

Evaluating the Differences between r/lifehacks and r/LifeProTips

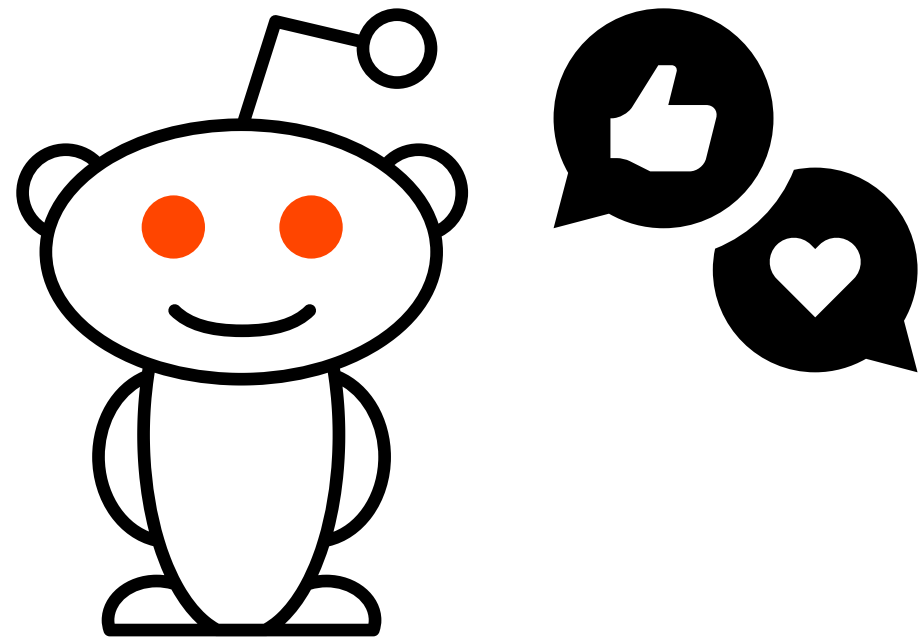
Franco Sanchez

Problem Statements

- Create a model that predicts if a Reddit post belongs to `r/lifehacks` or `r/LifeProTips`, based on post contents.
- Identify the top 10 keywords that distinguish both subreddits.

Background

Two subreddits dedicated to sharing advice that makes life easier.



r/lifehacks (LH)

"Lifehacks: Uncommon solutions to common problems."

r/LifeProTips (LPT)

"Tips that improve your life in one way or another."

Example

r/lifehacks top post from April.



20.8k



Posted by u/friphazeph 18 days ago



Optimal pizza placement for small ovens

i.imgur.com/MORhdq...



Posted by u/eljuan1161 17 days ago

 6  8  7  3

Miscellaneous

LPT If you feel tired and want to sleep with kids in the house. Tell them to wake you up in about 30 minutes so we can start cleaning the house and they will do literally anything to avoid waking you up.

 **90.2k**   **1132 Comments** ...

Example

r/LifeProTips top post from April.

Modeling Process

01

Data wrangling

- Gathered data using the Pushshift Reddit API
- Aggregated data into single data frame

02

Data preprocessing

- Dropped nulls and other erroneous data
- Aggregated into single clean data frame
- X and y train-test splits

03 - 04

Feature Engineering

- Vectorization
- Tokenization
- Lemmatization

03 - 04

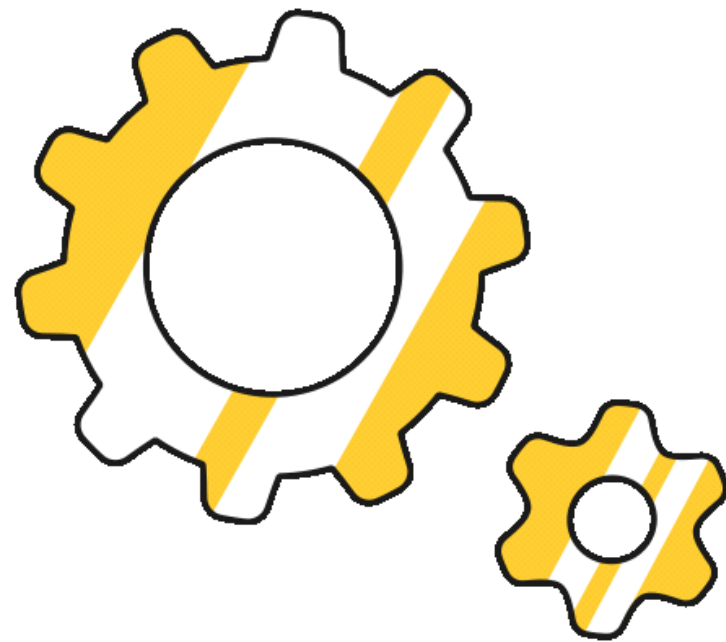
Modeling

- Classification models
- Tuning hyperparameters
- Performance

05

Insights

- Classification metrics
- Inferences



The Data

Feature Names	Descriptions
selftext	Body text of post
subreddit	r/lifehacks (1) or r/LifeProTips (0)
Rows: 2,476	LH: 0.241
Columns: 2	LPT: 0.759
Training split: 70% Test split: 30%	<i>*imbalanced data proportions</i>

The Models

Tried various models:

- Logistic Regression
- KNN
- Bagged Trees
- Random Forests
- Extra Trees
- Ada Boosting
- Ensembles

Baseline

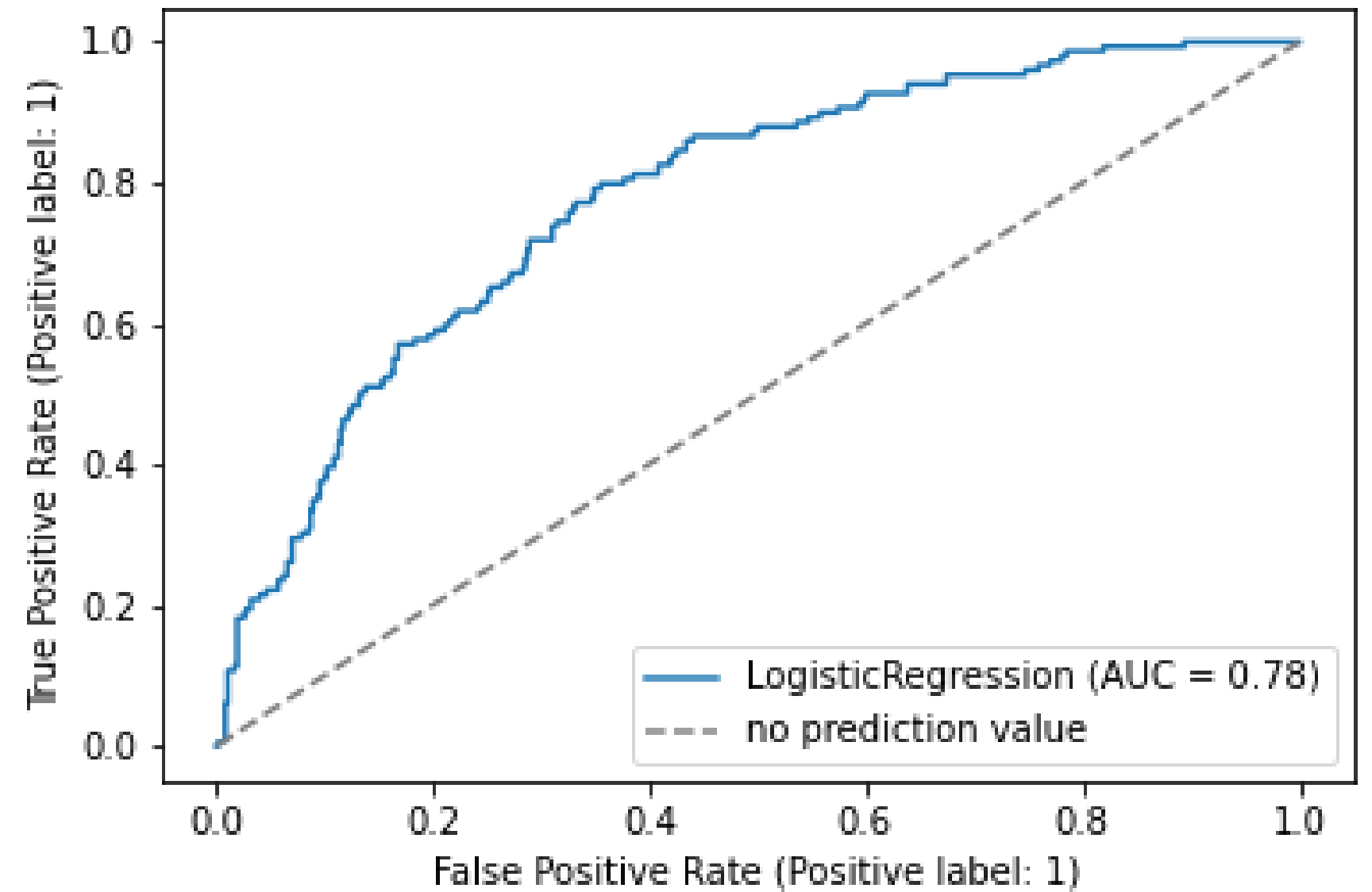
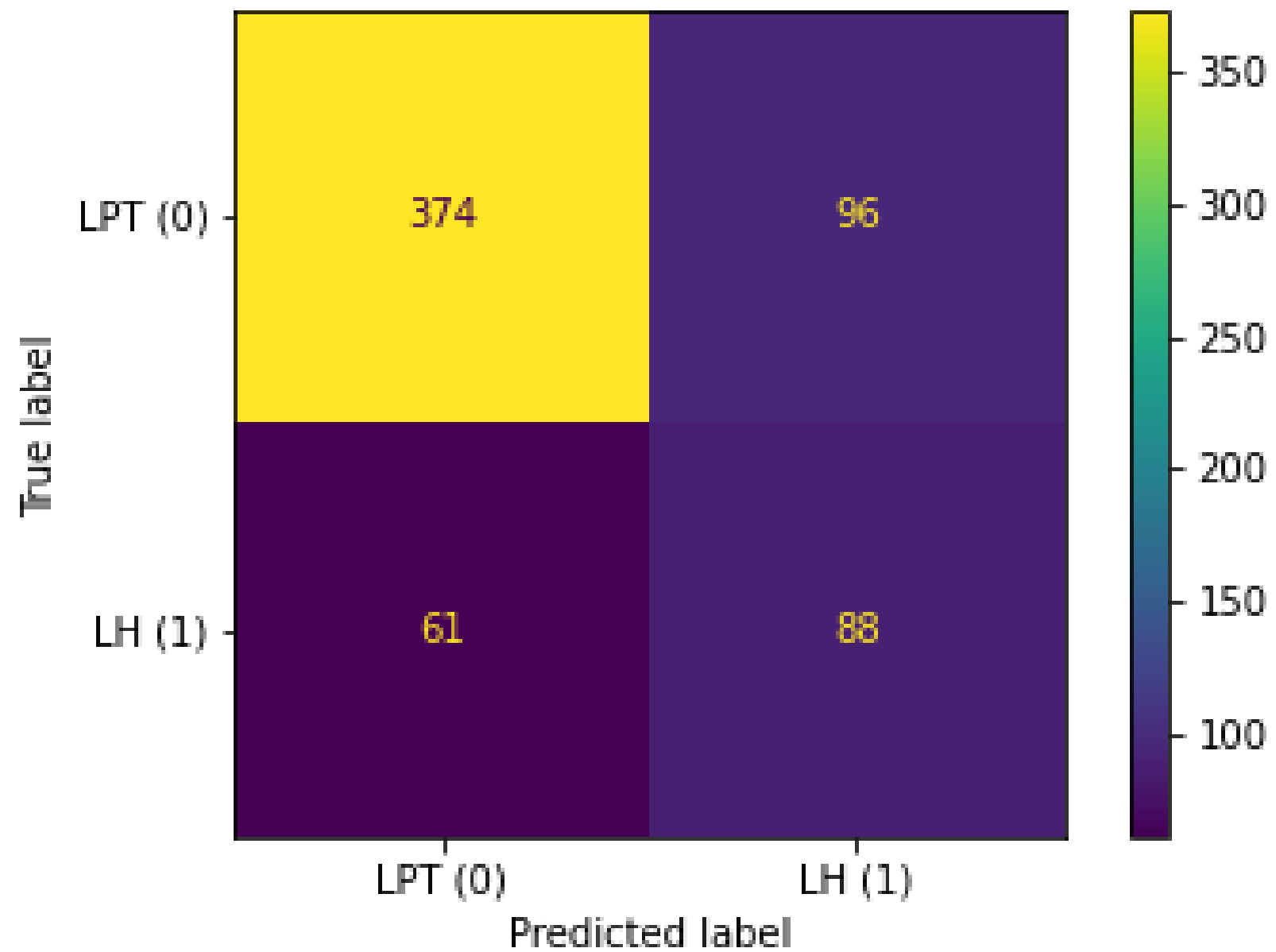
Accuracy	0.759
Recall	0
Specificity	1
Precision	Null (0)
F1 Score	0

Top Scores

Accuracy	0.746
Balanced Accuracy	0.693
Recall	0.591
Specificity	0.796
Precision	0.478
F1 Score	0.528

Best performing: **Logistics Regression**

Performances



Insights

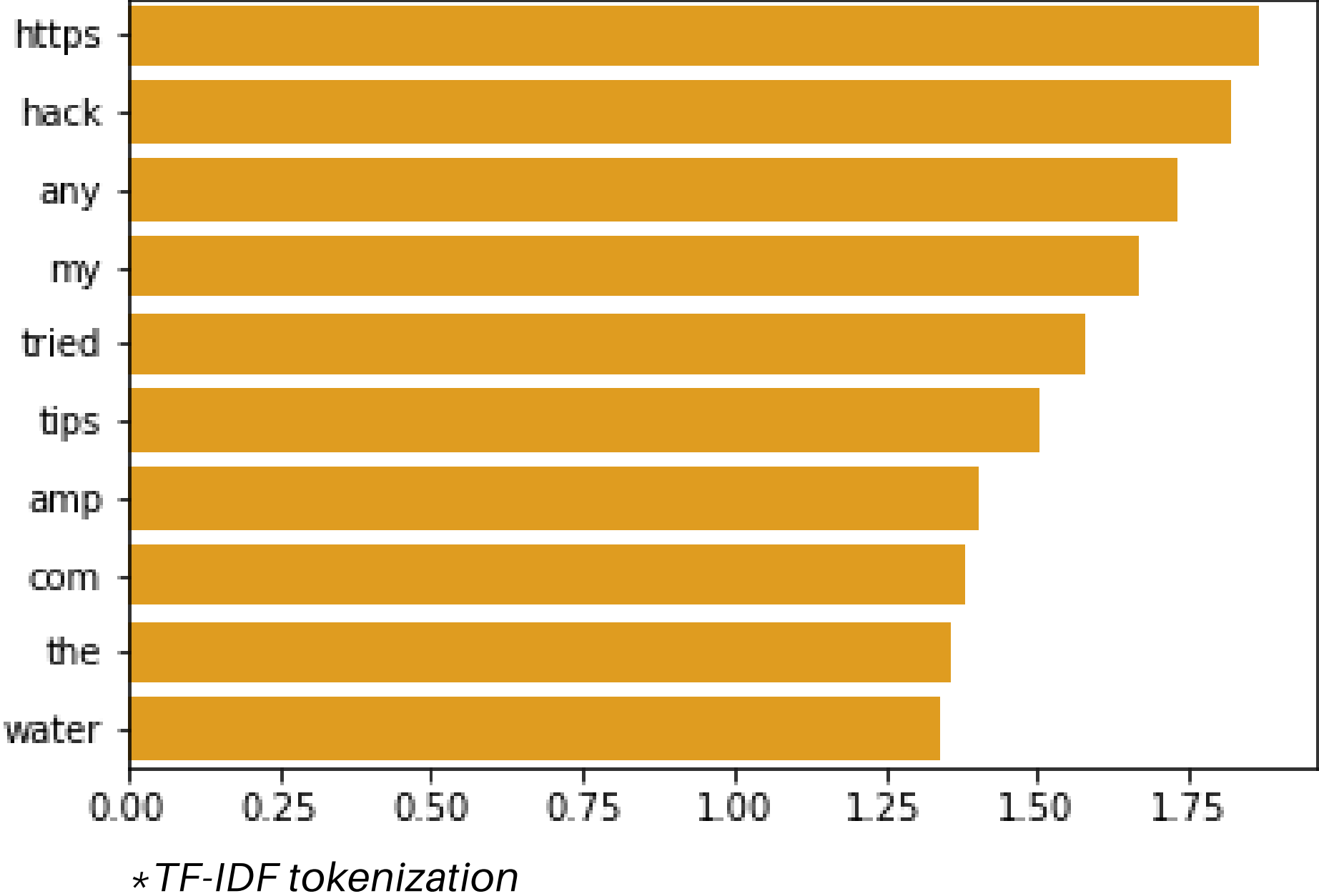
Classes

Label	Count	Proportion
True Positive	88	0.48
False Positive	96	0.52
False Negative	61	0.14
True Negative	374	0.86

Metrics Inference

Recall	The proportion of correctly predicted <code>r/lifehacks</code> posts over actual <code>r/lifehacks</code> posts
Precision	The proportion of correctly predicted <code>r/lifehacks</code> posts over all <code>r/lifehacks</code> predictions
F1 Score	Harmonic mean of recall and precision

Top 10 most impactful words in predicting LH



Discussion

Imbalanced data requires more feature engineering.

- Add more relevant stopwords
- Oversampling / undersampling
- Eliminate more erroneous posts (bots/spam posts)

Shift classification threshold to account for imbalanced data.

Wrangle more balanced data from the subreddits.

Thank you