



# R E D D I T N L P

PRESENTATION BY ROWAN SCHAEFER

# TL;DR

The goal of this analysis is to classify text posts as either belonging to one subreddit or another, using natural language processing.

Success will be evaluated using the model's accuracy score in correctly classifying the posts, with a second goal of trying many models to obtain the highest accuracy score possible.



# RESEARCH QUESTIONS



```
.stdfs=[]  
for key, val in range_gen:  
    i = query('learnsql',  
    print('pulled 100 rows')  
    listdfs.append(i)  
    time.sleep(10)
```



```
class 'pandas.core.frame.DataFrame'>  
Int64Index: 3396 entries, 0 to 98  
Data columns (total 12 columns):  
   Column      Non-Null Count  Dtype    
-----    
author      3396 non-null    object   
created_utc  3396 non-null    object   
id           3396 non-null    object   
is_created_from_ads_ui  1098 non-null    object   
num_comments 3396 non-null    int64    
score       3396 non-null    float64  
selftext    3388 non-null    object   
subreddit   3396 non-null    object   
title       3396 non-null    object   
total_awards_received  3396 non-null    int64    
upvote_ratio 2398 non-null    float64  
url         3396 non-null    object   
dtypes: float64(1), int64(4), object(7)  
memory usage: 344.9+ KB
```

- What preprocessing steps lead to highest accuracy in classifying subreddits?
- Which classification models and hyperparameters lead to highest accuracy and precision in classifying subreddits?

# SUBREDDITS

r/learnpython and r/learnsql are mostly filled with questions and resources for learners.

r/ learnpython

- 613,000 subscribers
- 3,396 rows of data

r/ learnsql

- 17,700 subscribers
- 3,407 rows of data
- Had to pull from API differently to avoid class imbalance.

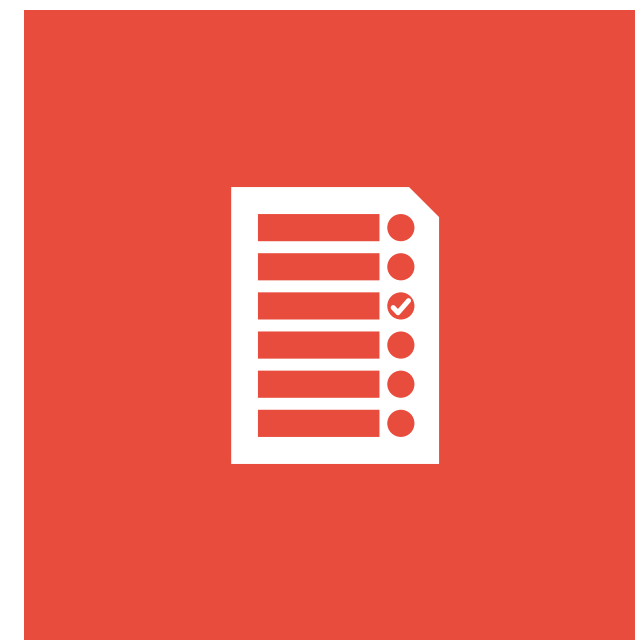
# THE PROCESS



find data



cleaning/EDA



preprocessing



modeling

# GETTING DATA

To get data, I used the `requests` library with reddit's API, **PushshiftAPI**. In order to automate pulls, I wrote a function that looped through "before" and "after" criteria for the date of the reddit posts, and set a sleep timer in between so I didn't overwhelm the API.

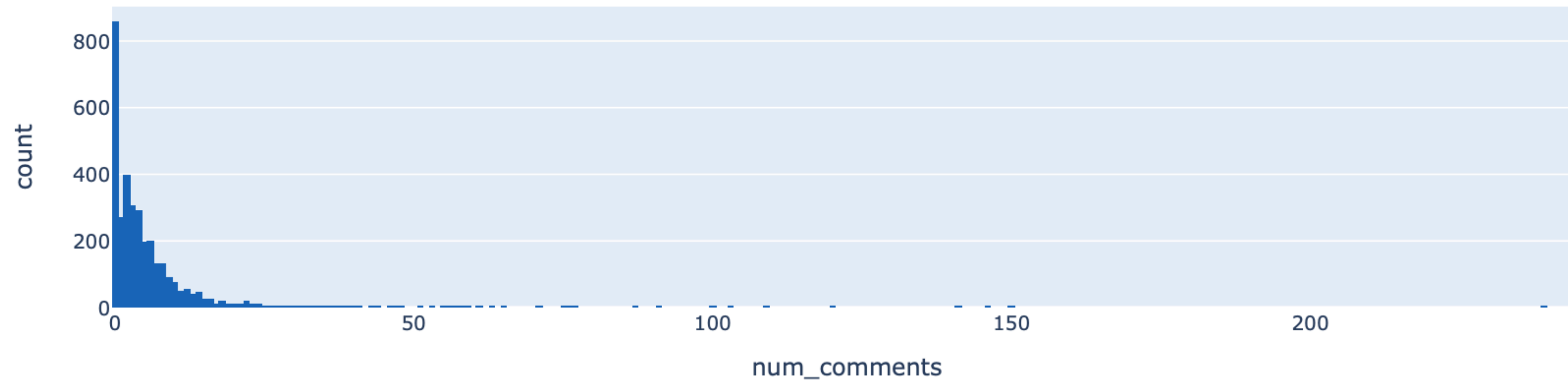
60 REQUESTS  
PER MINUTE

100 ROWS PER  
REQUEST

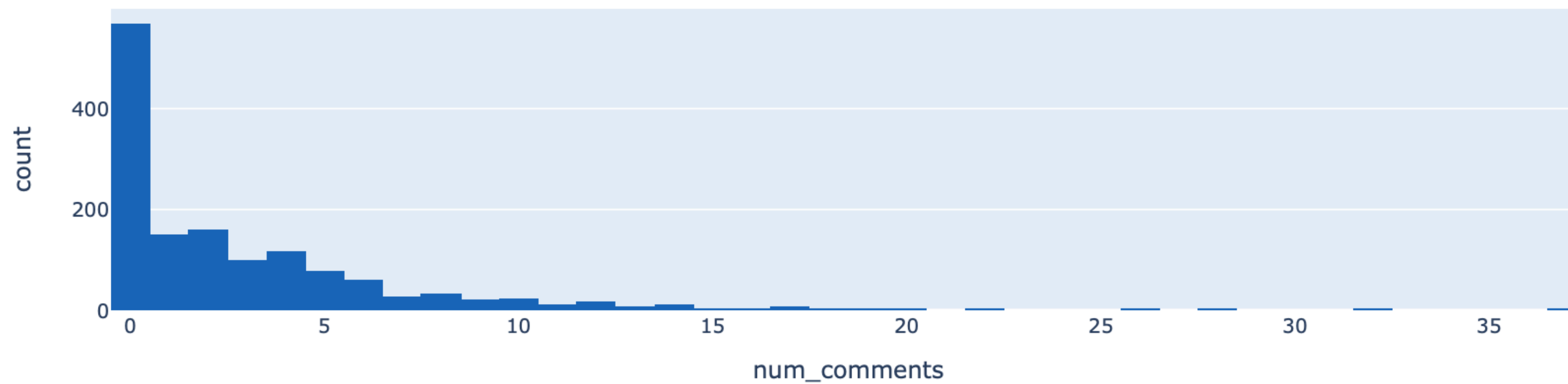
TOTAL ROWS:  
~7000

# EXPLORATORY ANALYSIS

Count of Comments per Post, r/learnpython

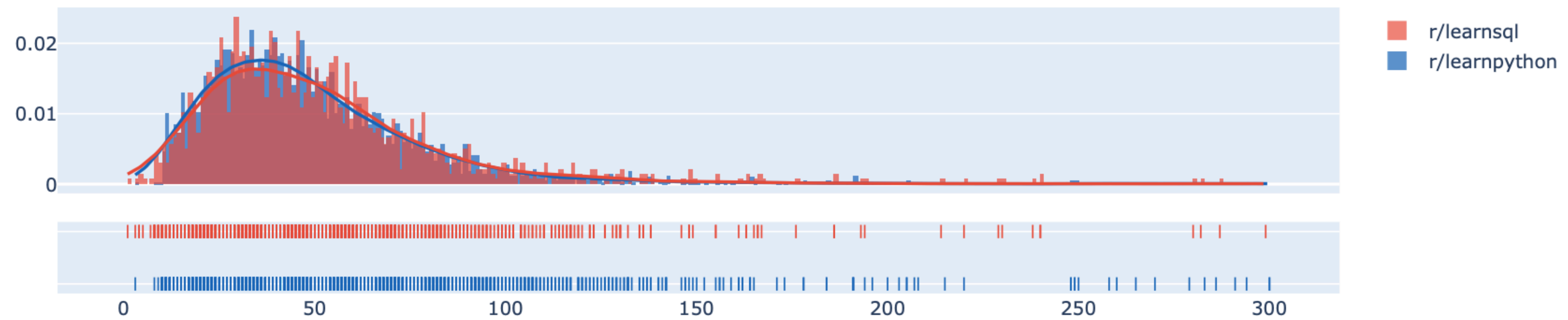


Count of Comments per Post, r/learnsql

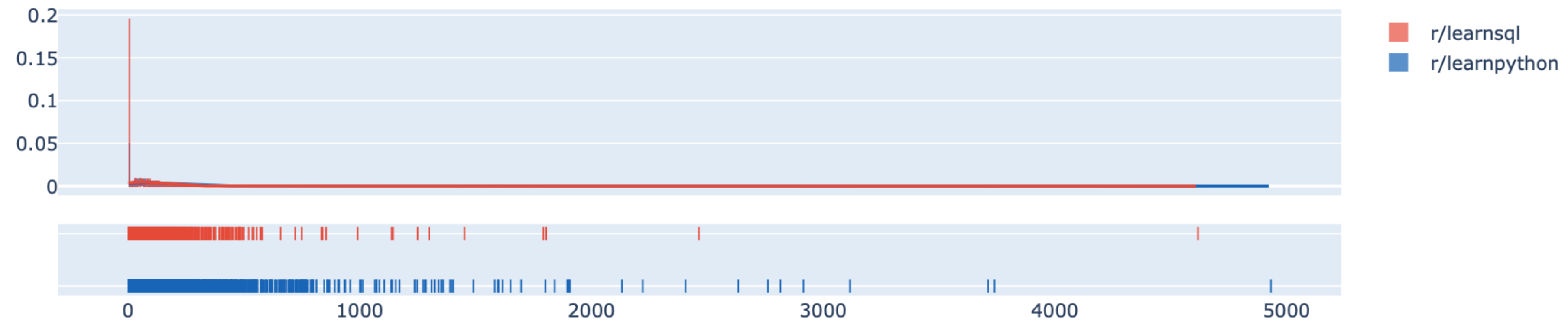


# EXPLORATORY ANALYSIS

Distributions of Title Length



Distributions of Word Count (selftext)





# PREPROCESSING

## AND SOME DATA CLEANING

- Removed 170 posts from user Andre380
- Removed ['removed'] selftext
- Removed numbers
- Removed giveaway words: SQL, postgre, python, snowflake, etc
- Models performed best with no bigrams

FREE  
Daily SQL  
Challenge  
#22

FREE  
Daily SQL  
Challenge  
#20

FREE  
Daily SQL  
Challenge  
#24

FREE  
Daily SQL  
Challenge  
#25

FREE  
Daily SQL  
Challenge  
#21

SQL Prep  
- Daily  
Challenge  
#95

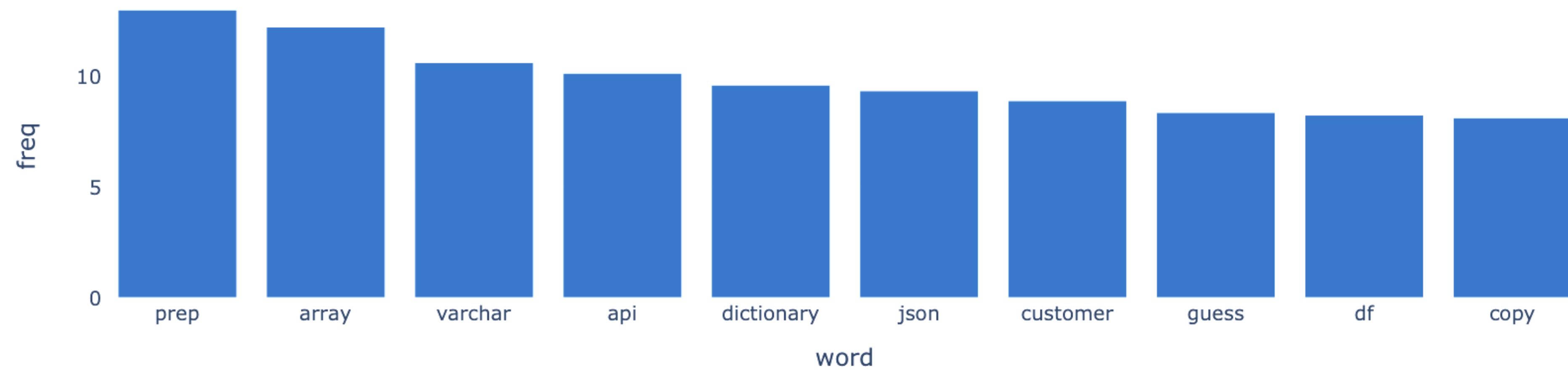
FREE  
Daily SQL  
Challenge  
#21

SQL Prep  
- Daily  
Challenge  
#101

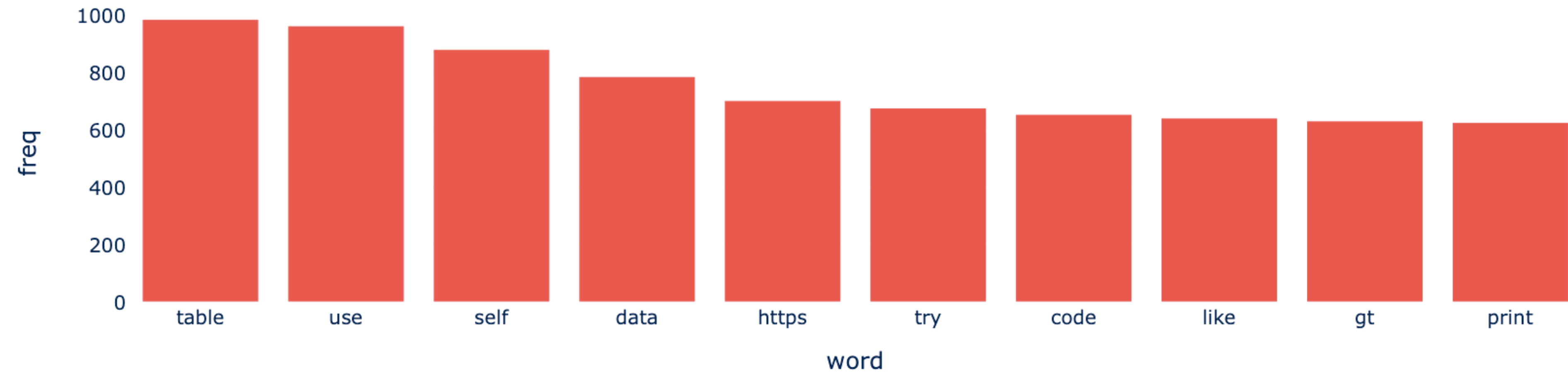
FREE  
Daily SQL  
Challenge  
#19

FREE  
Daily SQL  
Challenge  
#25

10 Most Frequent Words with TfidfVectorizer



10 Most Frequent Words with CountVectorizer



# BEST MODELS

1

>

2

>

3

Random Forest

Accuracy: **91.5%**

Precision: **91.5%**

Multinomial Bayes

Accuracy: **90.4%**

Precision: **93%**

Logistic Regression

Accuracy: **90.5%**

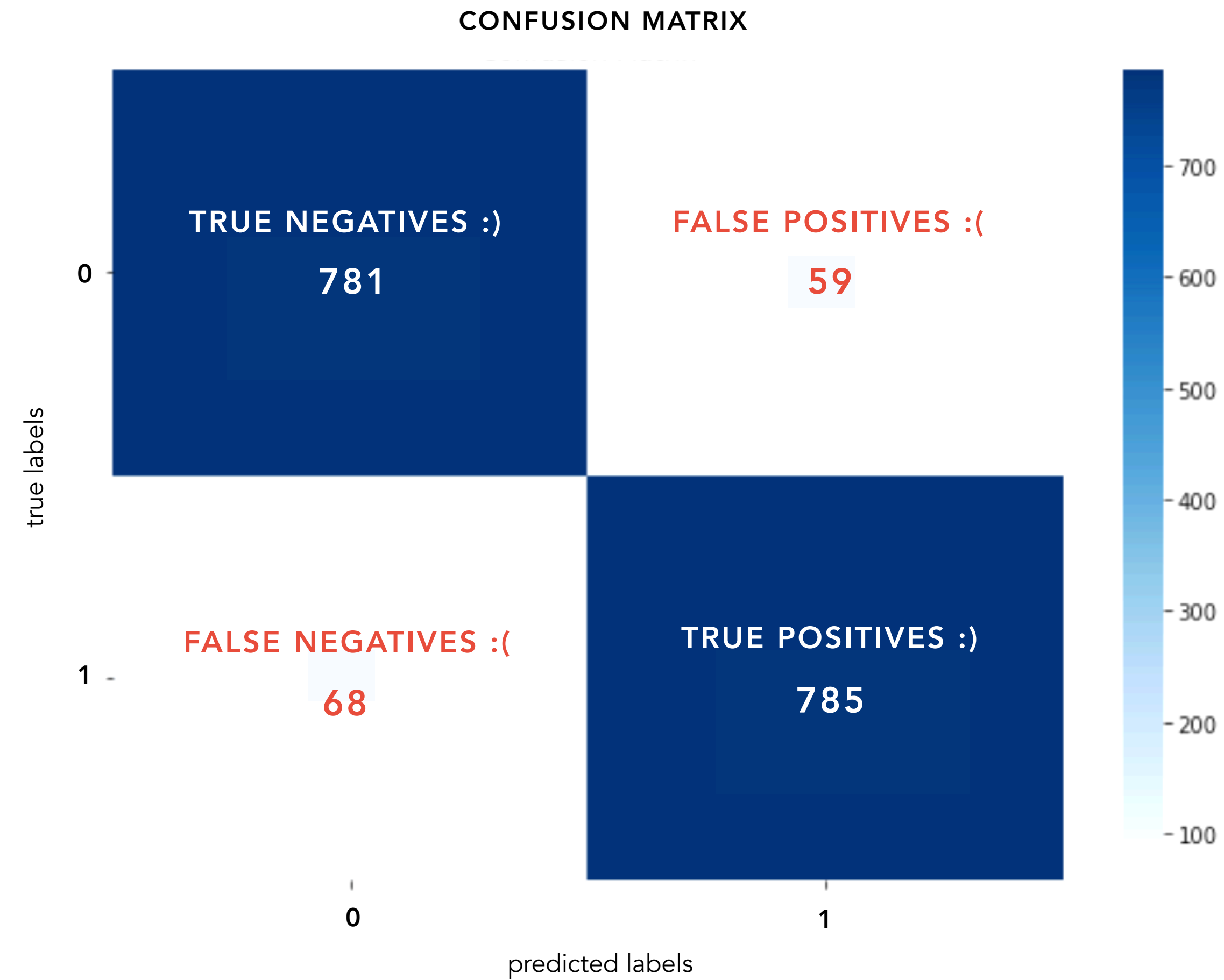
Precision: **90.9%**

# ENSEMBLE MODEL



Accuracy score: 92.5% / Precision score: 93%

# ENSEMBLE MODEL



# SUGGESTIONS

For someone else looking to solve a similar classification problem, I'd recommend the following:

- Try using boosting algorithms on highest-performing models
- Make sure one user isn't carrying the whole subreddit.
- Consider finding two subreddits that are more "evenly matched", or be careful about balancing classes.



ପ୍ର_ପ୍ର	
<hr/>	
5475	1

THANK YOU