# Predicting Subreddits

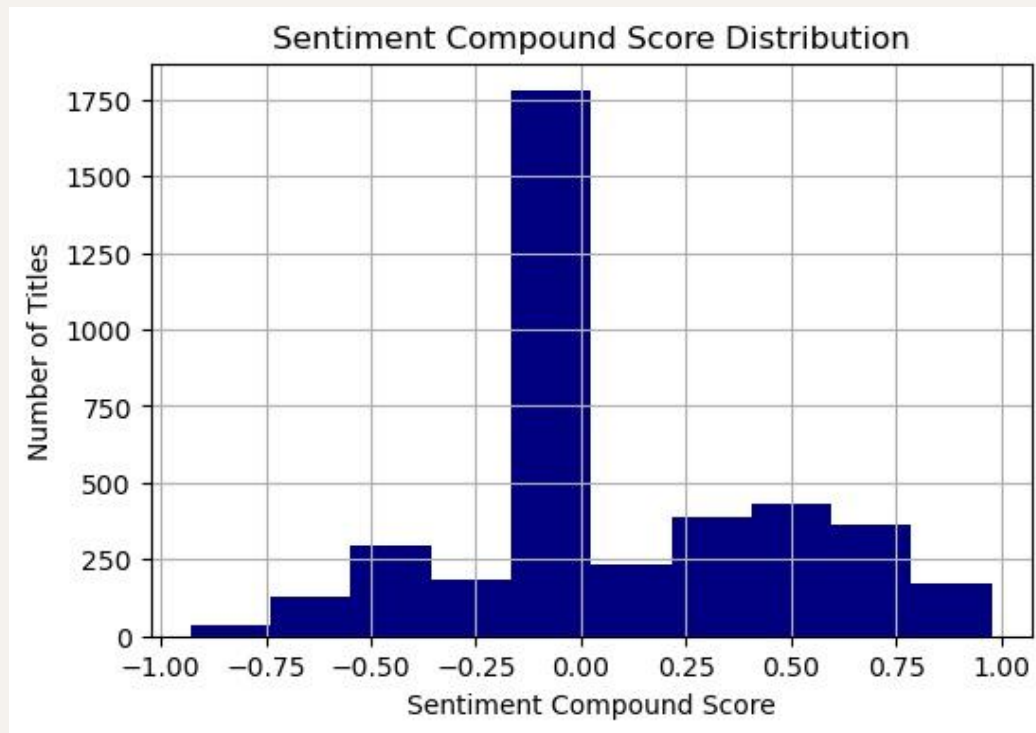March 4, 2022

# NBA or Chess?

# Overview

1. EDA
2. Modeling
3. Conclusions
4. Next Steps
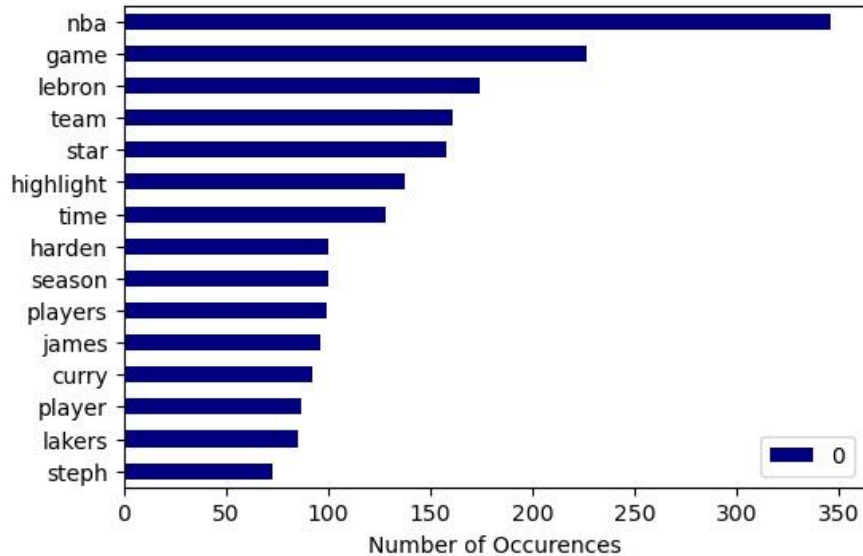
# EDA

## Engineered Features

- Title + Self-text

- Word Count
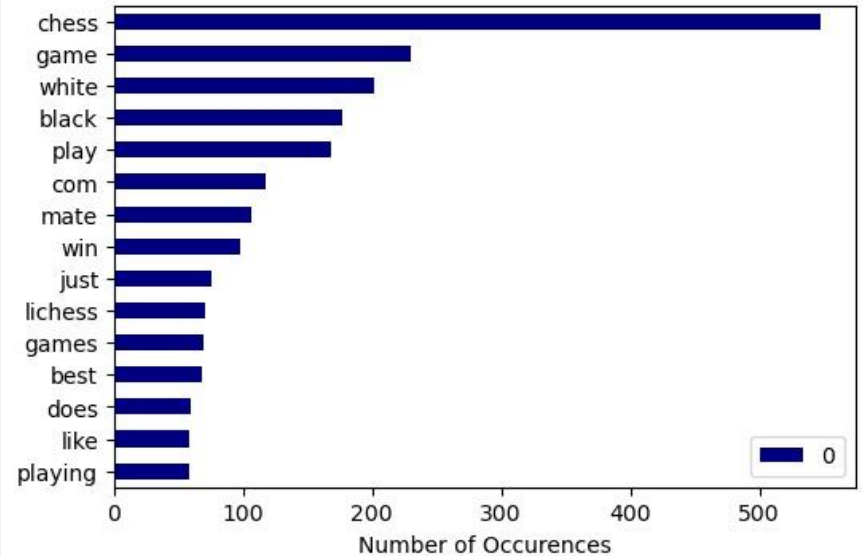
- Length (in characters)
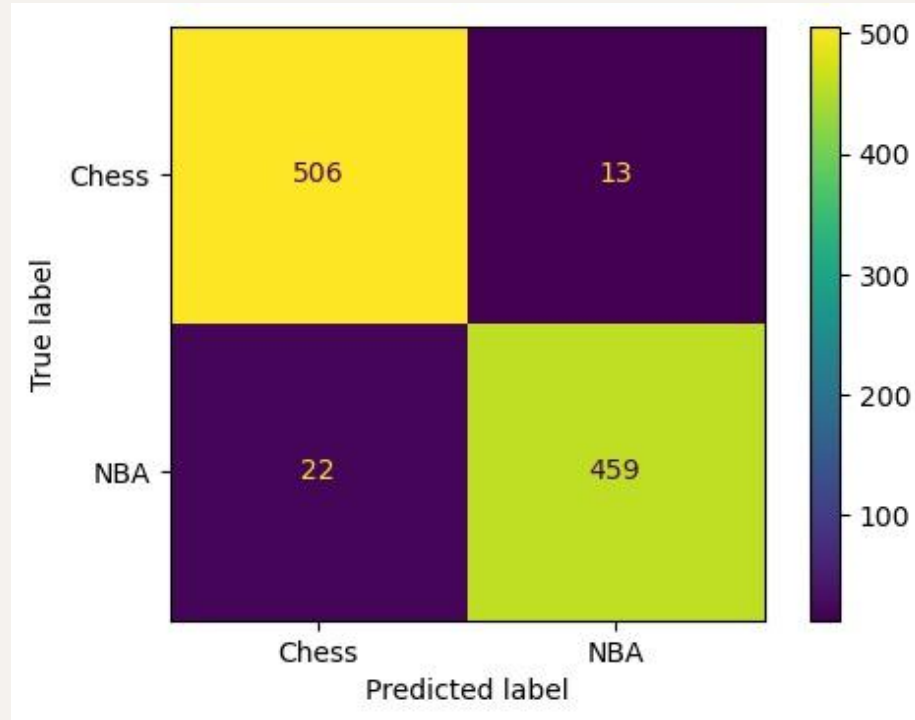
- Sentiment Analysis



Sentiment Compound Score Distribution

# EDA



Most Common Words in NBA Subreddit

Most Common Words in Chess Subreddit
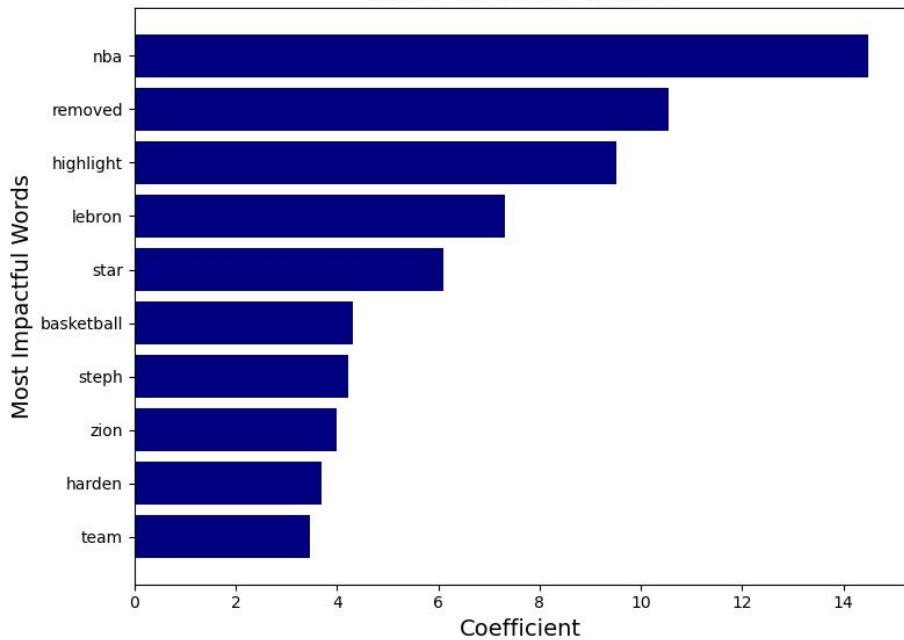
# Modeling

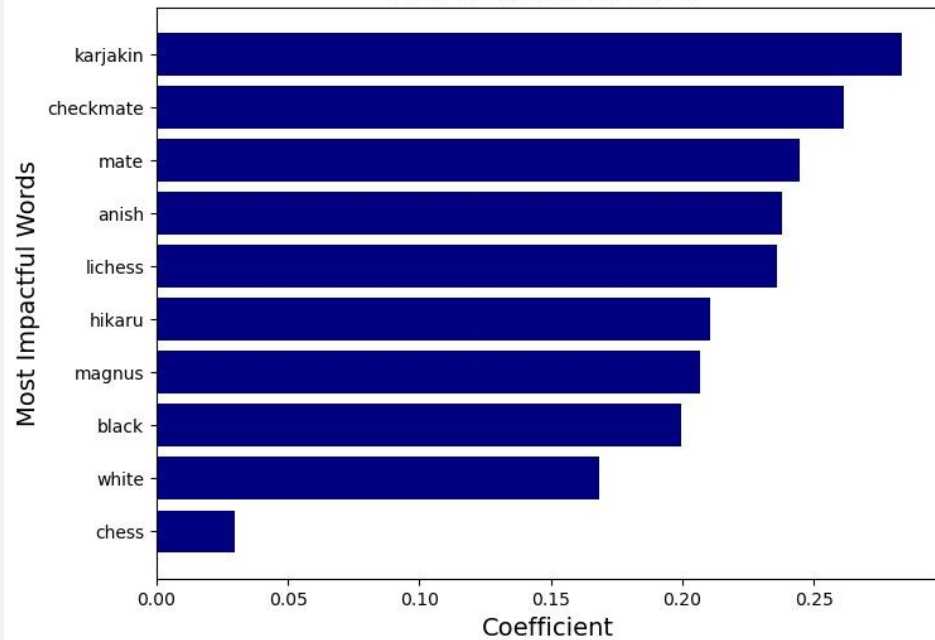| | Model 1 | Model 2 |
|---|---|---|
| Model Type | Logistic Regression | Random Forest |
| Text Analysis | CountVectorizer | TFIDF |
| Other Methodologies | Pipeline, FeatureUnion, GridSearchCV | Pipeline, GridSearchCV |
| Parameters | stop words, max features, min df, max df | stop words, max features, min df |
| Accuracy Score | 95.9% | 95.5% |

# Confusion Matrix

## Best Predictors For NBA

| Most Impactful Words | Coefficient |
| --- | --- |
| nba | |
| removed | |
| highlight | |
| lebron | |
| star | |
| basketball | |
| steph | |
| zion | |
| harden | |
| team | |

## Best Predictors for Chess

| Most Impactful Words | Coefficient |
| --- | --- |
| karjakin | |
| checkmate | |
| mate | |
| anish | |
| lichess | |
| hikaru | |
| magnus | |
| black | |
| white | |
| chess | |

# Conclusions & Next Steps

## Conclusions

- Yes, it was possible to create a model that beats the baseline model of 50% Accuracy.
- Consider the relationship between the two subreddits.

## Next Steps

- Additional CountVectorize tuning (different token patterns)
- More feature engineering
- Look into time patterns.