

Capstone Project

Identifying Offensive Tweets

Joel Cheung DSI-28



TABLE OF CONTENTS

1

INTRODUCTION

2

PROBLEM
STATEMENT

3

DATA SOURCES

4

DATA CLEANING

TABLE OF CONTENTS

5

EXPLORATORY
DATA ANALYSIS

6

MODELLING

7

MODELLING
RESULTS

8

LIMITATIONS &
FUTURE STEPS

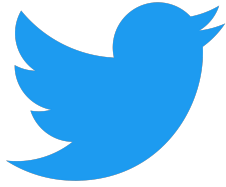
9

CONCLUSION



1

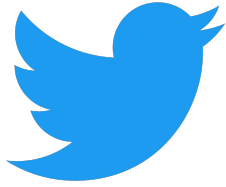
INTRODUCTION



TWITTER

- Is a micro-blogging social media site
- 217.5 million active users globally
- 500 million tweets everyday

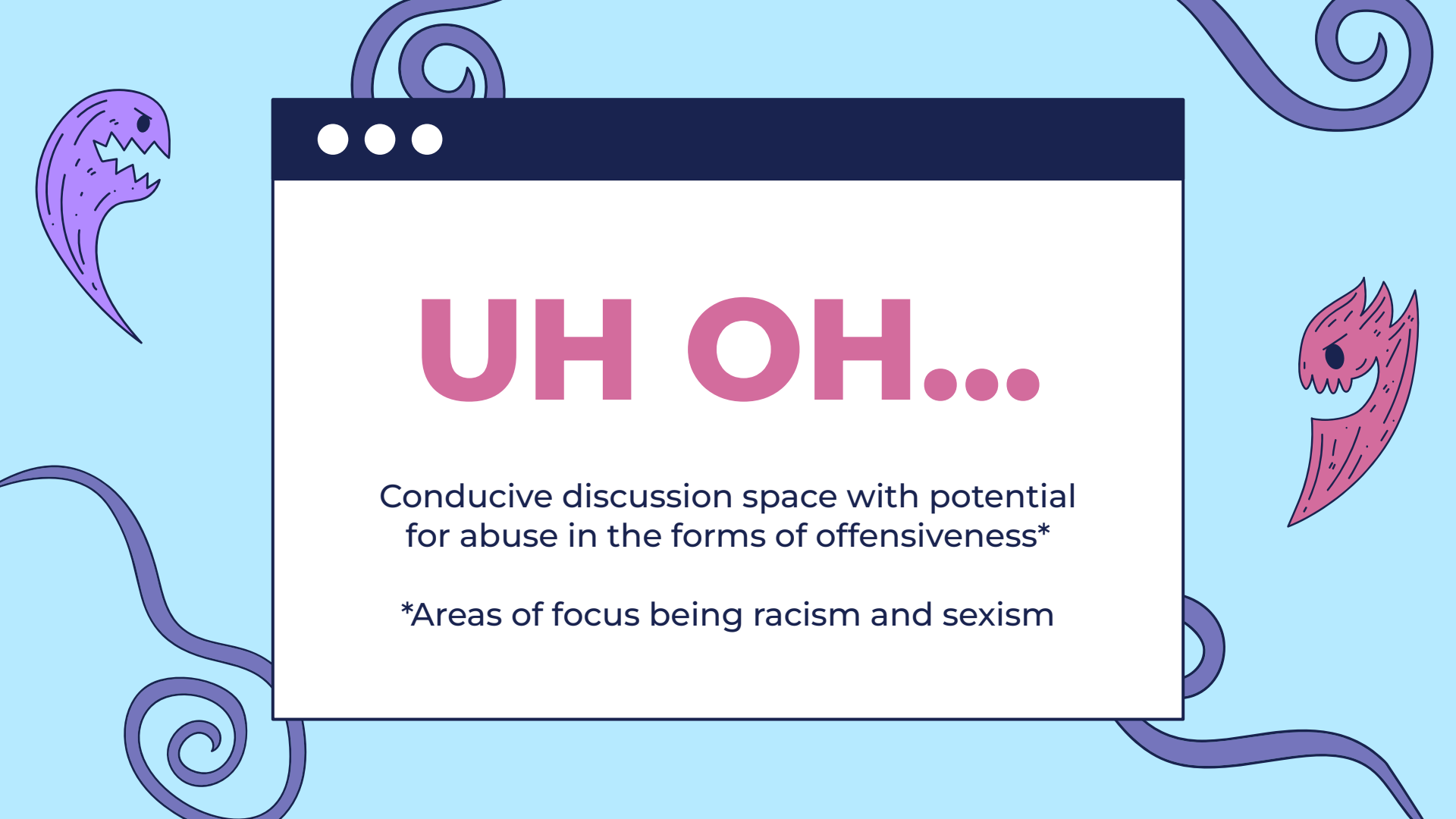




TWITTER

- Corporations - brand outreach
- Individuals - entertainment, news, discussions
- Freedom to discuss a wide range of topics (hashtags)



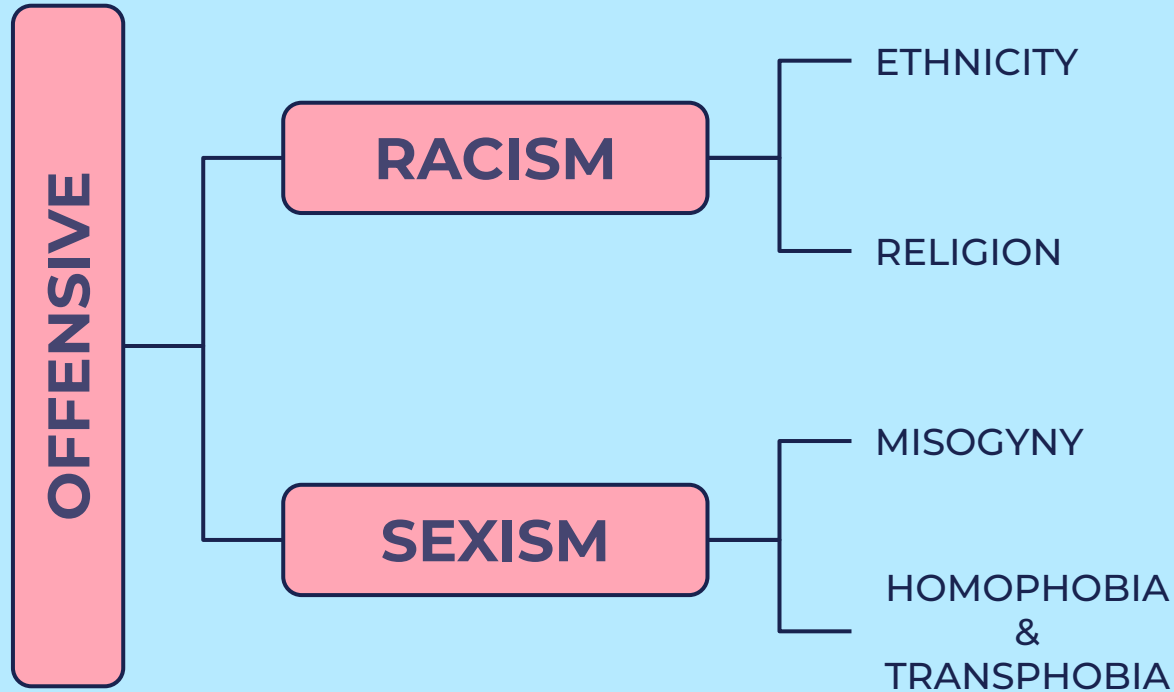


UH OH...

Conducive discussion space with potential
for abuse in the forms of offensiveness*

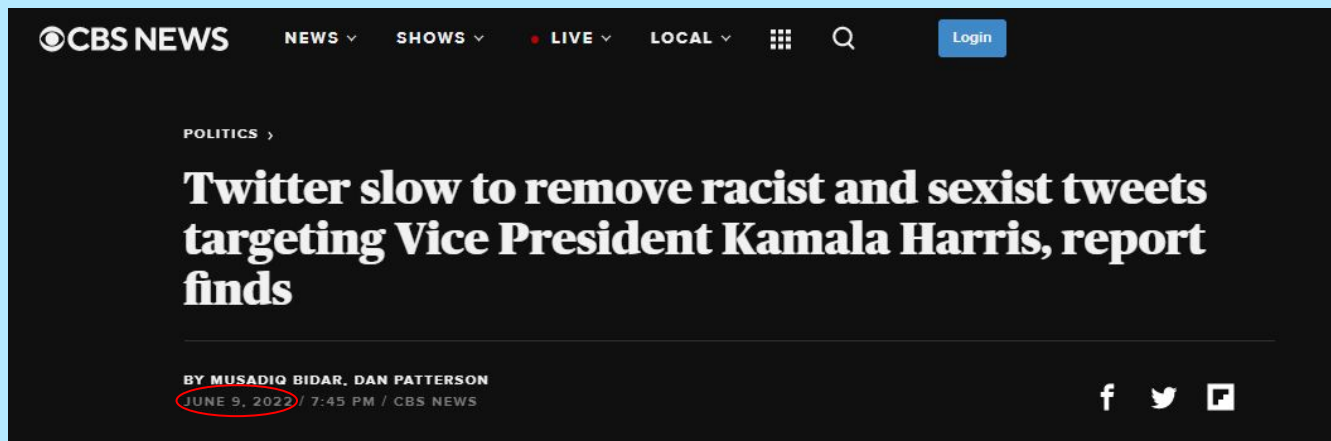
*Areas of focus being racism and sexism

RACISM AND SEXISM





TWITTER'S POLICY

Twitter is a social broadcast network that enables people and organizations to publicly share brief messages instantly around the world. This brings a variety of people with different voices, ideas, and perspectives. People are allowed to post content, including potentially inflammatory content, as long as they're not violating the [Twitter Rules](#). It's important to know that Twitter **does not screen content or remove potentially offensive content**.



The screenshot shows the CBS News website interface. At the top, the CBS News logo is on the left, and navigation links for NEWS, SHOWS, LIVE, and LOCAL are in the center. A search icon and a Login button are on the right. Below the navigation bar, the article is categorized under POLITICS. The main headline reads: "Twitter slow to remove racist and sexist tweets targeting Vice President Kamala Harris, report finds". Below the headline, the byline states "BY MUSADIQ BIDAR, DAN PATTERSON". The date and time are listed as "JUNE 9, 2022 / 7:45 PM / CBS NEWS". At the bottom right of the article preview, there are social media sharing icons for Facebook, Twitter, and a generic share icon.




CBS NEWS NEWS ▾ SHOWS ▾ LIVE ▾ LOCAL ▾   [Login](#)

POLITICS ▸

Twitter slow to remove racist and sexist tweets targeting Vice President Kamala Harris, report finds

BY MUSADIQ BIDAR, DAN PATTERSON

JUNE 9, 2022 / 7:45 PM / CBS NEWS

2

PROBLEM STATEMENT



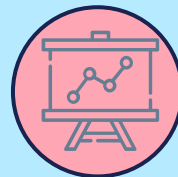


PROBLEM STATEMENT



BALANCING ACT

Maintaining a growing user base to attract corporations, while keeping a safe space for users

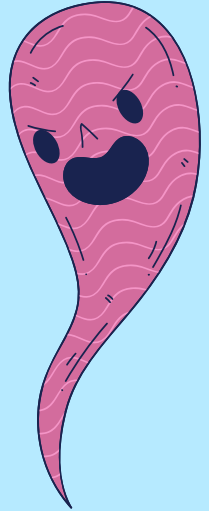
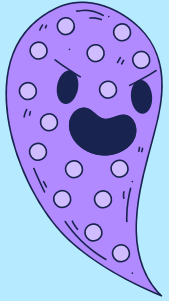


IDENTIFY OFFENSIVE TWEETS

A model to flag offensive tweets to point users to resources before they publish the tweet

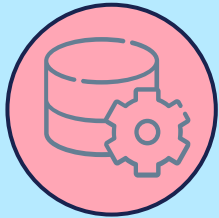
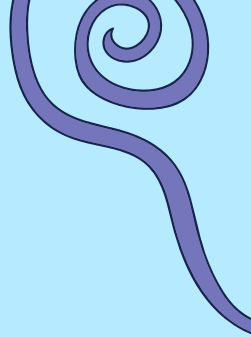
3

DATA SOURCES



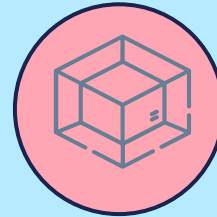


DATA SOURCES



TRAINING DATASET

kaggle:
Classified Tweets



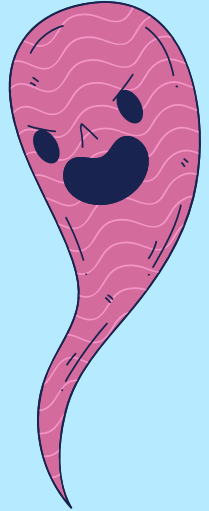
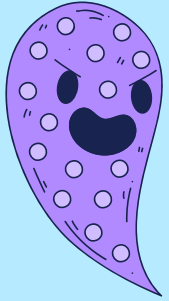
TEST DATASET

kaggle:
Cyberbullying
Classification challenge



4

DATA CLEANING



DATA CLEANING



- Remove emojis 😊
- Remove mentions (@username) and URLs
- Make texts lowercase
- Removing punctuations
- Removing stopwords
- Lemmatizing with Part-of-Speech (POS) tagging





5

EXPLORATORY DATA ANALYSIS

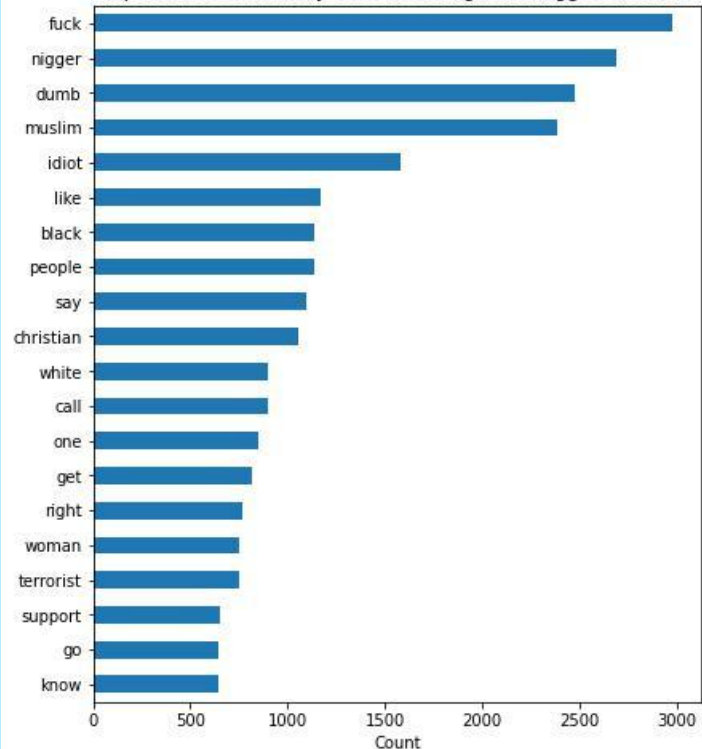


CONTENT WARNING!!

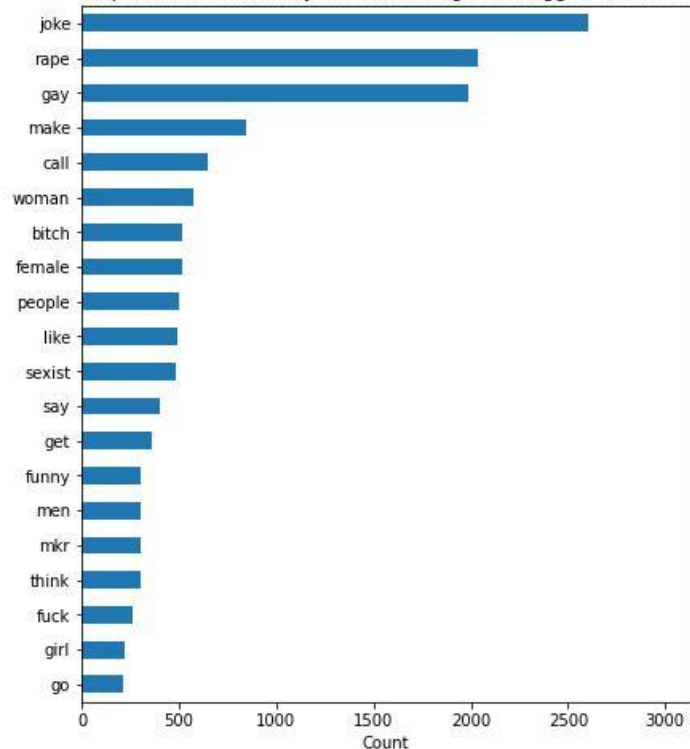
From this slide on, there will be mentions of sensitive topics such as religion, misogyny and other topics that may be offensive to many

UNIGRAM

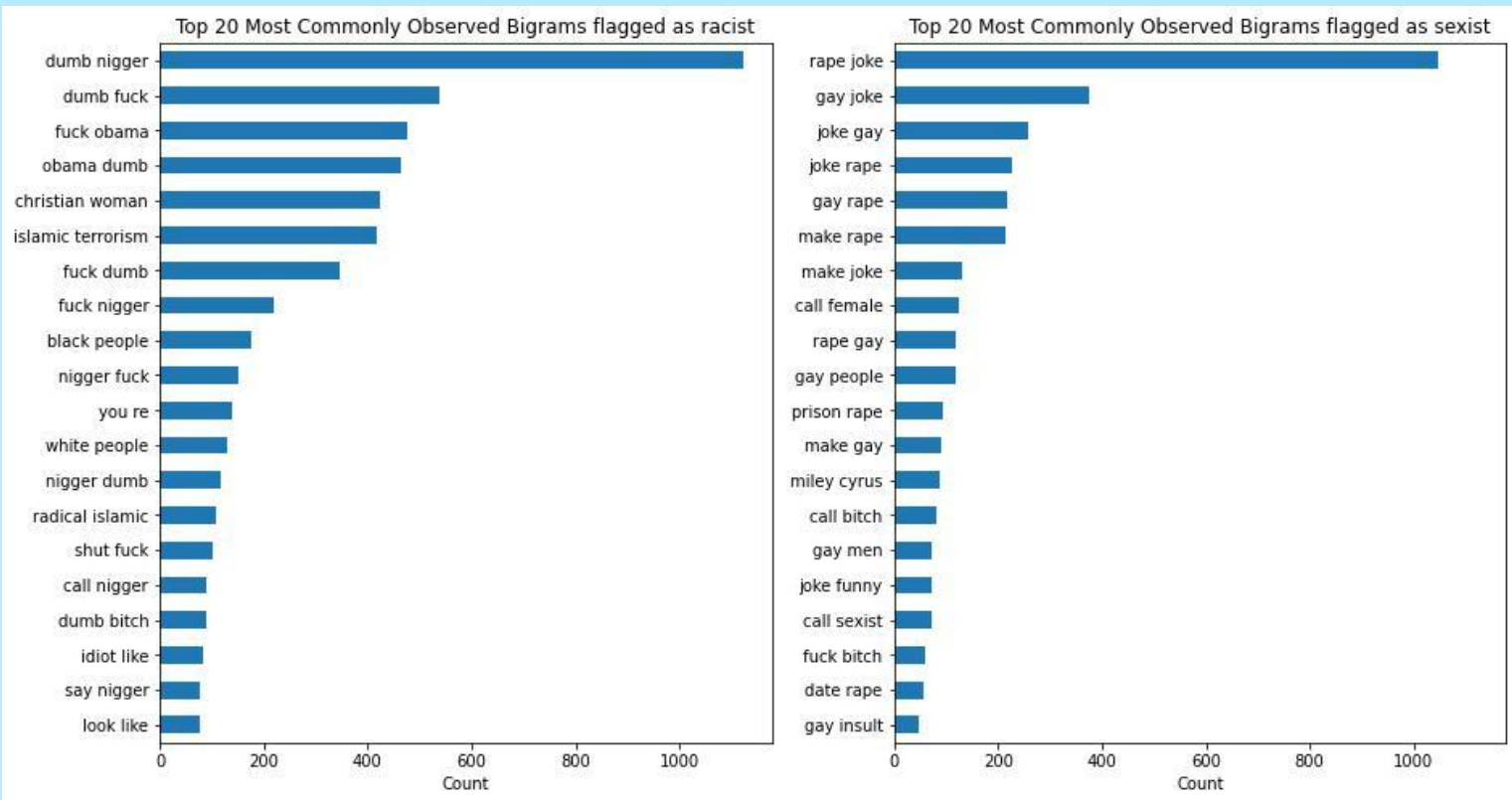
Top 20 Most Commonly Observed Unigrams flagged as racist



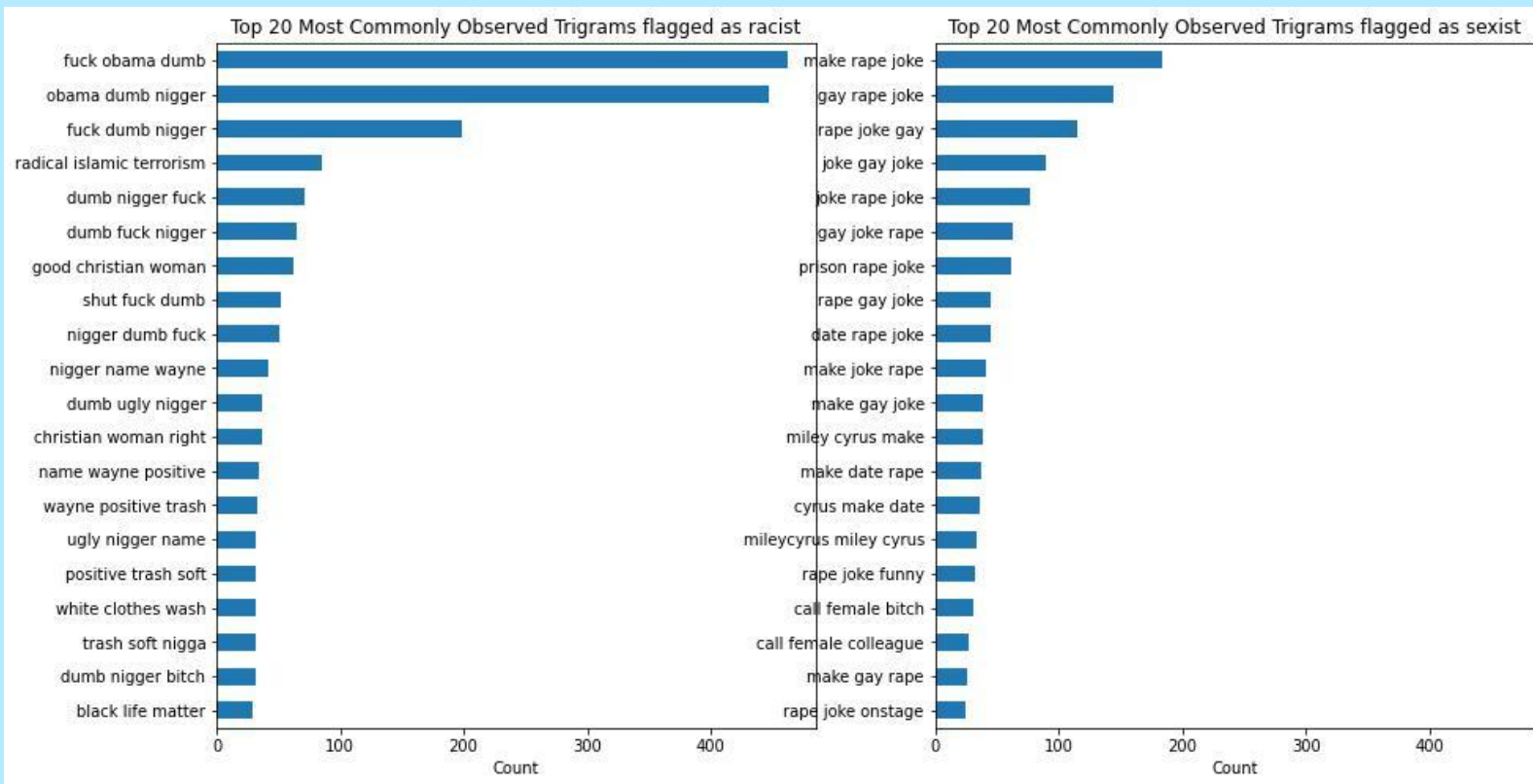
Top 20 Most Commonly Observed Unigrams flagged as sexist



BIGRAM



TRIGRAM



6

MODELLING



PERFORMANCE METRIC: F1-SCORE



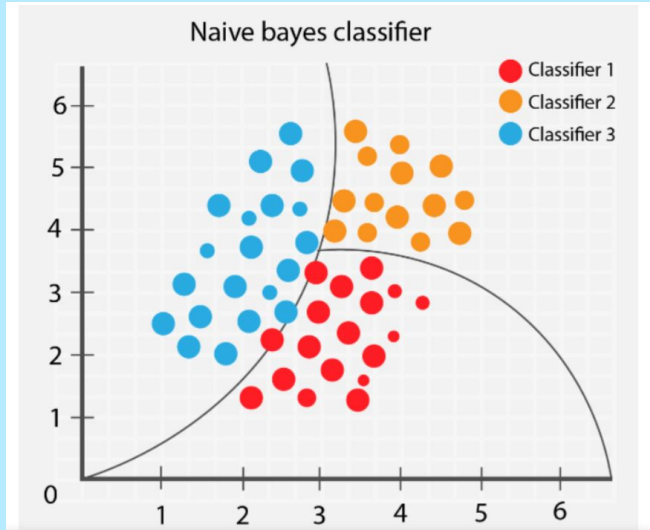
$$F_1\text{-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Precision: Of all positive predictions, how many are really positive?

Recall: Of all real positive cases, how many are predicted positive?



MULTINOMIAL NAIVE BAYES

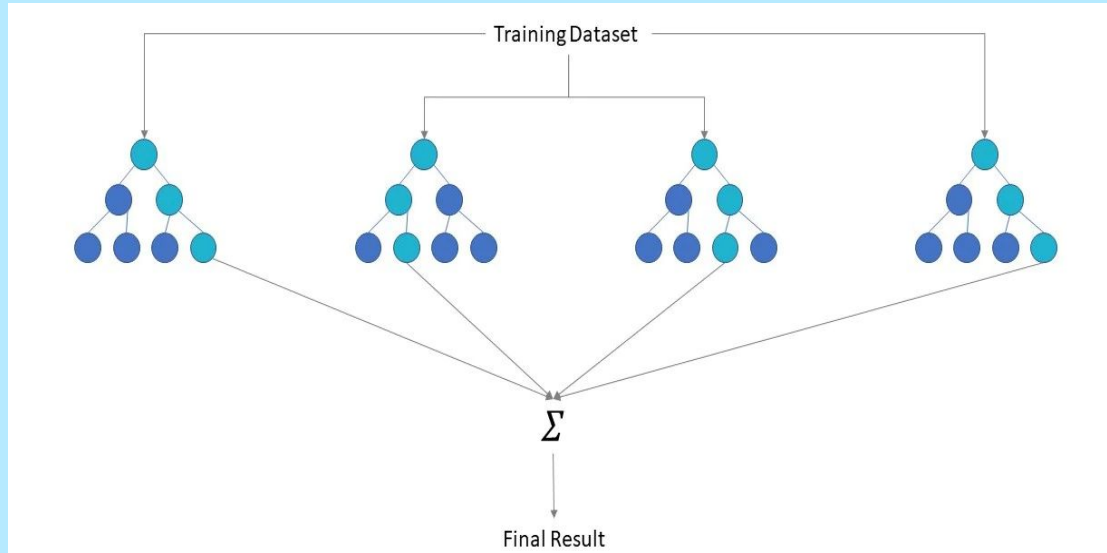


- Baseline Model
- Simple, easy to train model

$$\Pr(\text{Racist}|\text{text}) = \frac{\Pr(\text{text}|\text{Racist}) * \Pr(\text{Racist})}{\Pr(\text{text})}$$



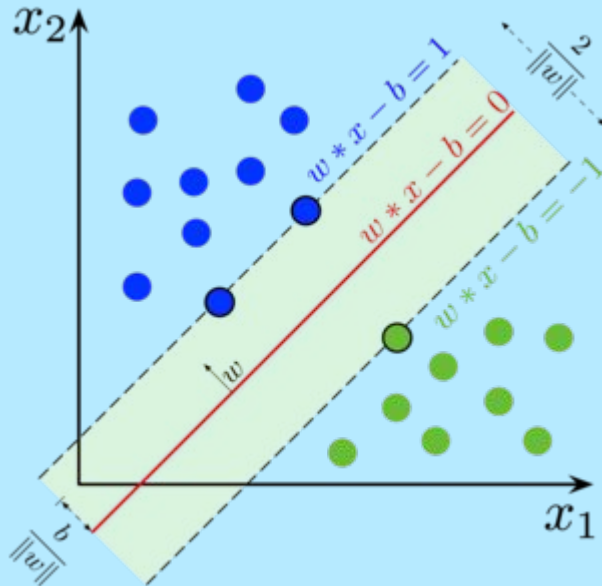
RANDOM FOREST



Algorithm trains many decision trees and returns the classification of the majority vote



SUPPORT VECTOR MACHINE



Using some representative points (support vectors), the algorithm aims to maximize the space between each category (margin)

BERT

- Bidirectional Encoder Representations from Transformers (BERT)
- Pre-trained model created by Google in 2018
- Trained on the entire English Wikipedia and BookCorpus by having it fill in masked words (blanks)
 - Total of 3.3 billion words
 - Useful for contextual word embeddings
 - Different representation for different contexts (despite being same word e.g. dog's bark vs tree bark)
- Takes into account words before and after a particular word
 - Truly bidirectional compared to LSTM



7

MODELLING RESULTS

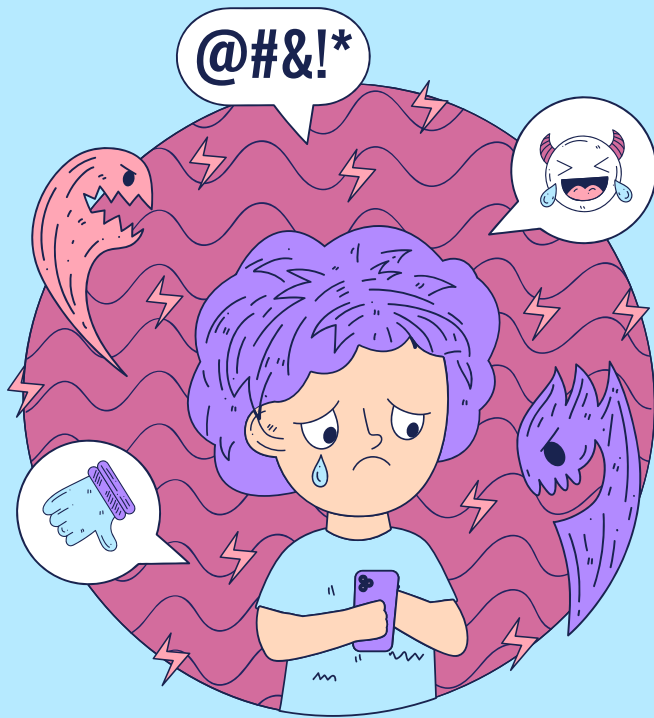


MODELLING RESULTS



Model	F1-Score		
	Train	Validation	Test
Multinomial NB	0.922	0.889	0.640
Random Forest	0.999	0.931	0.883
SVM	0.963	0.927	0.882
BERT	0.999	0.927	0.856





8

LIMITATIONS & FUTURE STEPS



LIMITATIONS



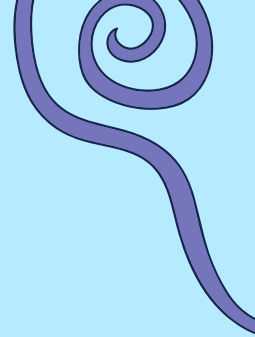
Dataset issues

Training dataset might not be comprehensive enough in terms of region and timeframe; no data dictionary available



Computing Constraints

Long training time, may not achieve true optimal hyperparameters





FUTURE STEPS



More data

More tweets from all over the globe in order to generalize better



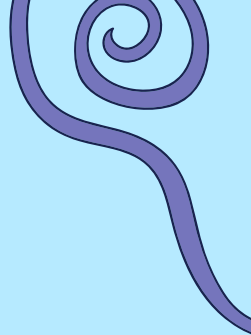
More types of data

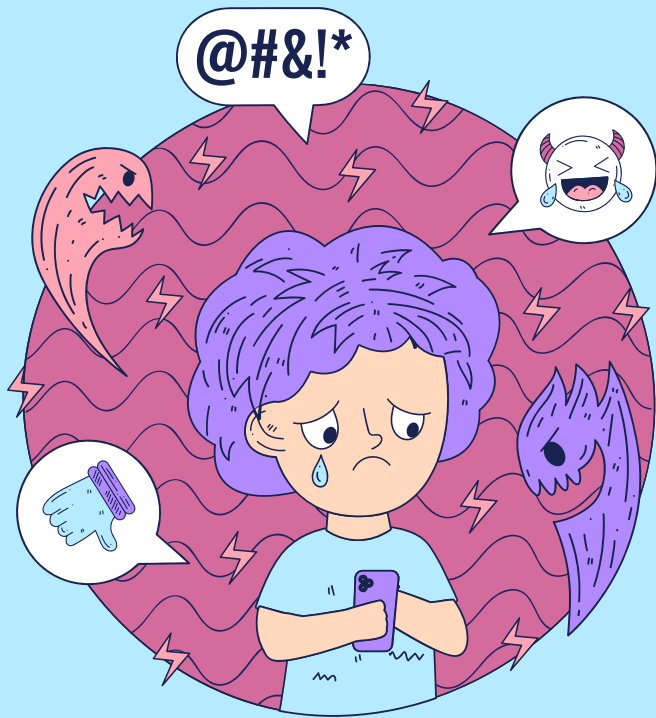
User information



Acquire better computing equipment

Allows for more granular GridSearch

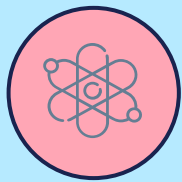




9

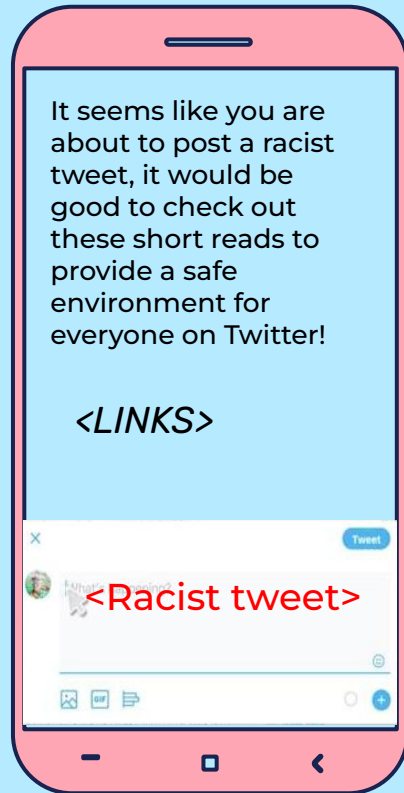
CONCLUSION

CONCLUSION



Model

Multinomial Naive
Bayes for **speed**,
Random Forest
for **better
predictions**





THANKS!

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

Please keep this slide for attribution