

Twitter Sentiment Analysis

Project Overview

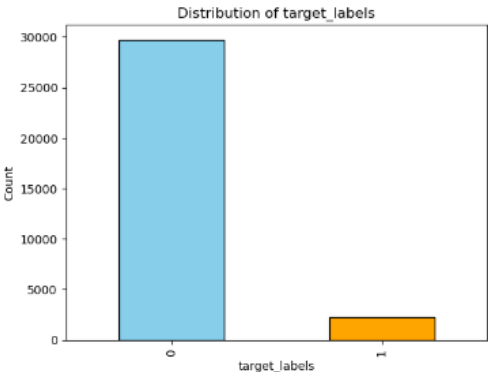
This project focuses on a Kaggle competition hosted by Twitter. The objective is to identify tweets that fall under the category of racism or sexism and potentially block them to reduce bullying and negative tweets as much as possible.

Exploratory-Data-Analysis

```
Data columns (total 3 columns):
#   Column  Non-Null Count  Dtype
---  -
0    id      31962 non-null    int64
1   label    31962 non-null    int64
2   tweet    31962 non-null    object
dtypes: int64(2), object(1)
```

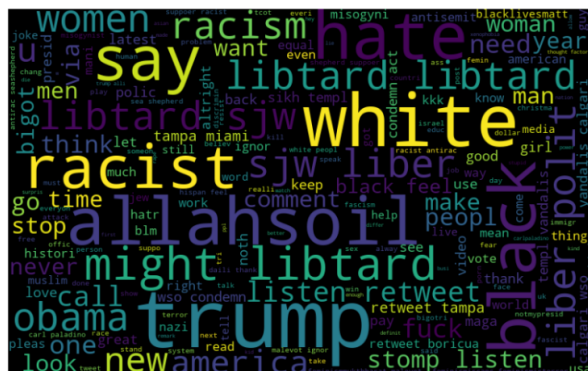
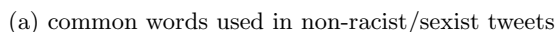
id	label	tweet
24972	0	be n #healthy do not speak in the hearing of...
24469	0	#late #ff to my #gamedev #indiedev #indiegam...
15099	0	after tonight, no more basketball or football.
27073	0	outrageously busy. we've even sold our bunting...
8386	0	ó¼~¥ó¼□□ó¼~□ó¼ #daughter @user just got #gra...

- 1. We observe that each row contains one tweet and a label indicating whether they are fall under the cateogry explained above. Thus, this is a **BINARY CLASSIFICATION PROBLEM**.
- 2. As part of basic preprocessing, we drop any rows that are duplicates or contain null values.
- 3. Next, we check for **class imbalance**, an important consideration in classification tasks. To address class imbalance, we assign a higher penalty when the model misclassifies a minority class during training. This approach encourages the model to pay attention to minority classes, even if the overall accuracy is high.



As part of the preprocessing steps, I have performed the following transformations:

- ## Feature-Engineering

[illegible]

Now, we create a corpus of words from the tweets column of both training and testing data so as to avoid out of vocabulary issue as much as possible. Then, we select the top 1000 most frequent words as the dimensions of vectors to represent each sentence. We will do a comparison between the Bag of Words and TF-IDF methods.

We will now combine these 1000 columns with the 20 features we created earlier (10 most common labels from each category) making a total of 1020 final features

Model Training & Evaluation

We will use a Logistic Regression classifier. The `f1_score` we obtain after training is approximately 0.544 using the Bag of Words (BOW) method and 0.559 using the TF-IDF method.