# PERSONAL FINANCE
# OR
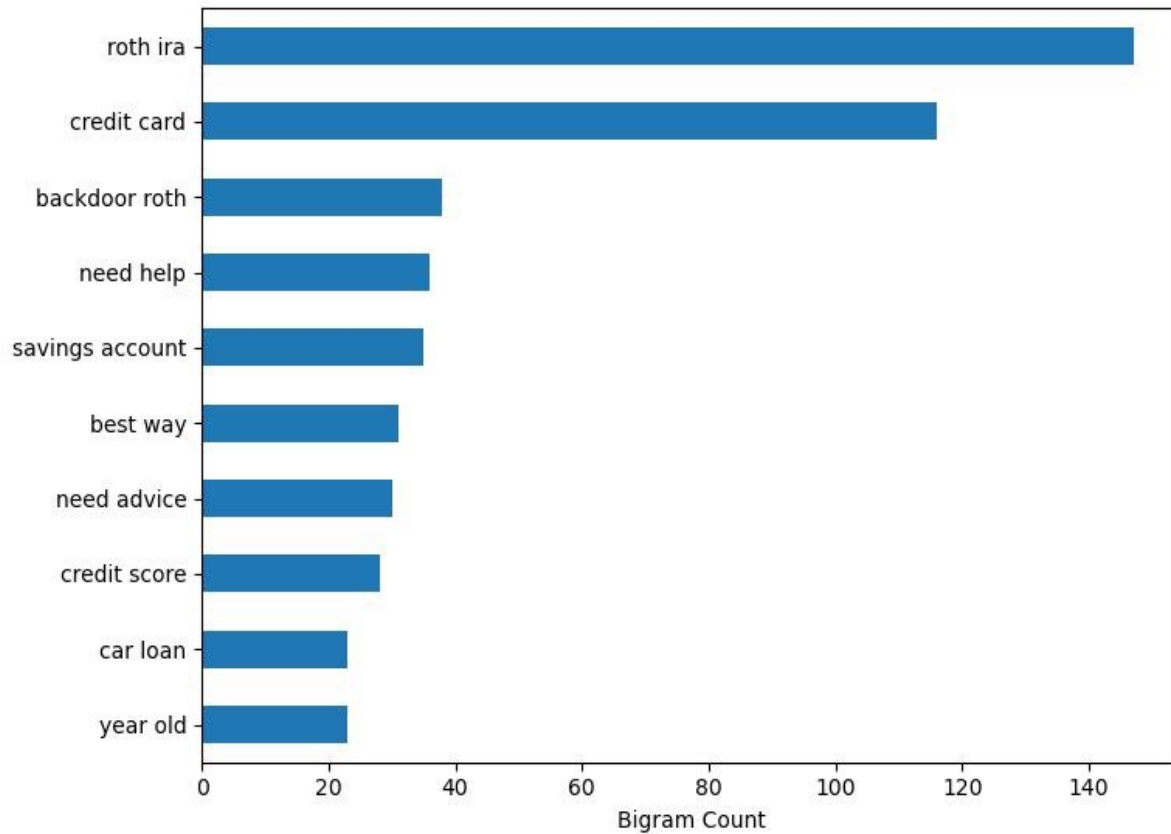# FINANCIAL INDEPENDENCE

by Osei Anom

# EXECUTIVE SUMMARY:

- collect unstuctured data, pre-process that data, then build a predictive classification model.

- subreddit's "Personal Finance" and "Financial Independence"

- Which classification model best predicts the subreddit a post will fall under

1 —————— Extracted Data using Pushshift API

2 —————— Exploratory Data Analysis

3 —————— Initial Modeling

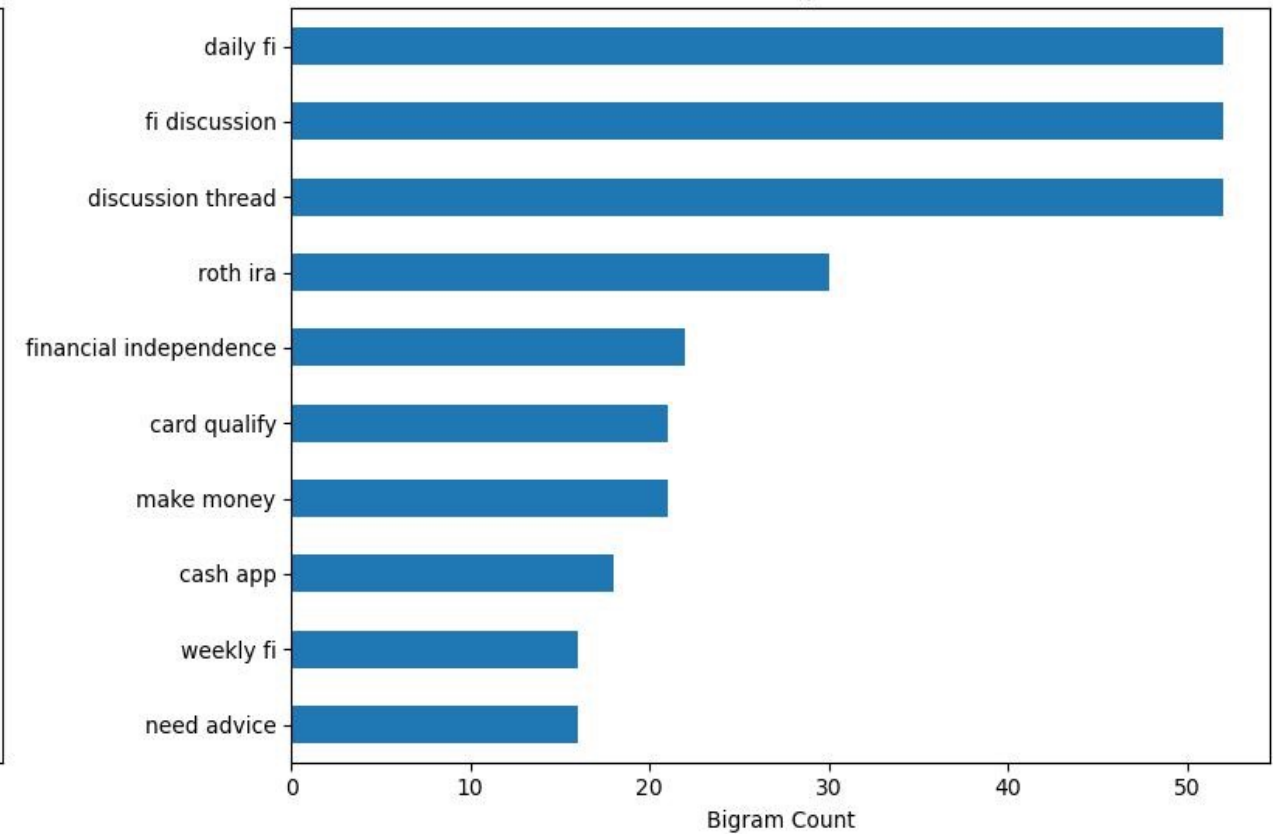4 —————— Secondary Modeling with additional data

\* Adding selftext with title

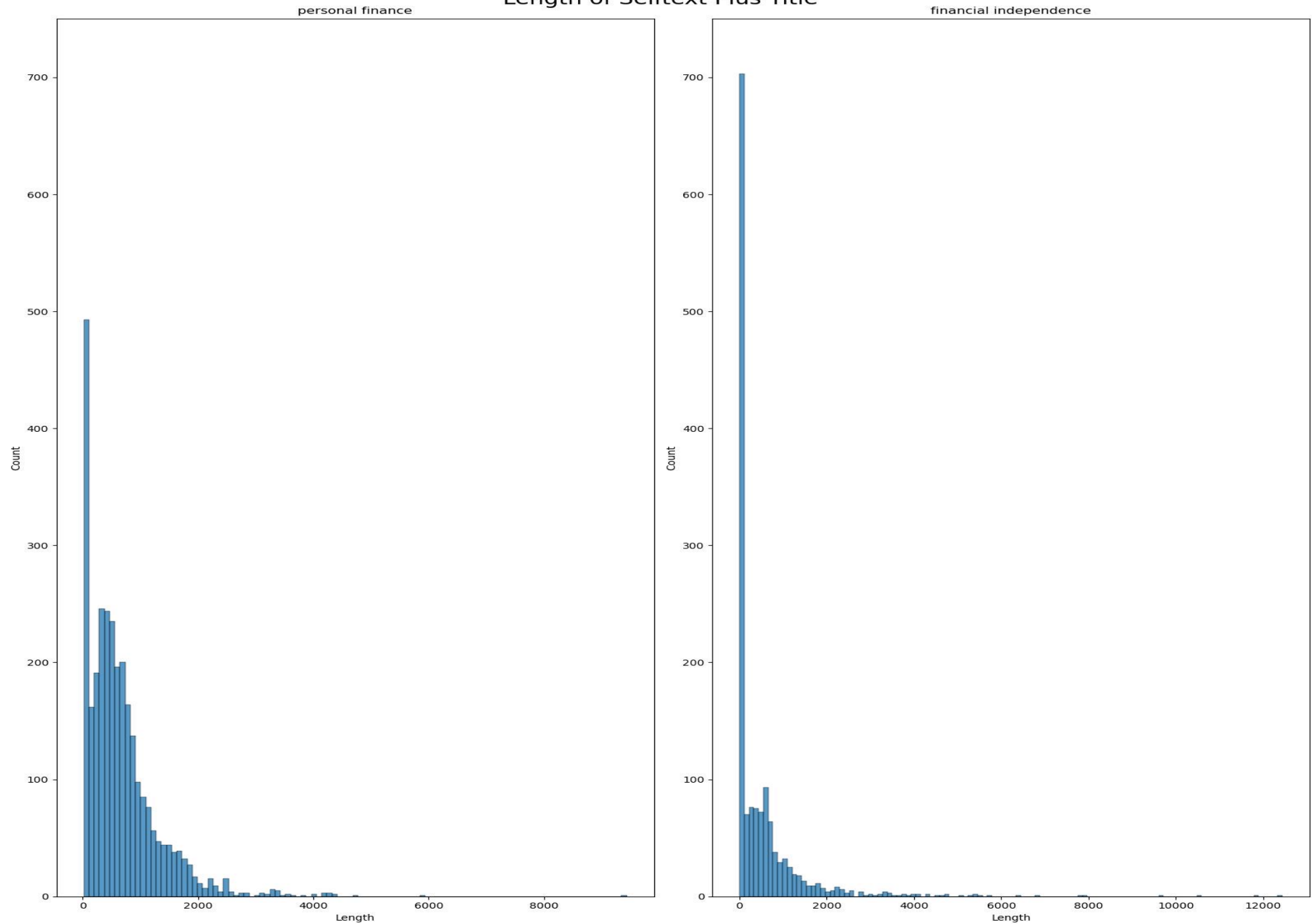PROCESS

# Top Bigrams by Subreddit



Personal Finance

| Bigram | Count |
|---|---|
| roth ira | ~147 |
| credit card | ~115 |
| backdoor roth | ~38 |
| need help | ~35 |
| savings account | ~34 |
| best way | ~31 |
| need advice | ~29 |
| credit score | ~27 |
| car loan | ~22 |
| year old | ~22 |

Financial Independence

| Bigram | Count |
|---|---|
| daily fi | ~52 |
| fi discussion | ~52 |
| discussion thread | ~52 |
| roth ira | ~30 |
| financial independence | ~22 |
| card qualify | ~21 |
| make money | ~21 |
| cash app | ~18 |
| weekly fi | ~16 |
| need advice | ~16 |

# Length of Selftext Plus Title
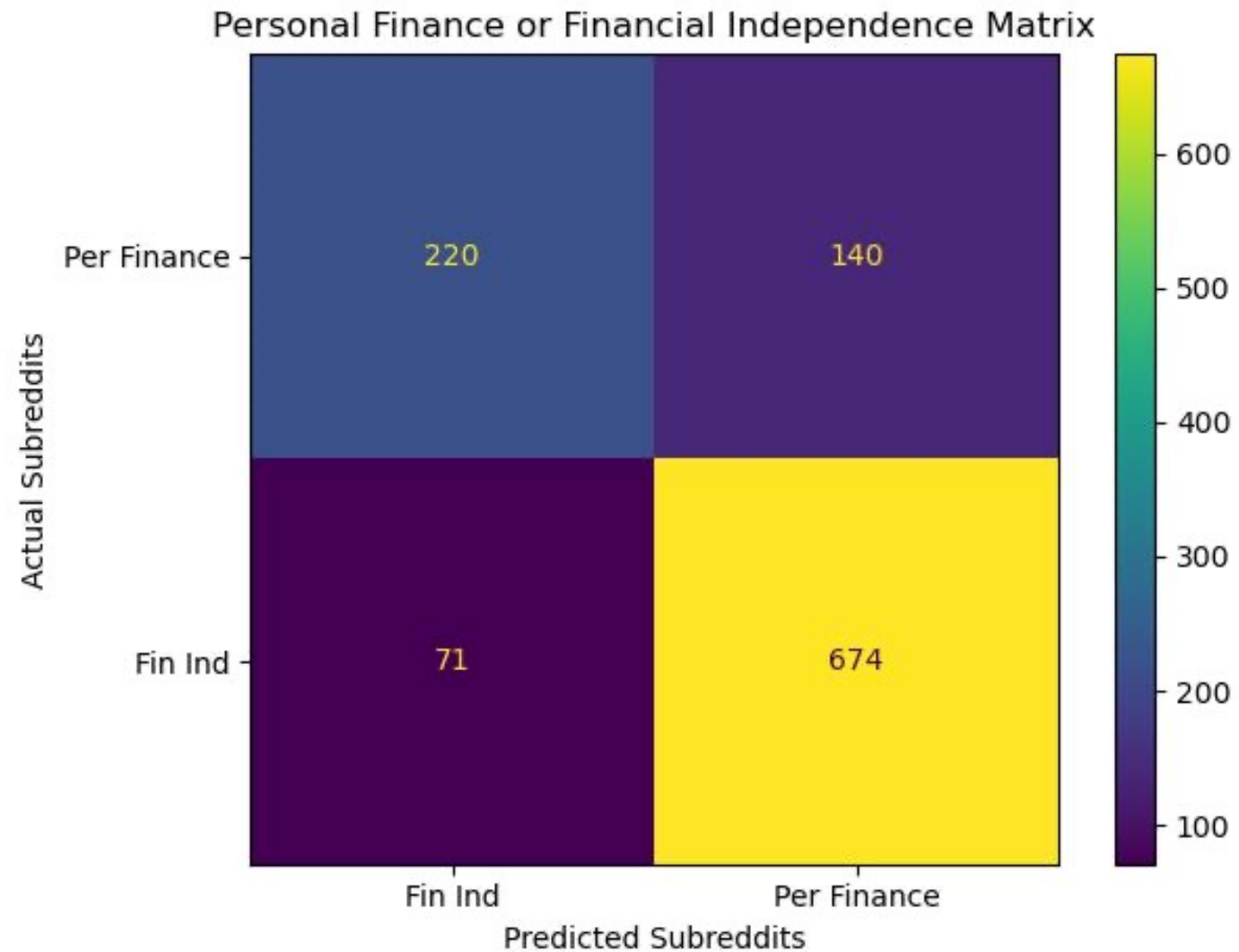
### personal finance

### financial independence

# METHODOLOGY, INFERENCES, ASSUMPTIONS:

- Used Tfid Vectorizer through gridsearch running Logistic Regression, Decision Tree, Multinomial Naive Bayes, KNeighbors, Random Forest, AdaBoost, and finally Stacking

- Used F1 to determine performance between models

- Giving additional data would help increase test accuracy and F1 score

# BEST PERFORMING MODELS
# (TEST SCORES)

|  | MN BAYES | RANDOM FOREST | ADA BOOST | STACKED |
|---|---|---|---|---|
| F1 | 84 % | 85 % | 84 % | 86 % |
| Accuracy | 76 % | 78 % | 78 % | 81 % |
| Misclassification | 24 % | 22 % | 22% | 19 % |

# EVALUATION

# SUMMARY & CONCLUSIONS

• All models were overfit

•TBD if including word count and word length in post and self text would help.

• The model over predicted personal finance subreddit (false positives), partially due to the imbalance of classes

• Would like to explore over / under sampling, Grid Searching with balanced accuracy,  synthetic data