

Classification Project

r/investing vs r/personalfinance

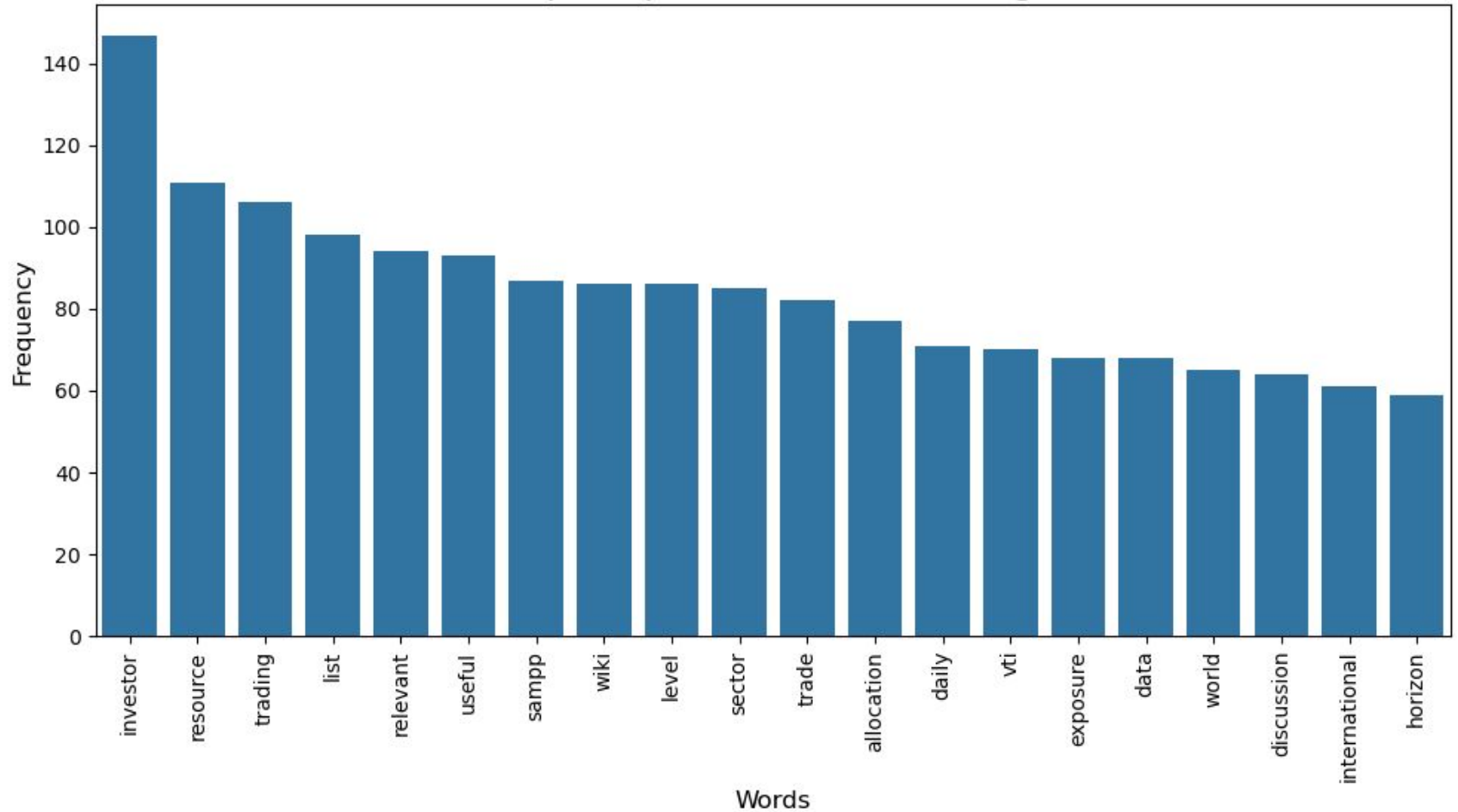
October 12, 2024

Problem Statement

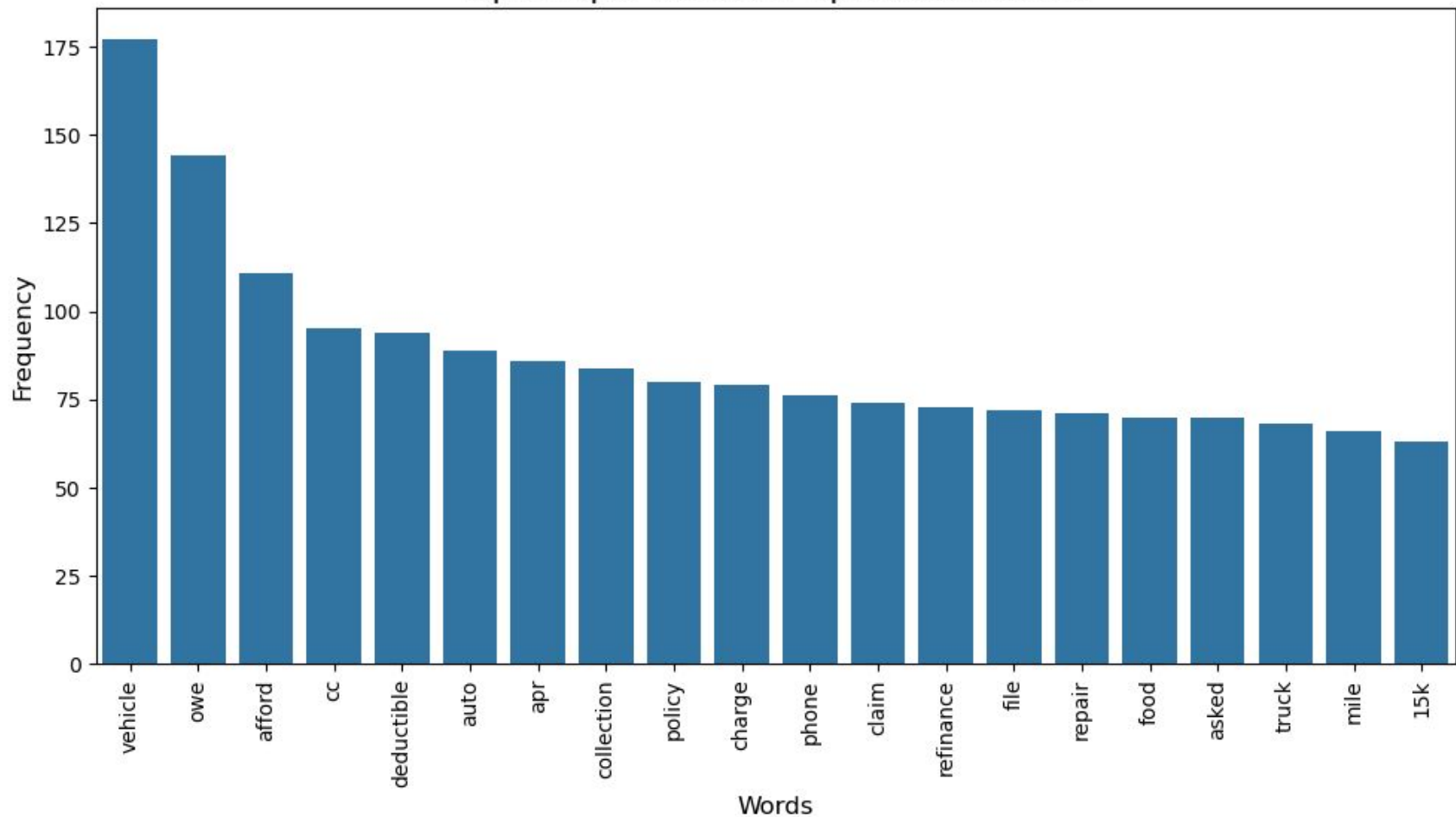
Develop a **classification model** that is able to distinguish between two types of users based on their Reddit posts:

- Users who are experienced investors (professionals)
 - Users who are interested in finance but are not heavily investing yet (amateurs).
-

Top Unique Words in r/investing



Top Unique Words in r/personalfinance



Logistic Regression + CountVectorizer: 0.84

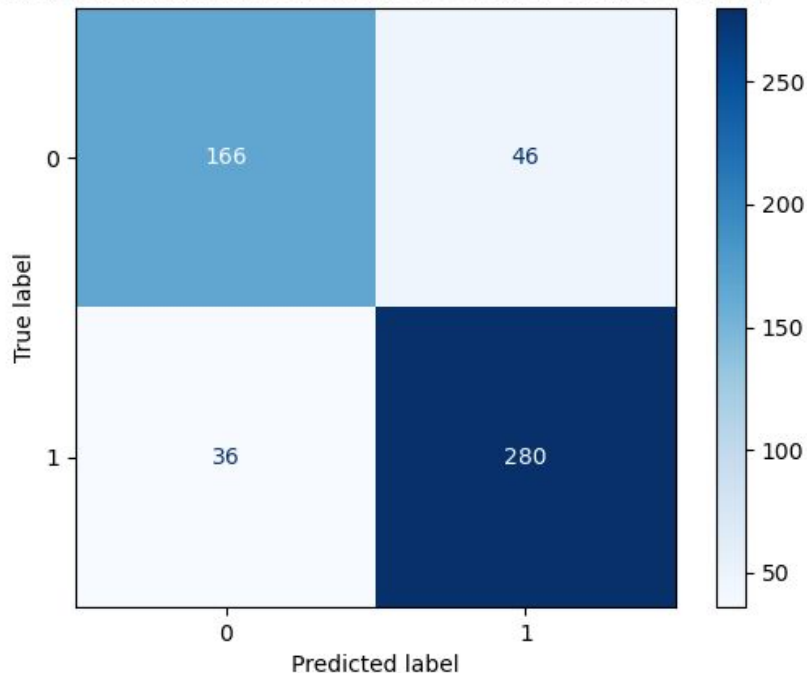
r/investing

- Precision: 0.82
- Recall: 0.78
- F1- score: 0.80

r/personalfinance

- Precision: 0.86
- Recall: 0.89
- F1- score: 0.87

Confusion Matrix for Logistic Regression + CountVectorizer



Logistic Regression + TF-IDF: 0.85

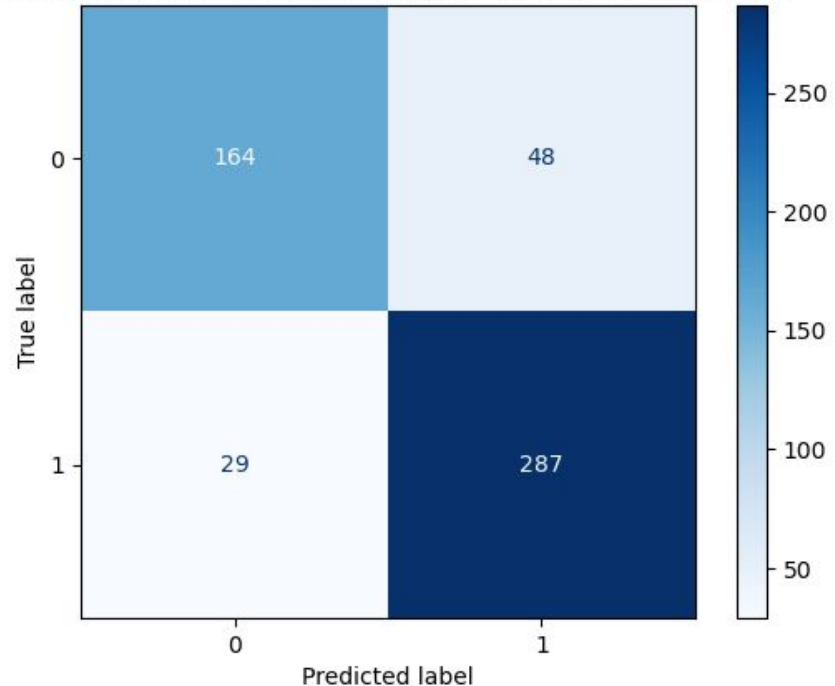
r/investing

- Precision: 0.85
- Recall: 0.77
- F1- score: 0.81

r/personalfinance

- Precision: 0.86
- Recall: 0.91
- F1- score: 0.88

Confusion Matrix for Logistic Regression + TF-IDFVectorizer



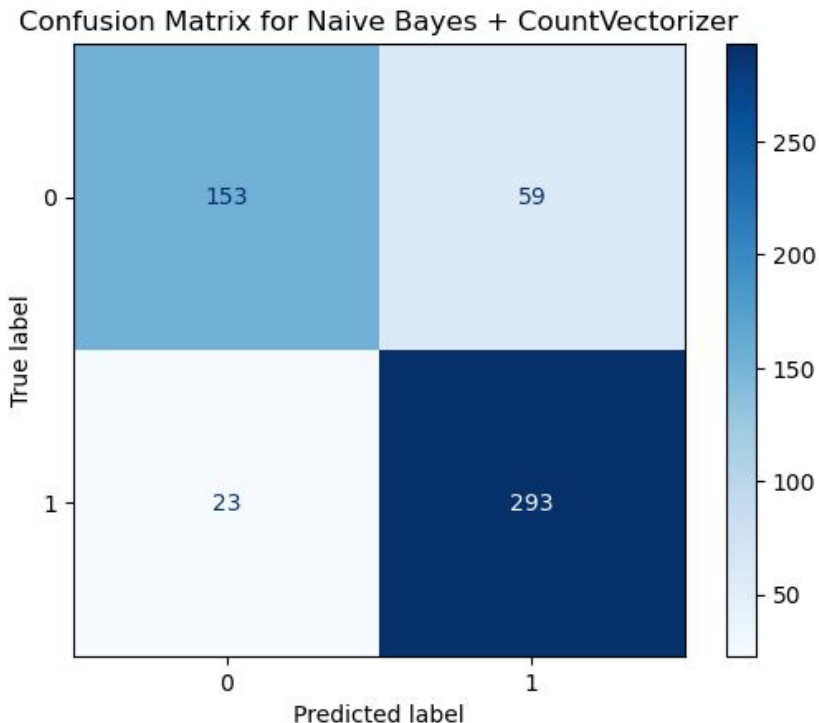
Naive Bayes + CountVectorizer: 0.85

r/investing

- Precision: 0.82
- Recall: 0.81
- F1- score: 0.81

r/personalfinance

- Precision: 0.87
- Recall: 0.88
- F1- score: 0.88



Naive Bayes + TF-IDF: 0.84

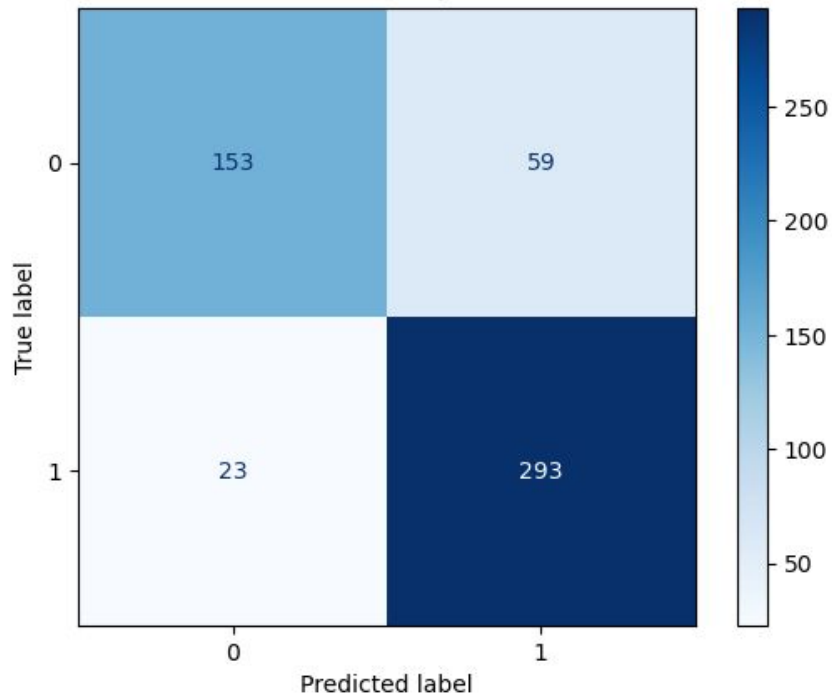
r/investing

- Precision: 0.87
- Recall: 0.72
- F1- score: 0.79

r/personalfinance

- Precision: 0.83
- Recall: 0.93
- F1- score: 0.88

Confusion Matrix for Naive Bayes + TF-IDFVectorizer



Best Model

Logistic Regression
with TF-IDF



Thank you!