# EDA Presentation and proposed modeling technique

**Advance NLP : Hate Speech detection using Transformers (Deep Learning)**
**Hiten Chadha**
**21.08.2022**

# Agenda

EDA Presentation

Proposed Modeling Technique

# Checking the Shape of Training and Test Data

```python
print("Training Set:"% training_data.columns, training_data.shape)
print("Test Set:"% testing_data.columns, testing_data.shape)
```

```
Training Set: (31962, 3)
Test Set: (17197, 2)
```

We have 31962 and 17197 tweets in the training and test data set respectively.

# Null Data

```
print('Train_Set -----')
print(training_data.isnull().sum())
print('Test_set -----')
print(testing_data.isnull().sum())
training_data.head()
```

```
Train_Set -----
id       0
label    0
tweet    0
dtype: int64
Test_set -----
id       0
tweet    0
dtype: int64
```

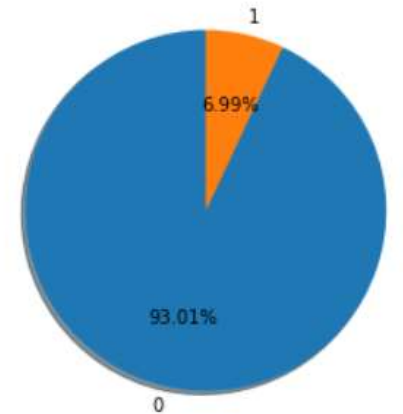| | id | label | tweet |
|---|---|---|---|
| 0 | 1 | 0 | @user when a father is dysfunctional and is s... |
| 1 | 2 | 0 | @user @user thanks for #lyft credit i can't us... |
| 2 | 3 | 0 | bihday your majesty |
| 3 | 4 | 0 | #model i love u take with u all the time in ... |
| 4 | 5 | 0 | factsguide: society now #motivation |

There are no null data in the datasets.

# Positive and Negative Tweets

```
training_data['label'].value_counts() #counting no of positives and negatives
```
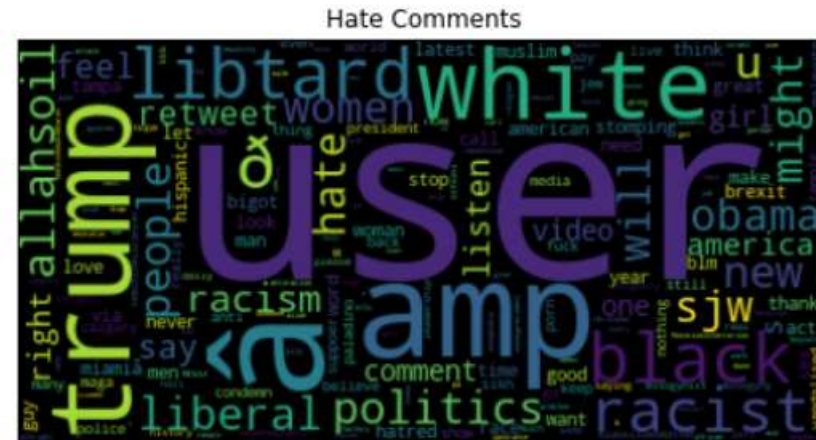
```
0    29720
1     2242
Name: label, dtype: int64
```

There are 2242 hate speech tweets (represented in yellow color in the given pie chart) in the training data and the rest contains no hate speech.
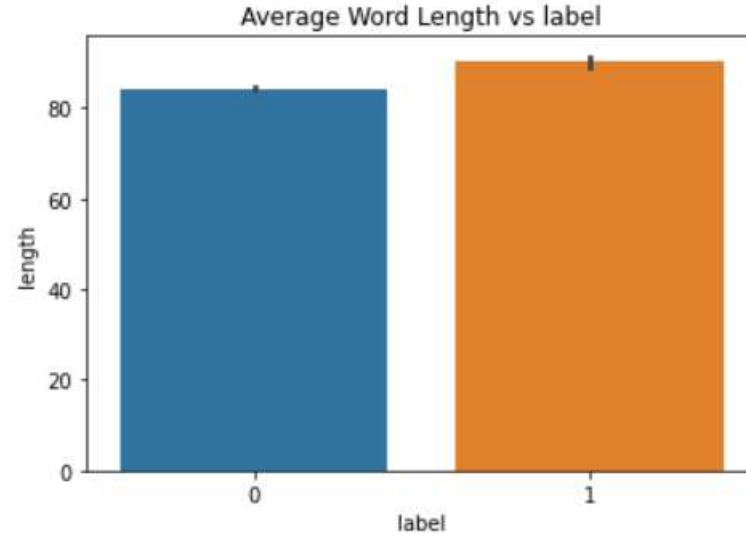


```
0    23783
1     1786
Name: label, dtype: int64
```
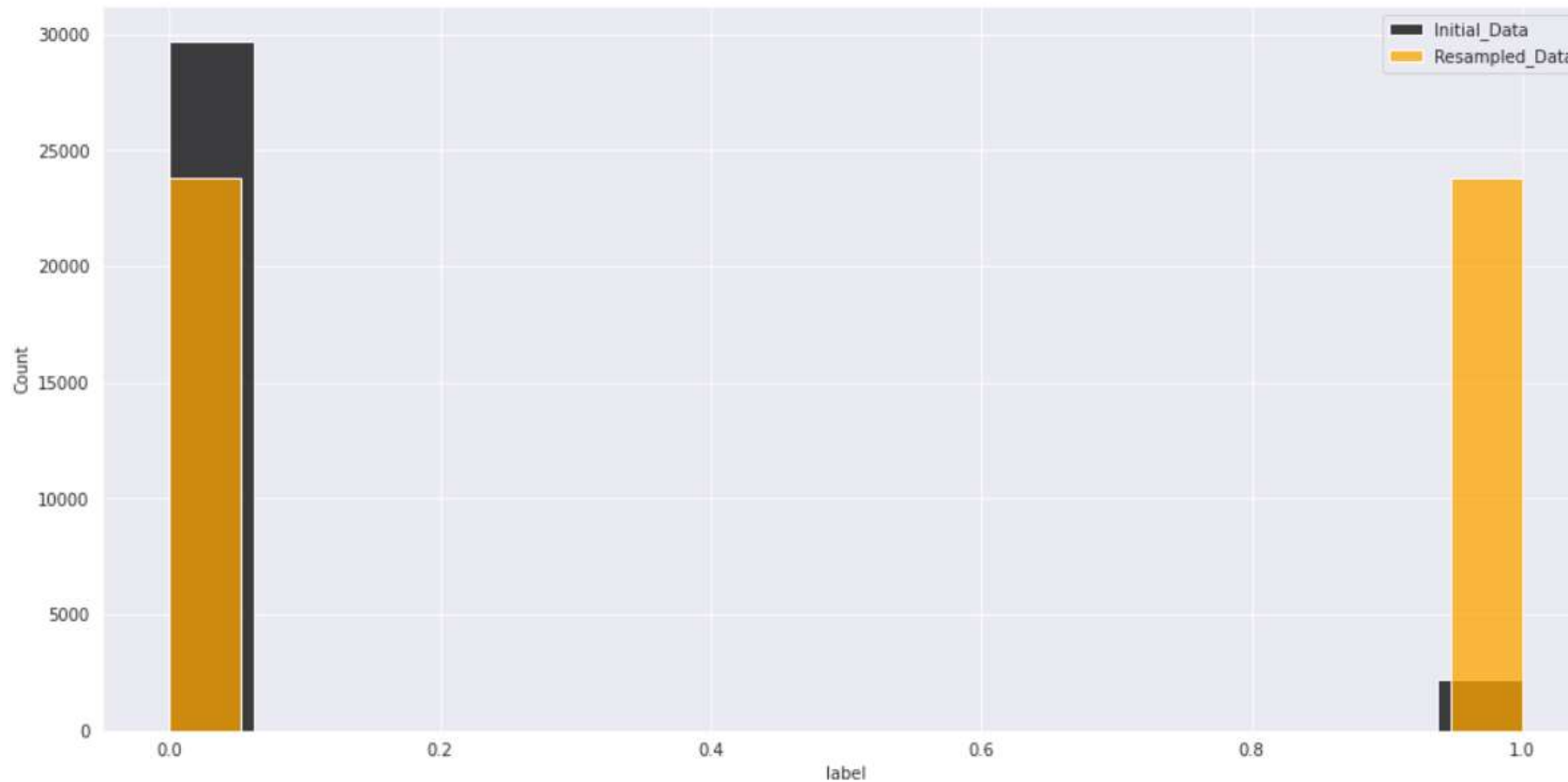
# Most Frequent Hate Words



Visual representation of most frequent hate words

# Average Word Lengths



Average word lengths for hate speech (orange) and non hate speech (blue) tweets.

# Undersampling and Overssampling results



Initial data(black) and after sampling data(orange) for hate and non hate words

# Recommended Models

- XGBClassifier

- LogisticRegression

- MultinomialNB

- SGDClassifier

- DecisionTreeClassifier

- RandomForestClassifier

- KNeighborsClassifier

- LinearSVC

- SVC

- BERT

- RoBERTa

# Thank You