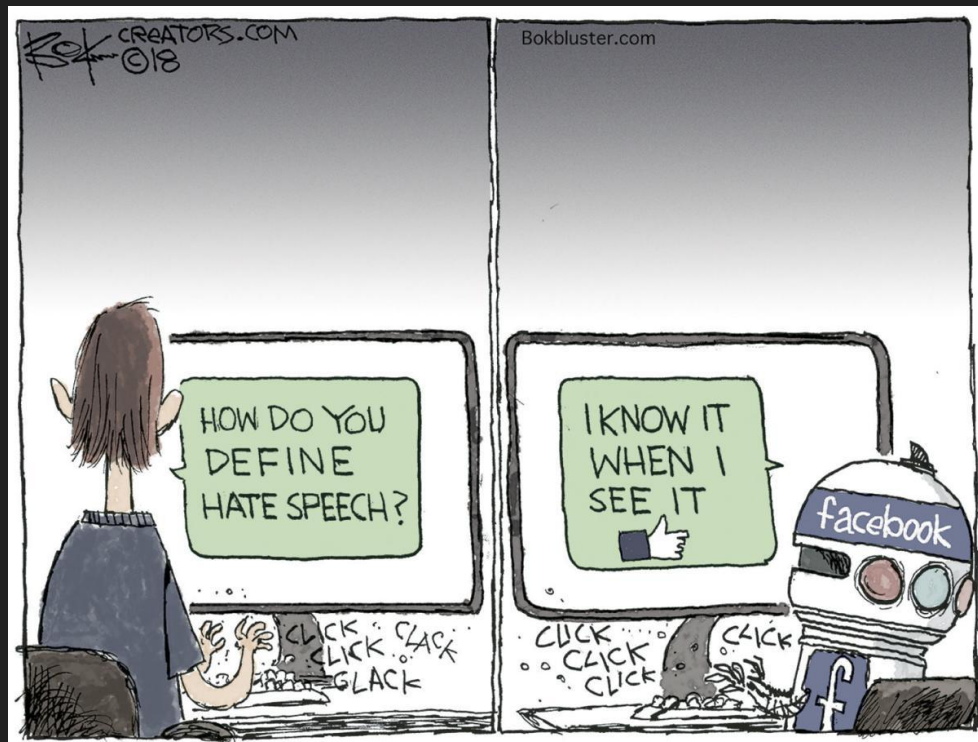# #hate speech recognition

# #overview
## The What

- Twitter - largest microblogging website in the world - 330 million global users

- Content is mostly unrestricted

- The Supreme Court has never created any category of speech that is defined by its hateful conduct.
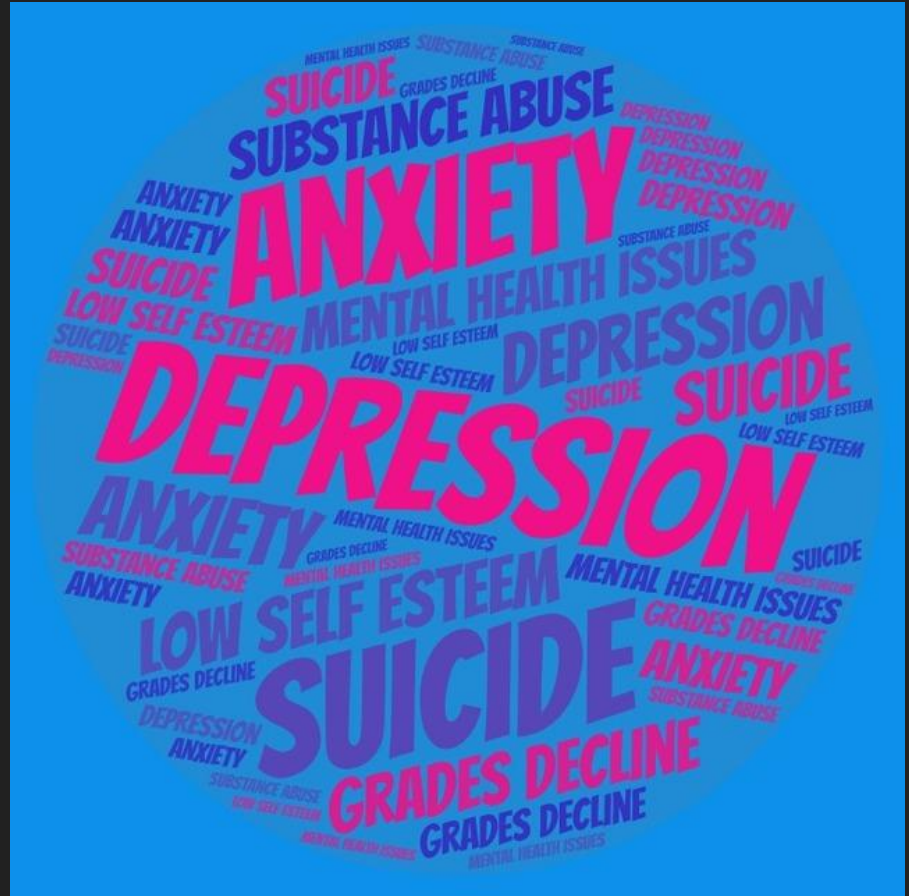
- So what is hate speech?

- Hate speech is defined as abusive or threatening speech that expresses prejudice against a particular group, especially on the basis of race, religion, or sexual orientation

- What is considered 'hate speech' is a fine line between the first amendment and harmful conduct unto others

- With this project, we hope to determine what can be recognized as hate speech on a dataset of over 32,000 tweets
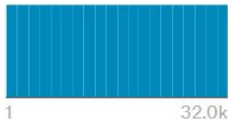
# #inspiration
## The Why

- 59% of Americans believe online hate and harassment make hate crimes more common

- 22% feel less safe in their community because of online hate

- 85% want the government to act by improving training and resources for police on cyber hate
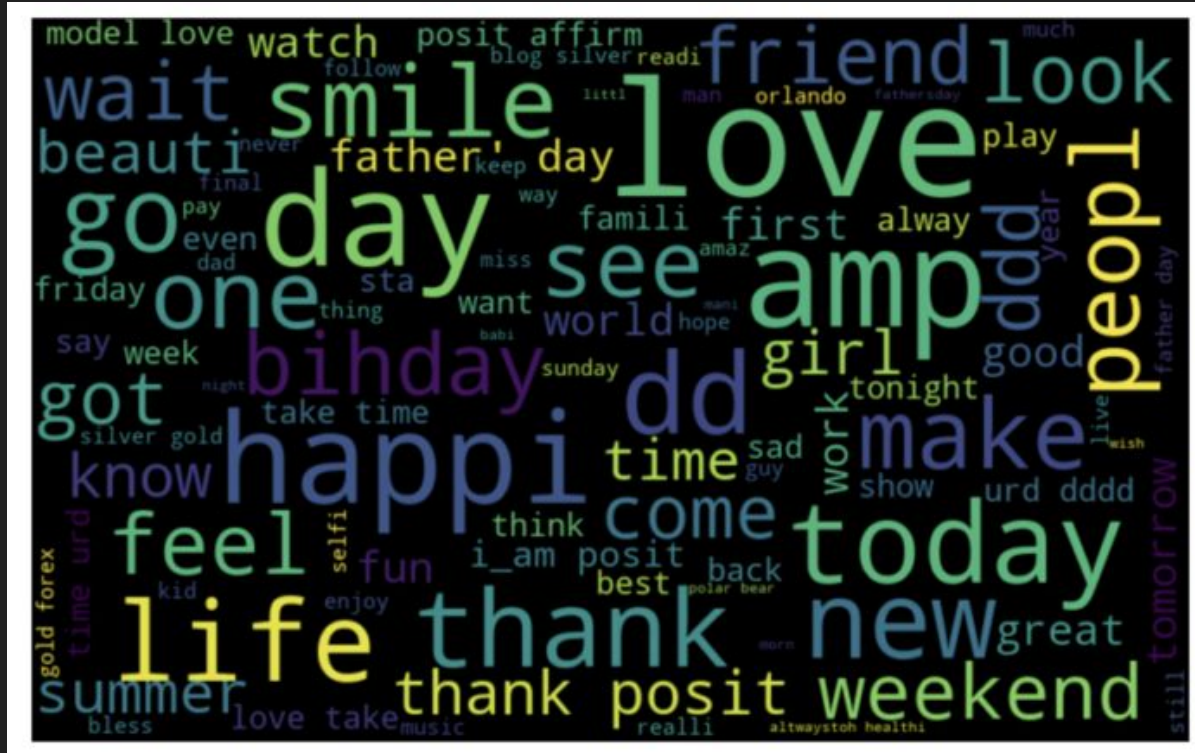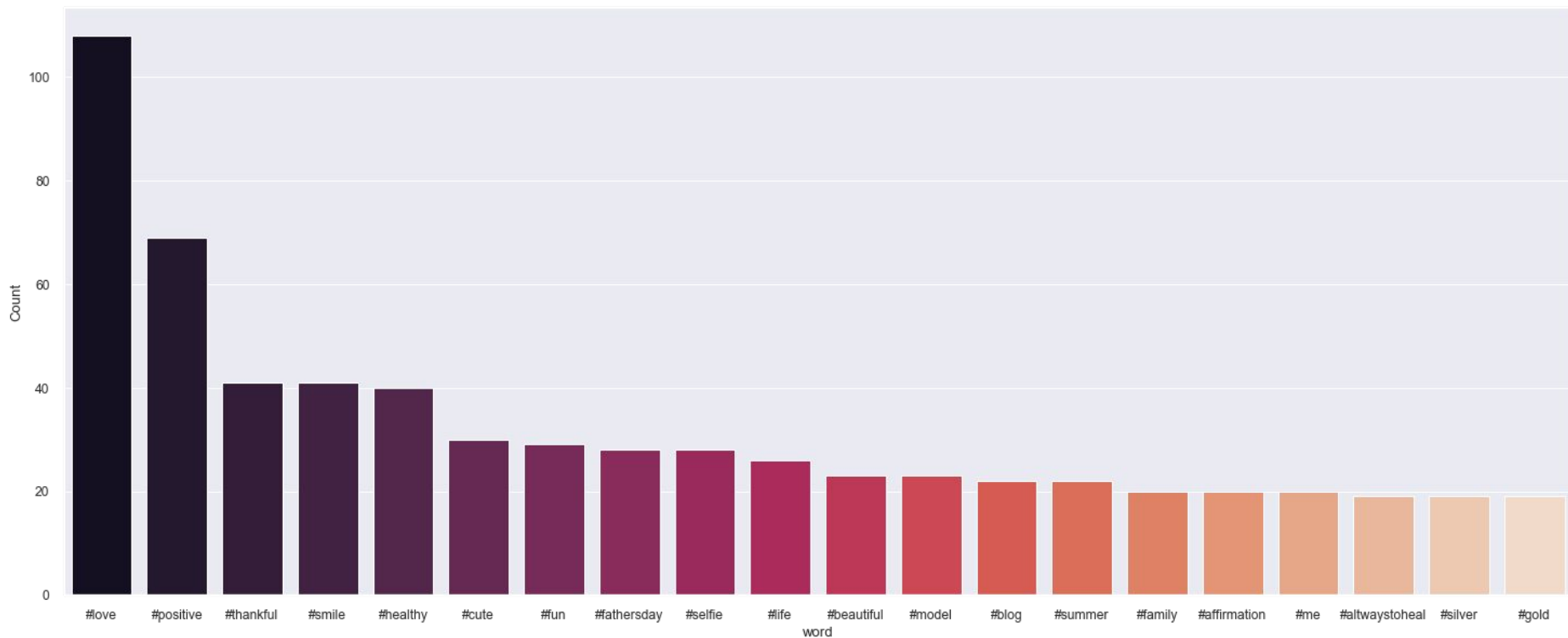
# #data

- 31935 records

- 29,695 Non-Hate Labeled Twitter Data (93%)
- 2240 Hate Labeled Twitter Data (7%)

- Source: https://www.kaggle.com/vkrahul/twitter-hate-speech
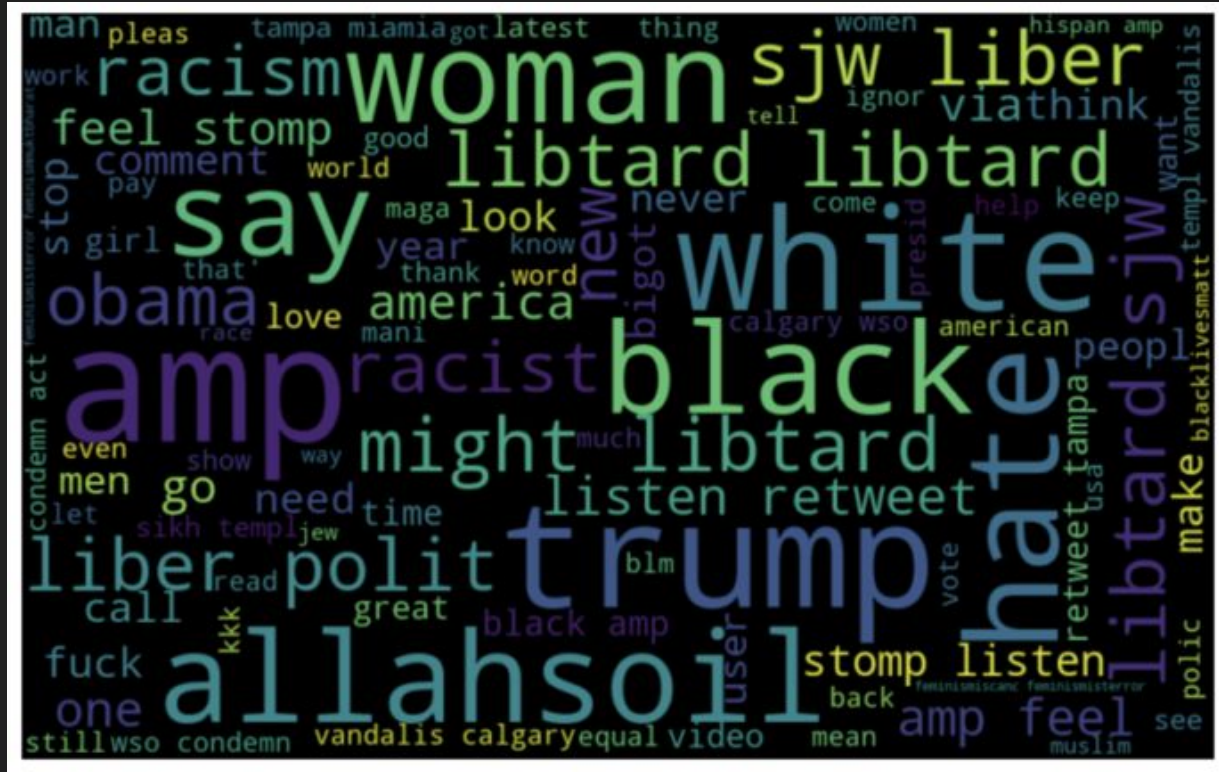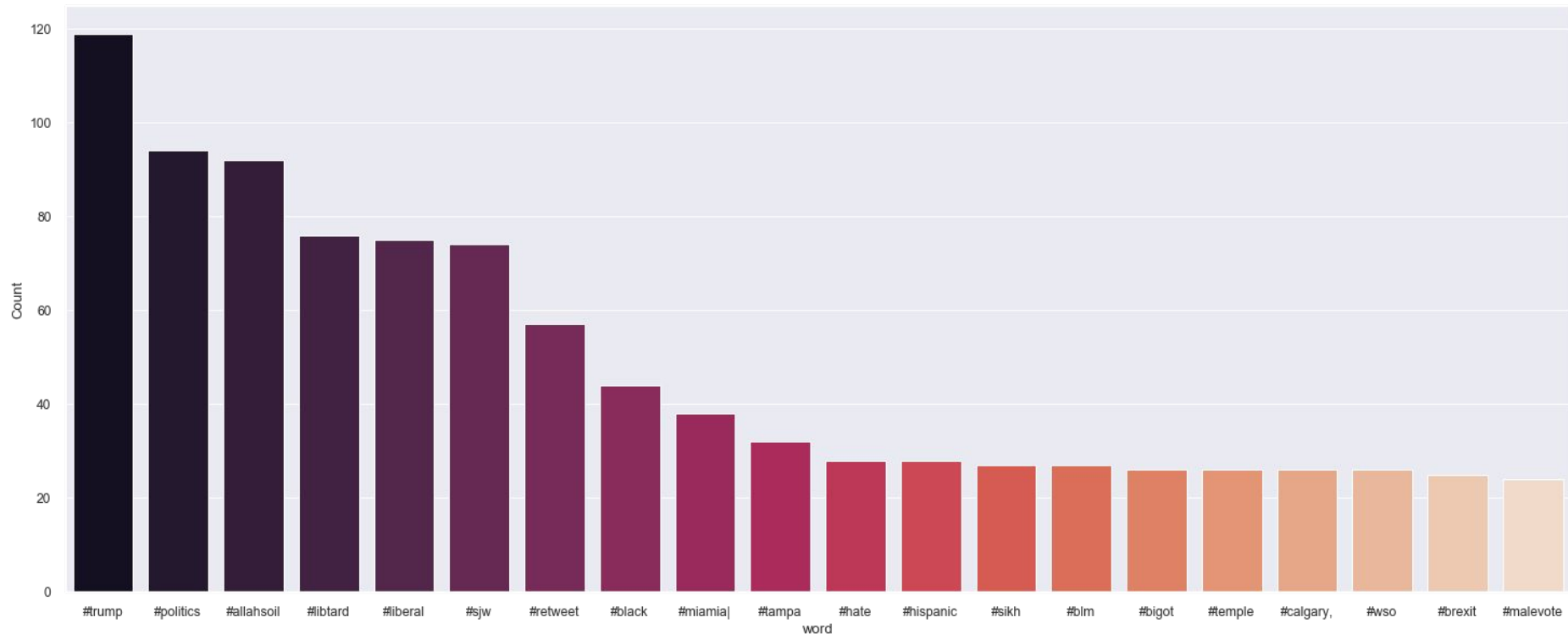
# #positive_word_cloud

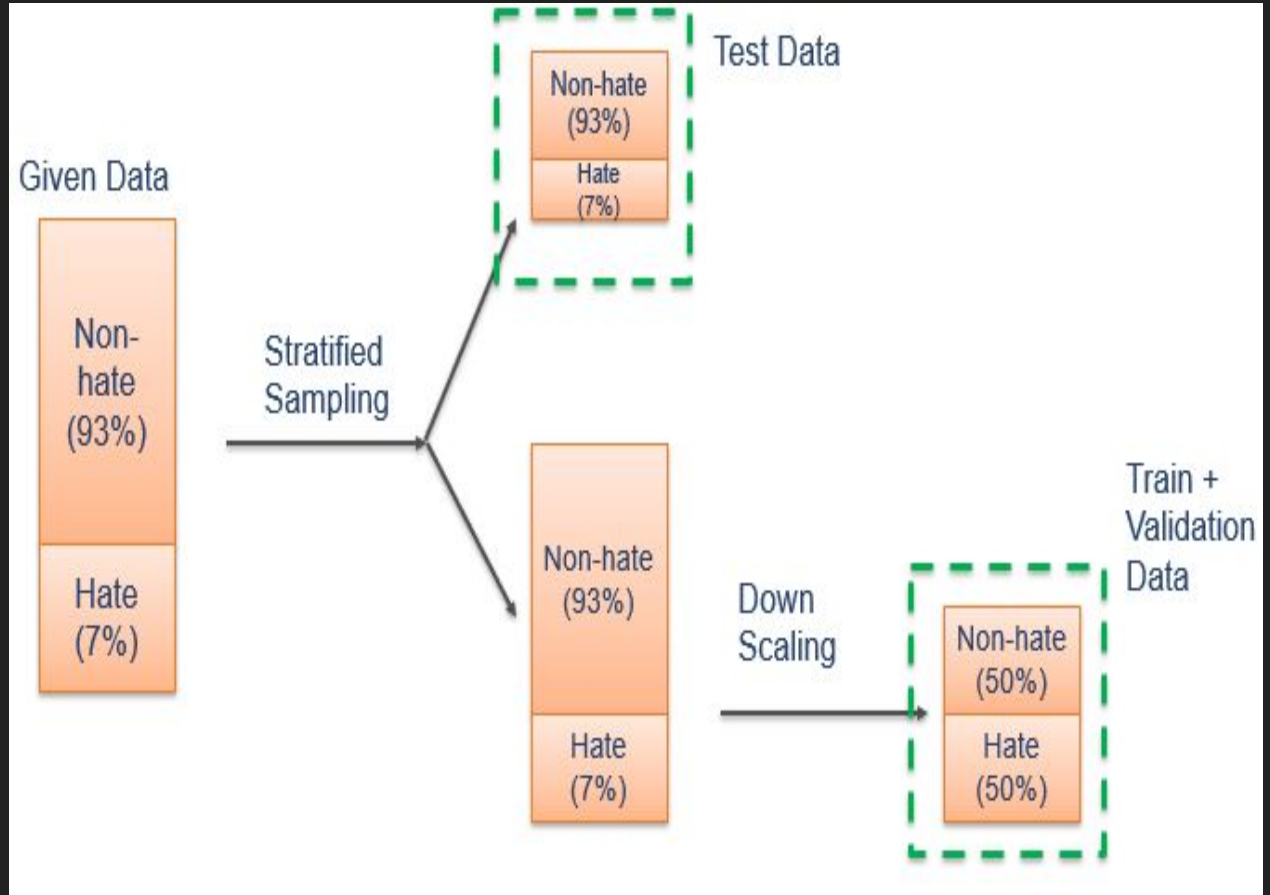# #top20_positive_hashtags

# #negative_word_cloud

# #top20_negative_hashtags

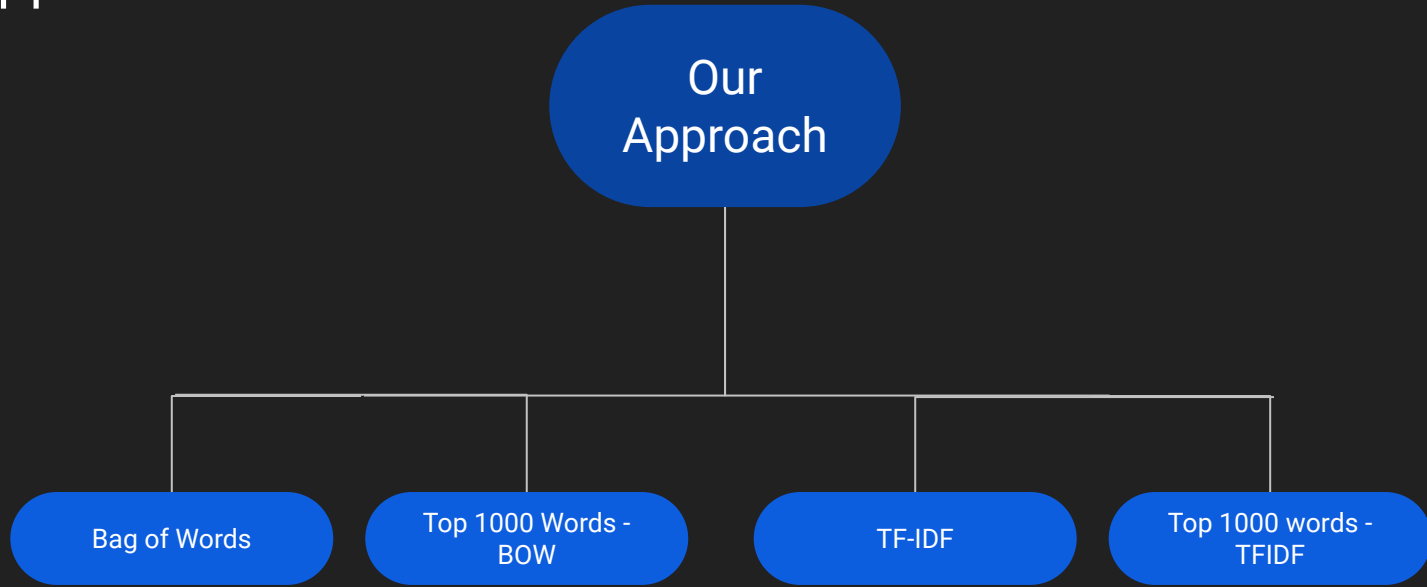# The #how

# #data architecture

- We had 93% of Non-Hate tweets and 7% Hate tweets
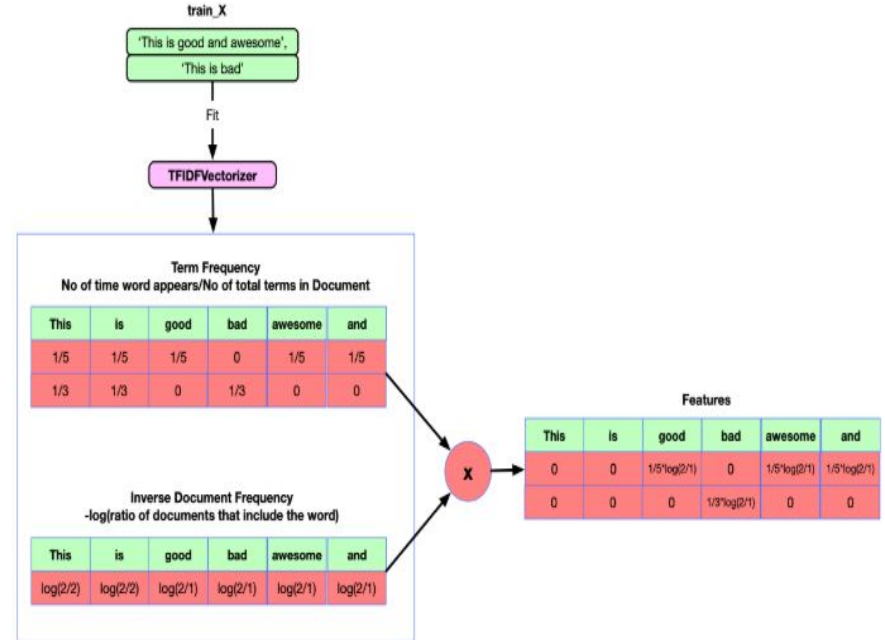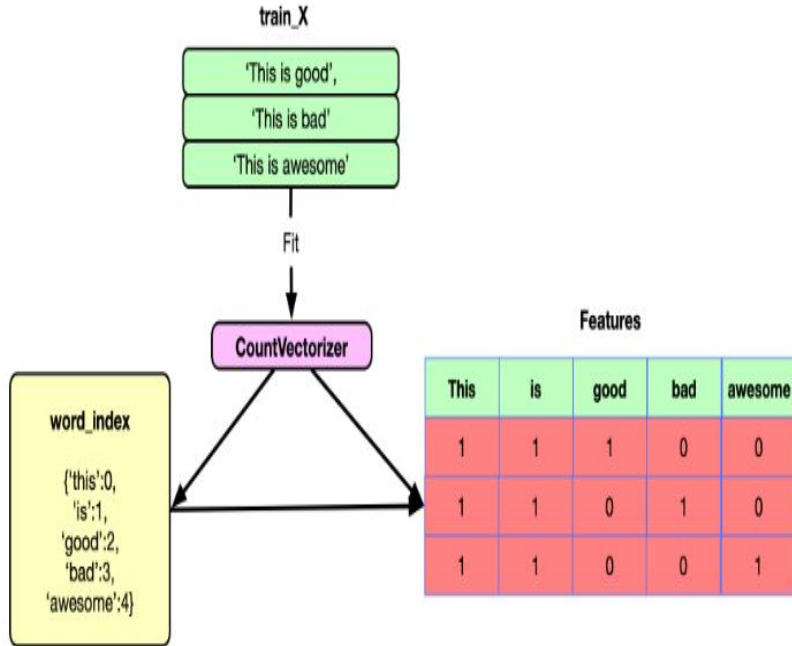
- We downscaled to remove the bias

# #our_approach

# #bag_of_words

# #term_frequency_inverse_docum ent _frequency

#data #cleaning

#LEMMATIZATION

#SLANG_WORDS

#STEMMING

#STOPWORDS

#GREEK_CHARACTERS

# #accuracy?

# #recall?

# #Bag_of_Words



Bag Of Words

| Model | Accuracy | Recall |
|---|---|---|
| Logistic | 0.88 | 0.86 |
| DecisionTree | 0.74 | 0.94 |
| RandomForest | 0.84 | 0.85 |
| NaiveBayes | 0.58 | 0.94 |
| GradientBoosting | 0.92 | 0.64 |

#TF_IDF



TF - IDF

# #1000_most_frequent_words



#TF_IDF

#Bag_Of_Words

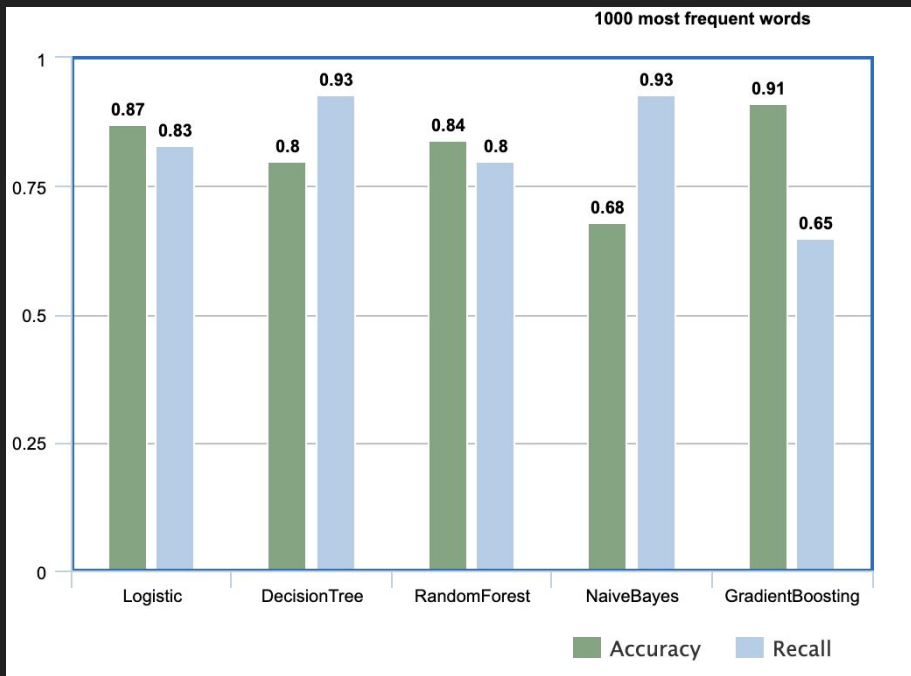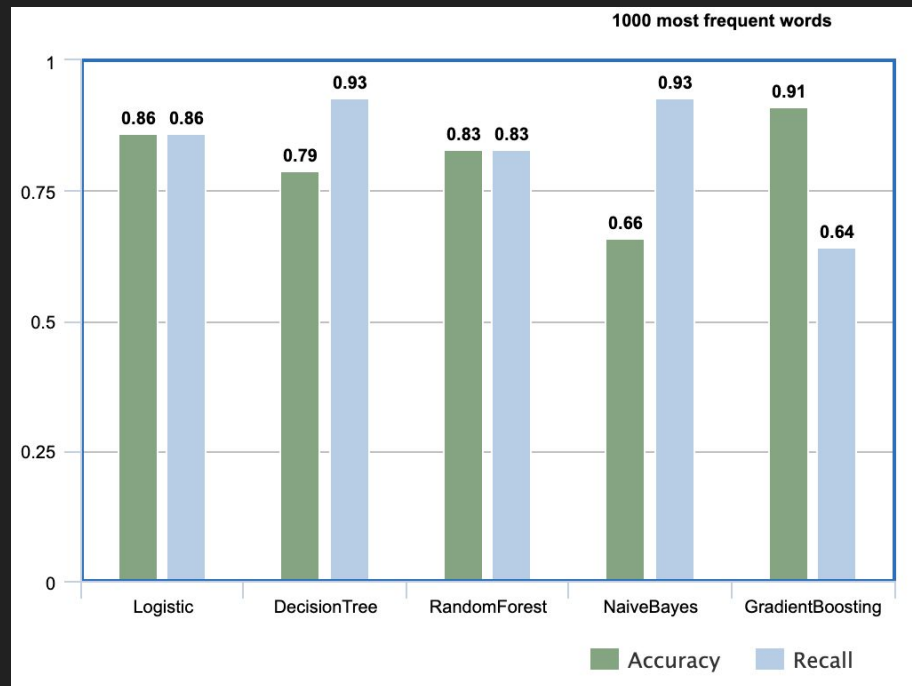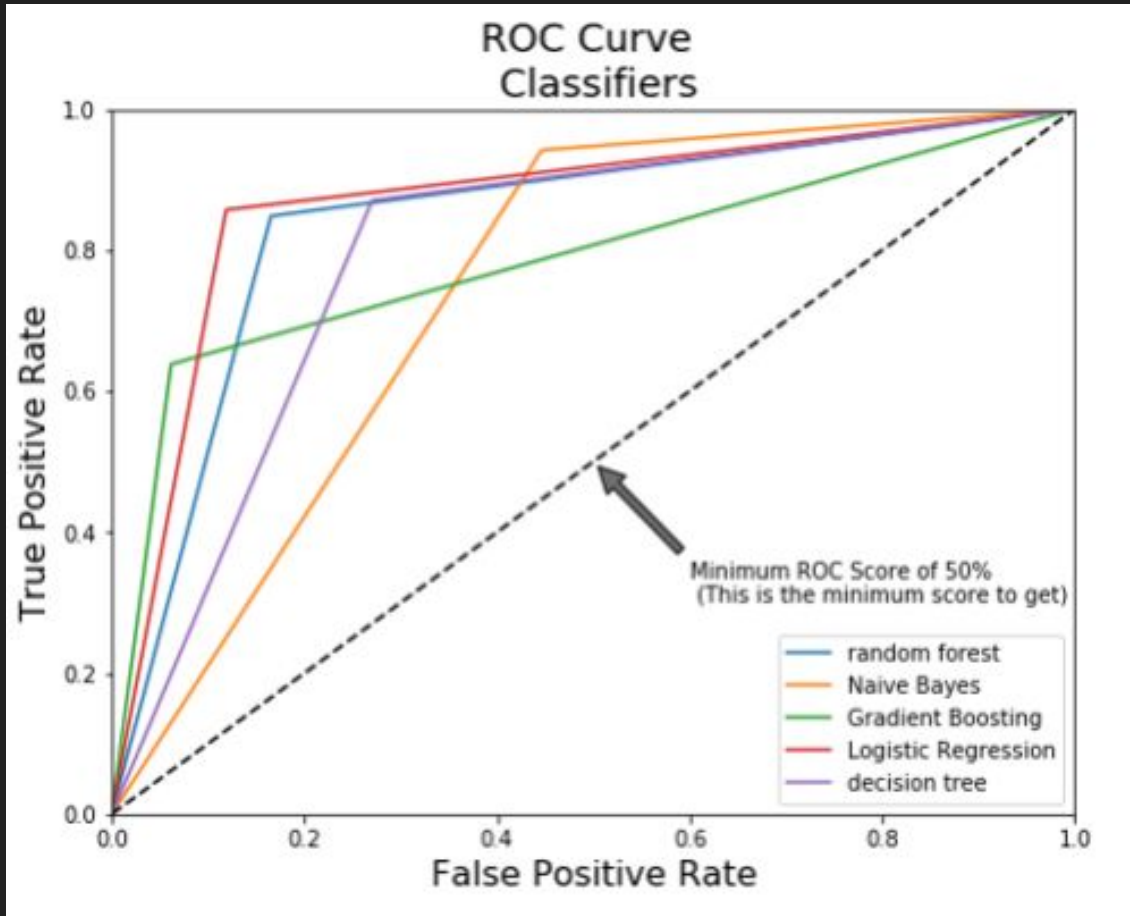# #ROC_curves_for_Bag_of_Words_approach



**AUC**

```
Logistic Regression Score:  0.8692322512282054
Naive Bayes Score:  0.7477859189866101
Random Forest Score:  0.8415043107600423
Decision Tree Score:  0.8000689360369907
Grad Boosting Score:  0.7883467994412869
```

# #key_insights

- Politics, race and sexual tweets form a major chunk of hate tweets
- Results obtained when we played around with an imbalanced data set were inaccurate. The model predicts new data to belong to the major class in the training set due to the skewed nature of the train data.
- The weighted accuracy of a classifier is not the only metric to be looked at while evaluating the performance of the model. Business context plays a vital role as well.

# #challenges #future_steps

- Change in attitude towards topics over time and historical context.

  *"The Nazi organization was great."*

- Automatic hate-speech detection is a closed-loop system:

  Hate verbiage is not restricted to text. Future prospects include recognising hate portrayed through pictures.

# #team



#adwait   #vedant   #sushrut   #sneha   #rishab   #tanvir   #chavi