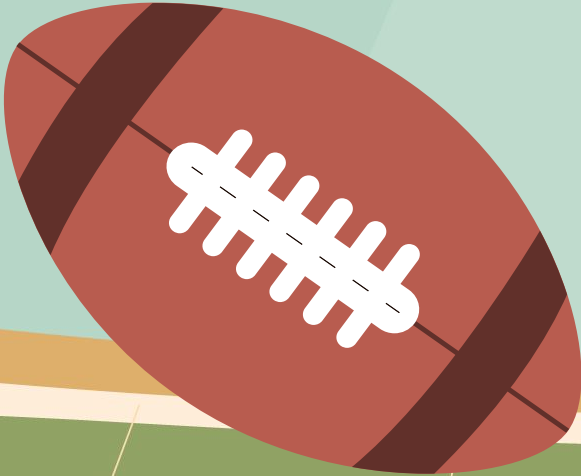


Subreddit Classification

Project 3

Football v Football

r/NFL



r/soccer



30

40

50

40

30

The Problem

- Do you want to know which football people are talking about?
- Do you want to know if you're in America or Europe?
- This project involves the creation of multiple models that classify reddit posts as belonging to either r/NFL or r/soccer
- Using data taken from the first two weeks of October 2022

30

40

50

40

30

Subreddit Comparison

r/NFL

- Subreddit for following the NFL specifically
- 3 million members
- Majority of posts are links to Twitter



r/soccer

- Subreddit for all soccer leagues
- 3.5 million members
- More match threads and highlights

30

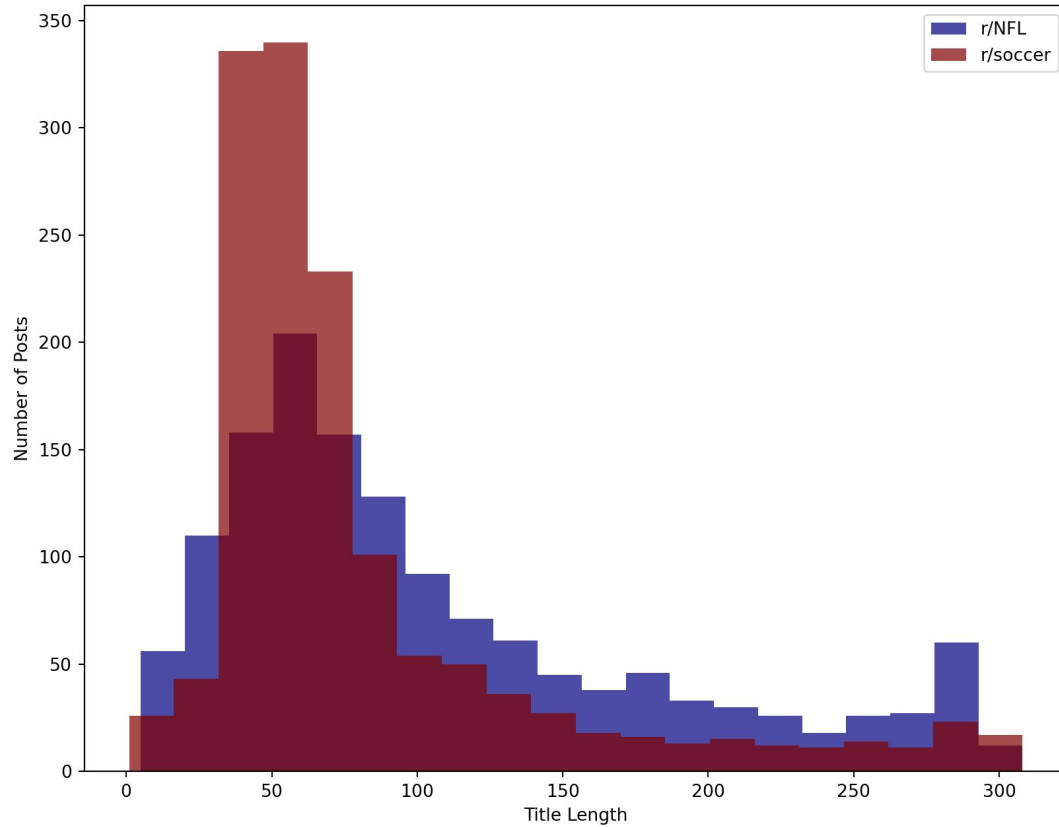
40

50

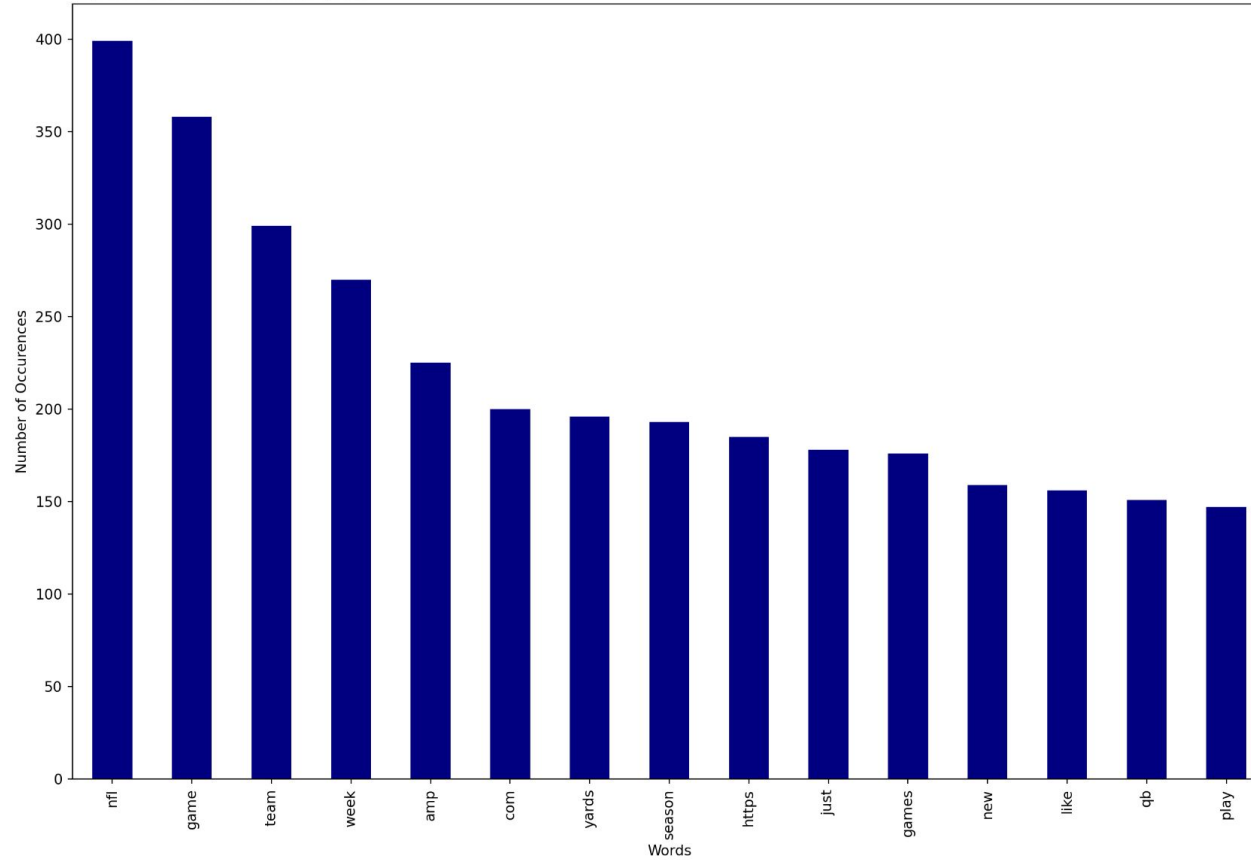
40

30

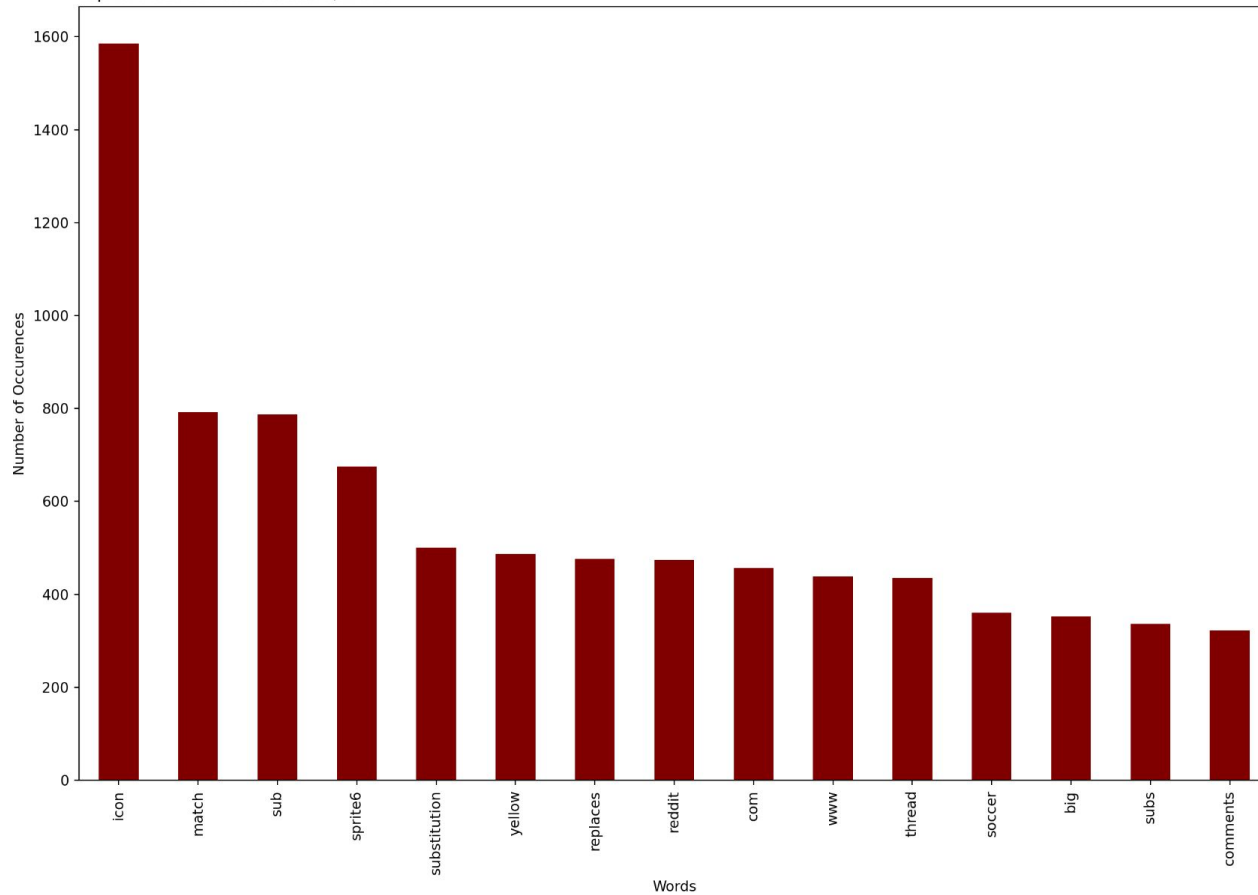
Distribution of Title Length for r/NFL v r/soccer



Top 15 Most Common Words r/NFL



Top 15 Most Common Words r/soccer



Icon?

- The most common word by far is 'icon'
- There are other confusing words like 'sprite6'



From Sports Illustrated

30

40


50

40


30


The Culprit Revealed


MATCH EVENTS | via ESPN^[2]


45'  Substitution, PSV Eindhoven. Luuk de Jong replaces Anwar El Ghazi.

55'  Goal! PSV Eindhoven 1, Arsenal 0. Joey Veerman (PSV Eindhoven) left footed shot from the centre of the box to the top right corner. Assisted by Luuk de Jong.

56'  Substitution, Arsenal. Thomas Partey replaces Albert Sambi Lokonga.

57'  Substitution, Arsenal. Bukayo Saka replaces Martin Ødegaard.

58'  Kieran Tierney (Arsenal) is shown the yellow card for a bad foul.

63'  Goal! PSV Eindhoven 2, Arsenal 0. Luuk de Jong (PSV Eindhoven) header from the centre of the box to the bottom right corner. Assisted by Cody Gakpo with a cross following a corner.

From reddit.com/r/soccer

30

40

50

40

30

Random Forest

- A simple model that still produces good results
- The random forest should have lower variance than other tree based models
- This model is 93.2% accurate on test data

30

40

50

40

30

Support Vector Machine

- Good fast model for generating predictions
- Works well with NLP data
- This model's accuracy is 91,7%

30

40

50

40

30

05

Conclusion

HOME

12:34

VISITOR

20

15

DOWN

TO GO

BALL ON

QTR

01

00

21

04

- Both models vastly outperform the baseline of 50/50
- Overall, the random forest model works the best

0

40

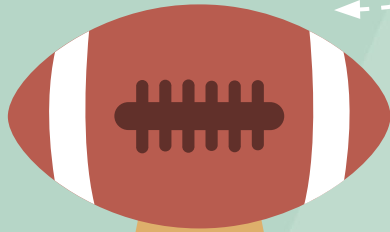
50

40

30

THANKS!

Do you have any questions?



CREDITS: This presentation template was created by
Slidesgo, including icons by Flaticon, and
infographics & images by Freepik
Please keep this slide for attribution

30

40

50

40

30