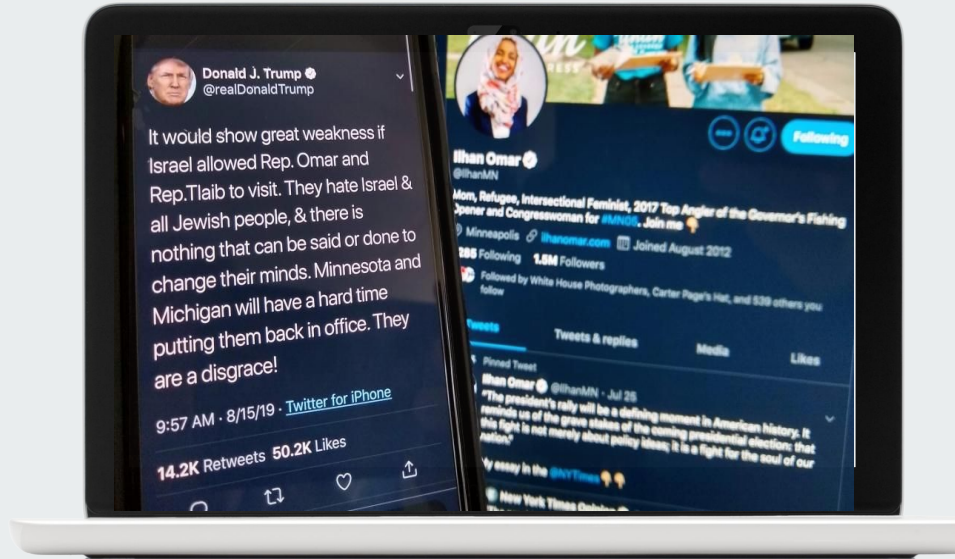# ``Share Love not Hate'' - Assessing Hate Speech Detection methods using TF-IDF and POS tagging approach on Twitter

Kshitija Hande
(1131778)

# Outline
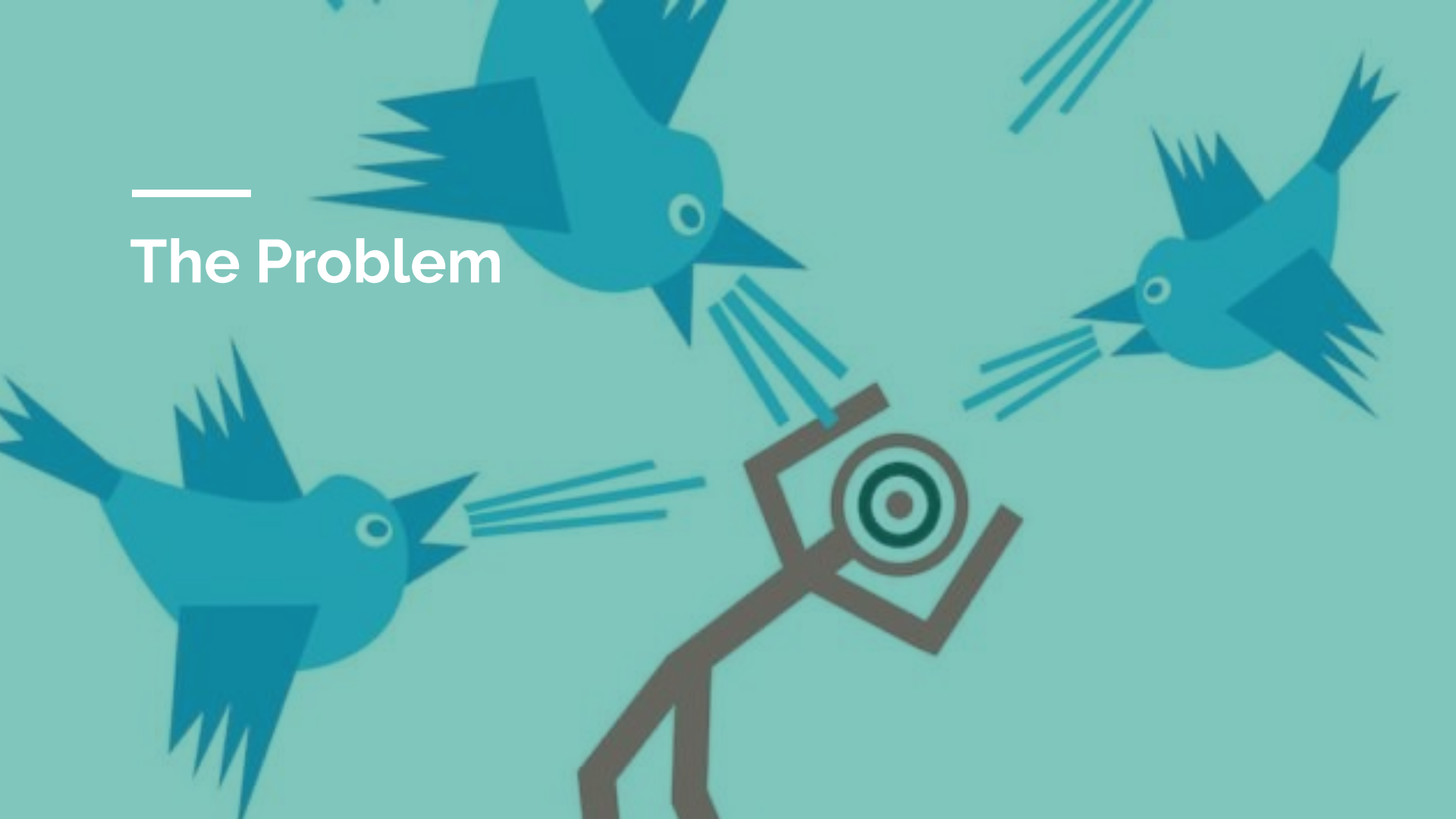
The Problem

# Why do we care?

Exponential increase in the use of the internet by people of different cultures, sexual orientation, ethnicities and educational backgrounds.

Differentiating hate speech and offensive language is a key challenge.

New York Times reported personal attacks motivated by bias and prejudice

# Problem statement

Classification of tweets on Twitter into three classes: hate speech, offensive language and neither. We perform experiments by leveraging TF-IDF and POS tagging features as input to machine learning models.

# Challenges

The style of social media especially hateful data is different

Uses abbreviations and spelling variations

Words cannot be taken literally.

# Problems with existing solutions

OOV words are impossible to train

Example : #BanPUBG

Issue with pre-trained embeddings

Splitting tweets to normalize input

Training on domain-specific data

# Methodology

# Data

01

| DATASET | COUNT |
|---|---|
| Train | 14,869 |
| Validation | 3,718 |
| Test | 6,196 |
| TOTAL | 24,783 |

# Data

02

| count |
| --- |
| hate_speech_count |
| offensive_language_count |
| neither_count |
| class |
| tweet |

VALUES IN THE CLASS COLUMN

0:"hate speech"

1:"offensive language"

2:"neither"

# Data

03

| Class labels | count |
|---|---|
| hateful | 840 |
| offensive | 11,571 |
| neither | 2,458 |
| Total | 14,869 |

TABLE I

CLASS DISTRIBUTION IN TRAINING SET

# **Data**

04

- Data Preprocessing :
  - Lowercase
  - Stemming
  - Lemmatization
- Data cleaning
  - Stopwords
  - URLs
  - Twitter Mentions
  - Retweet Symbols
  - Emojis

# Feature Extraction
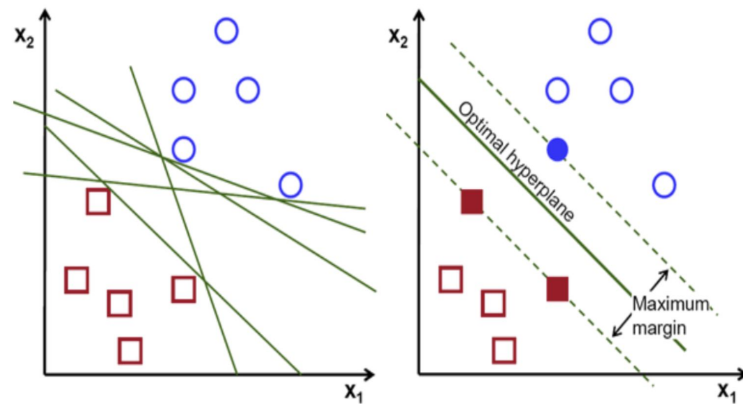
```
tfidf(d, t) =    tf(t) *  idf(d,t)
```

# Feature Selection - Logistic Regression

➔ L2 regularization

➔ Returns transformed version of input X

➔ This transformed input will be given as input to our classifier

# Model - Linear SVM

➔ Separate data points with a line across the hyperplane

➔ Goal is to maximize the margin and make a decision boundary

➔ One-versus-rest framework

➔ L2 regularization with 0.01 as regularization parameter

➔ Squared hinge loss function

# Results and Discussions

| Class labels | Precision | Recall | F1-score |
|---|---|---|---|
| hateful | 0.60 | 0.79 | 0.68 |
| offensive | 0.97 | 0.92 | 0.94 |
| neither | 0.81 | 0.93 | 0.87 |
| overall | | | **0.91** |

RESULTS FROM EXPERIMENT 2 USING TF-IDF AND POS TAGGING FOR FEATURE EXTRACTION

# Limitations

- Skewed dataset

- Hate comments without explicit toxic vocabulary

- Humorous content may be flagged due to use of controversial terms

- Overfitting can be addressed by using cross validation

# Conclusion

# What next?

➔ Build richer dictionary

➔ Experiment with deep neural network architectures

➔ Use larger dataset for training

# Thank You!

# Timeline

| MAY | JUN | TODAY | AUG | SEPT | OCT | NOV |
|-----|-----|-------|-----|------|-----|-----|
| Requirements gathering | User research | Wireframes | Review | Prototype | User testing | Dev hand-off |