# Reddit Post Classification

## Using NLP to Classify Text

By Eric Bayless
January 22, 2020

# The Problem

- Predict whether the given text of a submission came from the r/Fitness or r/Bodyweightfitness subreddit

# Data Acquisition

# Method

- Pushshift API

- Function

  - 100 posts per request

  - 10 seconds between

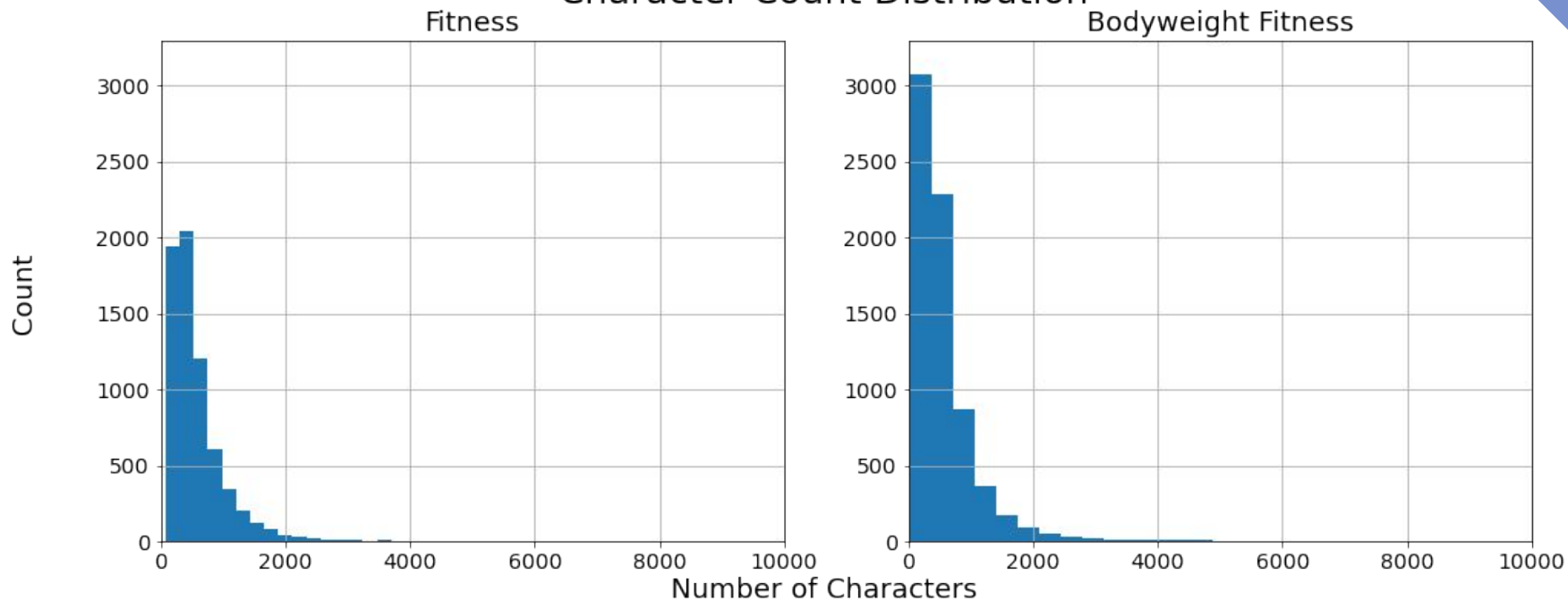  - Starting with most recent

# Data Cleaning

# Method

1. Drop null posts

2. Drop 'deleted' or 'removed'

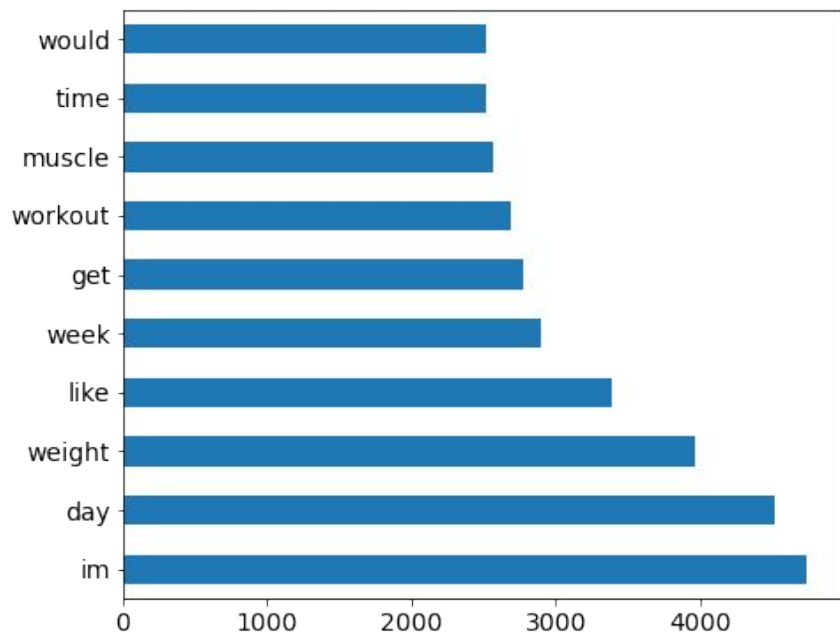3. Drop posts containing subreddit name

4. Drop posts shorter than title

Exploratory Analysis
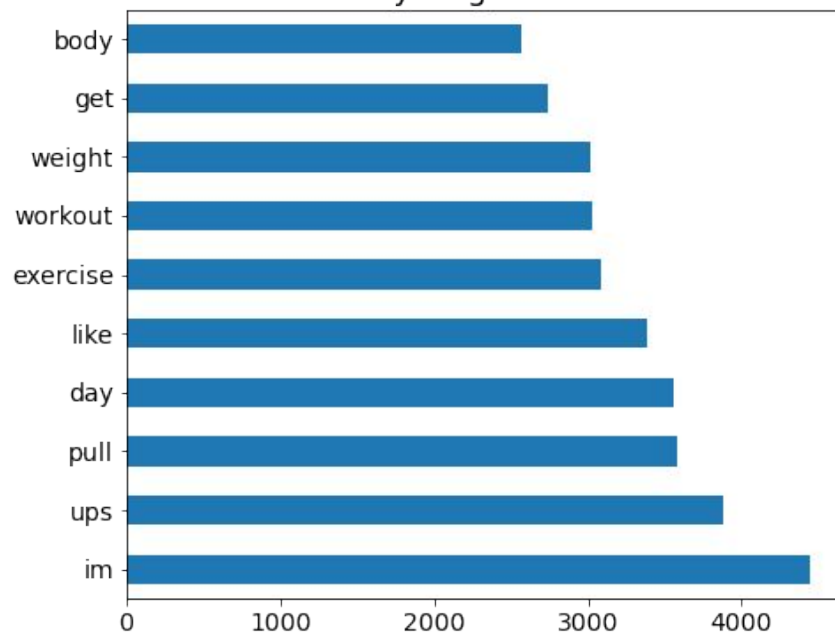
Character Count Distribution
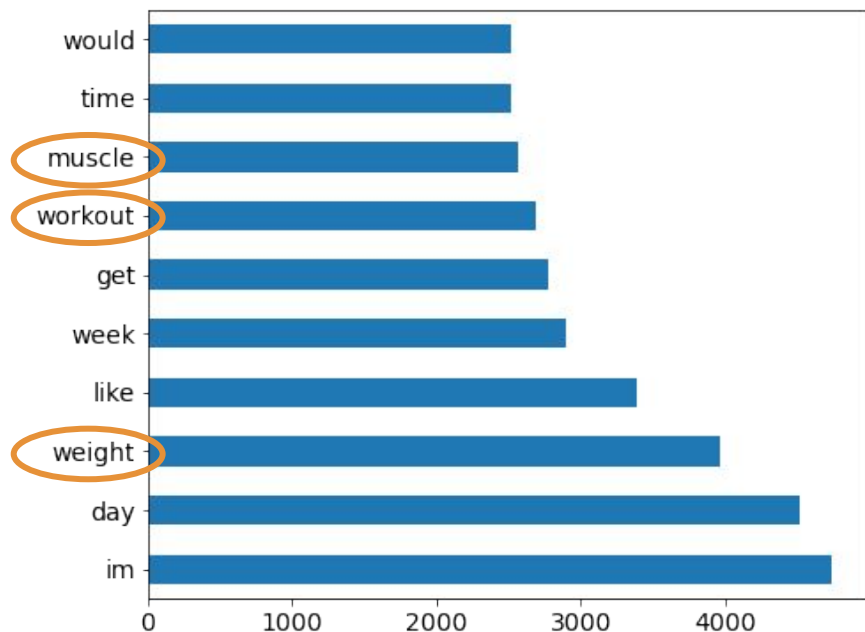
Top 10 Words

Fitness

| | Number of Occurances |
|---|---|
| would | 2500 |
| time | 2500 |
| muscle | 2600 |
| workout | 2700 |
| get | 2800 |
| week | 2900 |
| like | 3400 |
| weight | 3950 |
| day | 4500 |
| im | 4650 |

Bodyweight Fitness

| | Number of Occurances |
|---|---|
| body | 2600 |
| get | 2750 |
| weight | 3050 |
| workout | 3050 |
| exercise | 3100 |
| like | 3400 |
| day | 3550 |
| pull | 3600 |
| ups | 3850 |
| im | 4450 |

## Top 10 Words

### Fitness

| Word | Number of Occurances |
|------|---------------------|

would, time, muscle, workout, get, week, like, weight, day, im

### Bodyweight Fitness

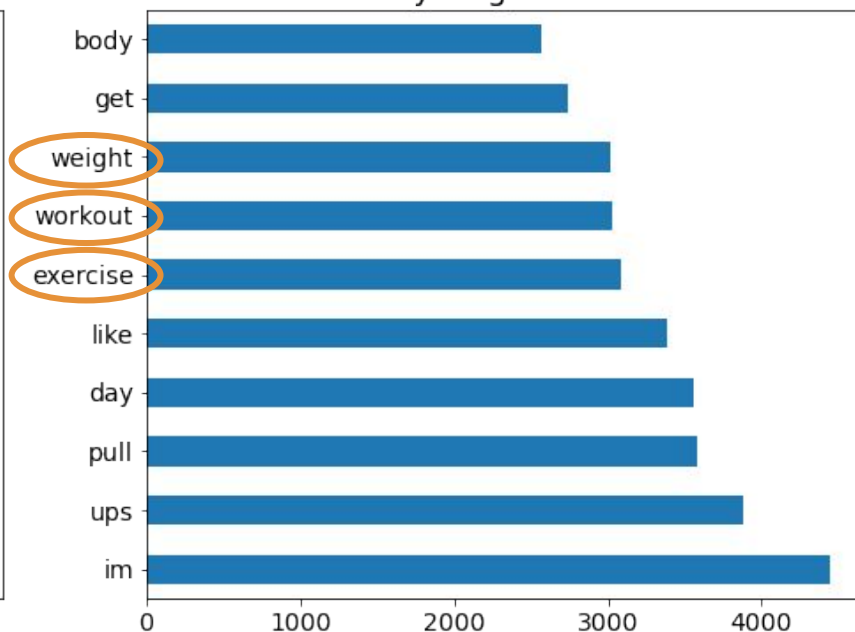body, get, weight, workout, exercise, like, day, pull, ups, im

Number of Occurances

Top 10 Words

Fitness

Bodyweight Fitness

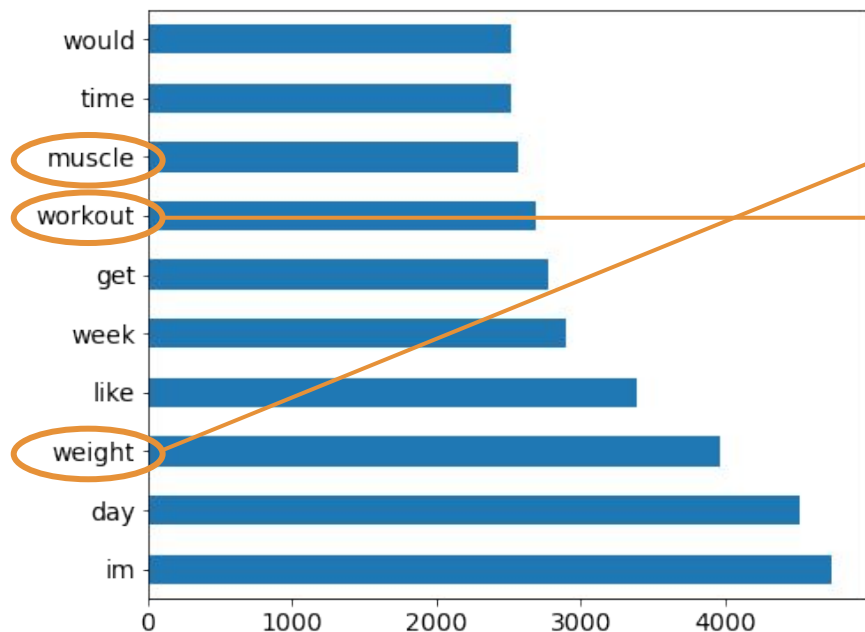Number of Occurances

# Top 10 Words

## Fitness

| | |
|---|---|
| would | |
| time | |
| muscle | |
| workout | |
| get | |
| week | |
| like | |
| weight | |
| day | |
| im | |

## Bodyweight Fitness

| | |
|---|---|
| body | |
| get | |
| weight | |
| workout | |
| exercise | |
| like | |
| day | |
| pull | |
| ups | |
| im | |

Number of Occurances

# Modeling

# Baseline

- Most Frequent

- r/Bodyweightfitness

- Accuracy:  51%

# Model Selection

- Multiple model types

- Accuracy 76-78%

- Custom preprocessing

# Best Model

- Logistic Regression

- Accuracy: 79%

- Interpretable coefficients
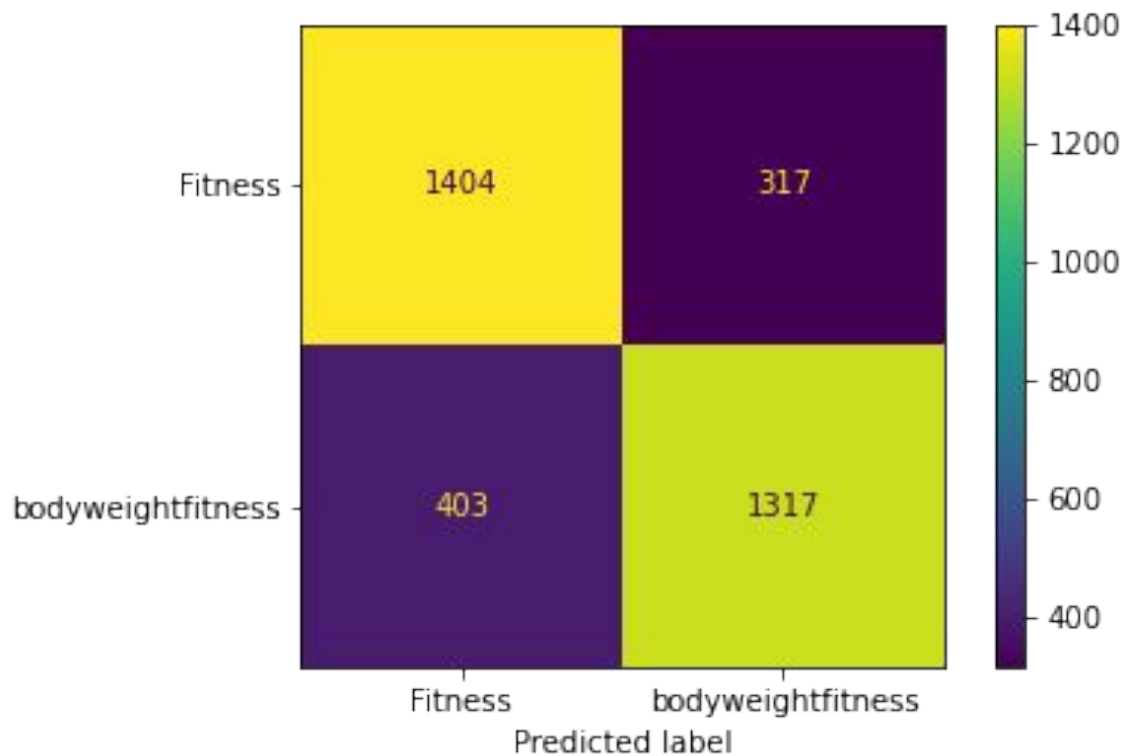
# Most Important Terms

## Bodyweight Fit.

1. rr
2. bodyweight
3. calisthenics
4. ring
5. covid
6. pull
7. dip
8. pushup
9. ups
10. pullups

## Fitness

1. Bench
2. Gym
3. Lifting
4. Lift
5. Deadlift
6. Machine
7. Fitness
8. Bike
9. Deadlifts
10. dumbbell

# Confusion Matrix

# Conclusions

# Results

- 79% Accuracy

- Similarity of subreddits may prevent further improvement

# Production Recommendation

- Logistic Regression
  - Highest accuracy
  - Interpretable

# Future Steps

- Custom Vocabulary

- Advanced Neural Network

- More Preprocessing

- More Data