



MALIGNANT COMMENTS **CLASSIFIER PROJECT**

Submitted by:

Kishan Verma

ACKNOWLEDGMENT

I would like to thank Flip Robo Technologies for providing me with the opportunity to work on this project from which I have learned a lot. I am also grateful to my SME, Mr. Shubham Yadav for his constant guidance. Some of the reference sources are as follows:

- [Stack Overflow](#)
- [Medium.com](#)
- [scikit-learn.org](#)
- [Python official documentation](#)

Contents

INTRODUCTION	1
BUSINESS PROBLEM FRAMING	1
CONCEPTUAL BACKGROUND OF THE DOMAIN PROBLEM	1
REVIEW OF LITERATURE.....	1
MOTIVATION FOR THE PROBLEM UNDERTAKEN.....	2
ANALYTICAL PROBLEM FRAMING	2
DATA SOURCES AND THEIR FORMATS.....	2
DATA PREPROCESSING DONE	3
HARDWARE AND SOFTWARE REQUIREMENTS AND TOOLS USED	5
MODEL/S DEVELOPMENT AND EVALUATION	6
IDENTIFICATION OF POSSIBLE PROBLEM-SOLVING APPROACHES (METHODS)	6
TESTING OF IDENTIFIED APPROACHES (ALGORITHMS)	6
RUN AND EVALUATE SELECTED MODELS	7
KEY METRICS FOR SUCCESS IN SOLVING PROBLEM UNDER CONSIDERATION	10
VISUALIZATIONS.....	11
INTERPRETATION OF THE RESULTS.....	14
CONCLUSION.....	15
KEY FINDINGS AND CONCLUSIONS OF THE STUDY.....	15
LEARNING OUTCOMES OF THE STUDY IN RESPECT OF DATA SCIENCE	15
LIMITATIONS OF THIS WORK AND SCOPE FOR FUTURE WORK.....	15

INTRODUCTION

BUSINESS PROBLEM FRAMING

The proliferation of social media enables people to express their opinions widely online. However, at the same time, this has resulted in the emergence of conflict and hate, making online environments uninviting for users. Although researchers have found that hate is a problem across multiple platforms, there is a lack of models for online hate detection.

Online hate, described as abusive language, aggression, cyberbullying, hatefulness and many others has been identified as a major threat on online social media platforms. Social media platforms are the most prominent grounds for such toxic behaviour.

There has been a remarkable increase in the cases of cyberbullying and trolls on various social media platforms. Many celebrities and influences are facing backlashes from people and have to come across hateful and offensive comments. This can take a toll on anyone and affect them mentally leading to depression, mental illness, self-hatred and suicidal thoughts.

Internet comments are bastions of hatred and vitriol. While online anonymity has provided a new outlet for aggression and hate speech, machine learning can be used to fight it. The problem we sought to solve was the tagging of internet comments that are aggressive towards other users. This means that insults to third parties such as celebrities will be tagged as unoffensive, but “u are an idiot” is clearly offensive.

Our goal is to build a prototype of online hate and abuse comment classifier which can be used to classify hate and offensive comments so that it can be controlled and restricted from spreading hatred and cyberbullying.

CONCEPTUAL BACKGROUND OF THE DOMAIN PROBLEM

REVIEW OF LITERATURE

In the past few years it's been seen that the cases related to social media hatred have increased exponentially. The social media is turning into a dark venomous pit for people now a days. Online hate is the result of difference in opinion, race, religion, occupation, nationality etc.

In social media the people spreading or involved in such kind of activities use filthy languages, aggression, images etc. to offend and gravely hurt the person on the other side. This is one of the major concerns now.

The result of such activities can be dangerous. It gives mental trauma to the victims making their lives miserable. People who are not well aware of mental health online hate or cyberbullying become life threatening for them. Such cases are also at rise. It is also taking its toll on religions. Each and every day we can see an incident of fighting between people of different communities or religions due to offensive social media posts.

Online hate, described as abusive language, aggression, cyberbullying, hatefulness, insults, personal attacks, provocation, racism, sexism, threats, or toxicity has been identified as a major

threat on online social media platforms. These kinds of activities must be checked for a better future.

MOTIVATION FOR THE PROBLEM UNDERTAKEN

The project is provided to me by FlipRobo as a part of the internship programme. The exposure to real world data and the opportunity to deploy my skillset in solving a real time problem has been the primary objective. However, the motivation for taking this project was that it is relatively a new field of research. Here we have many options but less concrete solutions. The main motivation is to build a prototype of online hate and abuse comment classifier which can be used to classify hate and offensive comments so that it can be controlled and restricted from spreading hatred and cyberbullying.

ANALYTICAL PROBLEM FRAMING

DATA SOURCES AND THEIR FORMATS

The data was provided by FlipRobo in CSV format. After loading the training dataset into Jupyter Notebook using Pandas and it can be seen that there are eight columns named as: id, comment_text, malignant, highly_malignant, rude, threat, abuse and loathe.

The description of each of the column is given below:

- **Malignant:** It is the Label column, which includes values 0 and 1, denoting if the comment is malignant or not.
- **Highly Malignant:** It denotes comments that are highly malignant and hurtful.
- **Rude:** It denotes comments that are very rude and offensive.
- **Threat:** It contains indication of the comments that are giving any threat to someone.
- **Abuse:** It is for comments that are abusive in nature.
- **Loathe:** It describes the comments which are hateful and loathing in nature.
- **ID:** It includes unique Ids associated with each comment text given.
- **Comment text:** This column contains the comments extracted from various social media platforms.

```
Features Present in the Dataset:  
Index(['id', 'comment_text', 'malignant', 'highly_malignant', 'rude', 'threat',  
      'abuse', 'loathe'],  
      dtype='object')
```

```
Total Number of Rows : 159571  
Total Number of Features : 8
```

```
Data Types of Features :  
id                object  
comment_text      object  
malignant         int64  
highly_malignant  int64  
rude              int64  
threat            int64  
abuse             int64  
loathe            int64  
dtype: object
```

DATA PREPROCESSING DONE

For Data pre-processing we did some data cleaning, where we used wordNet lemmatizer to clean the words and removed special characters using Regexp Tokenizer and filter the words by removing stop words and then used lemmatizers and joined and return the filtered words.

Used TFIDF vectorizer to convert those text into vectors, and split the data and into test and train and trained various Machine learning algorithms.

```

5]: 1 # Here I have made a function in which all the Data Cleaning steps like removing data which is not useful like
2 # email address, mobile numbers, removing punctuations, converting all the documents into lowercase,
3 # using Lemmatization technique, filtering documents using Stopwords, using POS tagging,
4 # all these type of data preprocessing steps are being performed with the help of the function defined below.
5
6 # function to filter using POS tagging..
7 def get_pos(pos_tag):
8     if pos_tag.startswith('J'):
9         return wordnet.ADJ
10    elif pos_tag.startswith('N'):
11        return wordnet.NOUN
12    elif pos_tag.startswith('R'):
13        return wordnet.ADV
14    else:
15        return wordnet.NOUN
16
17 # Function for data cleaning...
18 def Processed_data(comments):
19     # Replace email addresses with 'email'
20     comments=re.sub(r'^.+@[^\.\.]*\.[a-z]{2,}$',' ', comments)
21
22     # Replace 10 digit phone numbers (formats include paranthesis, spaces, no spaces, dashes) with 'phonenumber'
23     comments=re.sub(r'^\([?\d]{3}\)?[\s-]?[?\d]{3}[\s-]?[?\d]{4}$',' ', comments)
24
25     # getting only words(i.e removing all the special characters)
26     comments = re.sub(r'^\w',' ', comments)
27
28     # getting only words(i.e removing all the " - ")
29     comments = re.sub(r'[_]', ' ', comments)
30
31     # getting rid of unwanted characters(i.e remove all the single characters left)
32     comments=re.sub(r'\s+[a-zA-Z]\s+', ' ', comments)
33
34     # Removing extra whitespaces
35     comments=re.sub(r'\s+', ' ', comments, flags=re.I)
36
37     #converting all the letters of the review into lowercase
38     comments = comments.lower()
39
40     # splitting every words from the sentences
41     comments = comments.split()
42
43     # iterating through each words and checking if they are stopwords or not,
44     comments=[word for word in comments if not word in set(STOPWORDS)]
45
46     # remove empty tokens
47     comments = [text for text in comments if len(text) > 0]
48
49     # getting pos tag text
50     pos_tags = pos_tag(comments)
51
52     # considering words having Length more than 3only
53     comments = [text for text in comments if len(text) > 3]
54
55     # performing Lemmatization operation and passing the word in get_pos function to get filtered using POS ...
56     comments = [(WordNetLemmatizer().lemmatize(text[0], get_pos(text[1]))for text in pos_tags]
57
58     # considering words having Length more than 3 only
59     comments = [text for text in comments if len(text) > 3]
60     comments = ' '.join(comments)
61     return comments

```

```

5]: 1 df["clean_comment_text"] = df["comment_text"].apply(lambda x: Processed_data(x))

```

```

7]: 1 # checking the clean comments content..
2 df["clean_comment_text"]

```

```

In[17]: 0      explanation edits username hardcore metallica ...
1      match background colour seemingly stuck thanks...
2      trying edit constantly removing relevant infor...
3      real suggestion improvement wondered section s...
4      hero chance remember page
        ...
159566 second time asking view completely contradicts...
159567      ashamed horrible thing talk page
159568 spitzer there actual article prostitution ring...
159569      look like actually speedy version deleted look
159570 think understand came idea right away kind com...
Name: clean_comment_text, Length: 159571, dtype: object

```

HARDWARE AND SOFTWARE REQUIREMENTS AND TOOLS USED

HARDWARE:

HP Pavilion X360

SOFTWARE:

Jupyter Notebook (Anaconda 3) – Python 3

Microsoft Office 365 Package

LIBRARIES USED:

```
1 # Importing Libraries..
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 import string
7 import re
8 from collections import Counter
9 # packages from gensim
10 from gensim import corpora
11 from gensim.utils import simple_preprocess
12 from gensim.parsing.preprocessing import STOPWORDS
13
14 # packages from sklearn
15 from sklearn.feature_extraction.text import TfidfVectorizer
16
17 #packages from nltk
18 from nltk.corpus import wordnet
19 from nltk.stem import WordNetLemmatizer, SnowballStemmer
20 from nltk import pos_tag
21
22 import warnings
23 warnings.filterwarnings('ignore')
```

```
1 # Importing useful libraries for model training
2
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn.tree import DecisionTreeClassifier
6 from sklearn.linear_model import LogisticRegression,PassiveAggressiveClassifier
7
8 # Model selection Libraries...
9 from sklearn.model_selection import cross_val_score, cross_val_predict, train_test_split
10 from sklearn.model_selection import GridSearchCV
11
12
13 # Importing some metrics we can use to evaluate our model performance....
14 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
15 from sklearn.metrics import roc_auc_score, roc_curve, auc
16 from sklearn.metrics import precision_score, recall_score, f1_score
17 from sklearn.metrics import log_loss
18
19 # Creating instances for different Classifiers
20
21 LR=LogisticRegression()
22 MNB=MultinomialNB()
23 PAC=PassiveAggressiveClassifier()
24 DT=DecisionTreeClassifier()
25
```

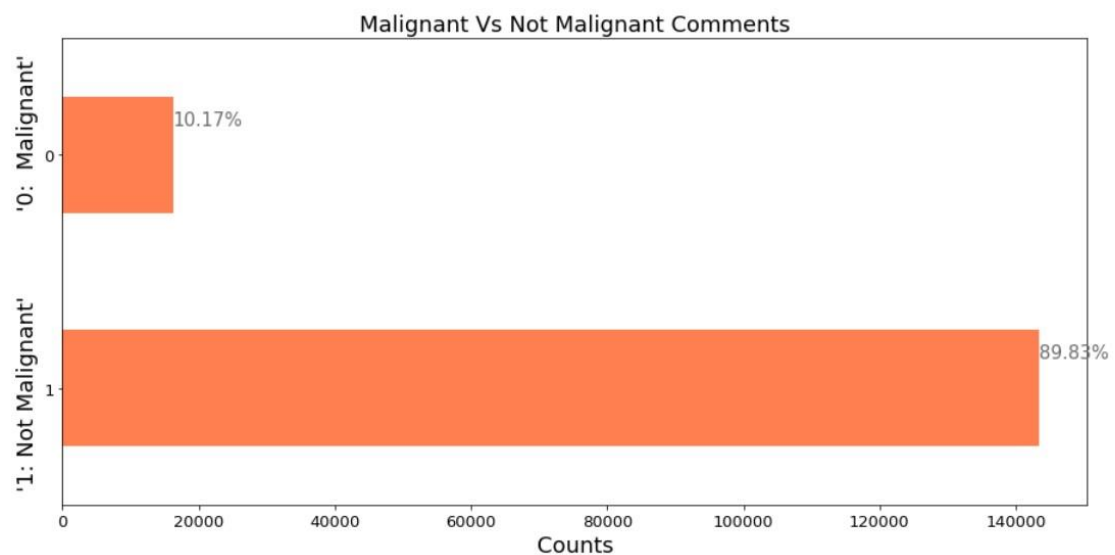

MODEL/S DEVELOPMENT AND EVALUATION

IDENTIFICATION OF POSSIBLE PROBLEM-SOLVING APPROACHES (METHODS)

```
1 # Tokenizing Documents..  
2 data=[]  
3 from nltk.tokenize import word_tokenize  
4 for j,i in enumerate(df['clean_comment_text']):  
5     a=word_tokenize(i,'english')  
6     data.append(a)
```

```
1 # Makin Word dictionary...  
2 dictionary = corpora.Dictionary(data)  
3 print(dictionary)
```

Dictionary(167144 unique tokens: ['closure', 'doll', 'edits', 'explanation', 'hardcore']...)



TESTING OF IDENTIFIED APPROACHES (ALGORITHMS)

The algorithms we used for the training and testing are as follows:-

- Logistic Regression
- Decision Tree Classifier
- Multinomial NB
- Passive Aggressive Classifier

RUN AND EVALUATE SELECTED MODELS

```
1 # Putting Scikit-Learn machine Learning Models in a List so that it can be used for further evaluation in Loop.
2 models=[]
3 models.append(('LogisticRegression',LR))
4 models.append(('MultinomialNB',MNB))
5 models.append(('PassiveAggressiveClassifier',PAC))
6 models.append(('DecisionTreeClassifier',DT))
```

```
1 # Function which will find best Random State and then calculate Maximum Accuracy Score corresponding to it
2 # and print accuracy score in one go.
3 def max_acc_score(clf,x,y):
4     max_acc_score=0
5     final_r_state=0
6     for r_state in range(42,100):
7         x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.30,random_state=r_state,stratify=y)
8         clf.fit(x_train,y_train)
9         y_pred=clf.predict(x_test)
10        acc_score=accuracy_score(y_test,y_pred)
11        if acc_score > max_acc_score:
12            max_acc_score=acc_score
13            final_r_state=r_state
14        print('Max Accuracy Score corresponding to Random State ', final_r_state, 'is:', max_acc_score)
15        print('\n')
16    return final_r_state
```

```
1 # Lists to store model name, Learning score, Accuracy score, cross_val_score, Auc Roc score .
2 Model=[]
3 Score=[]
4 Acc_score=[]
5 cvs=[]
6 rocscore=[]
7 lg_loss=[]
8 # For Loop to Calculate Accuracy Score, Cross Val Score, Classification Report, Confusion Matrix
9
10 for name,model in models:
11     print('*****',name,'*****')
12     print('\n')
13     Model.append(name)
14     print(model)
15     print('\n')
16
17     # Now here I am calling a function which will calculate the max accuracy score for each model
18     # and return best random state.
19     r_state=max_acc_score(model,x,y)
20     x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.30,random_state=r_state,stratify=y)
21     model.fit(x_train,y_train)
22     #.....Learning Score.....
23     score=model.score(x_train,y_train)
24     print('Learning Score : ',score)
25     Score.append(score*100)
26     y_pred=model.predict(x_test)
27     acc_score=accuracy_score(y_test,y_pred)
28     print('Accuracy Score : ',acc_score)
29     Acc_score.append(acc_score*100)
30
31     #.....Finding Cross_val_score.....
32     cv_score=cross_val_score(model,x,y,cv=10,scoring='roc_auc').mean()
33     print('Cross Val Score : ', cv_score)
34     cvs.append(cv_score*100)
35
```

```

36 #.....Roc auc score.....
37 false_positive_rate,true_positive_rate, thresholds=roc_curve(y_test,y_pred)
38 roc_auc=auc(false_positive_rate, true_positive_rate)
39 print('roc auc score : ', roc_auc)
40 rocscore.append(roc_auc*100)
41 print('\n')
42
43 loss = log_loss(y_test,y_pred)
44 print('Log loss : ', loss)
45 lg_loss.append(loss)
46 print('\n')
47
48 #.....Classification Report.....
49 print('Classification Report:\n',classification_report(y_test,y_pred))
50 print('\n')
51
52 print('Confusion Matrix:\n',confusion_matrix(y_test,y_pred))
53 print('\n')
54

```

***** LogisticRegression *****

LogisticRegression()

Max Accuracy Score corresponding to Random State 44 is: 0.9543574532085561

Learning Score : 0.9570721313530112
Accuracy Score : 0.9543574532085561
Cross Val Score : 0.9647437145554436
roc auc score : 0.7934277979300925

Log loss : 1.5764709396342023

Classification Report:				
	precision	recall	f1-score	support
0	0.94	0.59	0.72	4868
1	0.96	1.00	0.98	43004
accuracy			0.95	47872
macro avg	0.95	0.79	0.85	47872
weighted avg	0.95	0.95	0.95	47872

Confusion Matrix:
[[2879 1989]
[196 42808]]

***** MultinomialNB *****

MultinomialNB()

Max Accuracy Score corresponding to Random State 54 is: 0.9368524398395722

Learning Score : 0.9389072417837223
Accuracy Score : 0.9368524398395722
Cross Val Score : 0.9294912255583274
roc auc score : 0.6948768767912668

Log loss : 2.1810889674255955

Classification Report:

	precision	recall	f1-score	support
0	0.97	0.39	0.56	4868
1	0.94	1.00	0.97	43004
accuracy			0.94	47872
macro avg	0.95	0.69	0.76	47872
weighted avg	0.94	0.94	0.92	47872

Confusion Matrix:
[[1904 2964]
[59 42945]]

***** PassiveAggressiveClassifier *****

PassiveAggressiveClassifier()

Max Accuracy Score corresponding to Random State 58 is: 0.947735628342246

Learning Score : 0.9897313315249018
Accuracy Score : 0.9460018382352942
Cross Val Score : 0.9359602149598304
roc auc score : 0.8312225661376249

Log loss : 1.8650558733232439

Classification Report:

	precision	recall	f1-score	support
0	0.76	0.69	0.72	4868
1	0.96	0.98	0.97	43004
accuracy			0.95	47872
macro avg	0.86	0.83	0.85	47872
weighted avg	0.94	0.95	0.94	47872

Confusion Matrix:
[[3345 1523]
[1062 41942]]

```
***** DecisionTreeClassifier *****
```

```
DecisionTreeClassifier()
```

```
Max Accuracy Score corresponding to Random State 87 is: 0.941949364973262
```

```
Learning Score : 0.9982721420961692
Accuracy Score : 0.9420120320855615
Cross Val Score : 0.8346283576844614
roc auc score : 0.8343758433508737
```

```
Log loss : 2.0028579103214796
```

```
Classification Report:
              precision    recall  f1-score   support

     0           0.72       0.70       0.71         4868
     1           0.97       0.97       0.97        43004

   accuracy              0.94         47872
  macro avg           0.84       0.83       0.84         47872
 weighted avg           0.94       0.94       0.94         47872
```

```
Confusion Matrix:
[[ 3404 1464]
 [ 1312 41692]]
```

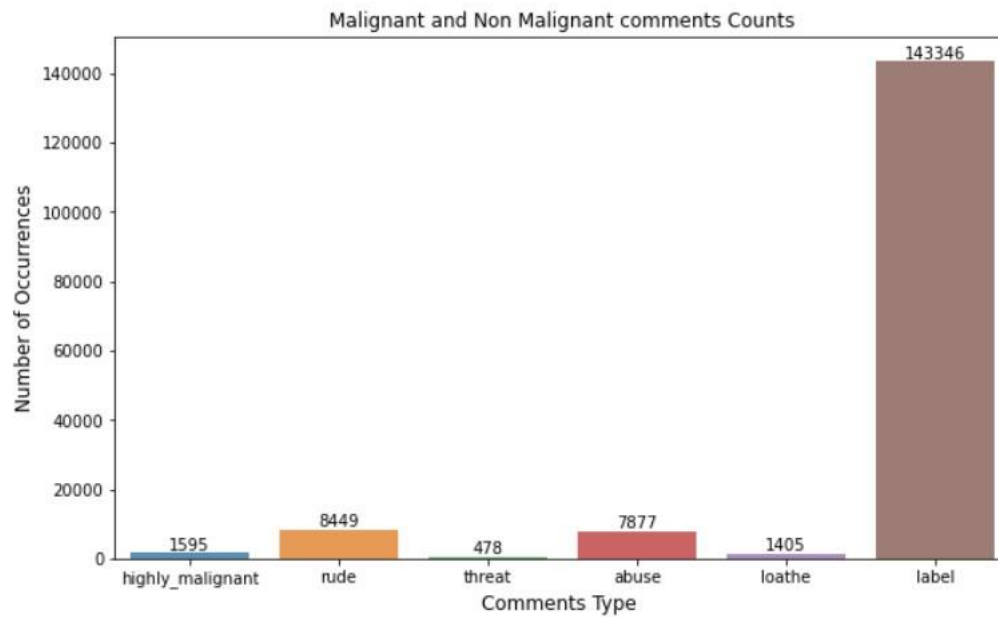
KEY METRICS FOR SUCCESS IN SOLVING PROBLEM UNDER CONSIDERATION

```
# Making a Dataframe comprises of Differnt Calculated Scores :
result=pd.DataFrame({'Model': Model, 'Learning Score': Score, 'Accuracy Score': Acc_score, 'Cross Val Score': cvs,
                    'Roc_Auc_curve': roc_score, 'Log_Loss': lg_loss})
result.style.background_gradient(cmap='YlGnBu')
```

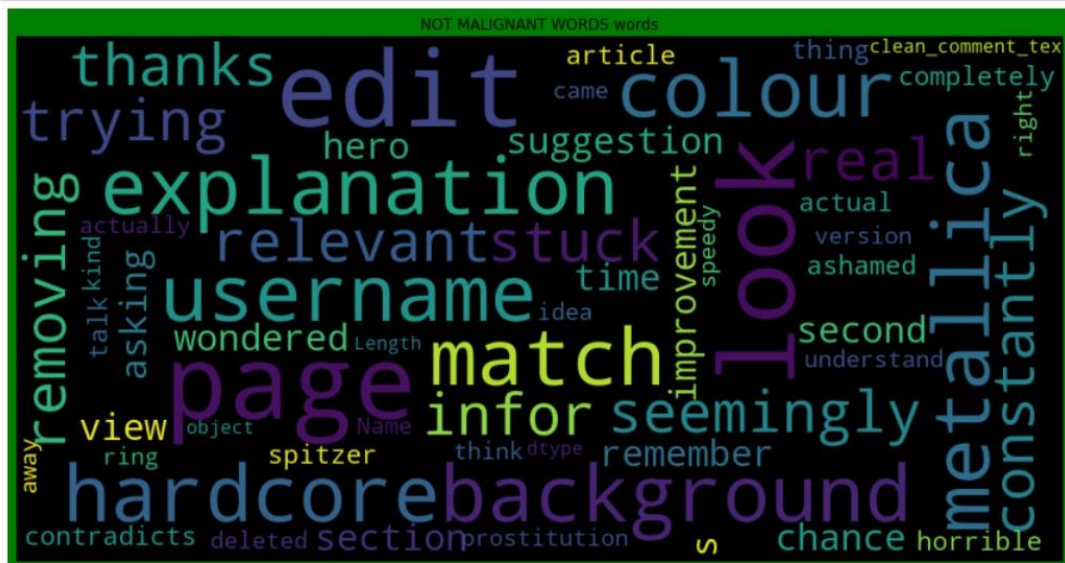
```
39]:
```

	Model	Learning Score	Accuracy Score	Cross Val Score	Roc_Auc_curve	Log_Loss
0	LogisticRegression	95.707213	95.435745	96.474371	79.342780	1.576471
1	MultinomialNB	93.890724	93.685244	92.949123	69.487688	2.181089
2	PassiveAggressiveClassifier	98.973133	94.600184	93.596021	83.122257	1.865056
3	DecisionTreeClassifier	99.827214	94.201203	83.462836	83.437584	2.002858

VISUALIZATIONS



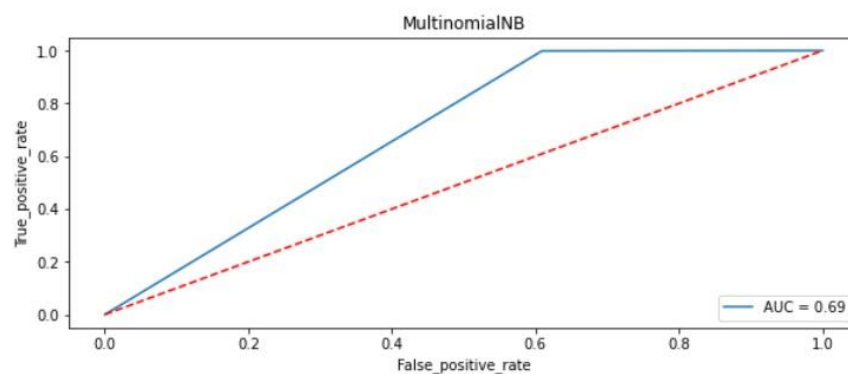
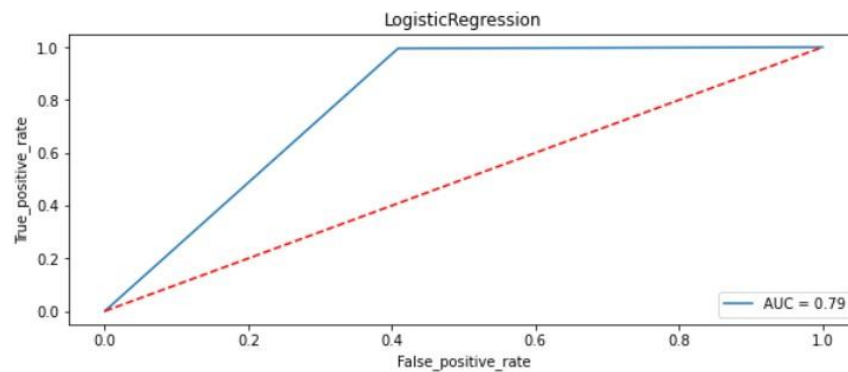
```
1 # NOT Malignant words in training data
2 Display_wordcloud(df['clean_comment_text'][df['label']==1], "NOT MALIGNANT WORDS")
```

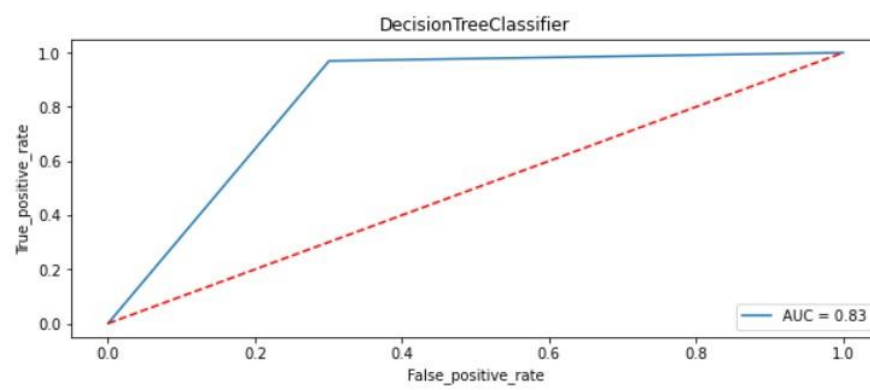
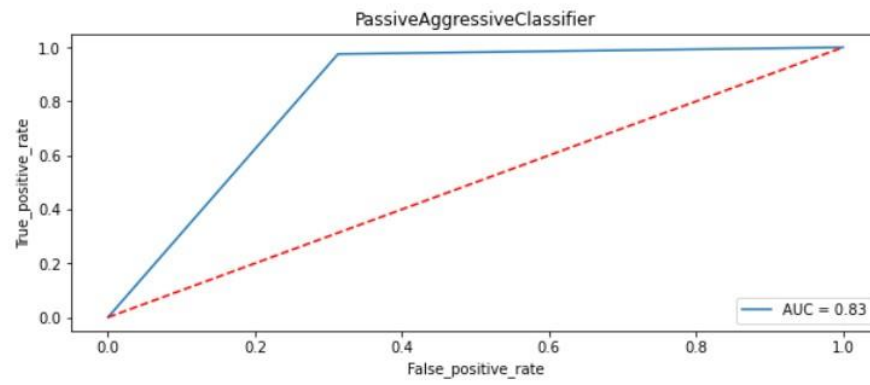


```

1 # Malignant words in training data
2 Display_wordcloud(df['clean_comment_text'][df['label']!=0], "MALIGNANT WORDS")

```





INTERPRETATION OF THE RESULTS

From the above visualization and matrices found that the Passive Aggressive Classifier performed the best AUC_ROC_SCORE.

```
1 # Using PassiveAggressiveClassifier for final model...
2 x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=57,test_size=.30,stratify=y)
3 PAC=PassiveAggressiveClassifier()
4 PAC.fit(x_train,y_train)
5 PAC.score(x_train,y_train)
6 PACpred=PAC.predict(x_test)
7 print('Accuracy Score:',accuracy_score(y_test,PACpred))
8 print('Log loss : ', log_loss(y_test,PACpred))
9 print('Confusion Matrix:',confusion_matrix(y_test,PACpred))
10 print('Classification Report:', '\n', classification_report(y_test,PACpred))
```

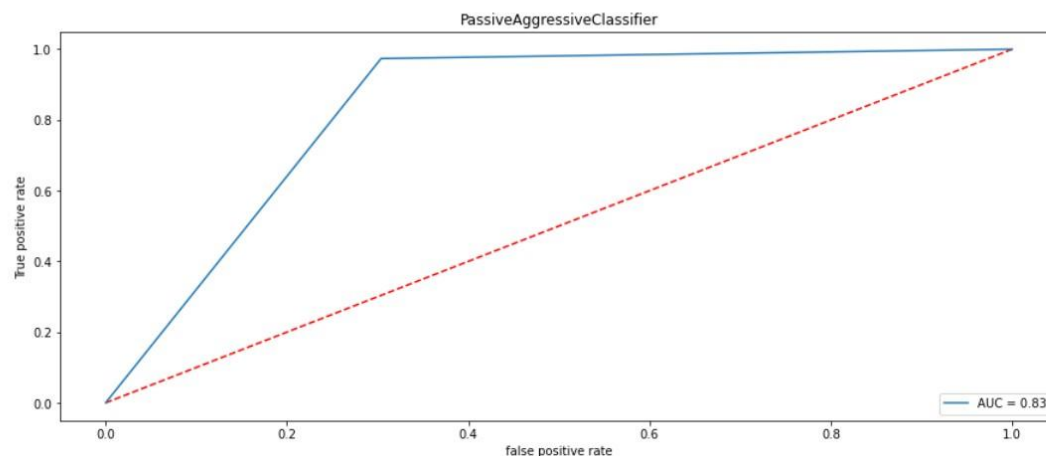
```
Accuracy Score: 0.9453960561497327
Log loss : 1.8859781103963864
Confusion Matrix: [[ 3389 1479]
 [1135 41869]]
Classification Report:
      precision    recall  f1-score   support

     0       0.75     0.70     0.72       4868
     1       0.97     0.97     0.97      43004

 accuracy      0.95      47872
 macro avg     0.86     0.83     0.85      47872
 weighted avg     0.94     0.95     0.94      47872
```

```
1 # Roc-Auc score
2 f,ax = plt.subplots(figsize = (15,6))
3 # Calculate fpr, tpr and thresholds
4 fpr, tpr, thresholds = roc_curve(y_test, PACpred)
5 ax.plot([0,1],[0,1], 'r--')
6 ax.plot(fpr,tpr,label='AUC = %.2f'% roc_auc_score(y_test, PACpred))
7 ax.legend(loc='lower right')
8 ax.set_xlabel('false positive rate')
9 ax.set_ylabel('True positive rate')
10 ax.set_title('PassiveAggressiveClassifier')
```

42]: Text(0.5, 1.0, 'PassiveAggressiveClassifier')



CONCLUSION

KEY FINDINGS AND CONCLUSIONS OF THE STUDY

Online hate, described as abusive language, aggression, cyberbullying, hatefulness and many others has been identified as a major threat on online social media platforms. Social media platforms are the most prominent grounds for such toxic behaviour.

From the above analysis the below mentioned results were achieved which depicts the chances and conditions of a comment being a hateful comment or a normal comment. With the increasing popularity of social media, more and more people consume feeds from social media and due differences they spread hate comments to instead of love and harmony. It has strong negative impacts on individual users and broader society.

LEARNING OUTCOMES OF THE STUDY IN RESPECT OF DATA SCIENCE

It is possible to classify the comments content into the required categories of Malignant and Non Malignant. However, using this kind of project an awareness can be created to know what is good and bad. It will help to stop spreading hatred among people.

LIMITATIONS OF THIS WORK AND SCOPE FOR FUTURE WORK

Machine Learning Algorithms like Decision Tree Classifier took enormous amount of time to build the model and Ensemble techniques were taking a lot more time thus I have not included Ensemble models.

Using Hyper-parameter tuning would have resulted in some more accuracy.

Every effort has been put on it for perfection but nothing is perfect and this project is of no exception. There are certain areas which can be enhanced. Comment detection is an emerging research area with few public datasets. So, a lot of works need to be done on this field.

THANK YOU