



**Can We Identify subreddits  
based on 1-3 word  
combinations?**

# The Subreddits



**JUST IN NO MIL**

2

3

4

# The Data

5

6

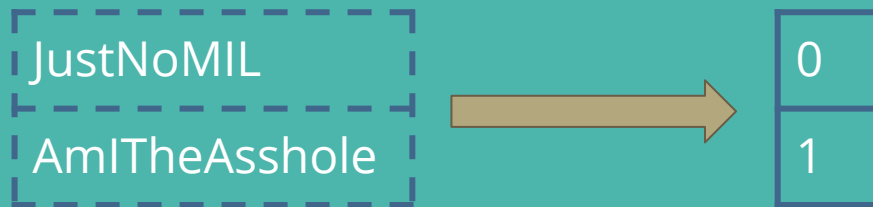
- Scrapped Reddit using PushShift API
- 9000 posts from r/AmITheAsshole
  - 7000 posts from r/JustNoMIL

7

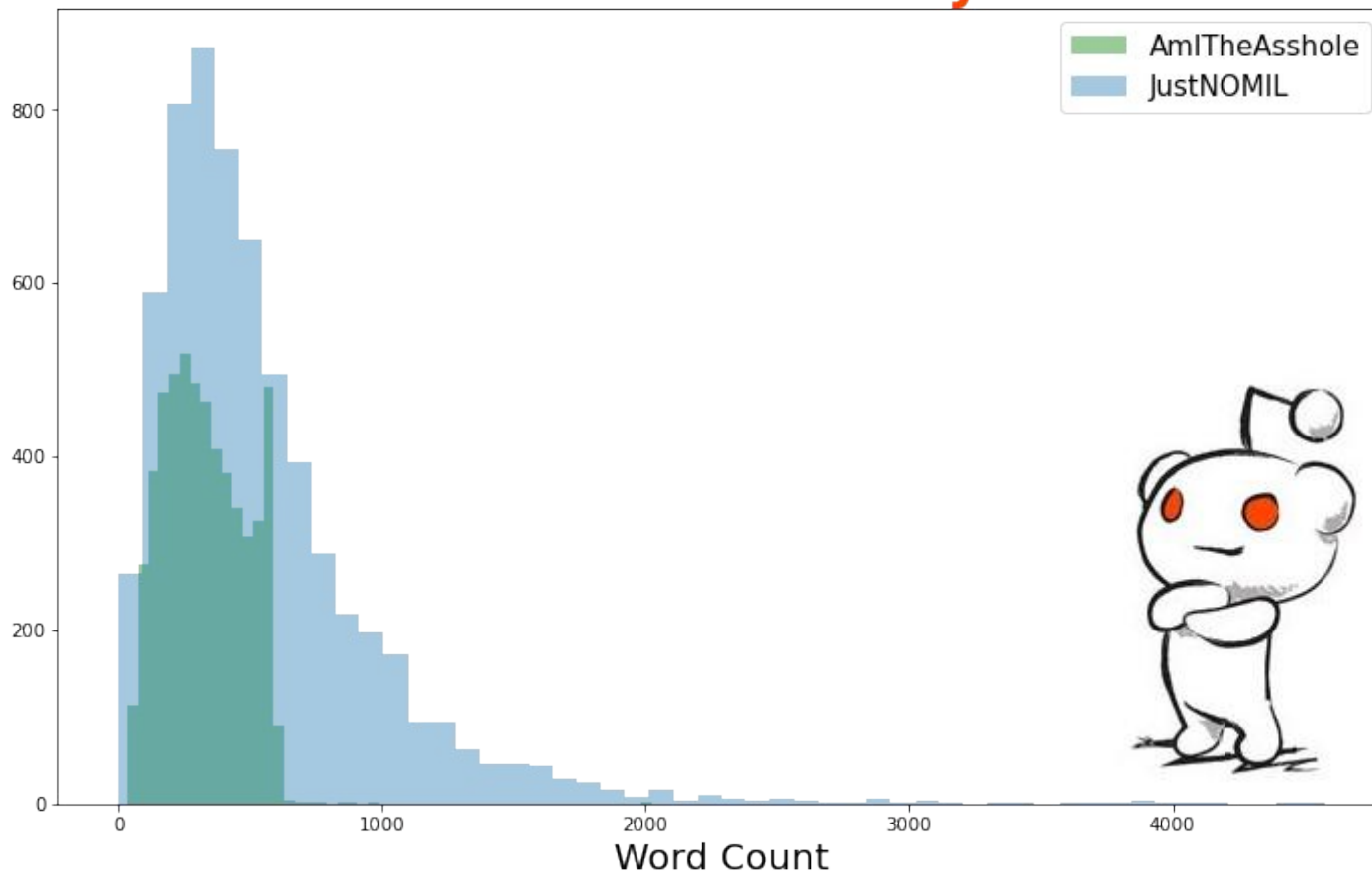
8

## Data Cleaning

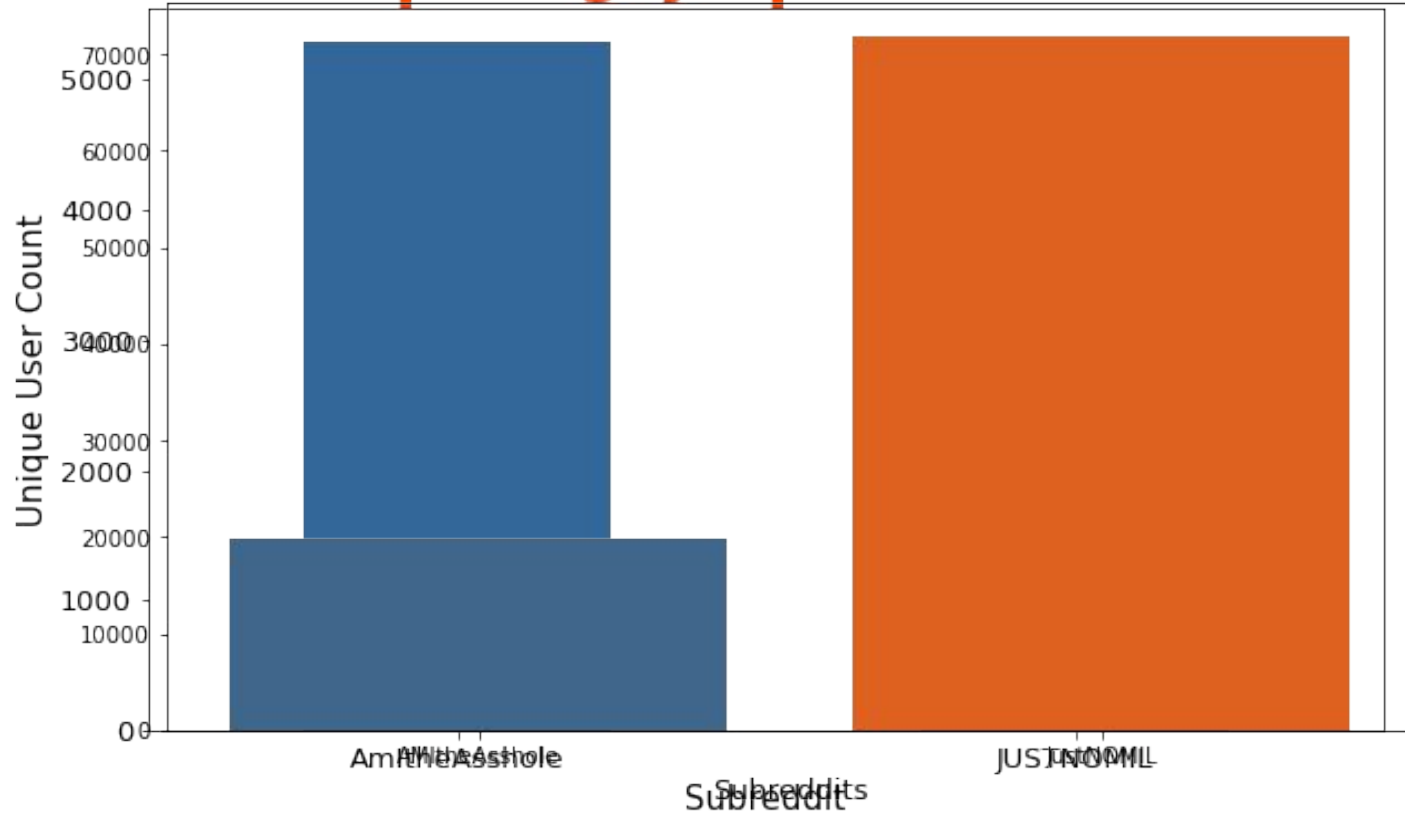
- [deleted] & [removed]
  - Duplicates
- AutoModerator posts
  - \n
  - &#x200B;



# Word Count Distribution by Subreddit



# How Angry Are The Users?

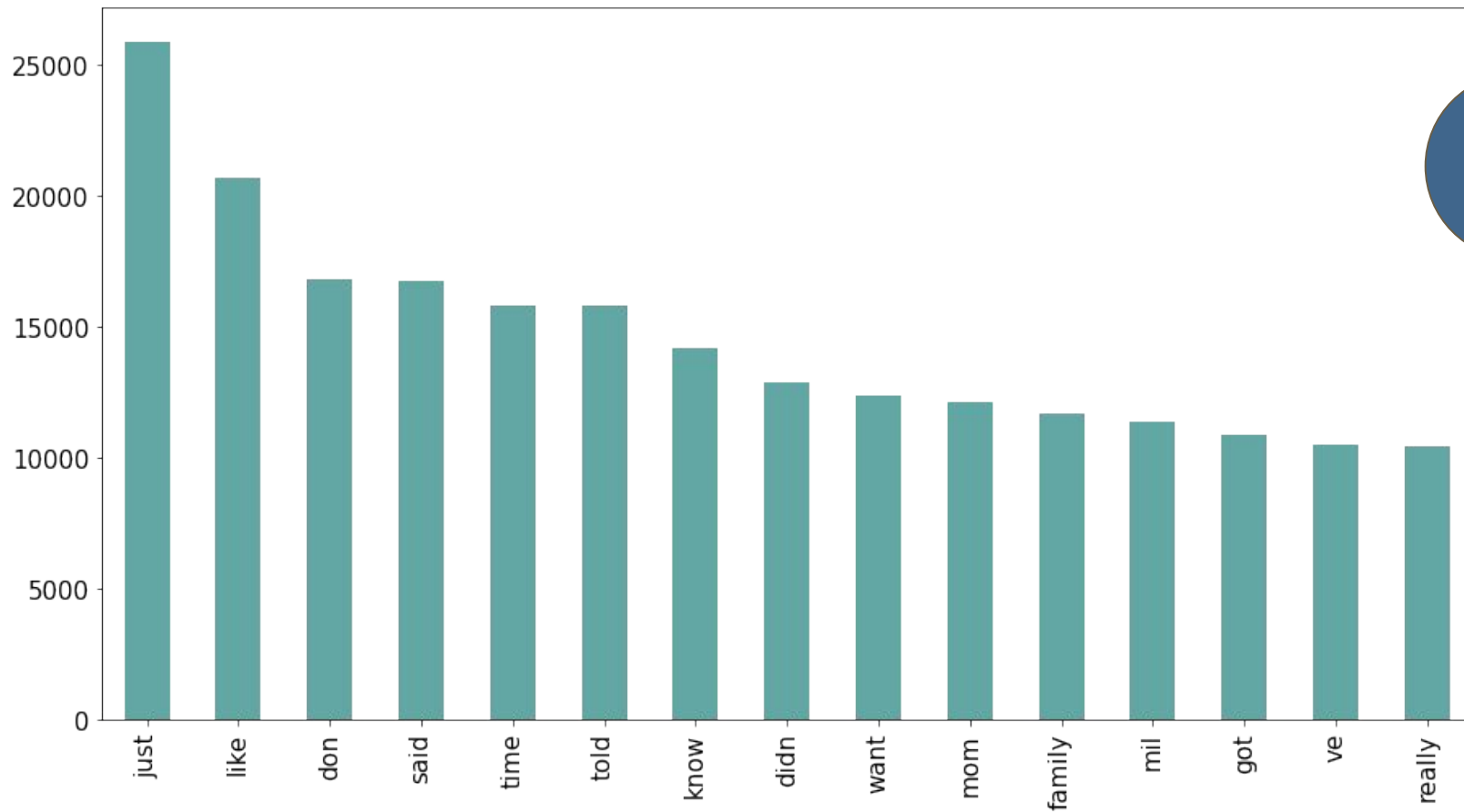


11

10



## Most Common Words in Titles & Posts



12



# Model Prep

- StopWords = 'english' + 'top 15'
- 'ref' for identifiers
- Combined title and selftext columns
- Target = 'sub-reddit'
- Train / Test split

just	justnomil	jnmil	mil	asshole
amita	amitah	ref wibta	wibtah	aita

# Models w. GridSearch Pipeline

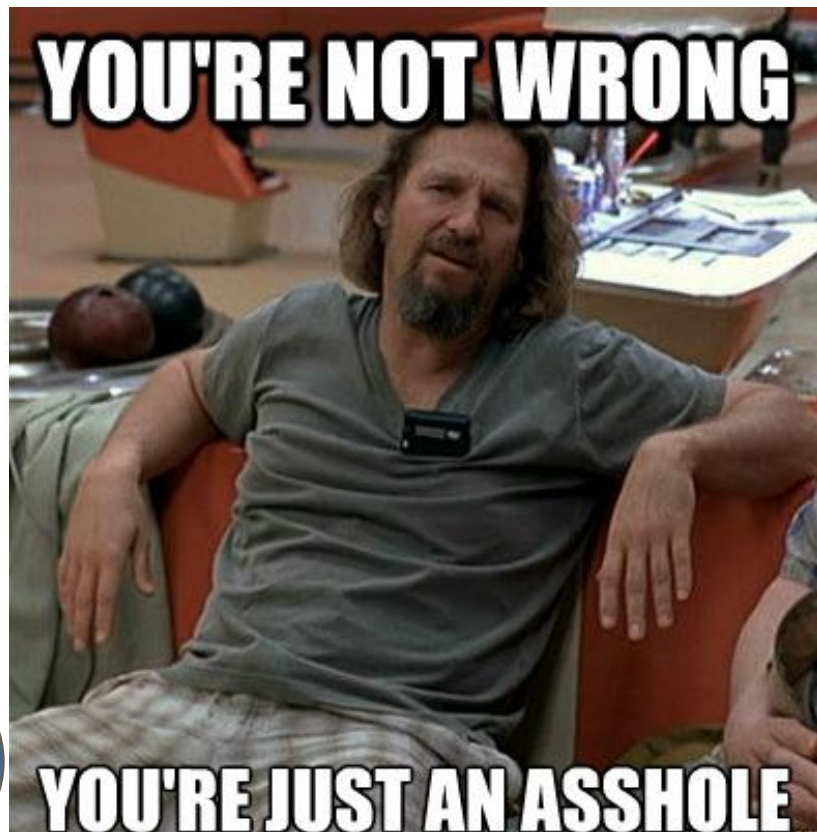
- Transformer
  - CountVectorizer
  - TfidfVectorizer
- Estimator
  - MultinomialNB
  - BernoulliNB
  - LogisticRegression
  - AdaBoost Classifier

16

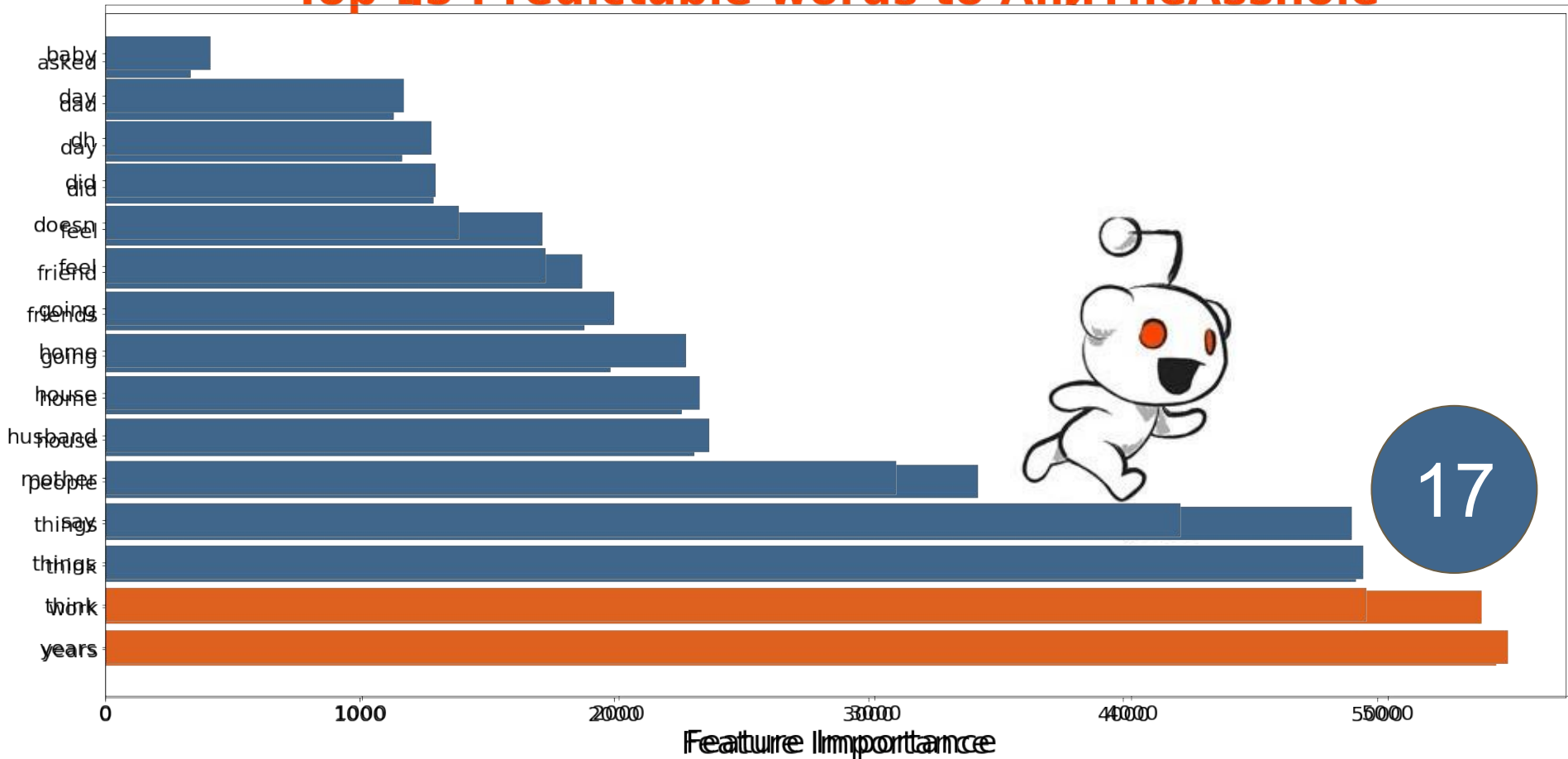
15

14

13



# Top 15 Predictable words to Just NOM!

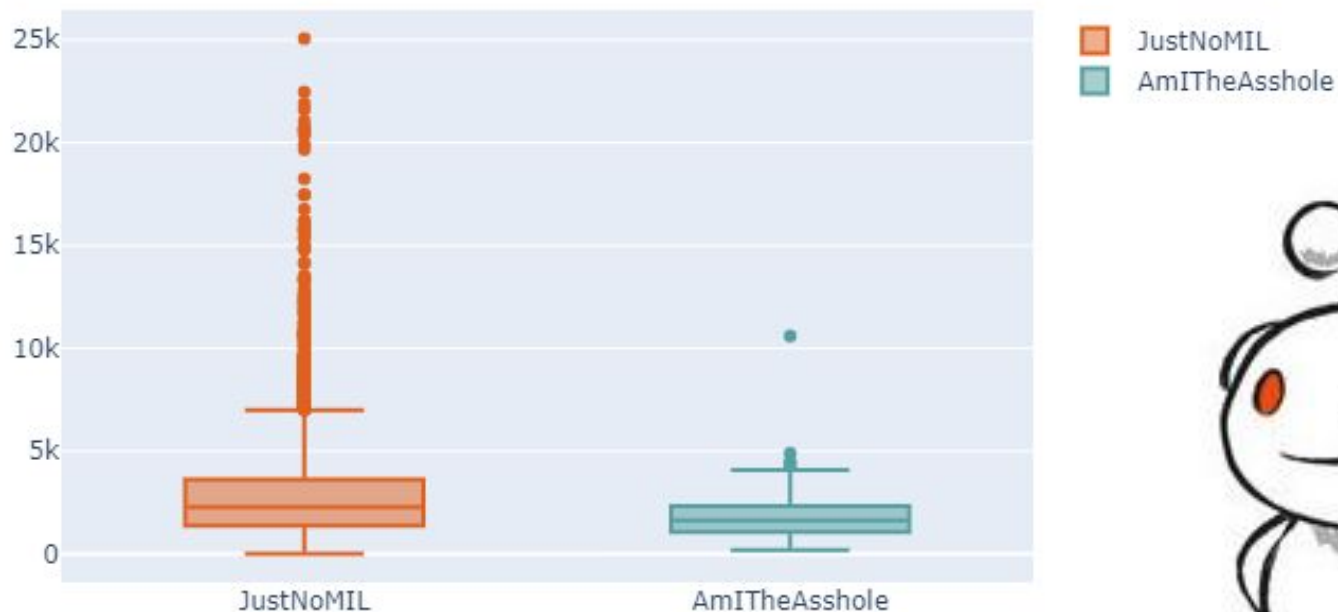


Train
Test
Specific



Class.
Vect
9.51%
7.59%
6.49%

## Post Length per Subreddit



THANK YOU!

