



Project Report on Malignant Comments Classifier



Submitted by:
Joshua Jonathan

Internship 31

ACKNOWLEDGEMENT

The following research papers helped me understand the topic of malignant comments, what are the different parameters which help define a comment as malignant, the process of building a malignant comments classifier & finally, helped me in my model building & predictions:

1. "Toxic Comment Classification"- Sara Zaheri, Jeff Leath, David Stroud

This paper presents a novel application of Natural Language Processing techniques to classify unstructured text into toxic and nontoxic categories. In the current century, social media has created many job opportunities and, at the same time, it has become a unique place for people to freely express their opinions. Meanwhile, among these users, there are some groups that are taking advantage of this framework and misuse this freedom to implement their toxic mindset (i.e insulting, verbal sexual harassment, threads, Obscene, etc.). The 2017 Youth Risk Behavior Surveillance System (Centers for Disease Control and Prevention) estimated that 14.9% of high school students were electronically bullied in the 12 months, prior to the survey. The primary result could be an Open Source model used by app developers in support of anti bullying efforts.

2. "Identification and Classification of Toxic Comment Using Machine Learning Methods" - P.Vidyullatha, Satya Narayan Padhya, Javvaji Geetha Priya, Kakarlapudi Srija, Sri Satyanjani Koppiseti

The increase in penetration of usage of internet services has increased exponentially in the past 4 months due to the ongoing pandemic, this has empowered an enormous number of dynamic new and old clients utilizing the web for different administration ranging from academic, entertainment, industrial, monitoring and the emergence of a new trend in the corporate life i.e work-from-home. Due to this sudden emergence of the crowd using the web, there has been an ascent in the number of mischievous persons too. Now it is the primary task of every online platform provider to keep the conversations constructive and inclusive. The best example can be referred to, can be twitter, a web-based media stage where people share their views. This platform has already drawn a lot of flak because of the spread of hate speech, insults, threats, and defamatory acts which becomes a challenge for many such online providers in regulating them. Thus, there is active research being conducted in the field of Toxic comment classification.

3. “A Machine Learning Approach to Comment Toxicity Classification” - Navoneel Chakraborty

Now-a-days, derogatory comments are often made by one another, not only in offline environments but also immensely in online environments like social networking websites and online communities. So, Identification combined with Prevention System in all social networking websites and applications, including all the communities, existing in the digital world is a necessity. In such a system, the Identification Block should identify any negative online behavior and should signal the Prevention Block to take action accordingly. This study aims to analyze any piece of text and detect different types of toxicity like obscenity, threats, insults and identity-based hatred. The labeled Wikipedia Comment Dataset prepared by Jigsaw is used for the purpose.

4. “An Automated Toxicity Classification on Social Media Using LSTM and Word Embedding” - Ahmad Alsharef, Karan Aggarwal, Sonia, Deepika Koundal, Hashem Alyami, and Darine Ameyed

The automated identification of toxicity in texts is a crucial area in text analysis since the social media world is replete with unfiltered content that ranges from mildly abusive to downright hateful. Researchers have found an unintended bias and unfairness caused by training datasets, which caused an inaccurate classification of toxic words in context. In this paper, several approaches for locating toxicity in texts are assessed and presented aiming to enhance the overall quality of text classification. General unsupervised methods were used depending on the state-of-art models and external embeddings to improve the accuracy while relieving bias and enhancing F1-score. Suggested approaches used a combination of long short-term memory (LSTM) deep learning model with Glove word embeddings and LSTM with word embeddings generated by the Bidirectional Encoder Representations from Transformers (BERT), respectively.

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 behavior.

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.

- **Conceptual Background of the Domain Problem**

In the past few years, it is seen that the cases related to social media hatred have increased exponentially. Social media is turning into a dark venomous pit for people nowadays. Online hate is the result of differences in opinion, race, religion, occupation, nationality etc. In social media the people spreading or involved in such kinds of activities use filthy language, 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.

• Review of Literature

Toxic comment classification is not clearly distinguishable from its related tasks. Besides looking at toxicity of online comments (Wulczyn et al., 2017; Georgakopoulos et al., 2018), related research includes the investigation of hate speech (Badjatiya et al., 2017; Burnap and Williams, 2016; Davidson et al., 2017; Gamback and Sikdar, 2017 Njagi et al, 2015 Schmidt and Wiegand, 2017; Vigna et al., 2017; Warner and Hirschberg, 2012), online harassment (Yin and Davison, 2009; Golbeck et al., 2017), abusive language (Mehdad and Tetreault, 2016; Park and Fung, 2017), cyberbullying (Dadvar et al., 2013; Dinakar et al., 2012; Hee et al., 2015; Zhong et al., 2016) and offensive language (Chen et al., 2012; Xiang et al., 2012). Each field uses different definitions for their classification, still similar methods can often be applied to different tasks. Besides traditional binary classification tasks, related work considers different aspects of toxic language, such as “racism” (Greevy and Smeaton, 2004; Waseem, 2016; Kwok and Wang, 2013) and “sexism” (Waseem and Hovy, 2016; Jha and Mamidi, 2017), or the severity of toxicity (Davidson et al., 2017; Sharma et al., 2018). These tasks are framed as multi-class problems, where each sample is labeled with exactly one class out of a set of multiple classes. Toxic comment identification is a supervised classification task and approached by either methods including manual feature engineering (Burnap and Williams, 2015; Mehdad and Tetreault, 2016; Waseem, 2016; Davidson et al., 2017; Nobata et al., 2016; Kennedy et al., 2017; Samghabadi et al., 2017; Robinson et al., 2018) or the use of (deep) neural networks (Ptaszynski et al., 2017; Pavlopoulos et al., 2017; Badjatiya et al., 2017; Vigna et al., 2017; Park and Fung, 2017; Gamback and Sikdar, 2017).

• **Motivation for the Problem Undertaken**

The main objective of this study is to investigate which method from a chosen set of machine learning techniques performs the best. So far, we have a range of publicly available models served through the Perspective API, including toxicity/malignant comments. But the current models still make errors, and they don't allow users to select which type of toxicity they are interested in finding. The project gives an insight to identify major factors that lead to cyberbullying and online abusive comments. The exposure to real world data and the opportunity to deploy my skills 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.

The main motivation was to classify the news in order to bring awareness and reduce unwanted chaos and make a good model which will help us to know such kinds of miscreants. 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.

• **ANALYTICAL PROBLEM FRAMING**

• **Mathematical/Analytical Modeling of the Problem**

We are provided with two different datasets. One for training and another one to test the efficiency of the model created using the training dataset. The training data provided here has both dependent and independent variables. As it is a multiclass problem it has 6 independent/target variables. Here the target variables are named "malignant", "highly malignant", "rude", "threat", "abuse" and "loathe". The label can be either 0 or 1, where 0 denotes a NO while 1 denotes a YES. There are various comments which have multiple labels. The first attribute is a unique ID associated with each comment. Clearly it is a binary classification problem as the target columns giving binary outputs and all independent variables have text so it is clear that it is a supervised machine learning problem where we can use the techniques of NLP and classification-based algorithms of Machine Learning. Here we will use NLP techniques like word tokenization, lemmatization, stemming and tfidf vectorizer then those processed data will be used to create best model using various classification based supervised ML algorithms like Logistic Regression, Multinomial NB, Extra Trees Classifier, XGBoost Classifier, LinearSVC, Decision Tree Classifier and Adaboost Classifier.

• Data Sources & their formats

Data set provided by Flip Robo was in the format of CSV (Comma Separated Values). The data set contains the training set, which has approximately 159571 samples and the test set which contains nearly 153164 samples.

• FEATURE INFORMATION:

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

• Data preprocessing done

- Importing necessary libraries and loading the train and test datasets.
- Checked some statistical information like shape, number of unique values present, info, null values, value counts, duplicated values etc.
- Checked for null values and did not find any null values. And removed Id.
- Done feature engineering and created new columns viz label: which contain both good and bad comments which is the sum of all the labels, comment_length: which contains the length of comment text.
- Visualized each feature using seaborn and matplotlib libraries by plotting categorical plots like pie plot, count plot, distribution plot and word cloud for each label.
- Done text pre-processing techniques like Removing Punctuations and other special characters, Splitting the comments into individual words, Removing Stop Words, Stemming and Lemmatization.
- Then created a new column as clean_length after cleaning the data. All these steps were done on both train and test datasets.
- Checked correlation using heatmap.
- After getting cleaned data use TF-IDF vectorizer. It'll help to transform the text data to feature vectors which can be used as input in our modeling. It is a common algorithm to transform text into numbers. It measures the originality of a word by comparing the frequency of appearance of a word in a document with the number of documents the words appear in. Mathematically, $TF-IDF = TF(t*d) * IDF(t,d)$
- Balanced the data using Randomoversampler method.
- Proceeded with Model Building.

• **Data Inputs- Logic- Output Relationships**

The dataset consists of a label and features. The features are independent and the label is dependent as the values of our independent variables change as our label varies.

- I checked the distribution of skewness using dist plots and used count plots to check the counts available in each column as a part of univariate analysis.
- To analyze the relation between features and label I have used many plotting techniques where I found some of the columns having strong relation with label.
- Got to know the sense of loud words in every label using word cloud which gives the words frequented in the labels.
- I have checked the correlation between the label and features using a heat map.

• **Libraries Used:**

- import pandas as pd
- import numpy as np
- import seaborn as sns
- import matplotlib.pyplot as plt
- import warnings; warnings.filterwarnings('ignore')
- import nltk
- import re
- import string
- from nltk.corpus import stopwords
- from wordcloud import WordCloud
- from nltk.tokenize import word_tokenize
- from nltk.stem import WordNetLemmatizer
- from sklearn.feature_extraction.text import TfidfVectorizer
- from wordcloud import WordCloud
- from sklearn.model_selection import train_test_split
- from imblearn.over_sampling import RandomOverSampler
- from collections import Counter
- from sklearn import metrics
- from sklearn.model_selection import cross_val_score
- from sklearn.model_selection import train_test_split

- from sklearn.metrics import classification_report, confusion_matrix
- from sklearn.metrics import roc_curve, accuracy_score, roc_auc_score, hamming_loss, log_loss.
- from sklearn.linear_model import LogisticRegression
- from sklearn.tree import DecisionTreeClassifier
- from sklearn.svm import LinearSVC
- from sklearn.naive_bayes import MultinomialNB
- from sklearn.ensemble import AdaBoostClassifier
- from xgboost import XGBClassifier
- from sklearn.ensemble import ExtraTreesClassifier
- from sklearn.model_selection import cross_val_score
- from sklearn import datasets
- from sklearn import model_selection
- from sklearn.metrics import plot_roc_curve
- from sklearn.model_selection import GridSearchCV
- import joblib

- **MODEL BUILDING & EVALUATION**

- **Identification of possible problem-solving approaches (methods)**

I have used both statistical and analytical approaches to solve the problem which mainly includes the pre-processing of the data and also used EDA techniques and heat map to check the correlation of independent and dependent features. Also, before building the model, I made sure that the input data was cleaned and scaled before it was fed into the machine learning models.

In this project there were 6 features which defines the type of comment like malignant, hate, abuse, threat, loathe but we created another feature named as “label” which is combined of all the above features and contains the labeled data into the format of 0 and 1 where 0 represents “NO” and 1 represents “Yes”.

I have used many classification algorithms and got the prediction results. By doing various evaluations I have selected the Logistic Regression Model as the best suitable algorithm to create our final model.

• Testing of Identified Approaches (Algorithms)

Since the target variable is categorical in nature, from this I can conclude that it is a classification type problem hence I have used the following classification algorithms. After the pre-processing and data cleaning I left with 10 columns including targets.

The algorithms used on training the data are as follows:

1. Logistic Regression Model
2. Decision Tree Classifier Model
3. Linear SVC Model
4. MultinomialNB Classifier Model
5. AdaBoost Classifier Model
6. Extreme Gradient Boosting Classifier (XGB) Model
7. Extra Trees Classifier Model

• Run and Evaluate selected models

• Classification Models:

Logistic Regression

```
from sklearn.linear_model import LogisticRegression

lg=LogisticRegression()
lg.fit(train_x,train_y)
pred_lg=lg.predict(x_test)

print("Accuracy score: ",accuracy_score(y_test,pred_lg))
print("Roc_auc_score: ",roc_auc_score(y_test,pred_lg))
print("Log loss : ",log_loss(y_test,pred_lg))
print("Hamming loss: ",hamming_loss(y_test,pred_lg))
print('Confusion matrix: \n',confusion_matrix(y_test,pred_lg))
print('Classification Report:\n ',classification_report(y_test,pred_lg))
```

```
Accuracy score: 0.9452707219251337
Roc_auc_score: 0.8943017052855605
Log loss : 1.8903122625429436
Hamming loss: 0.05472927807486631
Confusion matrix:
[[41210 1794]
 [ 826 4042]]
Classification Report:
              precision    recall  f1-score   support

     0       0.98         0.96         0.97         43004
     1       0.69         0.83         0.76          4868

   accuracy                   0.95         47872
  macro avg       0.84         0.89         0.86         47872
 weighted avg       0.95         0.95         0.95         47872
```

The Logistic Regression model gave us an Accuracy Score of 94.52 %.

Decision Tree Classifier

```
from sklearn.tree import DecisionTreeClassifier

dtc=DecisionTreeClassifier()
dtc.fit(train_x,train_y)
pred_dtc=dtc.predict(x_test)

print("Accuracy score: ",accuracy_score(y_test,pred_dtc))
print("Roc_auc_score: ",roc_auc_score(y_test,pred_dtc))
print("Log loss : ",log_loss(y_test,pred_dtc))
print("Hamming loss: ",hamming_loss(y_test,pred_dtc))
print('Confusion matrix: \n',confusion_matrix(y_test,pred_dtc))
print('Classification Report:\n ',classification_report(y_test,pred_dtc))
```

Accuracy score: 0.9280999331550802
Roc_auc_score: 0.838017972683667
Log loss : 2.4833754575676448
Hamming loss: 0.07190006684491979
Confusion matrix:
[[40901 2103]
 [1339 3529]]
Classification Report:

	precision	recall	f1-score	support
0	0.97	0.95	0.96	43004
1	0.63	0.72	0.67	4868
accuracy			0.93	47872
macro avg	0.80	0.84	0.82	47872
weighted avg	0.93	0.93	0.93	47872

The DecisionTree Classifier Model gave us an Accuracy Score of 92.88 %.

Linear SVC

```
from sklearn.svm import LinearSVC

lsvc=LinearSVC()
lsvc.fit(train_x,train_y)
pred_lsvc=lsvc.predict(x_test)

print("Accuracy score: ",accuracy_score(y_test,pred_lsvc))
print("Roc_auc_score: ",roc_auc_score(y_test,pred_lsvc))
print("Log loss : ",log_loss(y_test,pred_lsvc))
print("Hamming loss: ",hamming_loss(y_test,pred_lsvc))
print('Confusion matrix: \n',confusion_matrix(y_test,pred_lsvc))
print('Classification Report:\n ',classification_report(y_test,pred_lsvc))
```

Accuracy score: 0.9391084558823529
Roc_auc_score: 0.883949349994539
Log loss : 2.1031530493982213
Hamming loss: 0.06089154411764706
Confusion matrix:
[[40991 2013]
 [902 3966]]
Classification Report:

	precision	recall	f1-score	support
0	0.98	0.95	0.97	43004
1	0.66	0.81	0.73	4868
accuracy			0.94	47872
macro avg	0.82	0.88	0.85	47872
weighted avg	0.95	0.94	0.94	47872

The Linear SVC Model gave us an accuracy score of 93.96 %.

MultinomialNB Classifier

```
from sklearn.naive_bayes import MultinomialNB

mnb=MultinomialNB()
mnb.fit(train_x,train_y)
pred_mnb=mnb.predict(x_test)

print("Accuracy score: ",accuracy_score(y_test,pred_mnb))
print("Roc_auc_score: ",roc_auc_score(y_test,pred_mnb))
print("Log loss : ",log_loss(y_test,pred_mnb))
print("Hamming loss: ",hamming_loss(y_test,pred_mnb))
print('Confusion matrix: \n',confusion_matrix(y_test,pred_mnb))
print('Classification Report:\n ',classification_report(y_test,pred_mnb))
```

```
Accuracy score: 0.9112215909090909
Roc_auc_score: 0.8868266596820368
Log loss : 3.0663569153001258
Hamming loss: 0.08877840909090909
Confusion matrix:
[[39454 3550]
 [ 700 4168]]
Classification Report:
              precision    recall  f1-score   support

     0       0.98         0.92         0.95         43004
     1       0.54         0.86         0.66          4868

 accuracy          0.91         0.91         0.91         47872
 macro avg         0.76         0.89         0.81         47872
 weighted avg         0.94         0.91         0.92         47872
```

The MultinomialNB Classifier Model gave us an accuracy score of 91.08 %.

Ada Boost Classifier

```
from sklearn.ensemble import AdaBoostClassifier

abc=AdaBoostClassifier()
abc.fit(train_x,train_y)
pred_abc=abc.predict(x_test)

print("Accuracy score: ",accuracy_score(y_test,pred_abc))
print("Roc_auc_score: ",roc_auc_score(y_test,pred_abc))
print("Log loss : ",log_loss(y_test,pred_abc))
print("Hamming loss: ",hamming_loss(y_test,pred_abc))
print('Confusion matrix: \n',confusion_matrix(y_test,pred_abc))
print('Classification Report:\n ',classification_report(y_test,pred_abc))
```

```
Accuracy score: 0.9260319184491979
Roc_auc_score: 0.8083573869454118
Log loss : 2.554798580671883
Hamming loss: 0.07396808155080214
Confusion matrix:
[[41115 1889]
 [ 1652 3216]]
Classification Report:
              precision    recall  f1-score   support

     0       0.96         0.96         0.96         43004
     1       0.63         0.66         0.64          4868

 accuracy          0.93         0.93         0.93         47872
 macro avg         0.80         0.81         0.80         47872
 weighted avg         0.93         0.93         0.93         47872
```

The AdaBoost Classifier Model gave us an accuracy score of 92.61 %.

XGBoost Classifier

```
from xgboost import XGBClassifier

xgb=XGBClassifier()
xgb.fit(train_x,train_y)
pred_xgb=xgb.predict(x_test)

print("Accuracy score: ",accuracy_score(y_test,pred_xgb))
print("Roc_auc_score: ",roc_auc_score(y_test,pred_xgb))
print("Log loss : ",log_loss(y_test,pred_xgb))
print("Hamming loss: ",hamming_loss(y_test,pred_xgb))
print('Confusion matrix: \n',confusion_matrix(y_test,pred_xgb))
print('Classification Report:\n ',classification_report(y_test,pred_xgb))
```

Accuracy score: 0.9490516377005348
Roc_auc_score: 0.8523211366246949
Log loss : 1.759712950632643
Hamming loss: 0.05094836229946524
Confusion matrix:
[[41875 1129]
 [1310 3558]]
Classification Report:

	precision	recall	f1-score	support
0	0.97	0.97	0.97	43004
1	0.76	0.73	0.74	4868
accuracy			0.95	47872
macro avg	0.86	0.85	0.86	47872
weighted avg	0.95	0.95	0.95	47872

The XGBoost Classifier Model gave us an accuracy score of 94.90 %.
Extra Trees Classifier

```
from sklearn.ensemble import ExtraTreesClassifier

etc=ExtraTreesClassifier()
etc.fit(train_x,train_y)
pred_etc=etc.predict(x_test)

print("Accuracy score: ",accuracy_score(y_test,pred_etc))
print("Roc_auc_score: ",roc_auc_score(y_test,pred_etc))
print("Log loss : ",log_loss(y_test,pred_etc))
print("Hamming loss: ",hamming_loss(y_test,pred_etc))
print('Confusion matrix: \n',confusion_matrix(y_test,pred_etc))
print('Classification Report:\n ',classification_report(y_test,pred_etc))
```

Accuracy score: 0.9548587901069518
Roc_auc_score: 0.8135633195192252
Log loss : 1.5591286687917738
Hamming loss: 0.04514120989304813
Confusion matrix:
[[42614 390]
 [1771 3097]]
Classification Report:

	precision	recall	f1-score	support
0	0.96	0.99	0.98	43004
1	0.89	0.64	0.74	4868
accuracy			0.95	47872
macro avg	0.92	0.81	0.86	47872
weighted avg	0.95	0.95	0.95	47872

The Extra Trees Classifier Model gave us an accuracy score of 95.30 %.

From the above Classification Models, the highest accuracy score belongs to the Extra Trees Classifier, followed by the XGBoost Classifier and Logistic Regression Model.

Next, the Linear SVC Model followed by the AdaBoost Classifier and the Decision Tree Classifier. Lastly the MultinomialNB Classifier Model.

Cross Validation Scores

```
from sklearn.model_selection import cross_val_score
```

```
scr_lg=cross_val_score(lg,x,y,cv=5,scoring='accuracy')  
print("Cross validation score of this model is: ",scr_lg.mean())
```

Cross validation score of this model is: 0.9558441024494281

```
scr_dtc=cross_val_score(dtc,x,y,cv=5,scoring='accuracy')  
print("Cross validation score of this model is: ",scr_dtc.mean())
```

Cross validation score of this model is: 0.9420320736843184

```
scr_lsvc=cross_val_score(lsvc,x,y,cv=5,scoring='accuracy')  
print("Cross validation score of this model is: ",scr_lsvc.mean())
```

Cross validation score of this model is: 0.959547784553488

```
scr_mnb=cross_val_score(mnb,x,y,cv=5,scoring='accuracy')  
print("Cross validation score of this model is: ",scr_mnb.mean())
```

Cross validation score of this model is: 0.946982848930392

```
scr_abc=cross_val_score(abc,x,y,cv=5,scoring='accuracy')  
print("Cross validation score of this model is: ",scr_abc.mean())
```

Cross validation score of this model is: 0.9456981559137111

```
scr_xgb=cross_val_score(xgb,x,y,cv=5,scoring='accuracy')  
print("Cross validation score of this model is: ",scr_xgb.mean())
```

Cross validation score of this model is: 0.9532559184423365

```
scr_etc=cross_val_score(etc,x,y,cv=5,scoring='accuracy')  
print("Cross validation score of this model is: ",scr_etc.mean())
```

Cross validation score of this model is: 0.9572792002665083

From the above Cross Validation Scores, the highest CV score belongs to the LinearSVC model, followed by the Extra Trees Classifier & Logistic Regression Model. Next the XGBoost Classifier model, the MultinomialNB Classifier and the AdaBoost Classifier Model.

Lastly, the Decision Tree Classifier.

• Hyper Parameter Tuning:

Since the Accuracy score and the cross validation score of the Logistic Regression Model are good and the AUC score is the highest among others, we shall consider this model for hyper parameter tuning. We shall use GridSearchCV for hyper parameter tuning.

```
from sklearn.model_selection import GridSearchCV
```

```
parameters={
    'C': [0.2,0.3,0.4],
    'penalty': ['l1', 'l2'],
    'solver':['newton-cg','lbfgs'],
    'multi_class':['auto','ovr']}
grid_lg = GridSearchCV(lg, param_grid = parameters, cv = 4, scoring='accuracy')
```

```
grid_lg.fit(train_x,train_y)
```

```
GridSearchCV(cv=4, estimator=LogisticRegression(),
             param_grid={'C': [0.2, 0.3, 0.4], 'multi_class': ['auto', 'ovr'],
                         'penalty': ['l1', 'l2'],
                         'solver': ['newton-cg', 'lbfgs']},
             scoring='accuracy')
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
grid_lg.best_params_
```

```
{'C': 0.4, 'multi_class': 'auto', 'penalty': 'l2', 'solver': 'lbfgs'}
```

```
Final_Model= LogisticRegression(C=0.4,penalty='l2',solver='lbfgs',multi_class='auto')
Final_Model.fit(train_x,train_y)
pred = Final_Model.predict(x_test)

print("Accuracy score: ",accuracy_score(y_test,pred))
print("Roc_auc_score: ",roc_auc_score(y_test,pred))
print("Log loss : ",log_loss(y_test,pred))
print("Hamming loss: ",hamming_loss(y_test,pred))
print('Confusion matrix: \n',confusion_matrix(y_test,pred))
print('Classification Report:\n ',classification_report(y_test,pred))
```

```
Accuracy score: 0.945855614973262
Roc_auc_score: 0.8923501373857027
Log loss : 1.8701098870892066
Hamming loss: 0.05414438502673797
Confusion matrix:
[[41263 1741]
 [ 851 4017]]
Classification Report:

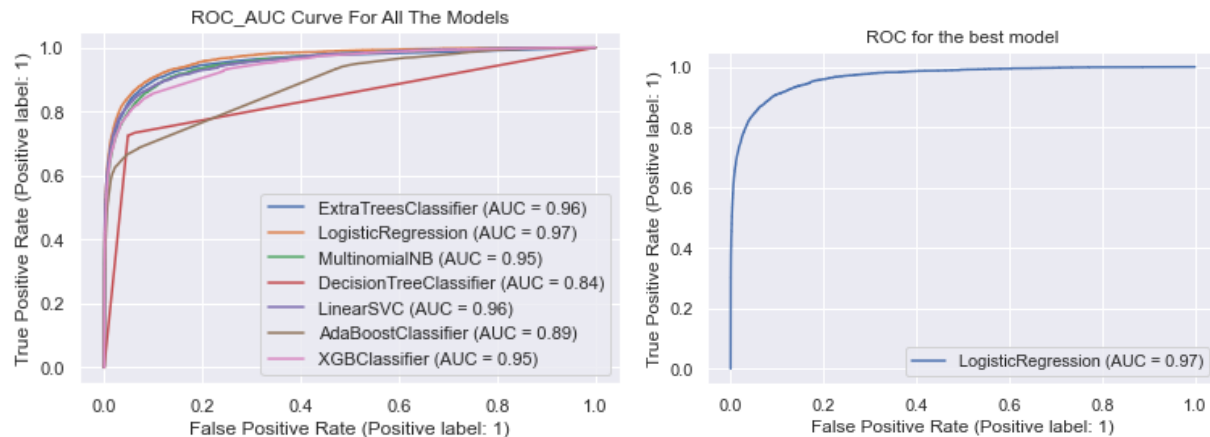
```

	precision	recall	f1-score	support
0	0.98	0.96	0.97	43004
1	0.70	0.83	0.76	4868
accuracy			0.95	47872
macro avg	0.84	0.89	0.86	47872
weighted avg	0.95	0.95	0.95	47872

- After multiple tries with hyper parameter tuning, the highest accuracy score obtained was 94.56 %.

• ROC-AUC Curve:

I have generated the ROC Curve for all the models and for the best model and compared it with AUC. The AUC score for my final model was 97%.



• Saving the final model and predicting the saved model

```
# Saving the model using .pkl
import joblib
joblib.dump(Final_Model, "Malignant_Comments_Classifier.pkl")
```

```
['Malignant_Comments_Classifier.pkl']
```

```
# Loading the final model
model = joblib.load('Malignant_Comments_Classifier.pkl')
```

```
# Predicting the values for test data after loading trained model
Predictions = model.predict(x1)
Predictions
```

```
array([0, 0, 0, ..., 0, 0, 0])
```

```
# Adding the predicted values to test dataframe
test_data['Predicted_Values'] = Predictions
test_data
```

	id	comment_text	comment_length	clean_length	Predicted_Values
0	00001cee341fdb12	yo bitch ja rule is more succesful then you ll...	367	358	0
1	0000247867823ef7	from rfc the title is fine as it is imo	50	39	0
2	00013b17ad220c46	sources zawe ashton on lapland	54	30	0
3	00017563c3f7919a	if you have a look back at the source the info...	205	200	0
4	00017695ad8997eb	i don t anonymously edit articles at all	41	40	0
...
153159	ffcd0960ee309b5	i totally agree this stuff is nothing but too ...	60	55	0
153160	fffd7a9a6eb32c16	throw from out field to home plate does it get...	198	182	0
153161	fffd9a9e8d6fafa9e	okinotorishima categories i see your changes a...	423	365	0
153162	fffe8f1340a79fc2	one of the founding nations of the eu germany ...	502	460	0
153163	ffffce3fb183ee80	stop already your bullshit is not welcome here...	141	124	0

- **Key Metrics for success in solving problem under consideration**

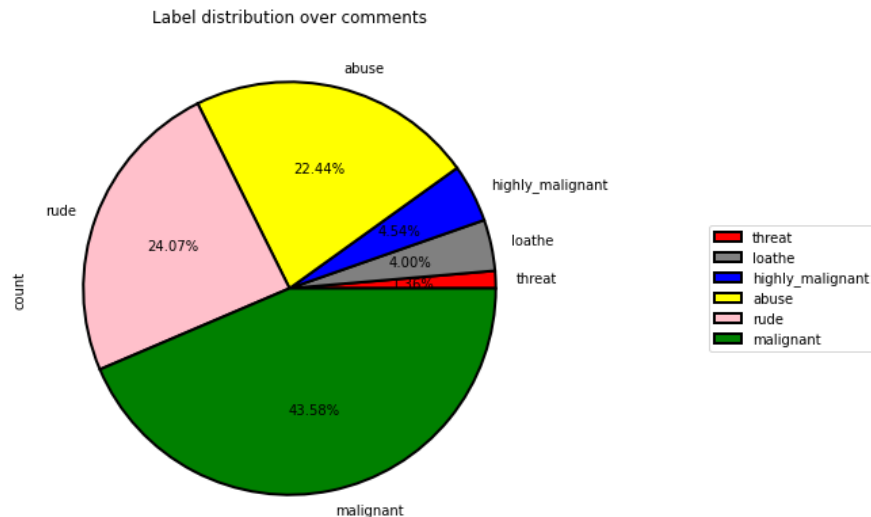
I have used the following metrics for evaluation:

- Accuracy score, which is used when the True Positives and True negatives are more important.
- Confusion Matrix, which is a matrix used to determine the performance of the classification models for a given set of test data.
- Classification report, which is used to measure the quality of predictions from a classification algorithm.
- Cross Validation Score, which is a technique for evaluating ML models by training several ML models on subsets of the available input data and evaluating them on the complementary subset of the data.
- Precision is the degree to which repeated measurements under the same conditions are unchanged. It is the amount of information that is conveyed by a value. It refers to the data that is correctly classified by the classification algorithm.
- Recall is how many of the true positives were recalled (found). Recall refers to the percentage of data that is relevant to the class.
- F1 Score is used to express the performance of the machine learning model (or classifier). It gives the combined information about the precision and recall of a model.
- This means a high F1-score indicates a high value for both recall and precision.
- Roc Auc Curve, the Receiver Operator Characteristic (ROC) curve, is an evaluation metric for binary classification problems. It is a probability curve that plots the TPR against FPR at various threshold values. The Area Under Curve (AUC) is the measure of the ability of a classifier to distinguish between classes and is used as a summary of the ROC curve.
- Log Loss is the most important classification metric based on probabilities. Log Loss is the negative average of the log of corrected predicted probabilities for each instance.
- Hamming Loss is the fraction of wrong labels to the total number of labels. In multi-class classification, hamming loss is calculated as the hamming distance between y_{true} and y_{pred} . In multi-label classification, hamming loss penalizes only the individual labels.

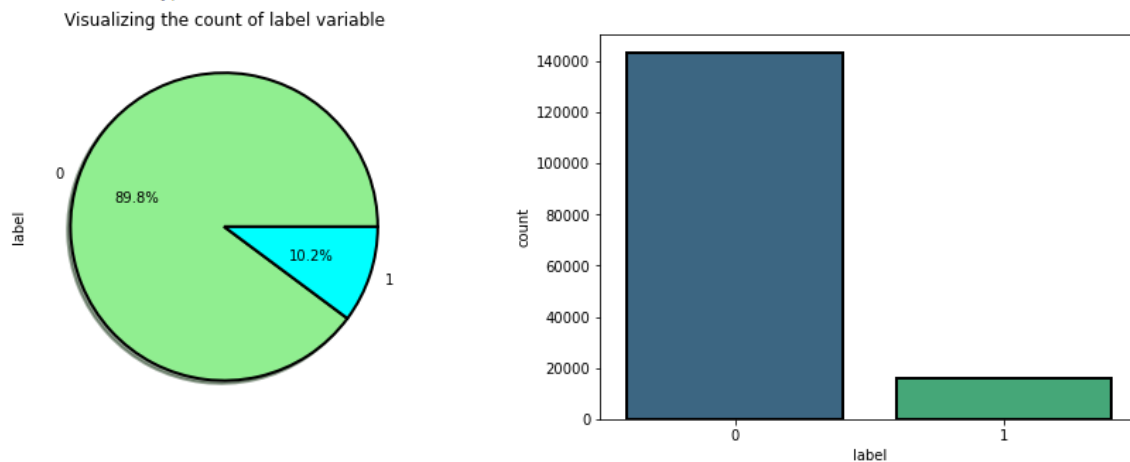
• Visualizations

I used pandas profiling to get the over viewed visualization on the pre-processed data. I have used various plots for visualization such as pie charts, count plots, distribution plots to help understand the features. I have used word clouds to get the sense of loud words in the labels.

• PLOTS:

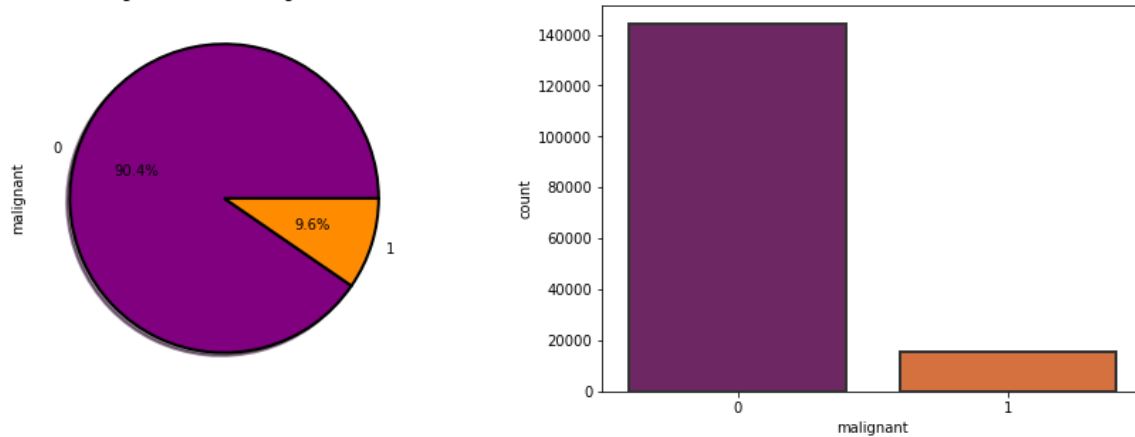


From the pie chart we can notice approximately 43.58 % of the comments are malignant, 24.07 % Of the comments are rude and 22.44 % are of abuse. The count of malignant comments is high compared to other types of comments and the count of threat comments are very less.



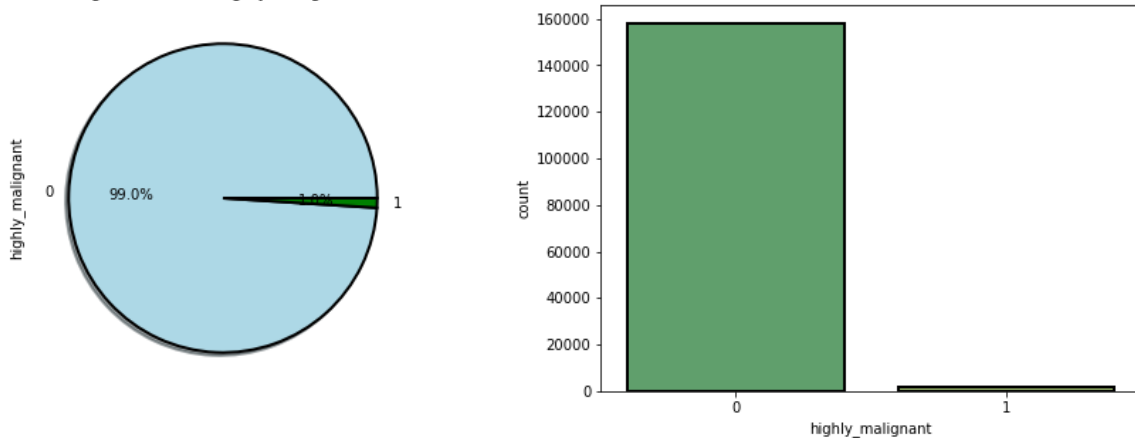
From the above plots we can observe the count of negative comments are high compared to the non negative comments. Here around 89.8% of the comments are turned out to be negative comments and only 10.2 % of them are considered to be positive or neutral comments. We can also observe the data imbalance issue here, we need to balance the data.

Visualizing the count of malignant variable



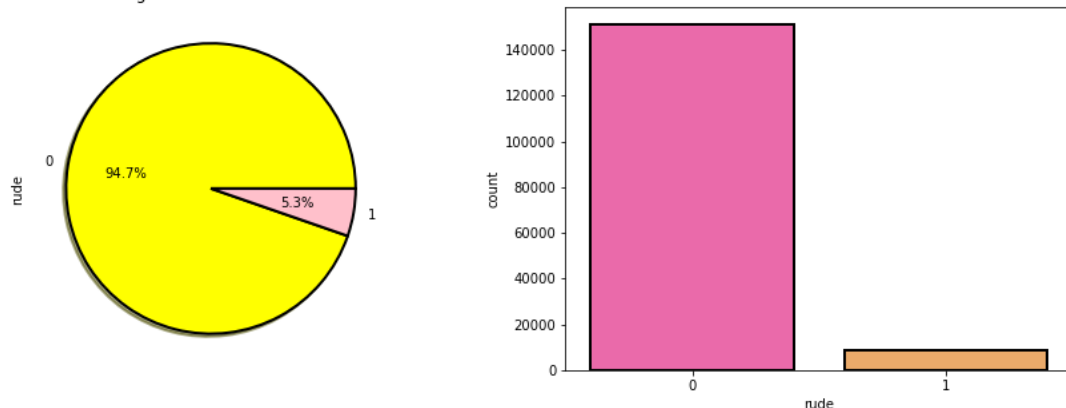
From the above plots we can observe the count of malignant comments is high compared to non malignant comments. That is around 90.4 % of the comments are malignant and only 9.6 % of the comments are good.

Visualizing the count of highly malignant variable



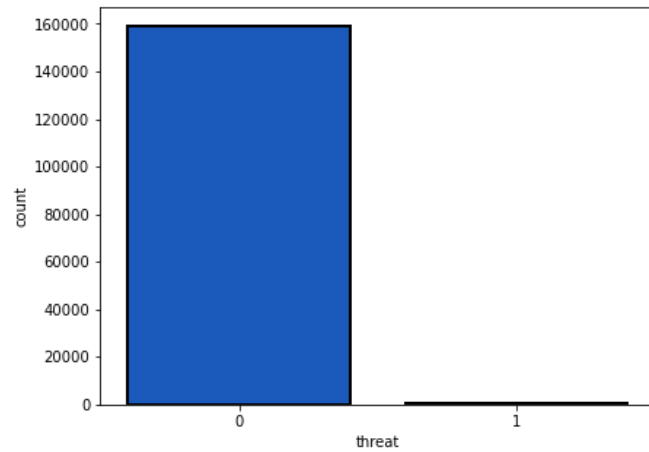
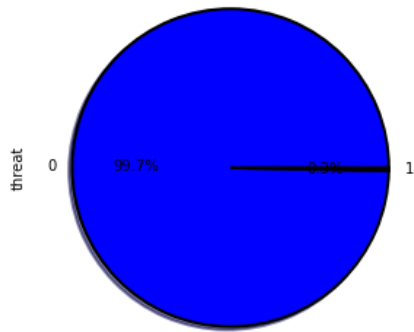
From the plot we can observe that the count of highly malignant comments is very high, which is about 99 % and only 1 % of the comments are normal.

Visualizing the count of rude variable



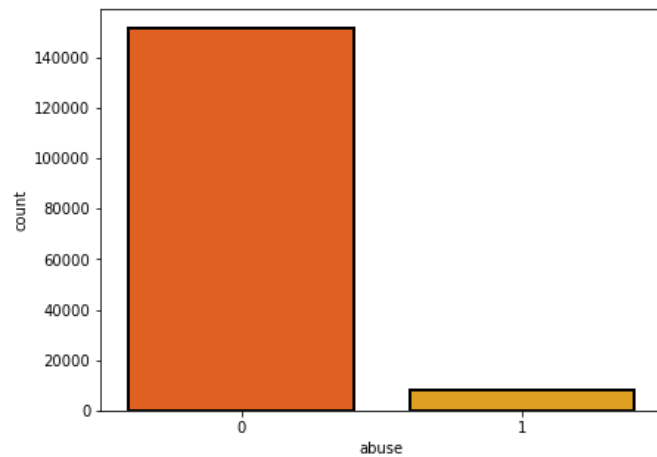
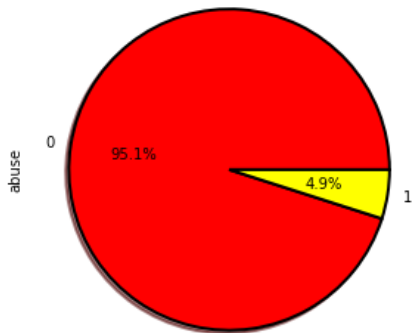
The number of rude comments is high compared to normal comments. Around 94.7 % of the comments fall into rude and remaining comments are considered to be normal comments.

Visualizing the count of threat variable



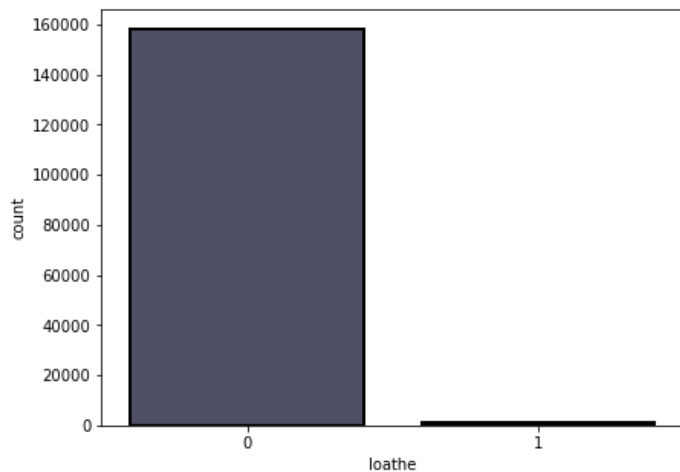
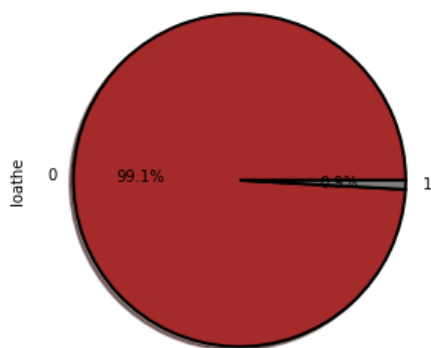
In the above visualization also, 99.7 % of the comments are threatening and only 0.3 % of the comments look normal.

Visualizing the count of abuse variable

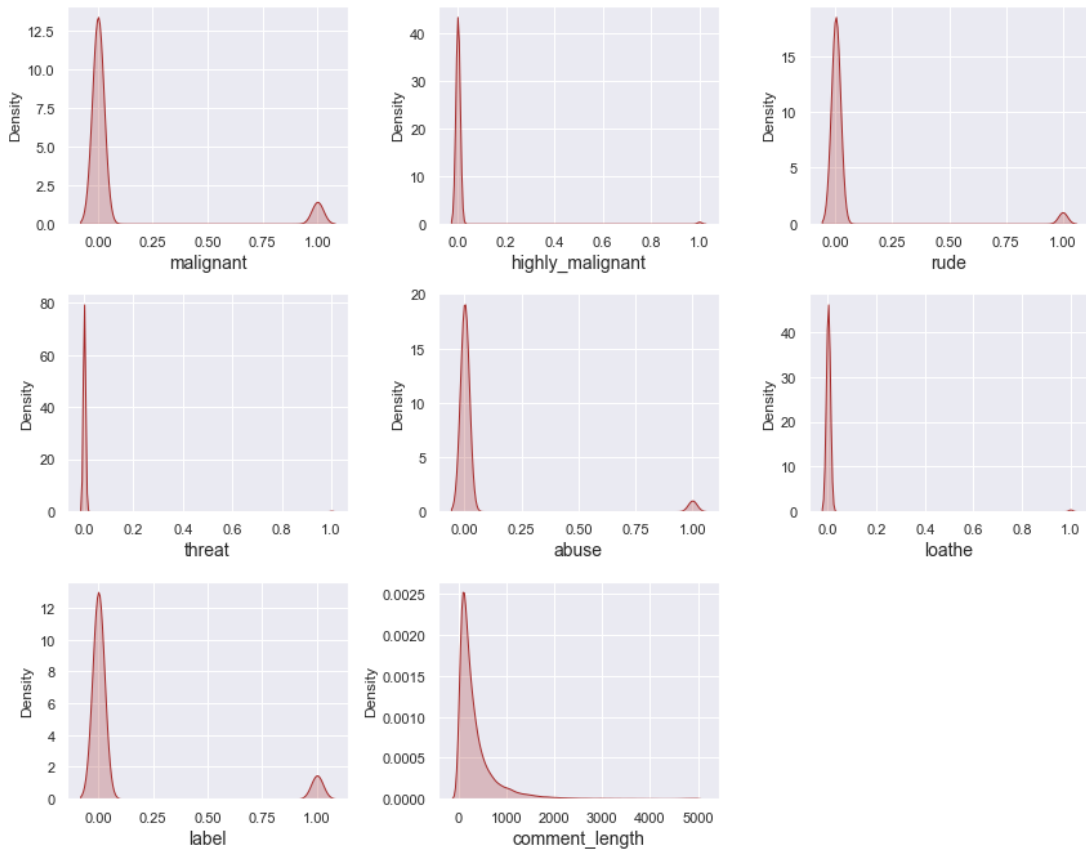


The count of abusing type comments is high which is 95.1 % and only 4.9 % of the comments are normal.

Visualizing the count of loathe variable



The count of loathe is high (99.1 %) compared to normal (0.9 %) text comments.



From the distribution plots we can notice that all the columns are skewed to the right except the comment_label column. Since all the columns are categorical in nature there is no need to remove skewness and outliers in any of the columns.

WORD CLOUDS:



Pre-processing: The dataset should be cleaned and scaled to build the ML models to get good predictions. I have performed a few NLP text processing steps which I have already mentioned in the pre-processing steps where all the important features are present in the dataset and ready for model building.

Model building: After cleaning and processing data, I performed a train test split to build the model. I have built multiple classification models to get the accurate accuracy score, and evaluation metrics like precision, recall, confusion matrix, f1 score, log loss, hamming loss. I got the Logistic Regression Model as the best model which gives a 94.52% accuracy score. I checked the cross-validation score ensuring there will be no overfitting. After tuning the best model, I got a 94.56% accuracy score. Finally, I saved my final model and got the good predictions results for the test dataset.

• CONCLUSION

• Key Findings and Conclusions of the Study

- 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 instead of love and harmony. It has strong negative impacts on individual users and broader society.
- From this dataset we were able to understand the impact of various malignant comments, the different types of malignant comments and how to identify them.
- In this study, we have used multiple machine learning models to predict malignant comments and identify them.
- We have gone through the data analysis by performing feature engineering, finding the relation between features and the label through visualizations. And got the important features and we used these features to predict the defaulters' rate by building ML models.
- After training the model we checked CV score to overcome the overfitting issue.
- Performed hyper parameter tuning, on the best model. We have also got good prediction results.

- **Learning Outcomes of the Study in respect of Data Science**

While working on this project I learned many things and gained new techniques and ways to deal with uncleaned text data. Found how to deal with multiple target features. Tools used for visualizations give a better understanding of the dataset. We have used a lot of algorithms and find that in the classification problem where we have only two labels, the Logistic Regression Model gives better results compared to others.

- **Limitations of this work and Scope for Future Work**

- **LIMITATIONS:**

This project was amazing to work on, it creates new ideas to think about but there were some limitations in this project like unbalanced dataset. There are certain areas which can be enhanced.

- **FUTURE WORK:**

- The potential future work for this project can focus on performance and error analysis of the model as lots of comments are misclassified into the hate category.
- Previous work has achieved success using various algorithms on data in the English language but in future, we can consider having data in regional languages.

Thank You!