# So, Reddit, AITA?
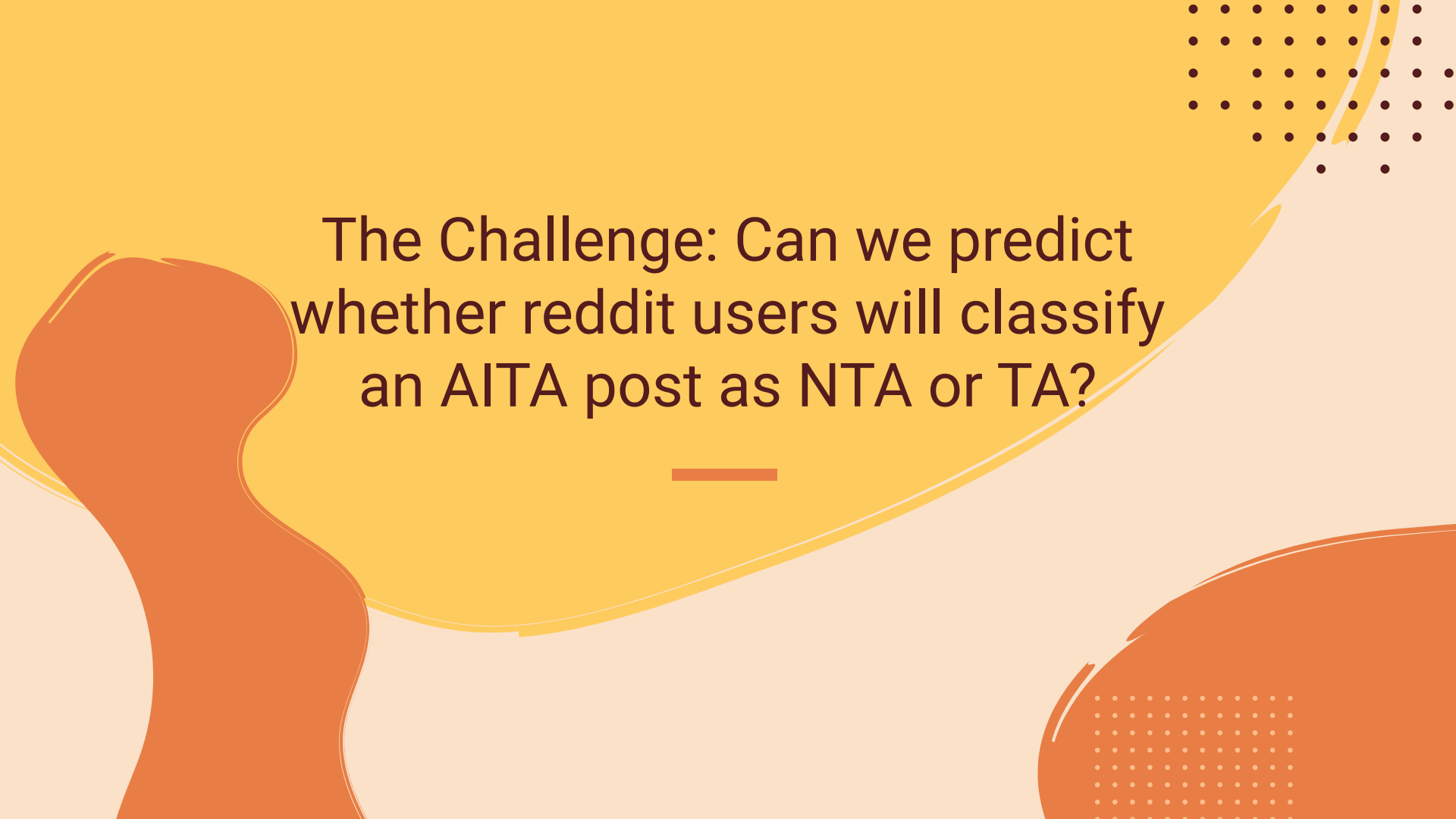
Please excuse the profanity in this presentation

# Background

- Social media
- Post, comment, upvote
- Sub-communities on reddit = Subreddits
  - AITA

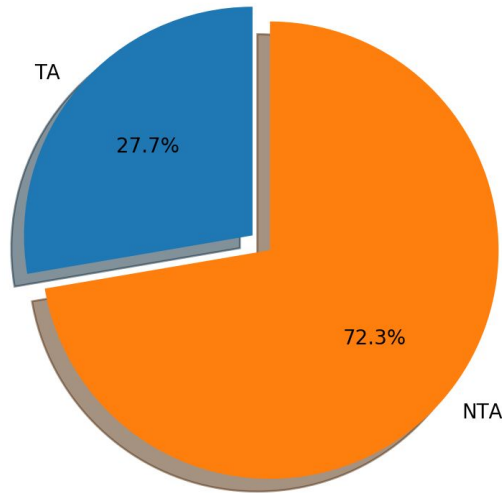The Challenge: Can we predict whether reddit users will classify an AITA post as NTA or TA?

# The Data

- Pushshift.io API
  - Title
  - Description
  - Flair
    - moderator-tagged short description, unique to the subreddit
    - AITA: NTA, TA, ESH, NAH

# Class Distribution



TA 27.7%

NTA 72.3%

# Data Engineering

## By Hand
- Description length
- Number of '!', '?', ", '...'
- Number of words in all caps
- Number of curse words
- Number of first versus third person pronouns

## Using Libraries
- Sentiment analysis of titles & descriptions
- TF-IDF → LDA
  - Latent Dirichlet Allocation AKA unsupervised categorization

NLTK

TextBlob

scikit learn

GENSIM
topic modelling for humans

# Modeling

First, we SMOTE-d to correct for class imbalance! Then, we tried vanilla versions of…

**Logistic Regression**

TP: 60%
TN:  38%
FP: 40%

**KNN**

TP: 60%
TN: 49%
FP: 40%

**Random Forest**
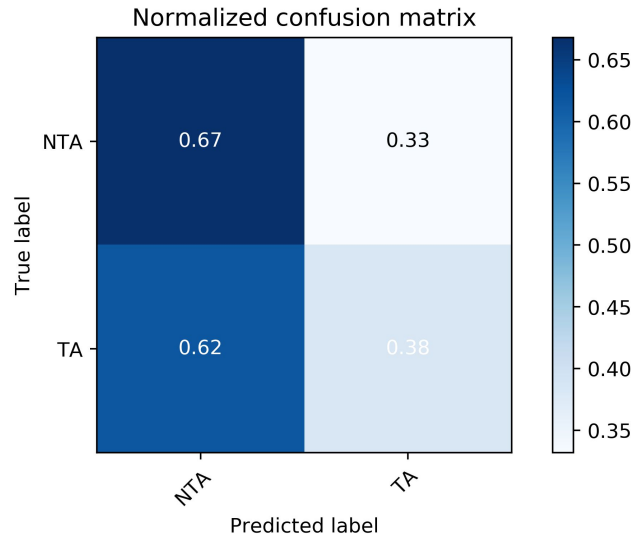
TP: 73%
TN:  26%
FP: 27%

**SVM**

TP: 54%
TN: 46%
FP: 46%

**XGBoost**

TP: 36%
TN:  66%
FP: 64%

**Naive Bayes**

TP: 69%
TN: 31%
FP: 31%

Normalized confusion matrix

# Final Model: Random Forest

- Highest accuracy ()
- One of the lowest false positives
- The lowest false negative rates

Overall, all models had a 50% or higher false negative rate :(

# Feature Importance

Top features:
- Curse count
- Quote count
- Afinn title sentiment
- Question mark count
- VADER description sentiment
- Length
- LDA category



Feature importances