

DataGlacier NLP Interns Group

MEMBER DETAILS

- Member 1:
 - Name – Gabriel Aluede
 - Email – ghabyaluede@gmail.com
 - Country – United Kingdom
 - College – Teesside University
 - Specialization – NLP
- Member 2:
 - Name – Yosha Udaya Shriyan
 - Email – yosha.shriyan@gmail.com
 - Country – United Kingdom
 - College – Royal Holloway University of London
 - Specialization – NLP

PROBLEM DESCRIPTION

The term hate speech is understood as any type of verbal, written or behavioural communication that attacks or uses derogatory or discriminatory language against a person or group based on what they are, in other words, based on their religion, ethnicity, nationality, race, colour, ancestry, sex or another identity factor. In this problem, we will take you through a hate speech detection model with Machine Learning and Python.

IDENTIFYING PROBLEMS IN DATA

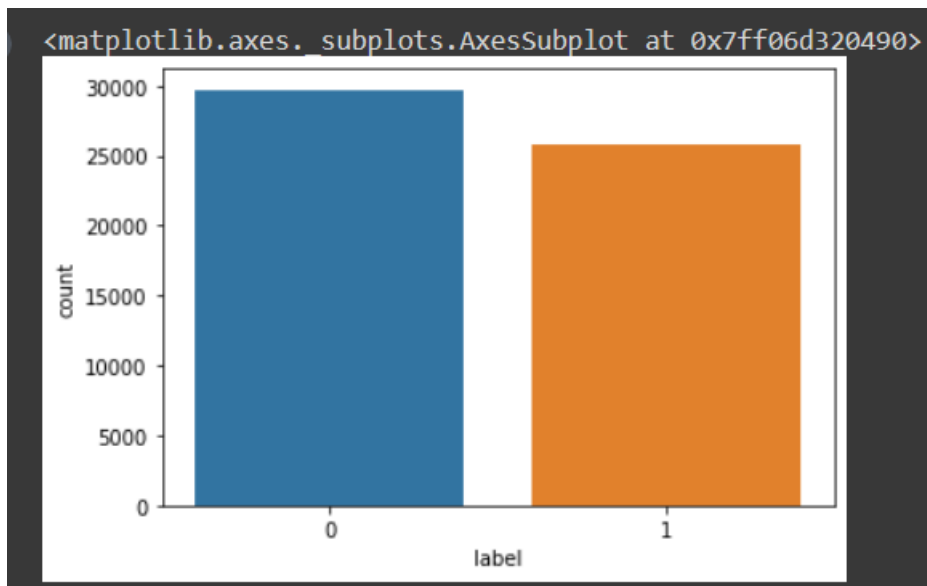
Pre-process data using the library ‘preprocessor’

```
import preprocessor as p
REPLACE_NO_SPACE = re.compile("[^a-zA-Z]http\S+\S*RT|cc#\S+@user\S+!\S+")
REPLACE_WITH_SPACE = re.compile("<br\S/><br\S/?>|(-)|(/)|(:).")
```

Defining function to preprocess each tweet and clean the data

```
def processed_tweets(train_data):
    subArr = []
    for line in train_data:
        #remove punctuation
        tweet = p.clean(line)
        tweet = REPLACE_NO_SPACE.sub("", tweet.lower())
        tweet = REPLACE_WITH_SPACE.sub(" ", tweet)
        subArr.append(tweet)
    return subArr
```


Finally, we merge the resampled data with the other majority up-sampled class. We get a more balanced data as the result.



GITHUB REPO LINK

<https://github.com/YoshaRHUL/DataGlacierProject>