

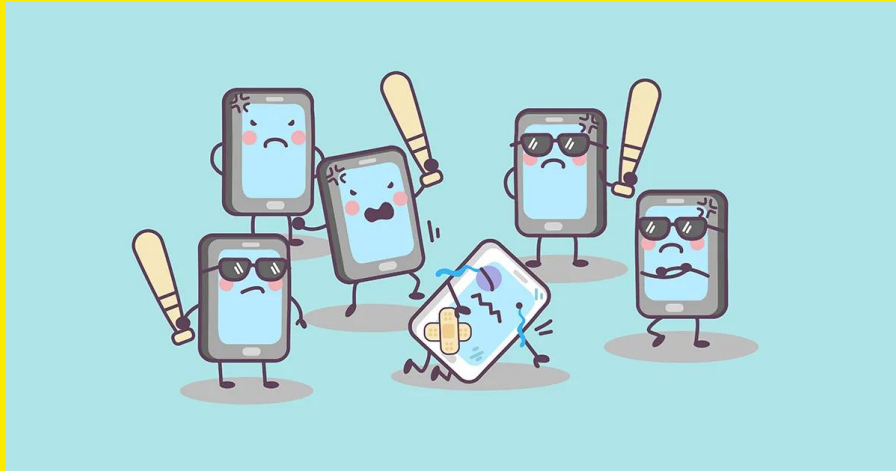
# **DABPT01**

# **Capstone**

# **Assignment**

**Final Presentation**

# THE PROBLEM



## Background

The online space can be a dangerous place to be, especially since anyone and everyone can partake in it

Toxic online comments are a result of online anonymity, and can be dangerous for unsuitable audiences if left unchecked for too long

## Project Objective

By way of a natural language processing (NLP) model, ascertain if any given online comment, could be classified as a 'problematic' / 'negative' one, or left alone without consequence

## Applications

Moderation of any website / application with:

- Free-text comment section
- Message boards
- Chatbots

# 6 FLAGS

- toxic
- severe\_toxic
- obscene
- threat
- insult
- identity\_hate

# 159,571

Total no. of social media comments analysed

# 10.16%

Proportion of comments with ≥1 toxicity flags

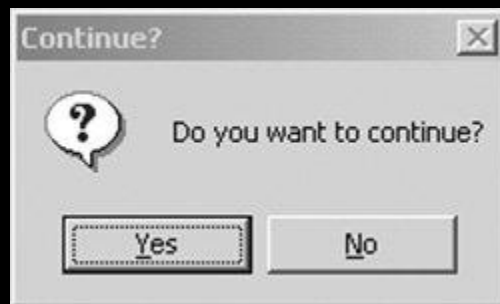
# DISCLAIMER!!!

Presentation material beyond this point, while constructed strictly for academic purposes, may be considered offensive, controversial, or triggering to some viewers.

This includes, but is not limited to, themes of:

- Violence
- Discrimination
- Sexually suggestive / explicit content
- Strong language

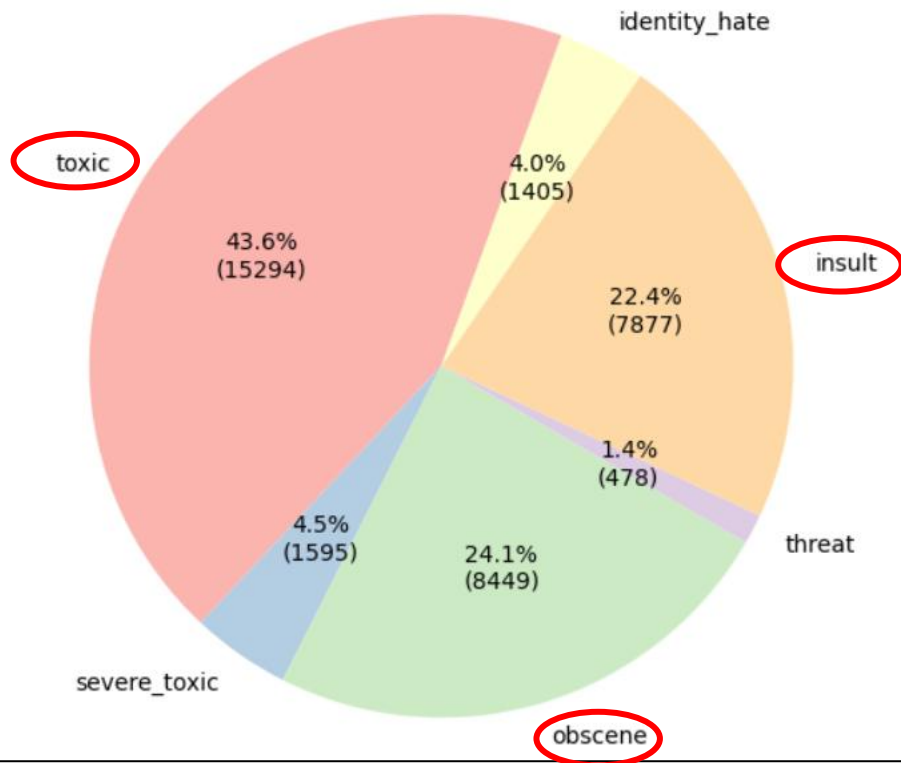
Viewer discretion is advised. If you are sensitive to such topics, you may wish to consider this before proceeding.



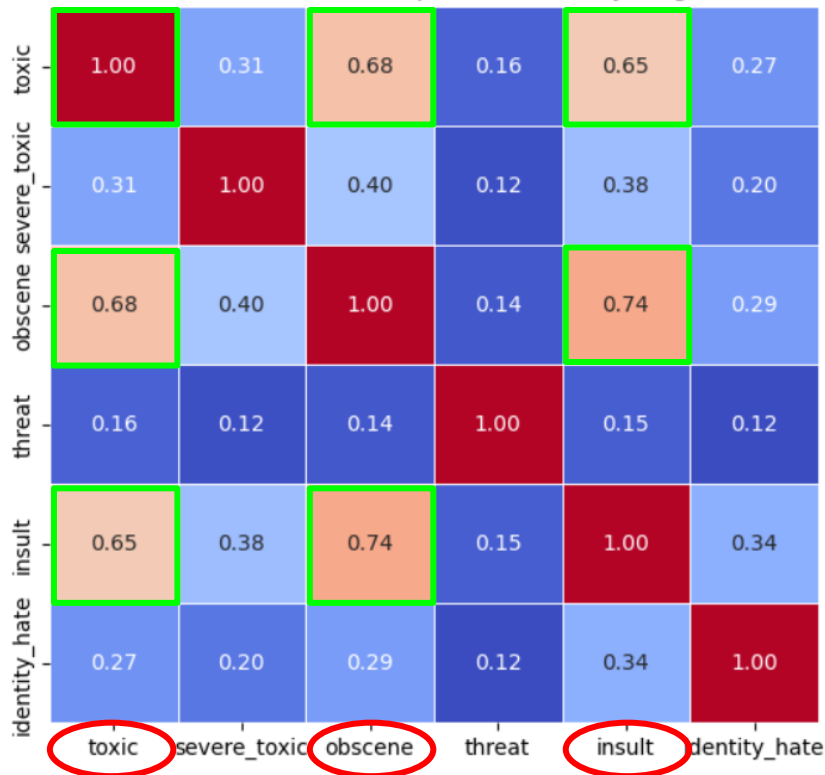
[illegible][illegible]

# EXPLORATORY DATA ANALYSIS (EDA)

Distribution of Toxicity Flags across Flagged Comments (Count)



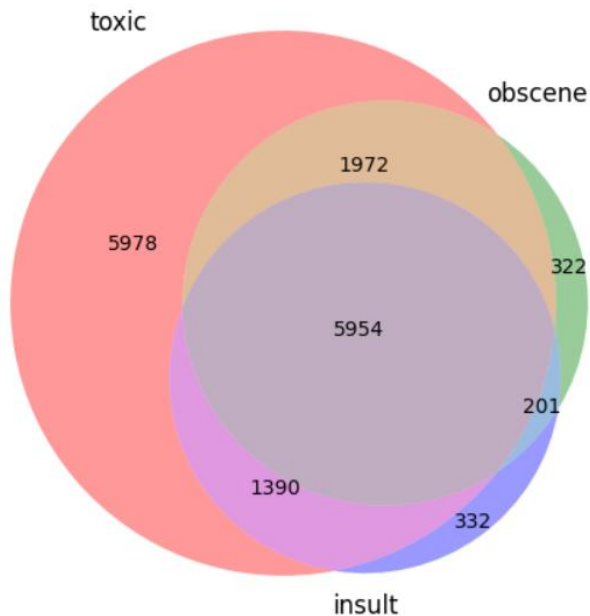
Correlation Heatmap across Toxicity Flags



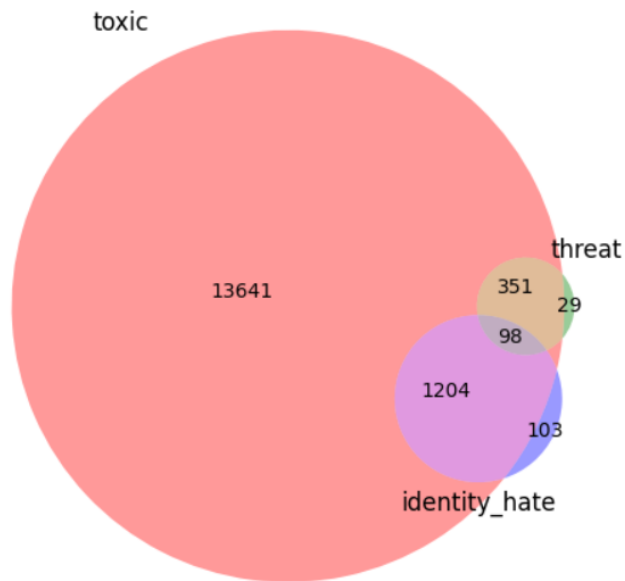


# HIGH DEGREE OF OVERLAP BETWEEN FLAGS?

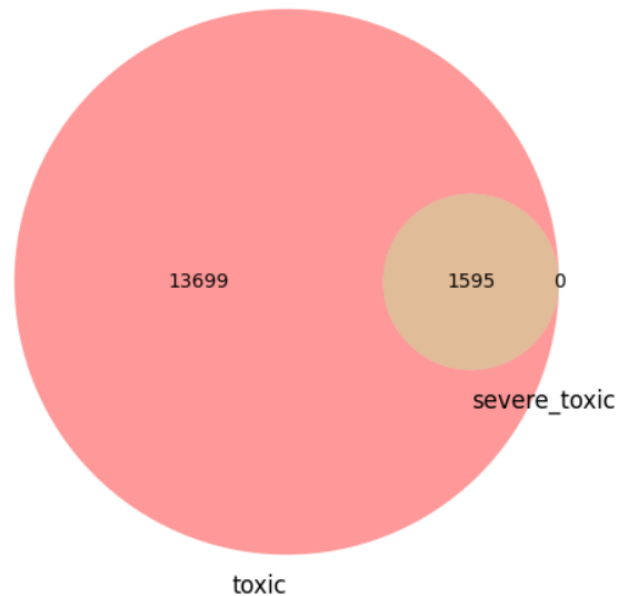
Venn diagram for 'toxic', 'obscene' and 'insult'



Venn diagram for 'toxic', 'threat' and 'identity\_hate'



Venn diagram for 'toxic' and 'severe\_toxic'



# Urgency of corrective action?

## Value of multi-label data?

- In real-world context, swift corrective action on problematic social media comments would be of higher importance
- Toxicity labels can be condensed into one single flag (eg. 'flagged') and analysed on binary basis
  - Problematic ('flagged' = 1)
  - Innocent ('flagged' = 0)



# 5,461,520

Total no. of tokens generated from 159,571 comments

**TRUE +VE**

*"i am going to pee on you"*

**TRUE -VE**

*"yep it was cut around \_ seconds however it has been released uncut now"*

**FALSE +VE**

*"burial where did he die where is he buried"*

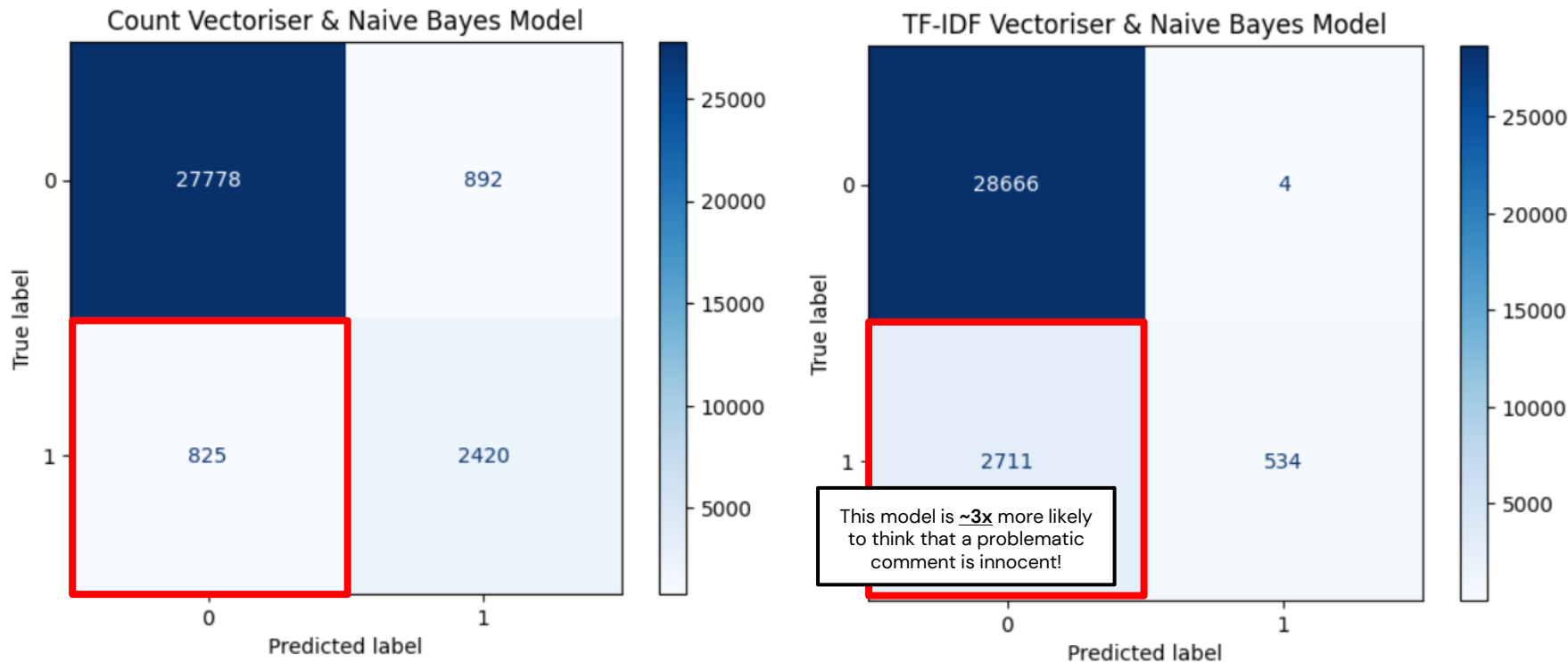
*"would care less good luck have fun enjoy your life"*

**FALSE -VE**

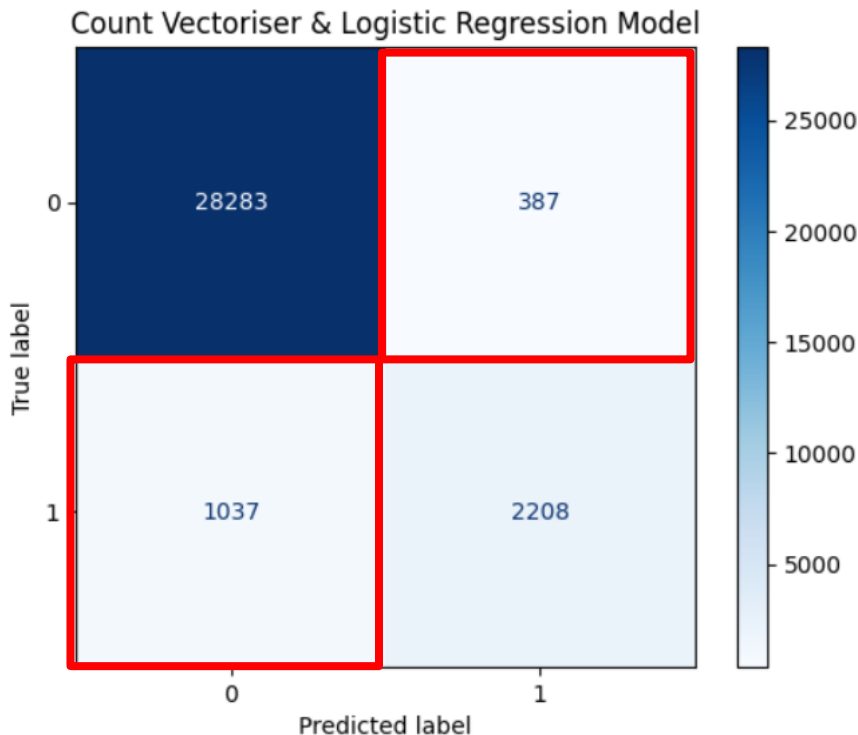
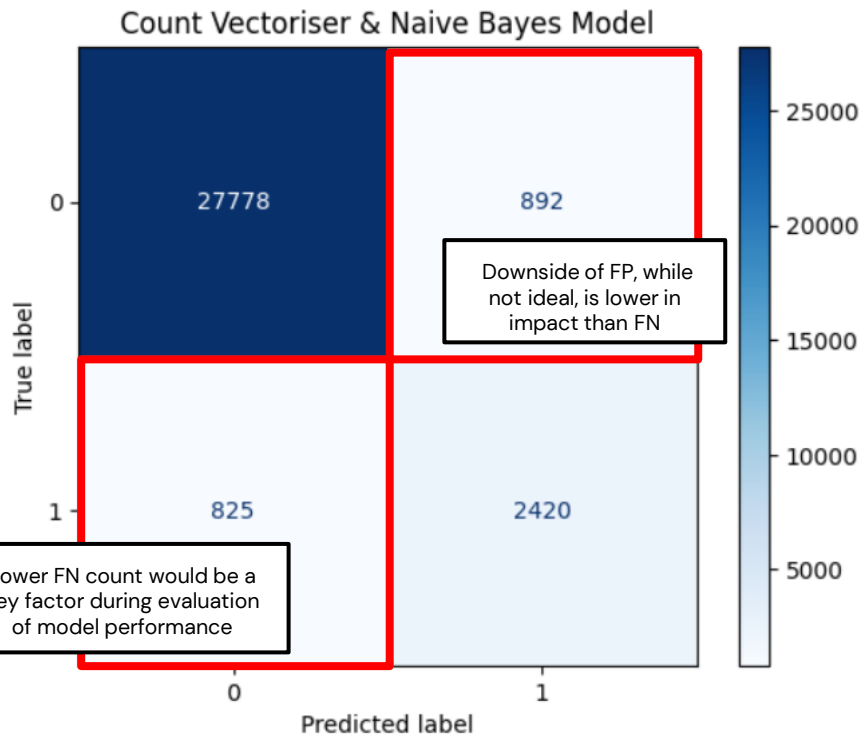
*"i am having my period"*

*"do it and i will cut you"*

# WHY IS FALSE -VE SIGNIFICANT HERE?



# BEST PERFORMING VECTOR-MODEL COMBOS



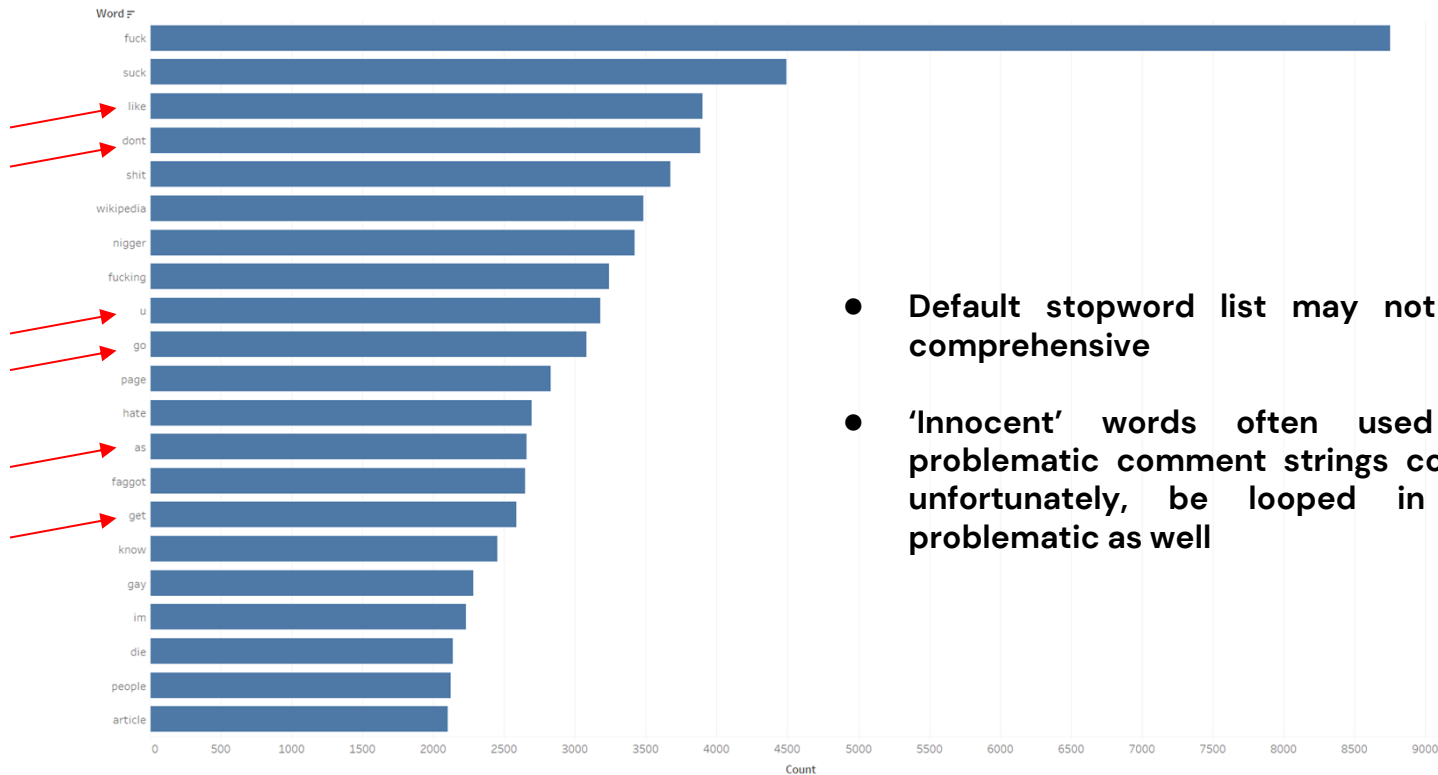
# VECTORISER-MODEL PERFORMANCE

Vectoriser	Model	Accuracy	Precision	Recall	F1-Score
Count	Dummy (Baseline)	82.07%	11.60%	11.53%	11.56%
Count	Logistic Regression	95.54%	85.09%	68.04%	<b>75.62%</b>
Count	Naive Bayes	94.62%	73.07%	<b>74.58%</b>	73.81%
TF-IDF	Dummy (Baseline)	81.54%	9.77%	9.89%	9.83%
TF-IDF	Logistic Regression	<b>95.60%</b>	92.01%	62.10%	74.15%
TF-IDF	Naive Bayes	91.49%	<b>99.26%</b>	16.46%	28.23%

Metrics above are measured against the vectoriser-model correctly labelling a ('flagged' = 1) comment

# ANALYSIS LIMITATIONS

Most frequent words in flagged comments



- Default stopwords list may not be comprehensive
- 'Innocent' words often used in problematic comment strings could, unfortunately, be looped in as problematic as well

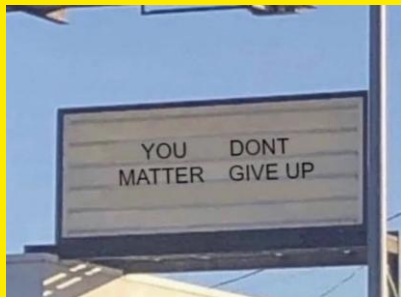
# ANALYSIS LIMITATIONS

## Absence of Context



Each word is analysed as is

## Unigram Analysis



Word order is ignored

## Creative Spelling

7H15 M3554G3  
53RV35 7O PR0V3  
H0W 0UR M1ND5 C4N  
D0 4M4Z1NG 7H1NG5!  
1MPR3551V3 7H1NG5!  
1N 7H3 B3G1NN1NG  
17 W45 H4RD BU7  
N0W, 0N 7H15 L1N3  
YOUR M1ND 1S  
R34D1NG 17  
4U70M471C4LLY  
W17H 0U7 3V3N  
7H1NK1NG 4B0U7 17,  
B3 PROUD! ONLY  
C3R741N P30PL3 C4N  
R3AD 7H15.

Social media users, while trying to circumvent rules and regulations, could still spell offensive words with similar alphanumeric characters

## Boundless Nature of Language



If target word was not fed into machine learning model during training, prediction will fail during test phase



# Some areas to look into as part of future work...

- Customised / updated stopwords library for text cleaning
- Look into more sophisticated models involving multi-label classification, where more information can be gleaned from predictions (eg. OneVsRestClassifier)
- Look into  $n$ -gram analysis (bi-gram, tri-gram)

**Digitalisation breeds anonymity.  
Anonymity emboldens recklessness.  
Recklessness begets suffering.**

**Be kind, even when faceless.**

**Thank You**