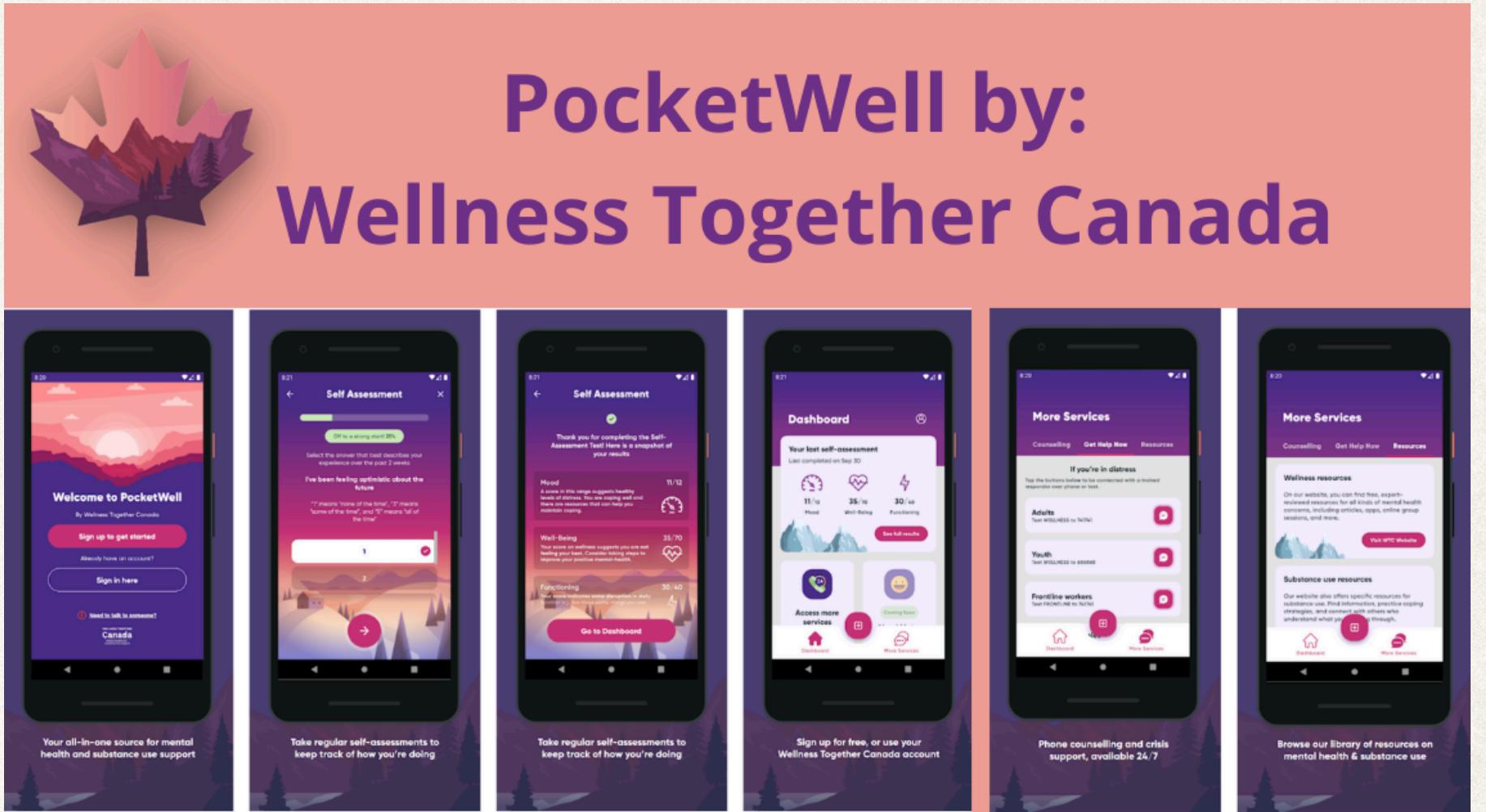
Uses of Social Media for Support

- ✿ 51.4% using social media for support
- ✿ Sense of community
- ✿ Raising awareness & combatting stigma
- ✿ Safe space for expression
- ✿ Coping and empowerment



Uses of Online Mental Health Apps & Resources for Support

- ✿ 32.6% using mental health apps for support
- ✿ Wellness Together Canada (Apr 2020) & PocketWell (Jan 2022) were launched by the Canadian Government
- ✿ 7 Cups has provided emotional support for over 53 million people since its launch in 2013



Natural Language Processing on Reddit and Twitter Data

- ❖ These social media platforms allow for accessible data collection of discussions
- ❖ Tokenizing important words and converting them into numerical vectors allows for a prediction of narratives
- ❖ Reddit: abuse vs. non-abuse
- ❖ Twitter: #WhyIStayed vs. #WhyILeft



Reddit Data



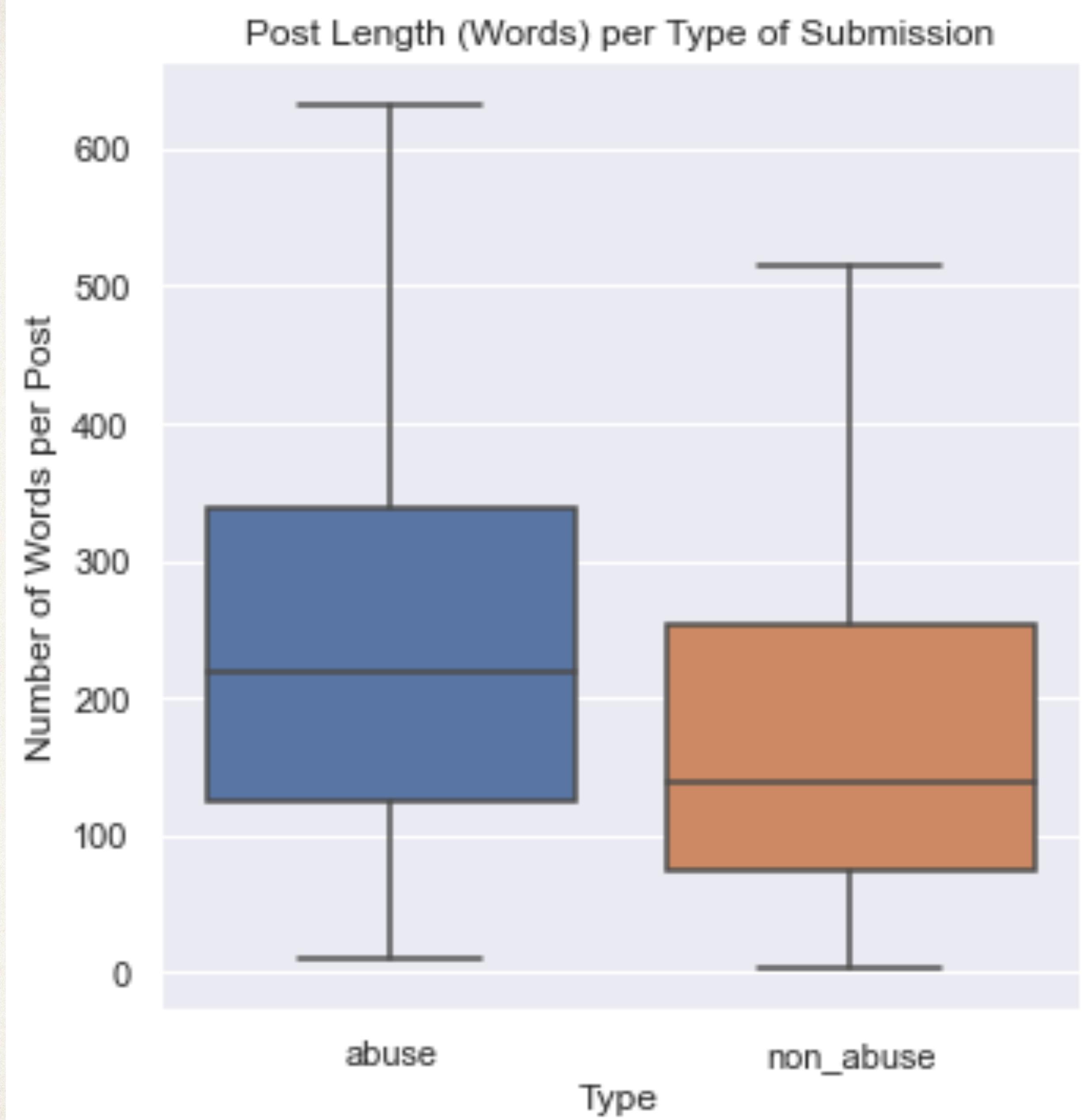
- ✿ Datasets by Nicholas Schrading: <http://www.nicschrading.com/data/>
- ✿ Shelved set of data for Reddit - used even submissions and comments dataset
- ✿ 552 abuse & 552 non-abuse posts, at least 1 comment under each submission
- ✿ Abuse: r/AbuseInterrupted, r/DomesticViolence, r/SurvivorsofAbuse
- ✿ Non-abuse: r/CasualConversation, r/Advice, r/Anxiety, r/Anger

Emotions Sensor Data

- ✿ <http://www.kaggle.com/datasets/iwilldoit/emotions-sensor-data-set>
- ✿ 1104 widely used words, corresponding scores of 7 emotions: disgust, surprise, neutral, anger, sad, happy, and fear
- ✿ Results to be used only as reference, complete dataset containing over 20,000 words should be used in future analyses

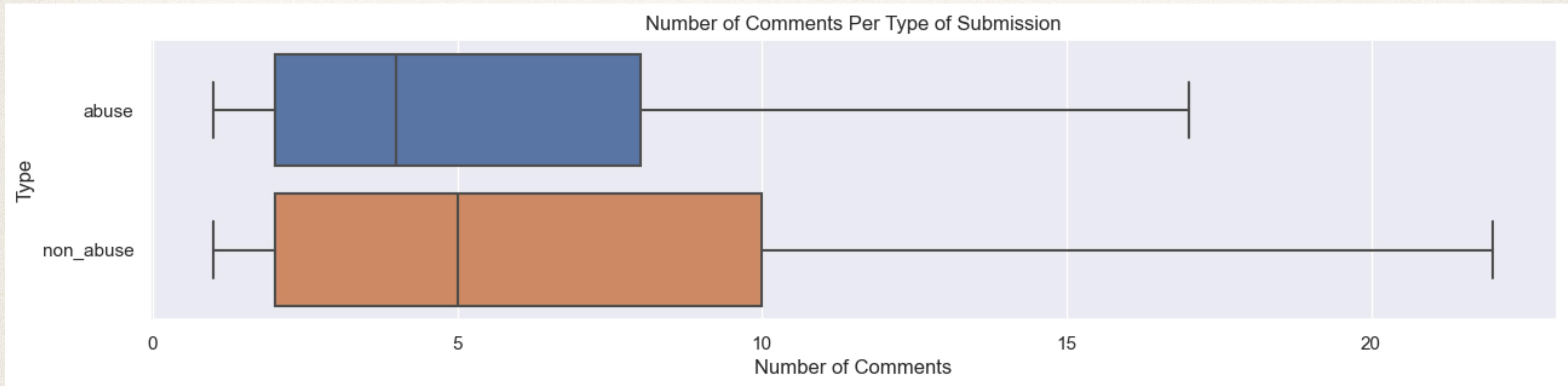
Reddit Data: Post Length

- Posts and comments related to abuse contain more characters and words ($p << 0.001$)
- Average word length in abuse posts and comments ($p << 0.001$)
- Non-abuse posts and comments were easier to read than abuse posts and comments ($p << 0.001$)
- More capital letters in non-abuse posts and comments ($p < 0.01$)



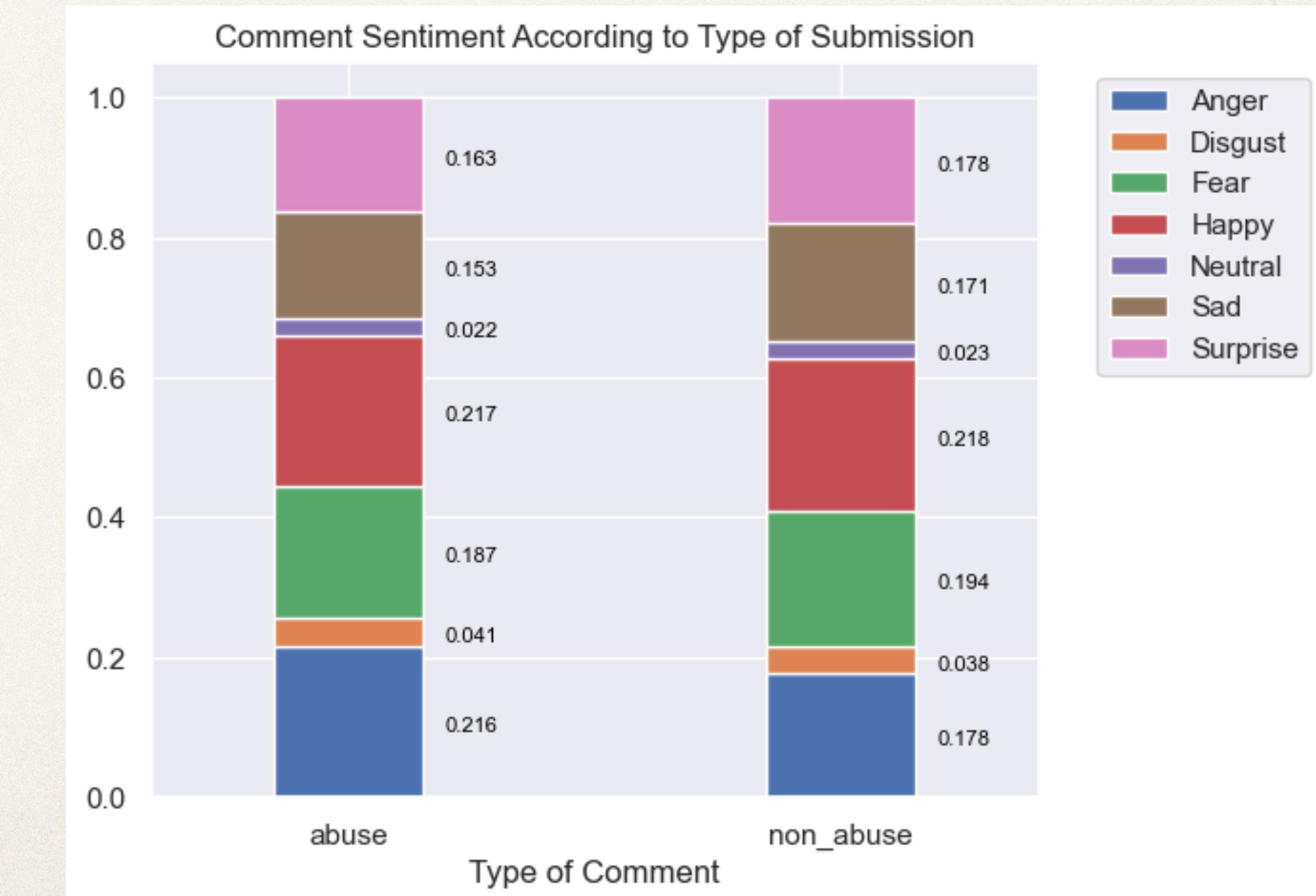
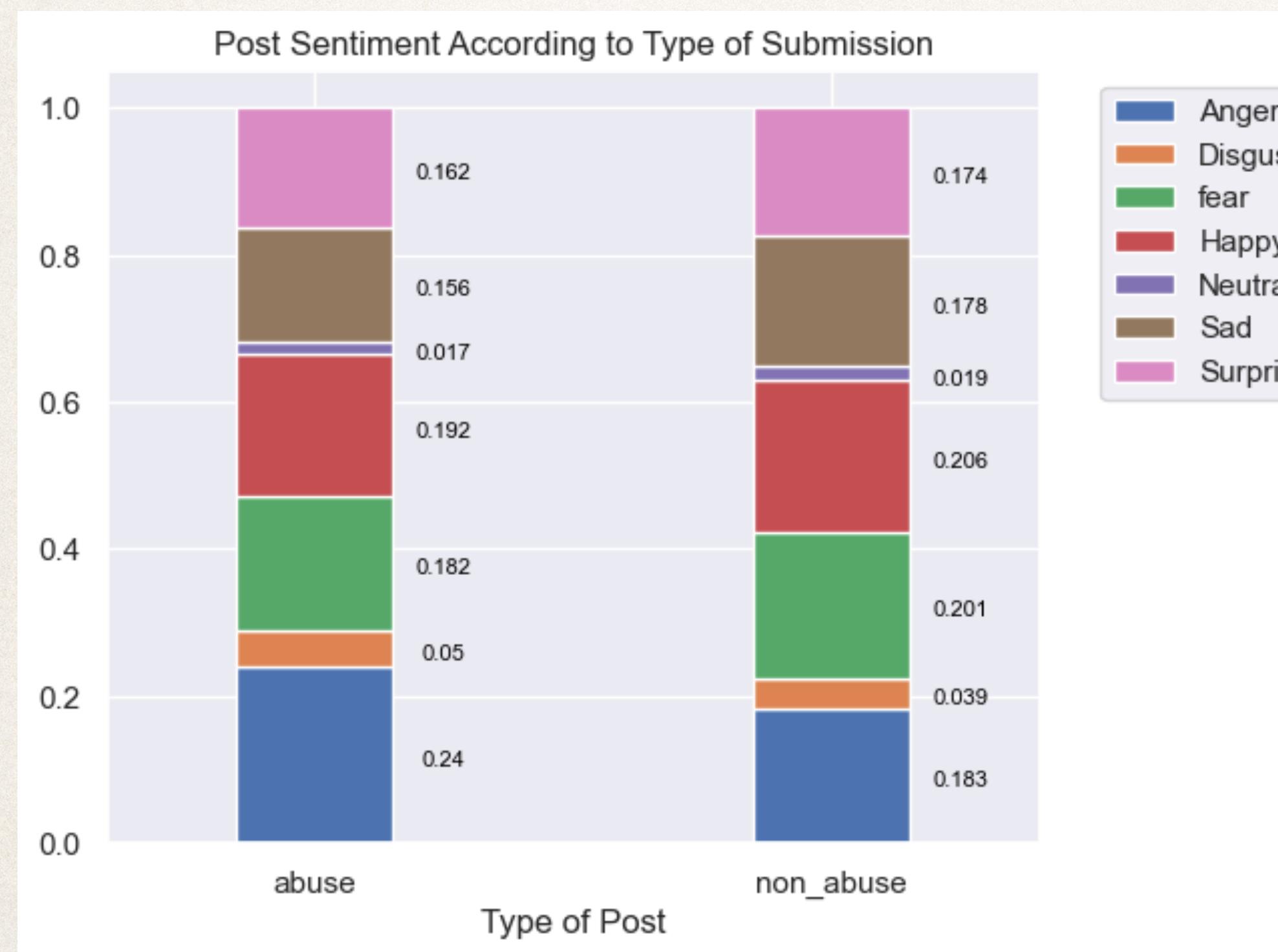
Reddit Data: Number of Comments per Post

- ✿ At least 1 comment per post
- ✿ Posts about abuse have less comments ($p << 0.001$)



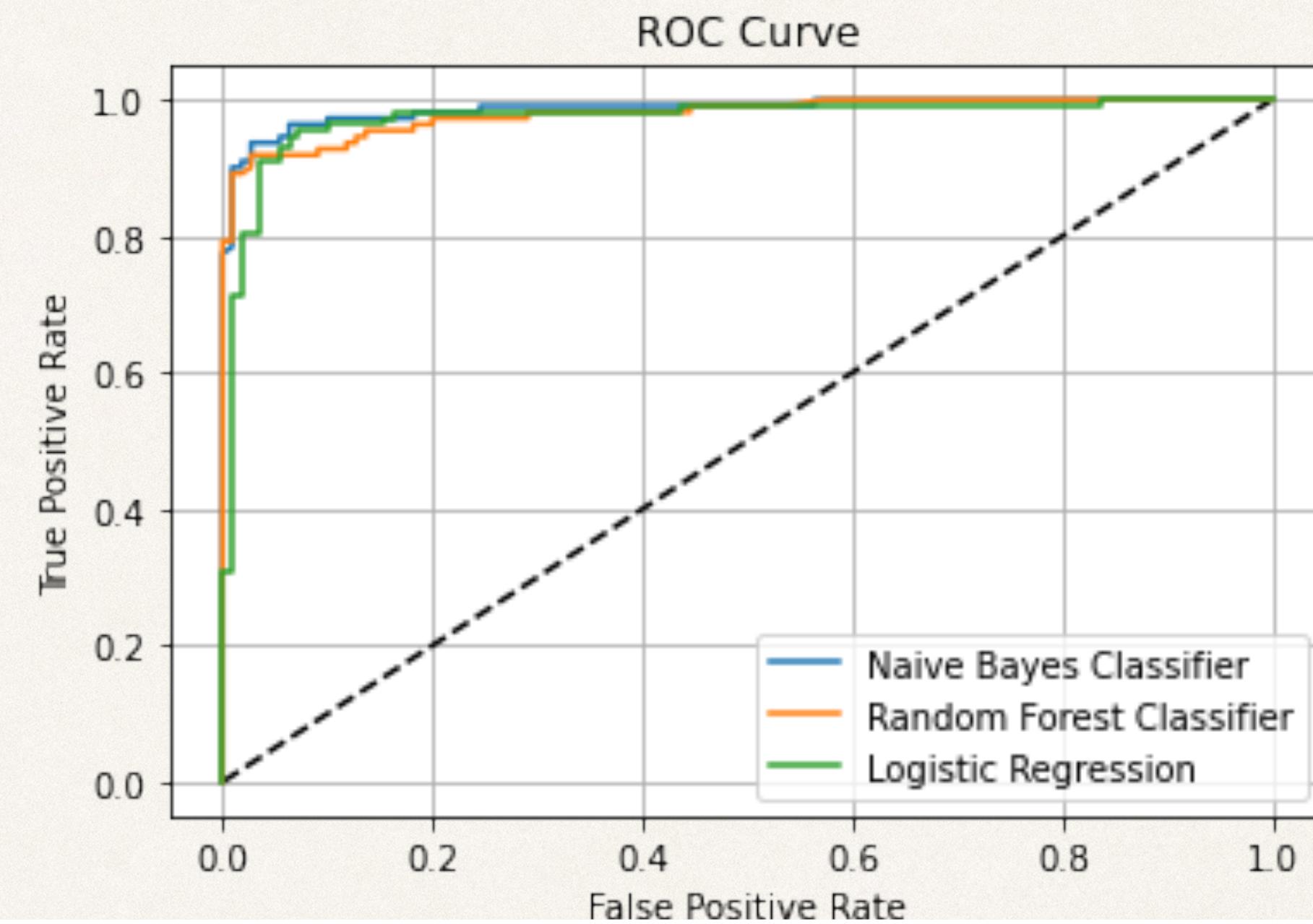
Reddit Data: Relative Frequencies of Top Sentiments

- ❖ Sentiment analysis using mean scores of emotions corresponding to tokens



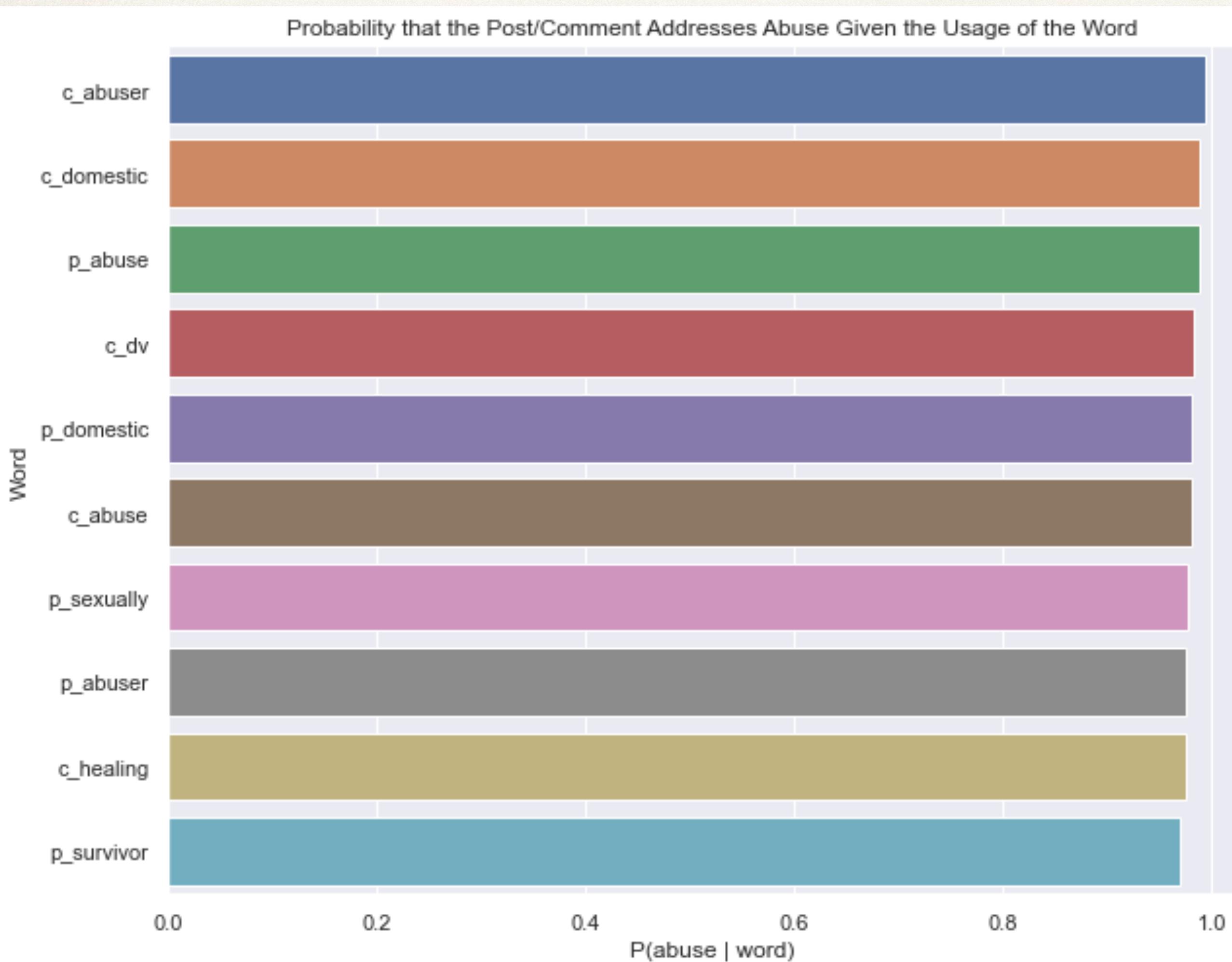
Reddit Data: Modeling

- 19518 tokens
- Numbers,
punctuations,
insignificant words
removed
- Standardized urls
- Lemmatization



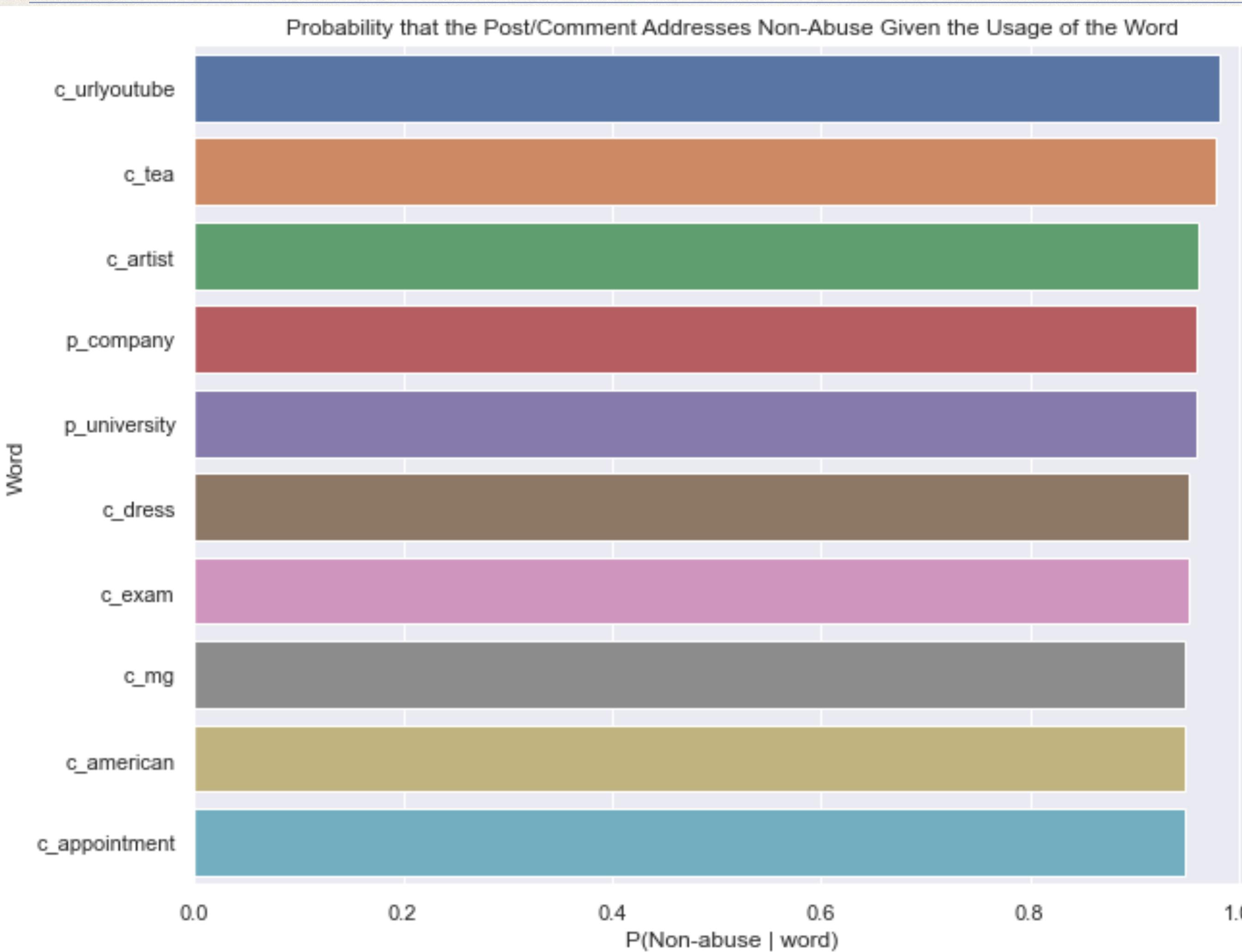
| Model | Best min_df for CountVectorizer | Best Parameters | ROC-AUC Score |
|--------------------------|---------------------------------|------------------------------------------------------------------------------------------|---------------|
| Naive Bayes Classifier | 1 | {'alpha': 1} | 0.987 |
| Random Forest Classifier | 14 | {'criterion': 'entropy', 'max_depth': None, 'max_features': 'log2', 'n_estimators': 700} | 0.979 |
| Logistic Regression | 1 | {'C': 0.1} | 0.973 |

Reddit Data: Feature Importances for Abuse Posts and Comments



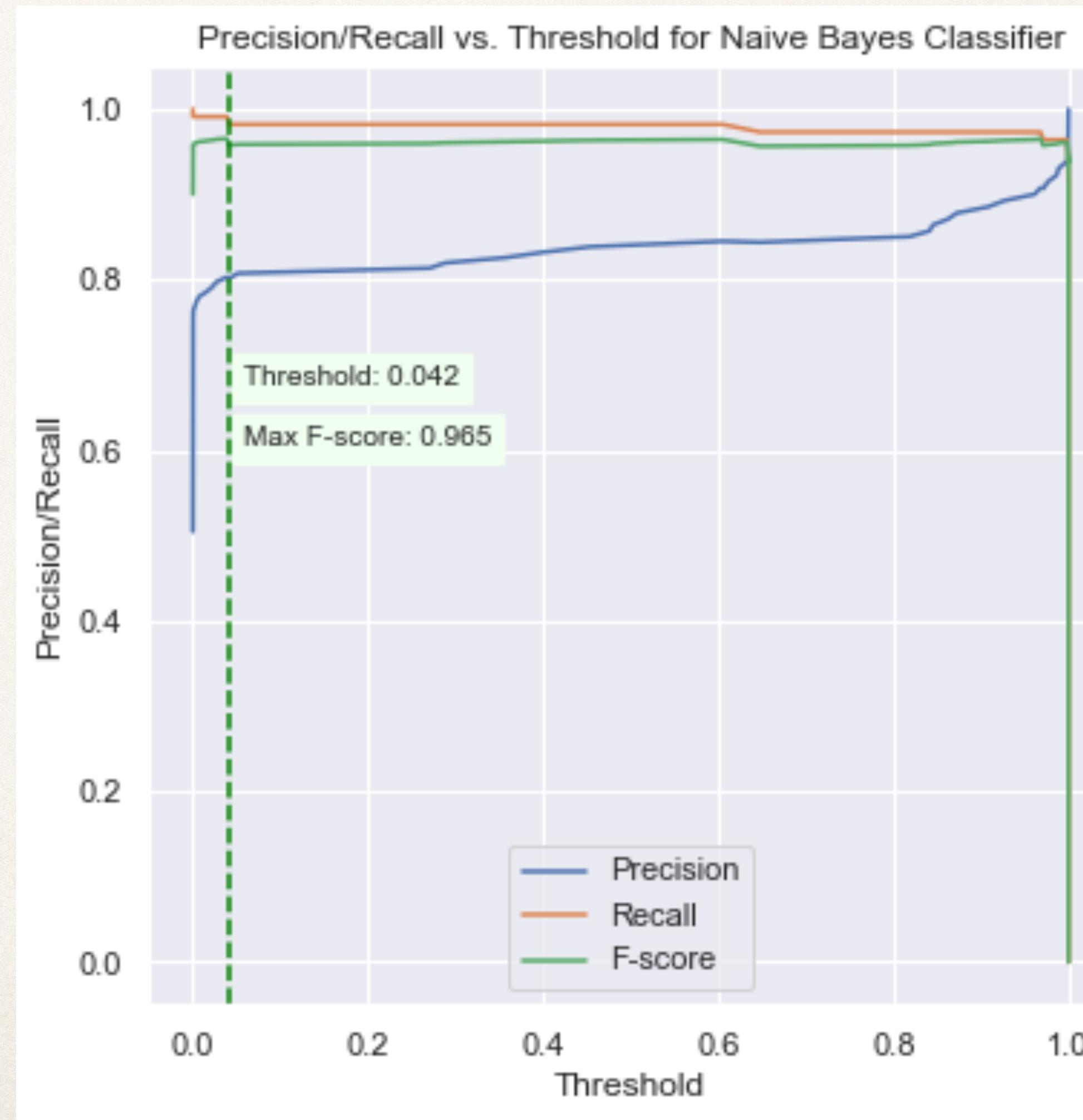
- ✿ “Depending on where you live, there may be a law in place that releases you of a lease if you’ve experienced DV, reporting or not” (Comment)
- ✿ “When you focus on your healing you will begin to feel happy and trustful of the world once again” (Comment)

Reddit Data: Feature Importances for Non-Abuse Posts and Comments



- ✿ “For a real song that always gives me feels I choose [this](<https://www.youtube.com/watch?v=9pkLDEEs20U>)” (Comment)
- ✿ “Zoloft was horrible for me for the first two weeks. I was extremely anxious and unable to sleep. It was just awful. In my case, it was because I had to increase the dosage quickly, all the way up to 150 mg in just over a week” (Comment)

Reddit Data: Thresholding



| | Precision | Recall |
|---|-----------|--------|
| 0 | 0.99 | 0.75 |
| 1 | 0.80 | 0.99 |

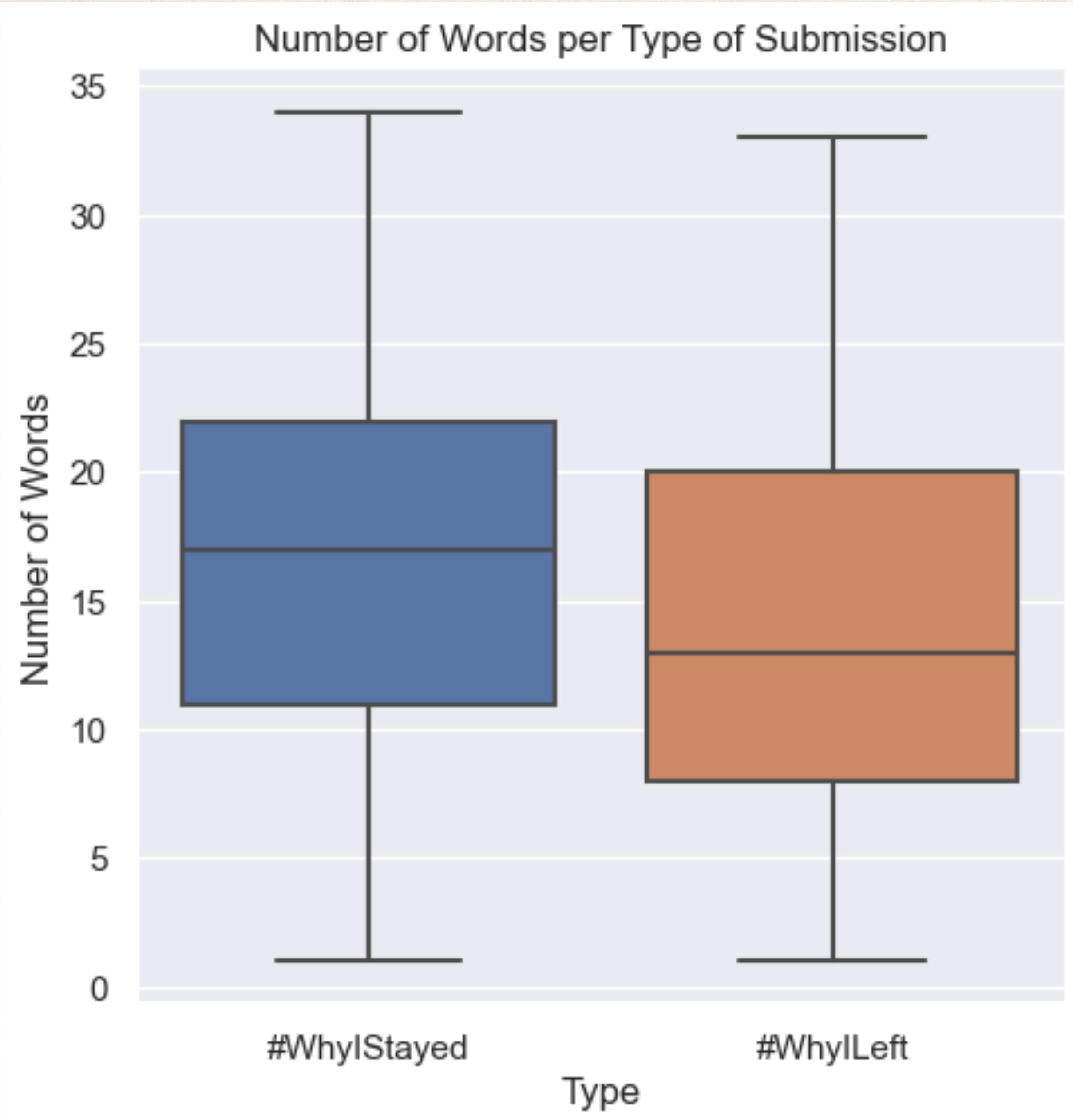
Twitter Data



- ✿ Twitter: #WhyIStayed vs. #WhyILeft - Ray Rice and Janay Palmer
- ✿ Final Twitter dataset: 13,731 #WhyIStayed & 5,821 #WhyILeft tweets totalling 19,552 tweets
- ✿ Tweets containing both hashtags were split, and emojis replaced, and mentions, hashtags, and redundant retweets deleted

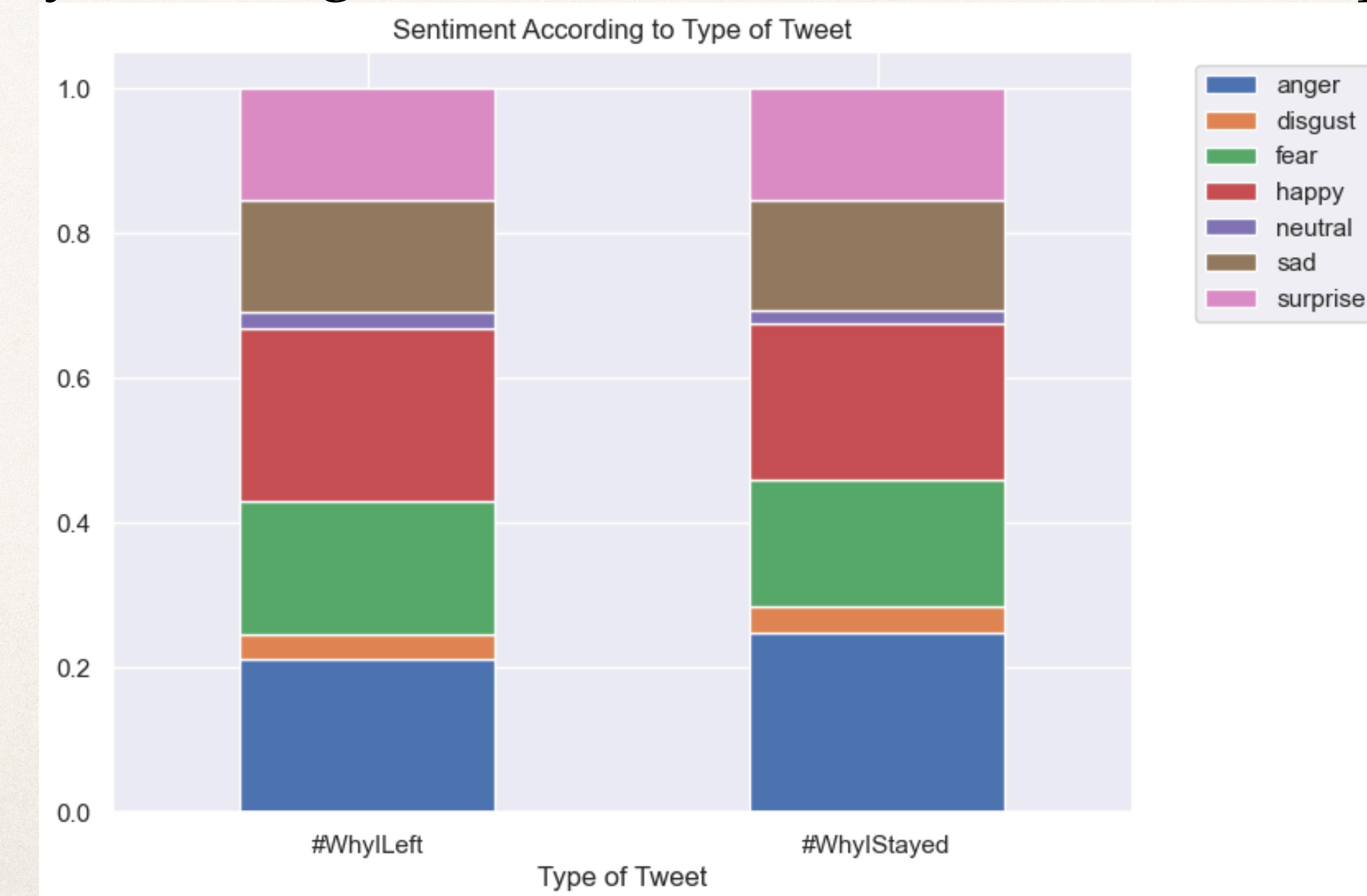
Twitter Data: Post Length

- #WhyIStayed tweets contained more words and tokens than #WhyILeft tweets
- #WhyILeft tweets were easier to read
- #WhyIStayed tweets contained more urls, mentions, and capital letters
- No significant difference between the counts of emojis



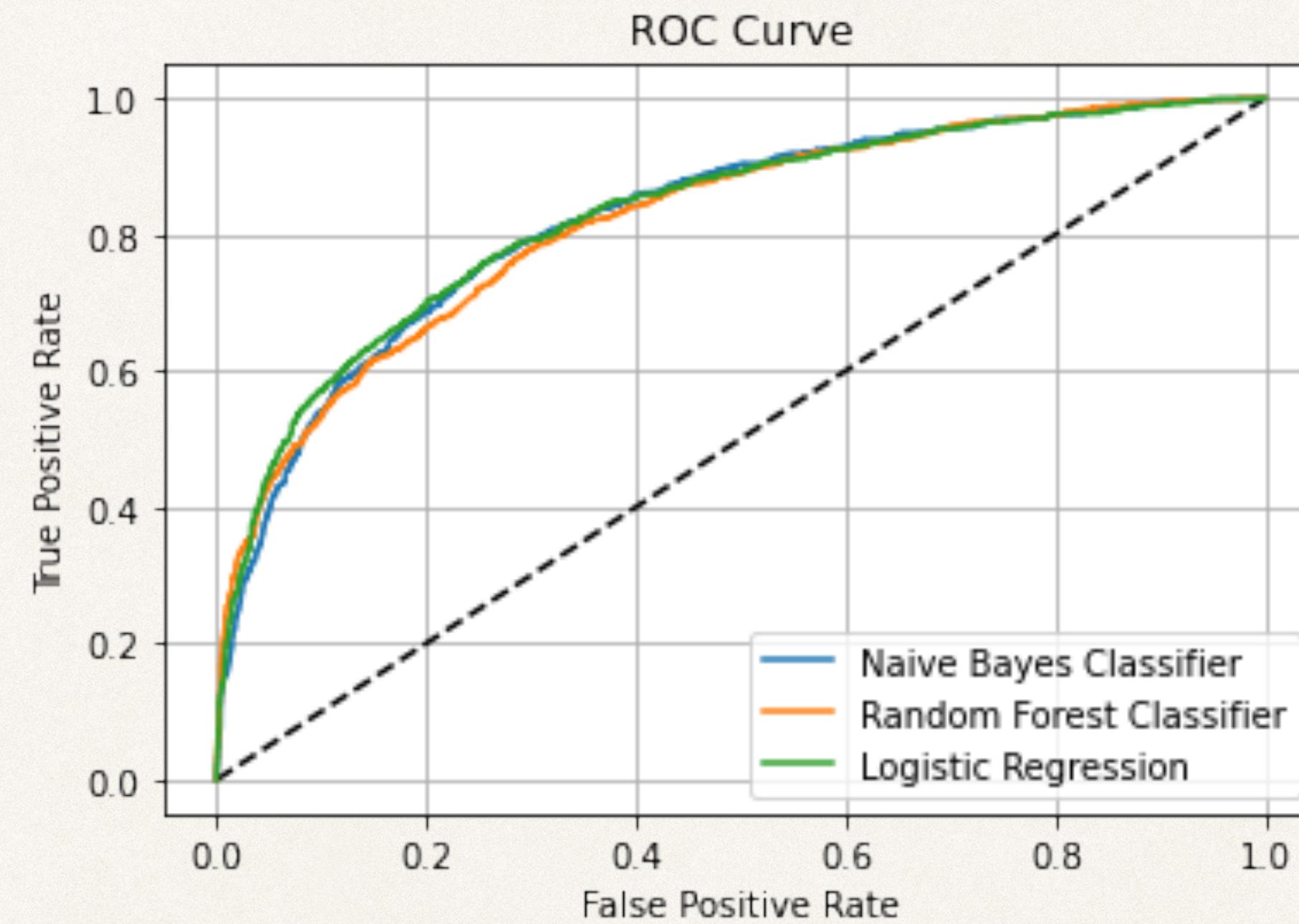
Twitter Data: Relative Frequencies of Top Sentiments

- ❖ Sentiment analysis using mean scores of emotions corresponding to tokens



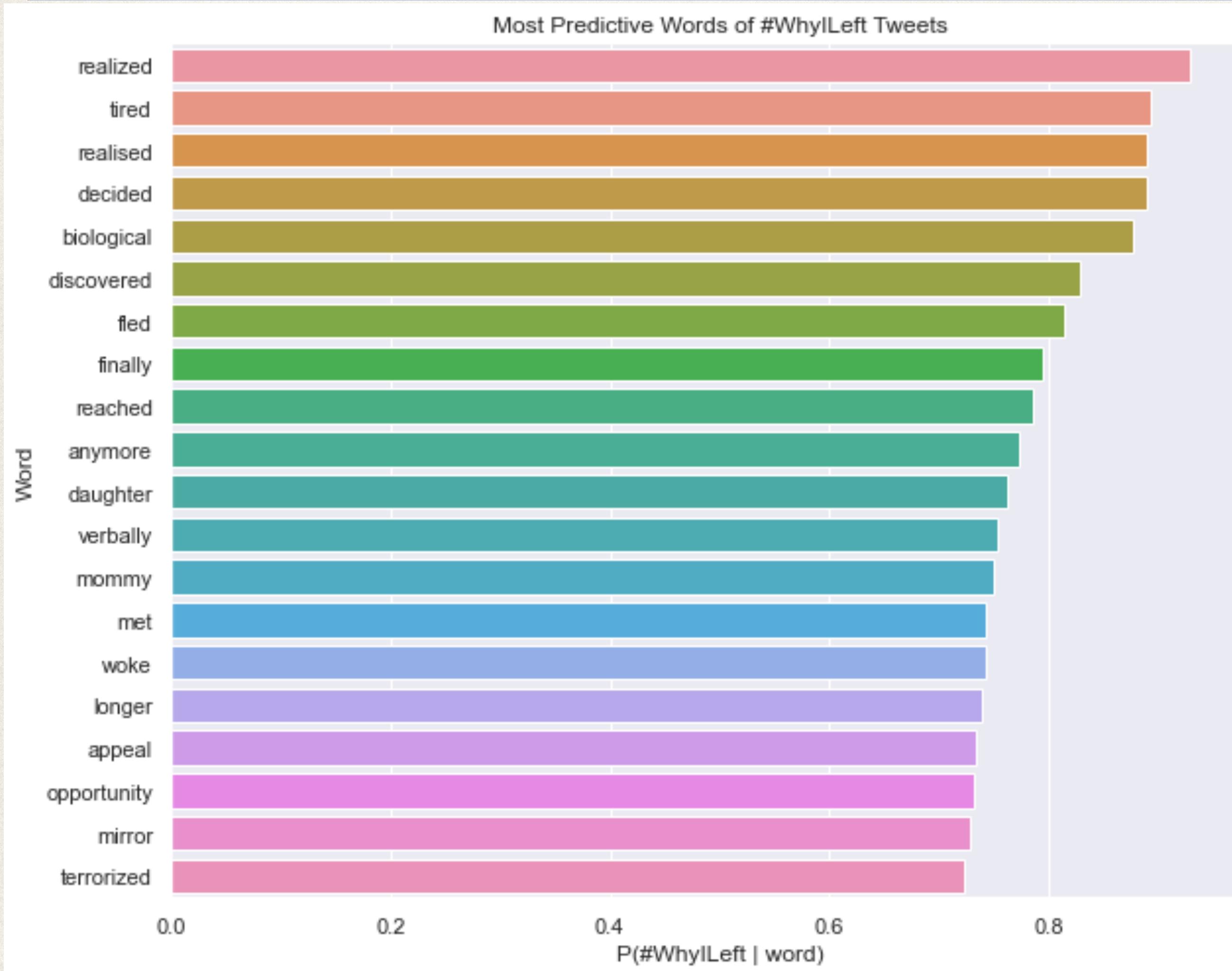
Twitter Data: Modeling

- ❖ 8,550 tokens



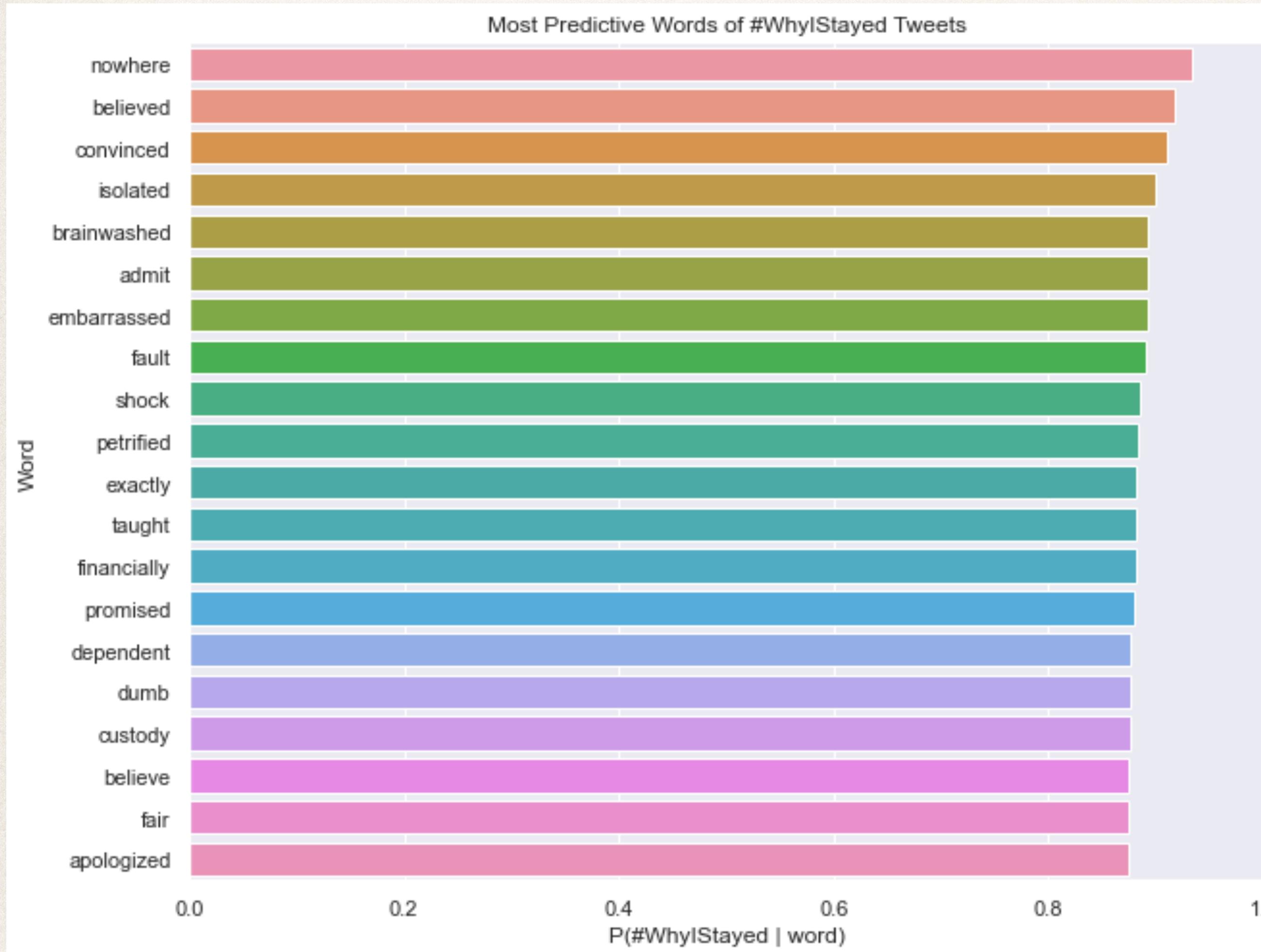
| Model | Best min_df for CountVectorizer | Best Parameters | ROC-AUC Score |
|----------------------------|---------------------------------|-----------------------------------------------------------------------------------------|---------------|
| Naive Bayes Classifier | 1 | {'alpha': 1} | 0.824 |
| Random Forest Classifier | 9 | {'criterion': 'entropy', 'max_depth': 30, 'max_features': 'log2', 'n_estimators': 1000} | 0.821 |
| Logistic Regression | 1 | {'C': 1} | 0.829 |

Twitter Data: Feature Importances for #WhyILeft



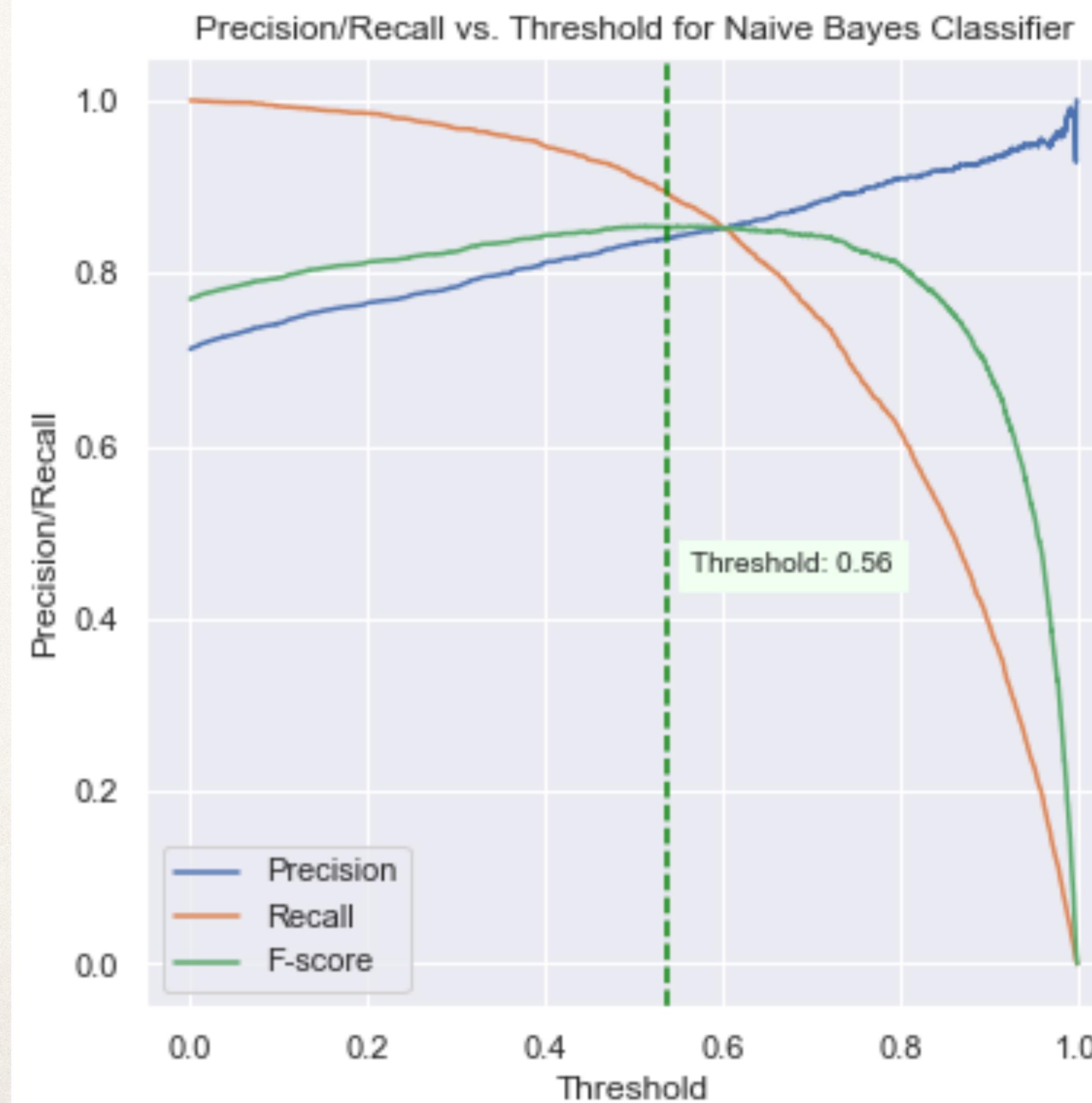
- ✿ “Because I **realized** someone had to actually love me, even if it was only myself”
- ✿ “An old friend **reached** out, showed me that I was not alone, and helped me find a way out”
- ✿ “My son refusing to go to school cos “daddy hits **mommy**” was the last straw. I finally saw it through his eyes. Near death. Twice”

Twitter Data: Feature Importances for #WhyStayed



- ✿ “I felt **dependent** on him. After you are told and treated like you are worthless repetitively enough you start to believe it”
- ✿ “Because I was a strong intelligent woman and I didn’t want people to know I was **dumb** enough to let it happen to me”
- ✿ “He also made me quit working, made me completely **dependent** on him. I realized all too late he had cleaned out on me **financially**”

Twitter Data: Thresholding



| | Precision | Recall |
|---|-----------|--------|
| 0 | 0.67 | 0.6 |
| 1 | 0.84 | 0.88 |

Conclusions

- ✿ Twitter analysis congruent with battered person syndrome
- ✿ Reddit Naive Bayes + Twitter Logistic Regression
- ✿ Suggestions for future analysis: filtering Reddit comments, larger emotions dataset, balanced Twitter dataset