

UNRESOLVED MYSTERIES VS FAN THEORIES SUBREDDITS

Binita Patel
Data Scientist



TABLE OF CONTENTS



01

**DATA SCIENCE
PROBLEM**

02

**PROCEDURE /
METHODOLOGY**

03

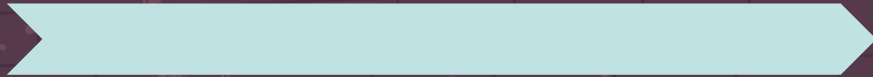
**PRIMARY
FINDINGS**

04

**NEXT STEPS /
CONCLUSIONS**

01

DATA SCIENCE PROBLEM



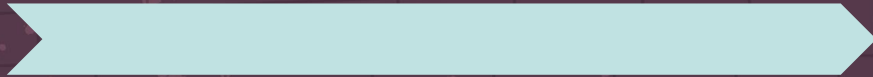
DATA SCIENCE PROBLEM

- Co-worker is about to quit job
- Goes through some of the subreddits and messes them up
- Build a classification model
- Focus on two subreddits :
UnresolvedMysteries and FanTheories.



02

PROCEDURE / METHODOLOGY



PROCEDURE/METHODOLOGY



DATA CLEANING

Cleaned dataset by removing Nulls



EDA

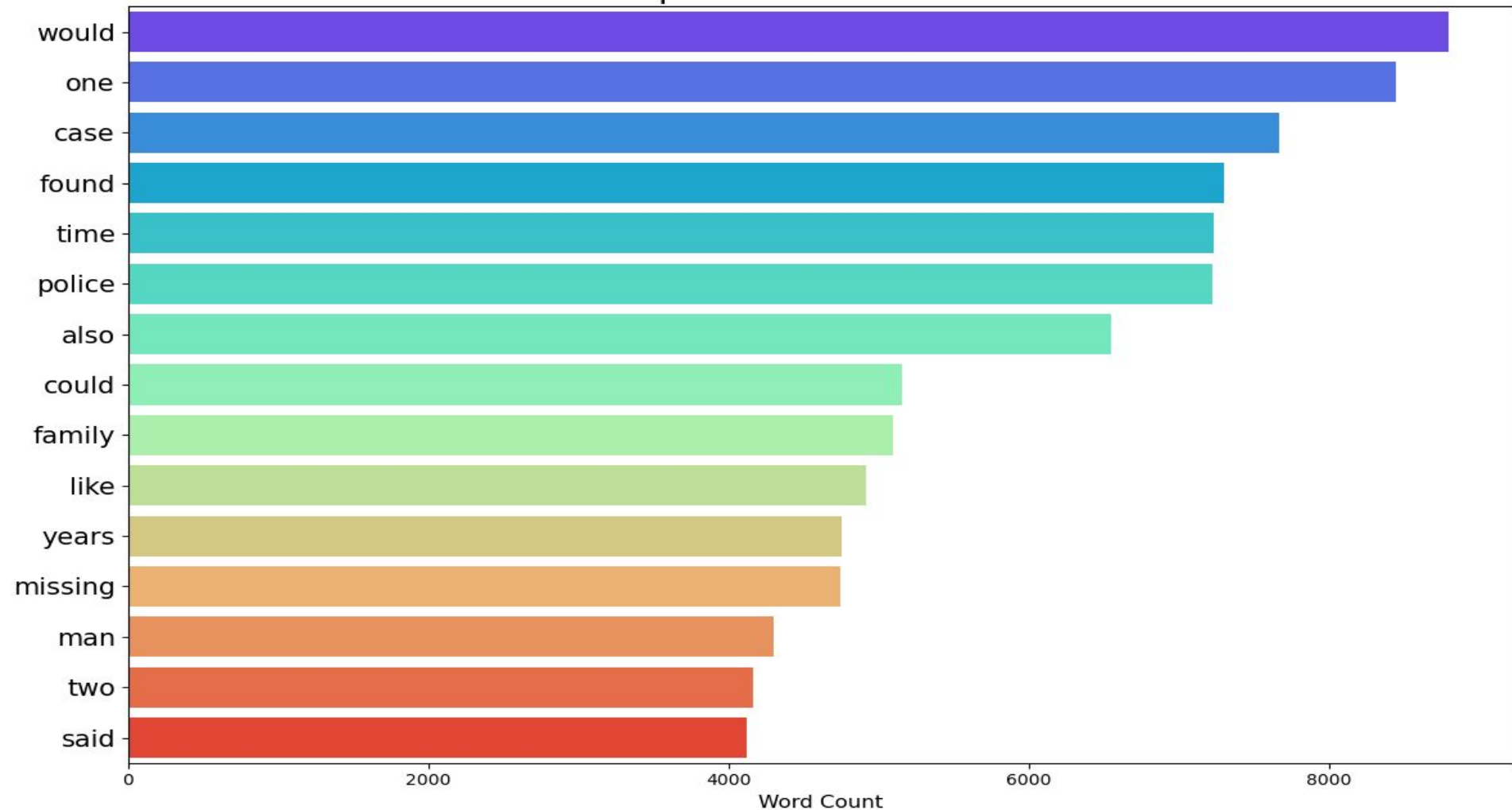
Looked at top words and did sentiment analysis



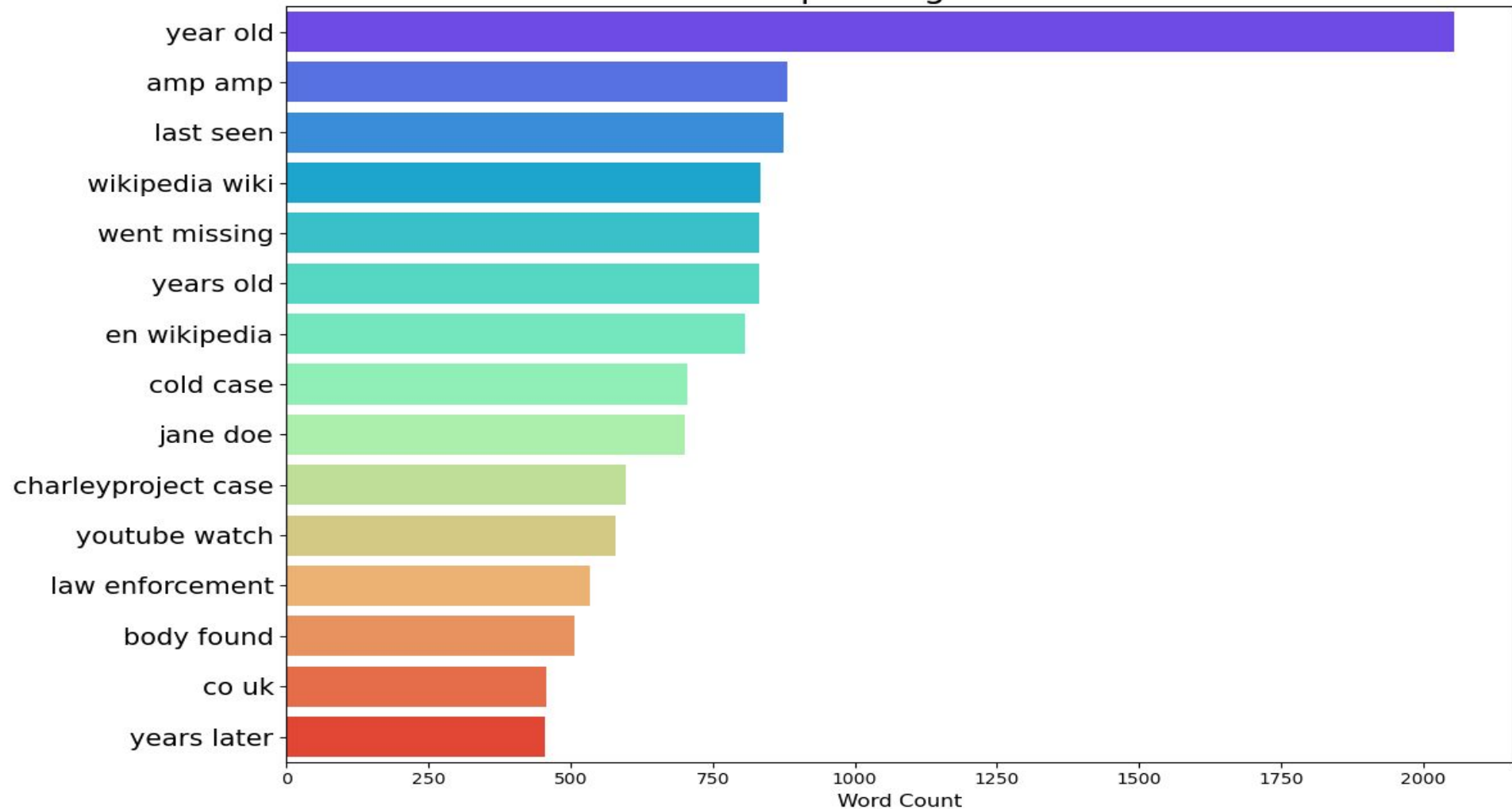
MODELING

Did preprocessing and ran 6 models

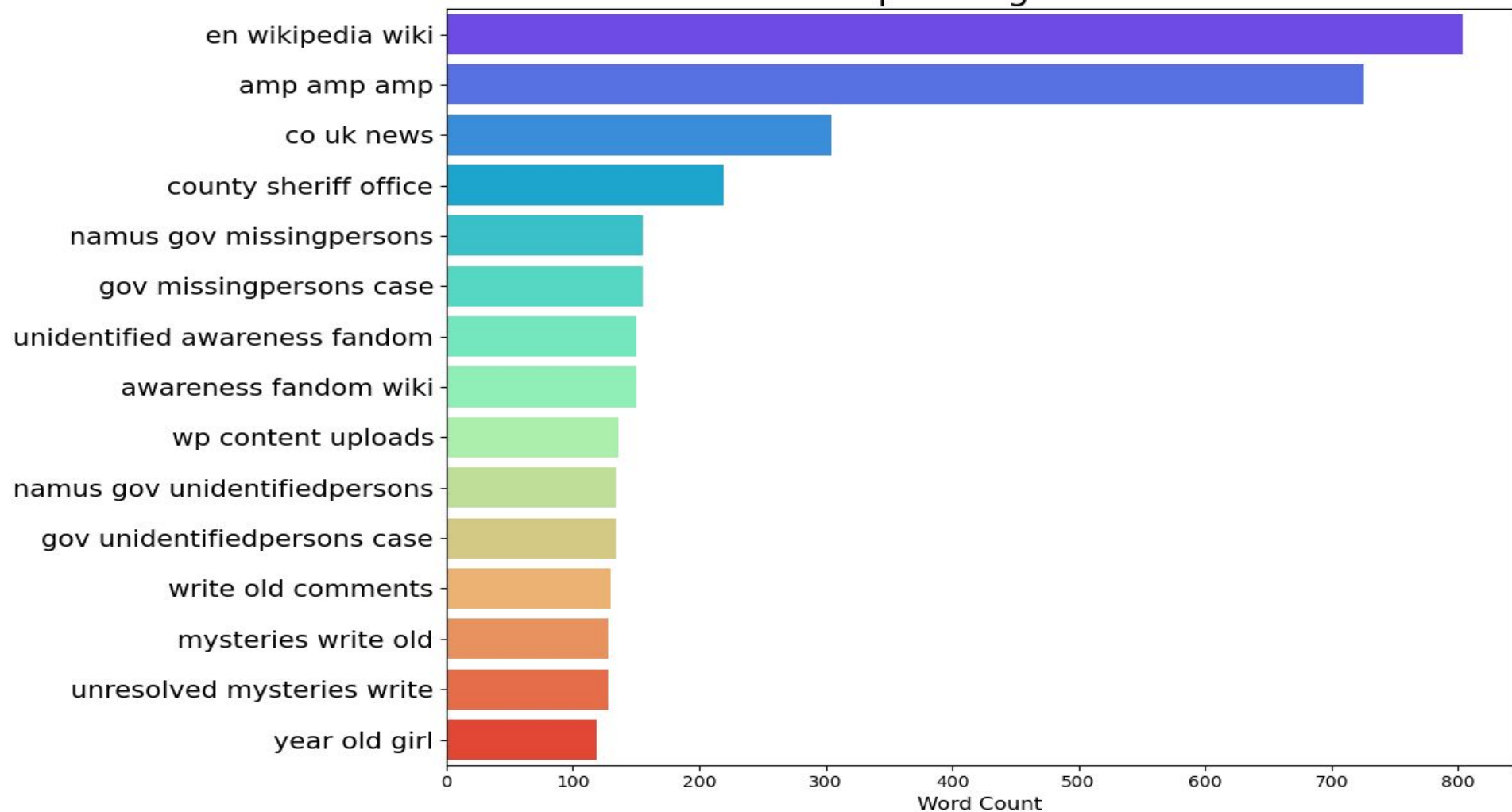
Top 15 common words



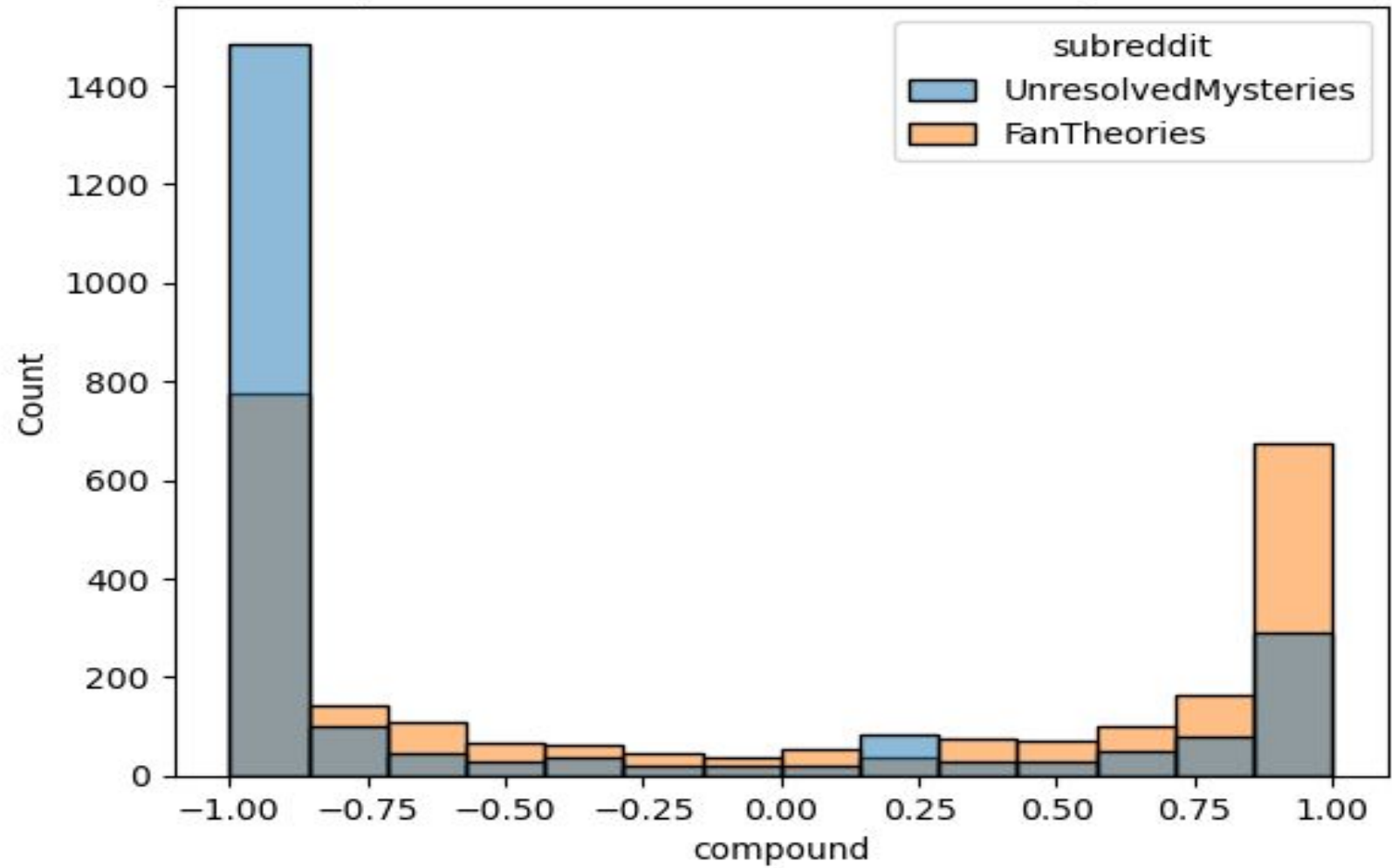
Top 15 bigrams



Top 15 trigrams



Distribution of Compound Score between UnresolvedMysteries and FanTheories



03

PRIMARY FINDINGS



MODELS

| Num | Vectorizer | Model | Best Score | Train Score | Test Score (Accuracy) | Best Parameters |
|-----|-----------------|--------------------------|------------|-------------|-----------------------|--|
| 1 | CountVectorizer | Multinomial Naïve Bayes | 0.9873 | 0.9969 | 0.9788 | {'cvec__max_features': None, 'cvec__min_df': 2, 'cvec__ngram_range': (2, 2)} |
| 2 | TfidfVectorizer | Multinomial Naïve Bayes | 0.9873 | 0.9975 | 0.9797 | {'tvec__max_features': None, 'tvec__min_df': 2, 'tvec__ngram_range': (2, 2)} |
| 3 | CountVectorizer | Random Forest Classifier | 0.9735 | 1.0 | 0.9755 | {'rf__max_depth': None, 'rf__min_samples_leaf': 1, 'rf__n_estimators': 150} |
| 4 | TfidfVectorizer | Random Forest Classifier | 0.9740 | 1.0 | 0.9772 | {'rf__max_depth': None, 'rf__min_samples_leaf': 1, 'rf__n_estimators': 500} |

04

NEXT STEPS / CONCLUSIONS



NEXT STEPS



REDDITS

Look at other similar
reddits



POST TITLES

Use post titles instead of
post content



MISPLACED BUTTON

Misplaced Subreddit Post
button



THANKS

Do you have any questions?

Binita Patel

binita.patel.bp@gmail.com

CREDITS: This presentation template was created by **Slidesgo**, and includes icons by **Flaticon**, and infographics & images by **Freepik**