

## **PROJECT : Hate-Speech-detection-using-Transformers-Deep-Learning**

**Group Name:** Hate Speech Detective

### **Members:**

**Name:** VARSHIT MANEPALLI

**Email:** varshitmanepalli1810@gmail.com

**Country:** INDIA

**College/Company:** Stevens Institute of Technology

**Specialization:** NLP

**Batch code :** LISUM32

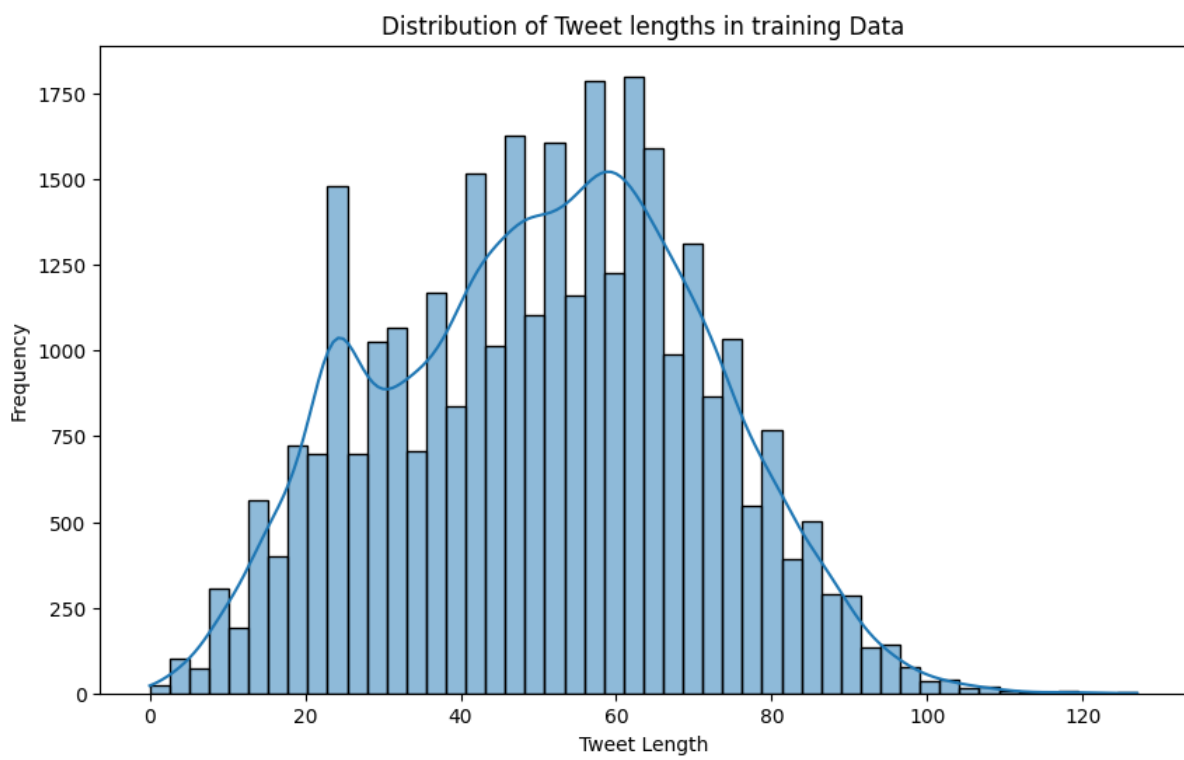
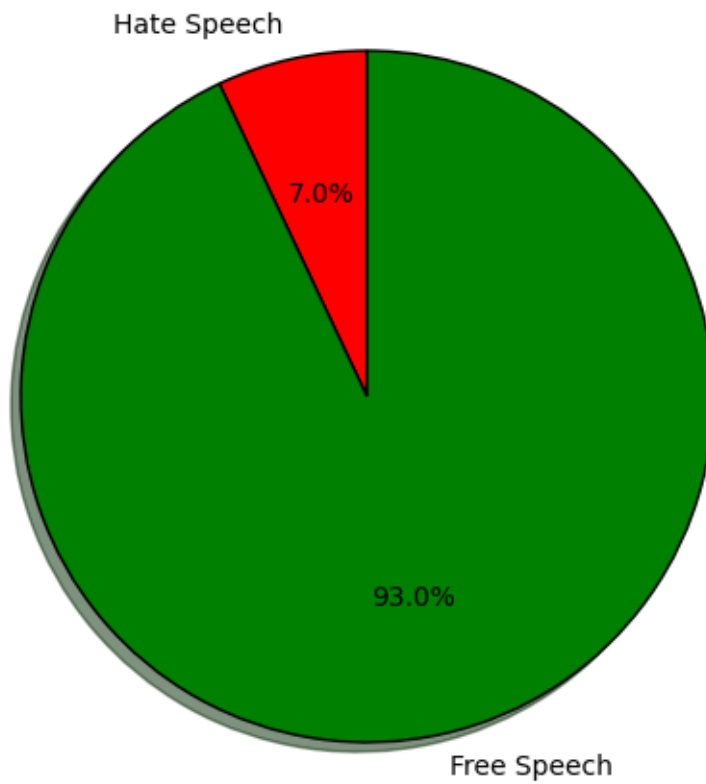
### **Problem Description**

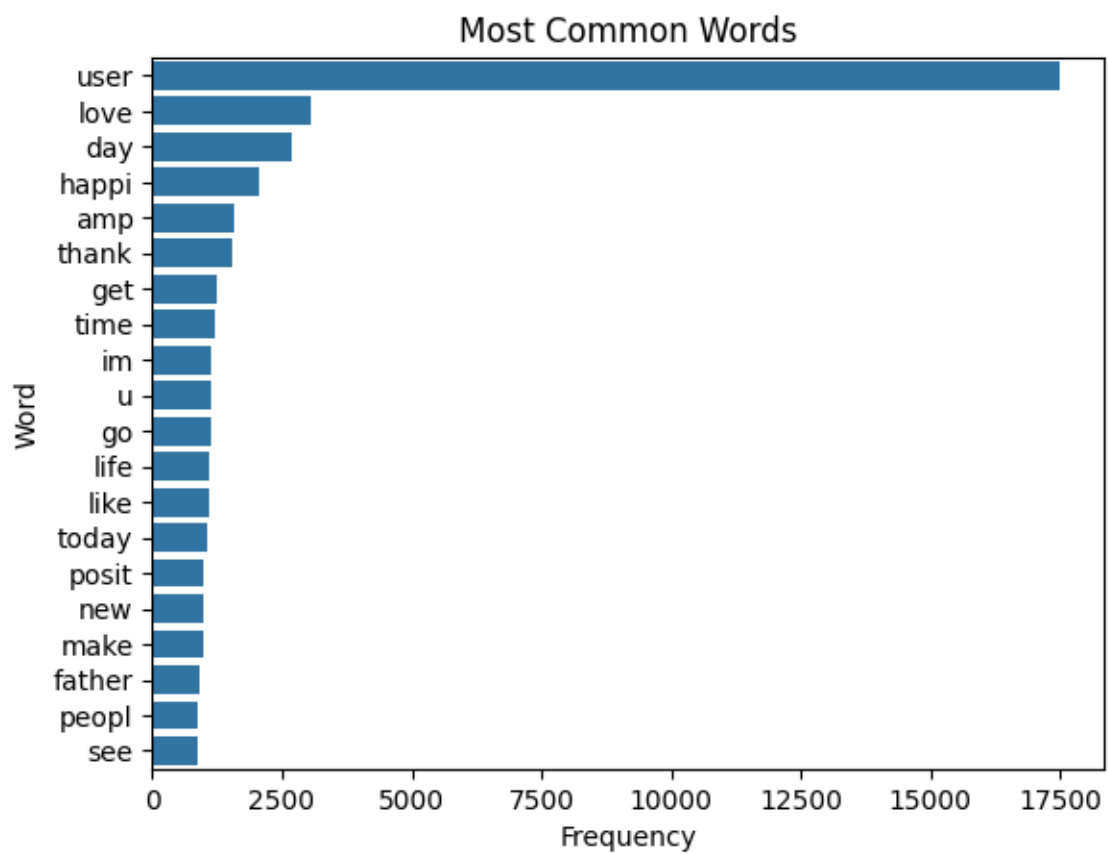
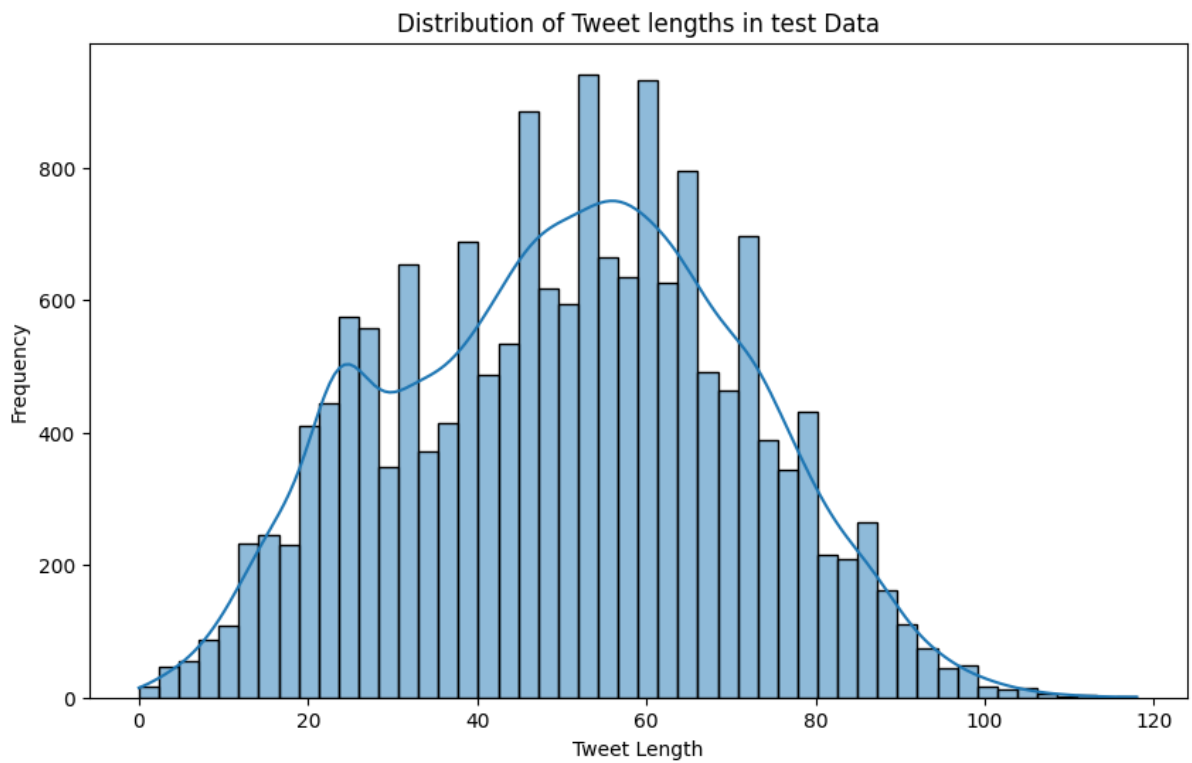
Hate speech detection aims to identify and classify statements that contain offensive, derogatory, or discriminatory language directed towards individuals or groups based on their identity factors such as religion, ethnicity, nationality, race, color, ancestry, sex, or other identity factors. This project involves developing a machine learning model to detect hate speech in Twitter tweets.

### **Business Understanding**

Hate speech can have serious consequences, including perpetuating discrimination, inciting violence, and causing psychological harm. Detecting hate speech on social media platforms like Twitter is crucial for maintaining a safe and inclusive online environment. By identifying and flagging hate speech, we can help prevent the spread of harmful content and protect vulnerable individuals and communities.

## EDA





## Word Cloud

### Word Cloud of Training Data



### Word Cloud of Test Data



## Class Distribution



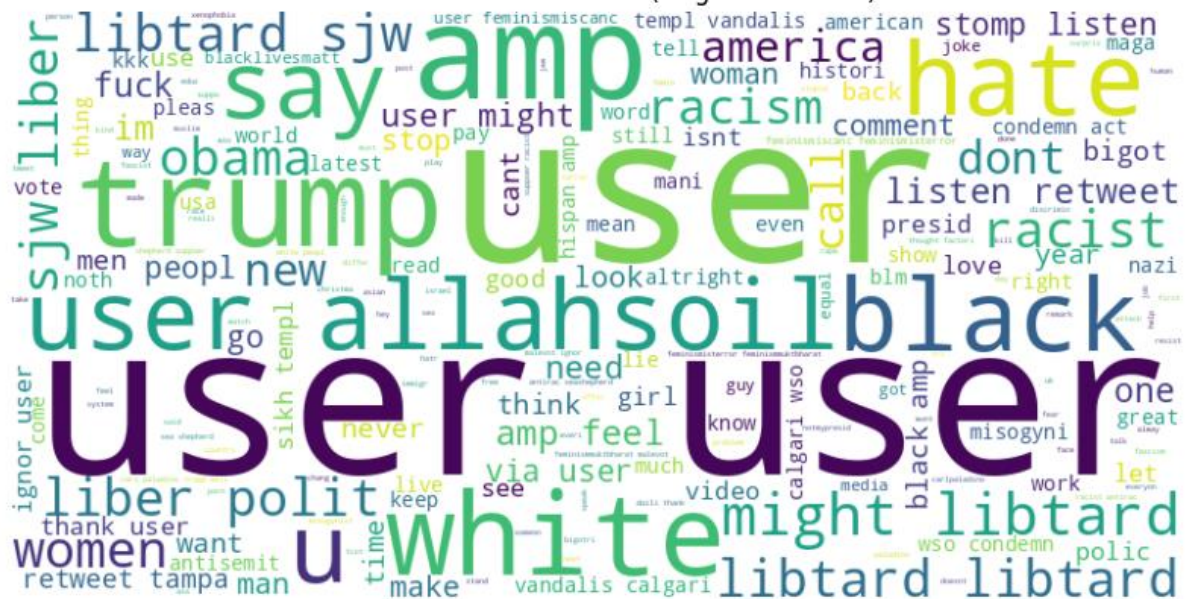


## Common Words By Class

### Common Words in class 0 (Non-Negative Tweets)



### Common Words in class 1 (Negative Tweets)



## Key Findings from EDA

- **Class Distribution** : The class distribution is imbalanced with significantly more non-hate speech tweets compared to hate speech tweets. This needs to be addressed in the model training phase.
- **Common Words** : Common words include generic terms like "user" indicating that further cleaning might be needed to remove such non-informative words.
- **Common Words** : Common words include generic terms like "user" and "thanks," indicating that further cleaning might be needed to remove such non-informative words.

## Recommendations

- We will use techniques such as oversampling (e.g., SMOTE) and under sampling to create balance in the dataset.
- We will do additional preprocessing steps to remove non-informative words and symbols to improve the quality of the features.
- We will be using the advanced models such as MLP classifier or Transformer based models such as BERT for this project.