# Subreddit Classification

Manjila Moktan
March 21, 2025

# BACKGROUND

- Declutter focuses on reducing clutter and organising spaces.

- Minimalism focuses on living with less.

# Problem Statement :

This project aims to scrap through two subreddits(Declutter and Minimalism) and apply Natural language Processing (NLP) techniques to build a supervised ML model that predicts if the post belongs to subreddit Declutter or Minimalism with high accuracy, despite their similarities and differences.
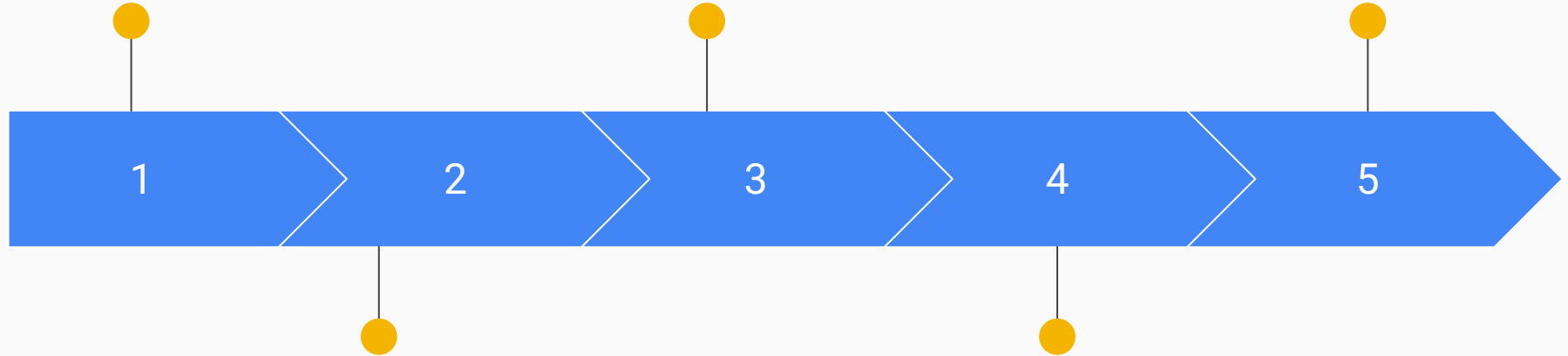
Data Collection:

reddit Praw (.new , .top)

Feature Extraction

-TF-IDF Vectorizer

Conclusion

1    2    3    4    5

Data Preprocessing:

-    Binarize the target
-    top 10 words
-    WordNetLemmatizer

Model Building
Evaluation

Top 20 Words in Declutter — Top 20 Words in Minimalism

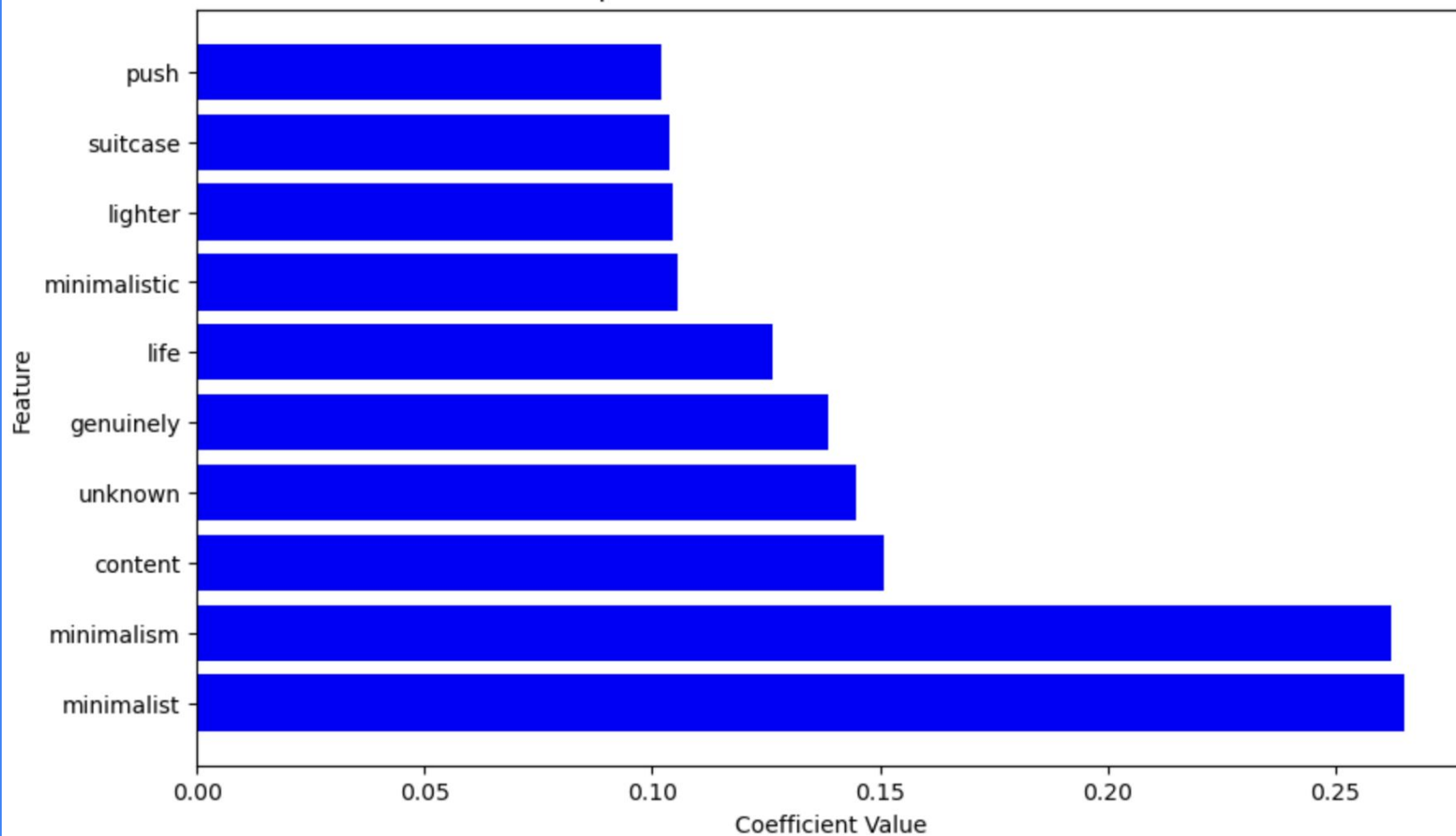Among the top 20 words from each subreddits in text; total of 40, only 26 words are unique!!!
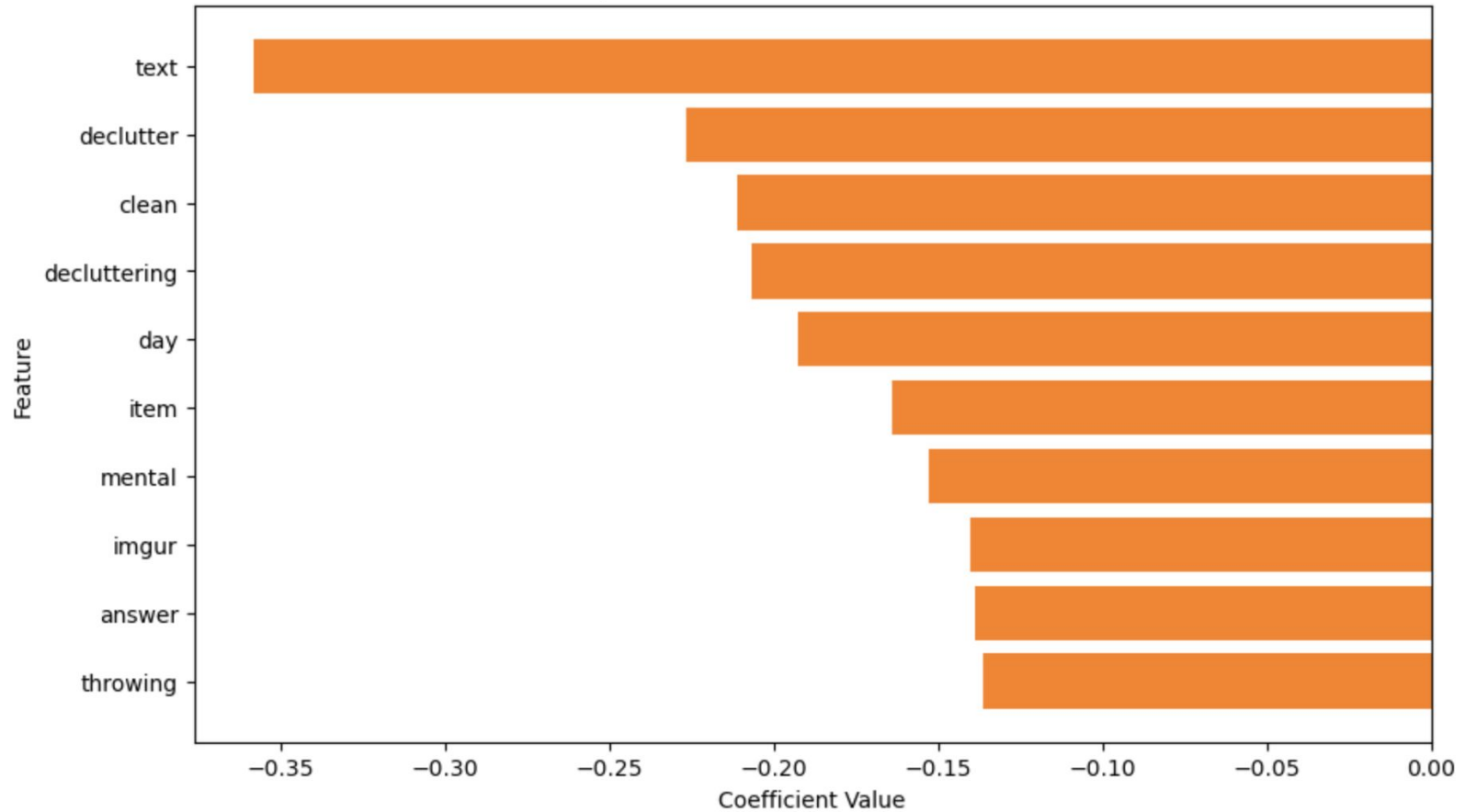
# Model Evaluation

- Logistic Regression Model

  (Accuracy = 77%)

|   | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.76      | 0.80   | 0.78     |
| 1 | 0.78      | 0.74   | 0.76     |

Top Features for Minimalism Subreddit

Top Features for Declutter Subreddit

# Conclusion And Recommendation

- Use of Logistic Regression
- Use aggressive stemming eg. PorterStemmer instead
- Exploring other techniques