



# **MALIGNANT COMMENT CLASSIFICATION PROJECT**

Submitted by:  
**SHALINI ROY**

*FlipRobo SME:*  
**KHUSHBOO GARG**

# ACKNOWLEDGMENT

I want to thank the Flip Robo Technologies team for giving me the chance to work with this dataset during my internship. It enabled me to develop my analytical abilities. The entire Data Trained team deserves a great thank you.

## **Bibliography:**

Reference used in this project:

- ❖ GitHub Notes & Repository.
- ❖ Hands-on Machine learning with scikit learn and tensor flow by Aurelien Geron.
- ❖ Various Kaggle and GitHub projects.
- ❖ Analytics Vidya's different papers on Data Science.
- ❖ SCIKIT Learn Library Documentation.
- ❖ Toxic Comment Classification by Nupur Baghel.

# INTRODUCTION

## **Business Problem Framing**

People can now freely express themselves online due to the growth of social media. However, concurrently, this has led to the emergence of conflict and hatred, making online spaces hostile for users. Online hatred has been identified as a significant hazard on social media websites, including abusive language, aggressiveness, cyberbullying, hatefulness, and many more. The most common environment for such harmful behaviour is social media platforms. 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**

It has been observed that the number of incidents involving online hatred has skyrocketed in recent years. Nowadays, social media is becoming a hole of darkness and poison for people. Differentiation in viewpoint, race, religion, occupation, nationality, and other factors are the cause of online hatred. People that engage in or disseminate these types of activities on social media sometimes use vulgar language, aggressive behaviour, offensive imagery, and other tactics to insult and seriously harm those on the opposing side. One of the main issues at the moment is this. On-line social media platforms have been highlighted as being particularly vulnerable to online hate, which is defined as abusive language, aggression, cyberbullying, hatefulness, insults, personal assaults, provocation, racism, sexism, threats, or toxicity. For a brighter future, these kinds of behaviours ought to be stopped.

## **Review of Literature**

Users now leave a lot of comments on various social networks, news websites, and forums. Certain comments are harmful or abusive. Since it is impractical to manually monitor so many comments, the majority of systems employ some form of machine learning models to automatically identify harmful content. In this study, we used machine learning techniques to conduct an in-depth analysis of the state-of-the-art in the classification of toxic comments. First, we looked into the papers' publication dates, locations, and levels of maturity. Each major study's data set, evaluation metric, machine learning techniques, types of toxicity, and comment language were all examined.

## **Motivation for the Problem Undertaken**

The project is provided to me by Flip Robo Technologies as a part of the internship programme (Internship Batch No-31). This problem is a real world dataset. The exposure of this data gives me the opportunity to locate my skills in solving a real time problem. It is the primary motivation to solve this problem. We have a lot of possibilities here, but not as many specific solutions. The primary aim is to create a prototype for an online hate and abuse comment classifier that can be used to categorise hateful and offensive comments in order to limit their ability to spread intolerance and cyberbullying. This study aims to analyse and predicting malignant comment when using Classification Model like Logistic Regression, Random Forest, Decision Tree, Gradient Boosting, Extra Tree, Ada Boost Classification algorithms. Thus, the purpose of this study is to grow the knowledge of Classification methods in machine learning fields. Those are the different factors to undertaken the problem for study purposes.

# Analytical Problem Framing

## Mathematical/ Analytical Modelling of the Problem

The goal of this project is to predict malignant comments in social media. We would use a classification method, which is a sort of supervised learning. Although performing a classification seems more logical given that we have between 5 and 6 classes to forecast. We will only do classification in this case. Filtering the words is necessary to avoid overfitting because the dataset only contains one feature.

This project have two big set of data one is training and another is testing. Throughout the project's classification phase, we would first eliminate email addresses, phone numbers, web addresses, spaces, and other terms with stops, in order to calculate the regularisation parameter. We also used TFID to transform the tokens from the train papers into vectors so that the machine could carry out additional processing in order to further enhance our models.

Here 6 different algorithm are used and final model is chosen by best AUC-ROC score and accuracy score.

## Data Sources and their formats

There are two set of data, training and testing. Training dataset has 159571 rows and 8 columns in ot her hand the testing dataset has 153164 rows and 2 columns. The model will train with the help of tra ining dataset. It has 6 integer datatype and 2 object datatype. All integer datatype are actually binary in nature.

Later the training dataset is divided into two parts, training and testing. After determine the proper model, the model is applied to predict the target variable for the test dataset.

```
In [16]: print('No. of Rows of train dataset : ',data.shape[0])  
         print('No. of Columns of train dataset : ',data.shape[1])
```

```
No. of Rows of train dataset : 159571  
No. of Columns of train dataset : 8
```

```
In [17]: data.columns.to_series().groupby(data.dtypes).groups
```

```
Out[17]: {int64: ['malignant', 'highly_malignant', 'rude', 'threat', 'abuse', 'loathe'], object: ['id', 'comment_text']}
```

```
In [19]: print('No. of Rows of test dataset : ',data_test.shape[0])  
         print('No. of Columns of test dataset : ',data_test.shape[1])
```

```
No. of Rows of test dataset : 153164  
No. of Columns of test dataset : 2
```

## Data Pre-processing Done:

### 1 Feature Engineering:

'--', 'null', 'NA', ' ' are not present in the training and testing dataset.

```
In [20]: data.isin(['--','null','NA',' ']).sum().any()
```

```
Out[20]: False
```

```
In [21]: data.isnull().sum()
```

```
Out[21]: id                0
comment_text             0
malignant                0
highly_malignant         0
rude                     0
threat                   0
abuse                    0
loathe                   0
dtype: int64
```

```
In [22]: data_test.isin(['--','null','NA',' ']).sum().any()
```

```
Out[22]: False
```

```
In [23]: data_test.isnull().sum()
```

```
Out[23]: id                0
comment_text             0
dtype: int64
```

### 2 Drop unnecessary columns:

Let's drop the unnecessary column 'id' from both dataset.

```
In [25]: data.drop('id',axis=1,inplace=True)
data_test.drop('id',axis=1,inplace=True)
```

### 3 Calculate length before cleaning of 'comment\_text' column of training and testing dataset:

Let's calculate the comment length before cleaning.

```
In [51]: data['length_before_cleaning'] = data['comment_text'].str.len()
data.head()
```

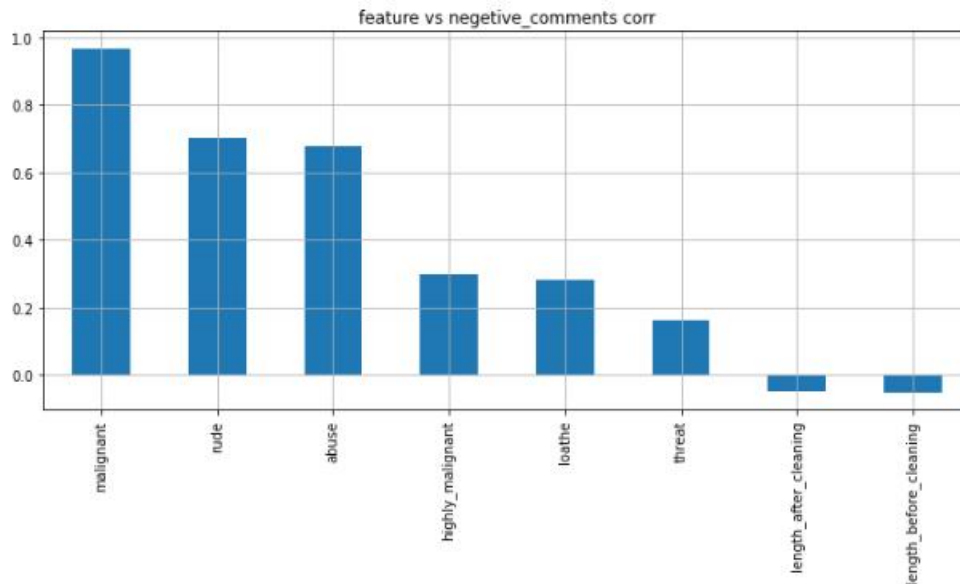
```
Out[51]:
```

	comment_text	malignant	highly_malignant	rude	threat	abuse	loathe	length_before_cleaning
0	Explanation\nWhy the edits made under my usern...	0	0	0	0	0	0	264
1	D'aww! He matches this background colour I'm s...	0	0	0	0	0	0	112
2	Hey man, I'm really not trying to edit war. It...	0	0	0	0	0	0	233
3	"\nMore!\nI can't make any real suggestions on ...	0	0	0	0	0	0	622
4	You, sir, are my hero. Any chance you remember...	0	0	0	0	0	0	67

## 4 Correlation:

Observations of the correlation:

1. Very obviously length before cleaning and length after cleaning are highly correlated with each other's.
2. Malignant and rude comments are highly correlated with target.
3. The two feature length before cleaning and length after cleaning are negatively correlated with target.



## 5 Make new column named negative\_comments:

Here, the comments are 6 different type. If anyone is present the comment tagged as a malignant comment (bad/ negative) comment. So let's make a new column named.

If it is 0= Good comment, 1= malignant comment.

```
In [74]: target_data = data[['malignant', 'highly_malignant', 'rude', 'threat', 'abuse', 'loathe', 'length_after_cleaning', 'length_before_cleaning']]

data['negative_comments'] = data[['malignant', 'highly_malignant', 'rude', 'threat', 'abuse', 'loathe', 'length_after_cleaning', 'length_before_cleaning']]
print(data['negative_comments'].value_counts())
data['negative_comments'] = data['negative_comments'] > 0
data['negative_comments'] = data['negative_comments'].astype(int)
print(data['negative_comments'].value_counts())
```

0	143346
1	6360
3	4209
2	3480
4	1760
5	385
6	31

Name: negative\_comments, dtype: int64

0	143346
1	16225

Name: negative\_comments, dtype: int64

### **Data Inputs- Logic- Output Relationships**

We can see in the correlation that every features are correlated with each other and also they are highly correlated with target variable label.

### **State the set of assumptions (if any) related to the problem under consideration**

No such assumptions are taken for this case.

### **Hardware and Software Requirements and Tools Used**

Processor: Intel(R) Core(TM) i3-5005U CPU @ 2.00GHz 2.00 GHz

RAM: 4.00 GB

System Type: 64-bit operating system, x64-based processor

Window: Windows 10 Pro

Anaconda – Jupyter Notebook

Libraries Used –

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')

from sklearn.model_selection import train_test_split
```

For Word-Cloud the following libraries are used.

```
#Importing Required Libraries
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
```

Except this, different libraries are used for machine learning model building from sklearn.

# Model/s Development and Evaluation

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

In this problem Classification-based machine learning algorithm like logistic regression can be used. Removed any excess spaces, changed the email addresses subject line to a phone number that is probably wise, etc. For building an appropriate ML model before implementing classification algorithms, data is split in training & test data using `train_test_split`. Then different statistical parameter like accuracy score, confusion matrix, classification report, precision, recall etc. are determined for every algorithm. Hyper parameter tuning is performed to get the accuracy score much higher and accurate than earlier.

Then the best model is chosen from 6 different algorithm.

## **Testing of Identified Approaches (Algorithms)**

Total 7 algorithms used for the training and testing are:

1. Logistic Regression
2. Decision Tree Classifier
3. Gradient Boosting Classifier
4. Random Forest Classifier
5. Extra Trees Classifier
6. Ada Boost Classifier

## **Key Metrics for success in solving problem under consideration:**

From metrics module of sklearn library import `classification_report`, `accuracy_score`, `confusion_matrix`, `classification_report` and `f1_score`. From `model_selection` also, we use `cross_val_score`. Those are the matrices use to validate the model's quality. Let's discuss every metrics shortly.

- ❖ **Classification report:** It is a performance evaluation metric in machine learning which is used to show the precision, recall, F1 Score, and support score of your trained classification model
- ❖ **Accuracy score:** It is used when the True Positives and True negatives are more important. Accuracy can be used when the class distribution is similar.
- ❖ **Confusion Matrix:** It is a table that is used in classification problems to assess where errors in the model were made. The rows represent the actual classes the outcomes should have been. While the columns represent the predictions we have made. Using this table it is easy to see which predictions are wrong.
- ❖ **Precision:** It can be seen as a measure of quality. If the precision is high, an algorithm returns more relevant results than irrelevant ones.
- ❖ **Recall:** The recall is calculated as the ratio between the numbers of Positive samples correctly classified as Positive to the total number of Positive samples.
- ❖ **F1 Score:**  $F1 = 2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$



### Run and Evaluate selected models

First find the best random state of train\_test\_split to get best accuracy. Here the random state is 88.

Then after splitting the data into 4 different part and check the shape of the data.

```
print('Training feature shape:',x_train.shape)
print('Training target shape:',y_train.shape)
print('Test feature shape:',x_test.shape)
print('Test target shape:',y_test.shape)
```

```
Training feature shape: (119678, 51153)
Training target shape: (119678,)
Test feature shape: (39893, 51153)
Test target shape: (39893,)
```

### *Logistic Regression:*

```
from sklearn.linear_model import LogisticRegression
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size = 0.25, random_st
log = LogisticRegression()
log.fit(x_train, y_train)
y_pred = log.predict(x_test)

print('accu score : ', accuracy_score(y_test, y_pred))
print('cof_mat:\n ', confusion_matrix(y_test, y_pred))
print('classification report:\n ', classification_report(y_test, y_pred))
print("-----")
print("-----")
print('training score : ', log.score(x_train, y_train))
print('testing score : ', log.score(x_test, y_test))
```

```
accu score : 0.957762013385807
cof_mat:
[[35763  144]
 [ 1541 2445]]
classification report:
              precision    recall  f1-score   support

     0       0.96       1.00       0.98       35907
     1       0.94       0.61       0.74        3986

 accuracy          0.96          0.96          0.96       39893
 macro avg          0.95          0.80          0.86       39893
 weighted avg       0.96          0.96          0.95       39893

-----
-----
training score : 0.9591236484566921
testing score : 0.957762013385807
```

In this way accuracy score is determined for each 6 different classification model.

## Decision Tree Classifier:

The accuracy score, confusion matrix and classification report after using Decision tree Classifier is as follows.

```
from sklearn.tree import DecisionTreeClassifier
clf = DecisionTreeClassifier()
clf.fit(x_train, y_train)

y_pred = clf.predict(x_test)

print('accu score : ', accuracy_score(y_test, y_pred))
print("\n")
print('cof_mat: ', confusion_matrix(y_test, y_pred))
print("\n")
print('classification report: \n\n', classification_report(y_test, y_pred))

print("-----")
print("-----")

print('training score : ', clf.score(x_train, y_train))
print('testing score : ', clf.score(x_test, y_test))
```

accu score : 0.9441255358082872

cof\_mat: [[34814 1093]  
[ 1136 2850]]

classification report:

	precision	recall	f1-score	support
0	0.97	0.97	0.97	35907
1	0.72	0.72	0.72	3986
accuracy			0.94	39893
macro avg	0.85	0.84	0.84	39893
weighted avg	0.94	0.94	0.94	39893

-----  
-----

training score : 0.9994735874596835  
testing score : 0.9441255358082872

## Gradient Boosting Classifier:

The accuracy score, confusion matrix and classification report after using Gradient Boosting Classifier is as follows.

```
accu score : 0.942070037349911
```

cof\_mat: [[35816 91]  
[ 2220 1766]]

classification report:

	precision	recall	f1-score	support
0	0.94	1.00	0.97	35907
1	0.95	0.44	0.60	3986
accuracy			0.94	39893
macro avg	0.95	0.72	0.79	39893
weighted avg	0.94	0.94	0.93	39893

-----  
-----

training score : 0.9393372215444776  
testing score : 0.942070037349911

### Random Forest Classifier:

The accuracy score, confusion matrix and classification report after using Random Forest Classifier is as follows.

```
accu score : 0.9590404331586995
```

```
cof_mat: [[35606  301]
 [ 1333 2653]]
```

```
classification report:                precision    recall  f1-score   support

     0       0.96       0.99       0.98       35907
     1       0.90       0.67       0.76       3986

 accuracy                   0.96       39893
 macro avg                  0.93       0.83       0.87       39893
 weighted avg              0.96       0.96       0.96       39893
```

```
-----
training score : 0.9994568759504671
testing score  : 0.9590404331586995
```

### Extra Trees Classifier:

The accuracy score, confusion matrix and classification report after using Extra Trees Classifier is as follows.

```
accu score : 0.9589401649412178
```

```
cof_mat: [[35552  355]
 [ 1283 2703]]
```

```
classification report:                precision    recall  f1-score   support

     0       0.97       0.99       0.98       35907
     1       0.88       0.68       0.77       3986

 accuracy                   0.96       39893
 macro avg                  0.92       0.83       0.87       39893
 weighted avg              0.96       0.96       0.96       39893
```

```
-----
training score : 0.9994735874596835
testing score  : 0.9589401649412178
```

## Ada Boost Classifier:

The accuracy score, confusion matrix and classification report after using Ada Boost Classifier is as follows.

```
accu score : 0.9477101245832602
```

```
cof_mat: [[35629  278]
 [ 1808  2178]]
```

```
classification report:                precision    recall  f1-score   support

      0      0.95      0.99      0.97    35907
      1      0.89      0.55      0.68    3986

 accuracy          0.95    39893
  macro avg       0.92    0.77    0.82    39893
  weighted avg    0.95    0.95    0.94    39893

-----
training score : 0.944818596567456
testing score  : 0.9477101245832602
```

As per 6 different model, for Random Forest the accuracy score is highest among all the models. But the difference between training and testing is large. But for Logistic Regression the difference between training and testing is very small as well as it also gives good accuracy.

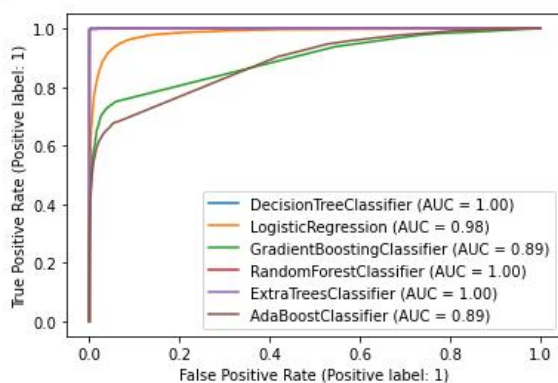
## AUC- ROC Curve:

```
from sklearn.metrics import plot_roc_curve

disp = plot_roc_curve(clf, x_train, y_train)

plot_roc_curve(log, x_train, y_train, ax=disp.ax_)
plot_roc_curve(gbdt, x_train, y_train, ax=disp.ax_)
plot_roc_curve(rf, x_train, y_train, ax=disp.ax_)
plot_roc_curve(etc, x_train, y_train, ax=disp.ax_)
plot_roc_curve(ada, x_train, y_train, ax=disp.ax_)

plt.show()
```

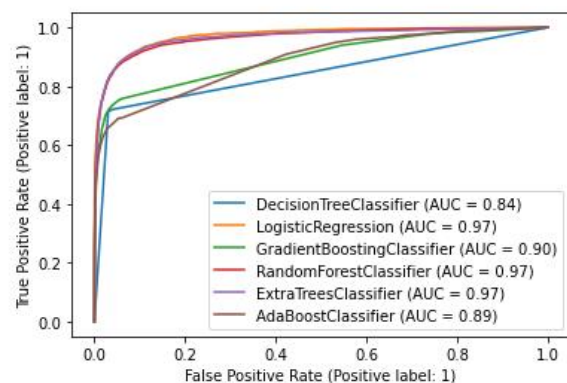


```
from sklearn.metrics import plot_roc_curve

disp = plot_roc_curve(clf, x_test, y_test)

plot_roc_curve(log, x_test, y_test, ax=disp.ax_)
plot_roc_curve(gbdt, x_test, y_test, ax=disp.ax_)
plot_roc_curve(rf, x_test, y_test, ax=disp.ax_)
plot_roc_curve(etc, x_test, y_test, ax=disp.ax_)
plot_roc_curve(ada, x_test, y_test, ax=disp.ax_)

plt.show()
```



## Hyper Parameter Tuning:

Here for Random Forest, Logistic Regression and Extra tree classifier, the AUC score is same. Here also we take Logistic Regression for final model. So it is the final model for this dataset.

```
from sklearn.model_selection import GridSearchCV
grid = dict(solver=['newton-cg', 'lbfgs', 'liblinear'],penalty=['l2','l1'], C=[1
grid_log = GridSearchCV(estimator=log, param_grid= grid,refit = True, verbose =
grid_log.fit(x_train, y_train)
print('best params : ', grid_log.best_params_)
```

Best params: {'C': 1.0, 'penalty': 'l1', 'solver': 'liblinear'}

Here accuracy score is slightly improved after using hyper parameter tuning. First, accuracy score was 0.957762013385807, but after applying hyper parameter tuning it is 0.9624495525530795.

## Final Model:

For final model the target variable is as follows after using Logistic Regression.

```
-----
accu score : 0.9624495525530795
cof_mat:

[[35625  282]
 [ 1216 2770]]
classification report:
              precision    recall  f1-score   support

     0       0.97       0.99       0.98       35907
     1       0.91       0.69       0.79       3986

   accuracy          0.96       39893
  macro avg       0.94       0.84       0.88       39893
weighted avg       0.96       0.96       0.96       39893

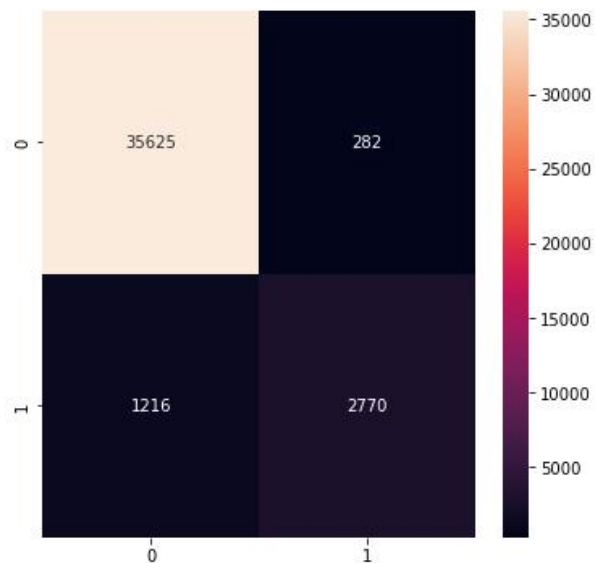
-----
-----
log loss: 1.2969521601797698
training score : 0.9624910175637962
testing score : 0.9624495525530795
```



### Confusion Matrix:

```
conf = confusion_matrix (y_test, y_pred)

fig , ax = plt.subplots(figsize=(6,6))
sns.heatmap(conf, annot = True, fmt = ".0f")
plt.show()
```



### Load the model:

Let's save the model using pickle for future use. Then see the actual and predicted value of 6 random sample.

```
import pickle
pickle.dump(grid_log_best, open("Malignant_Classification_model", "wb"))
load_Malignant_Classification_model= pickle.load(open("Malignant_Classification_m

y_pred = load_Malignant_Classification_model.predict(x_test)

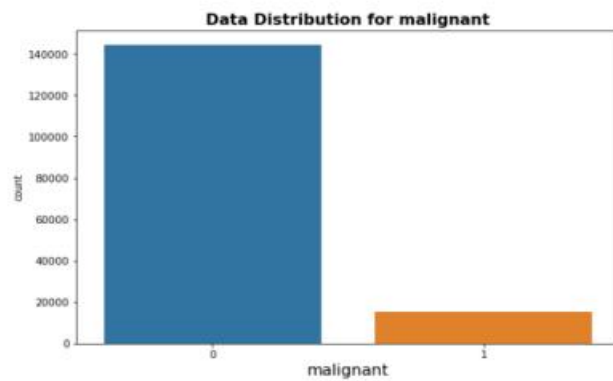
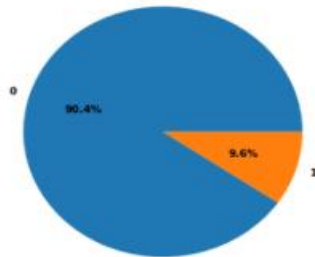
y_test = np.array(y_test)
data_prediction_by_model = pd.DataFrame()
data_prediction_by_model["Predicted Values"] = y_pred
data_prediction_by_model["Actual Values"] = y_test
data_prediction_by_model.sample(n=6)
```

	Predicted Values	Actual Values
38353	0	0
20720	0	0
27201	0	0
33658	1	1
372	0	0
11147	0	0

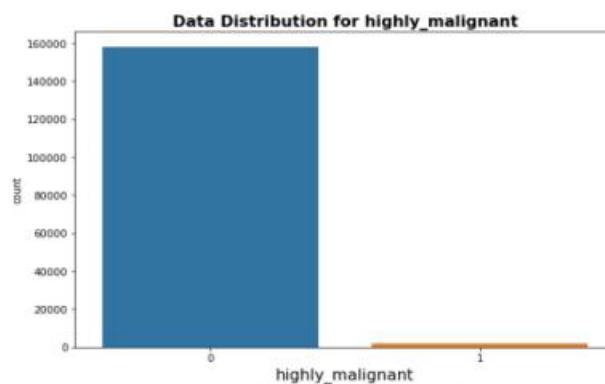
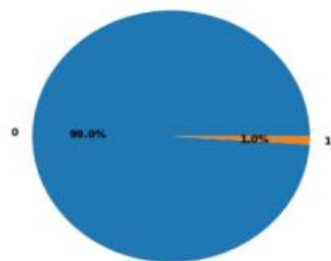
## Visualizations:

Let's start the observation exploration of feature analysis.

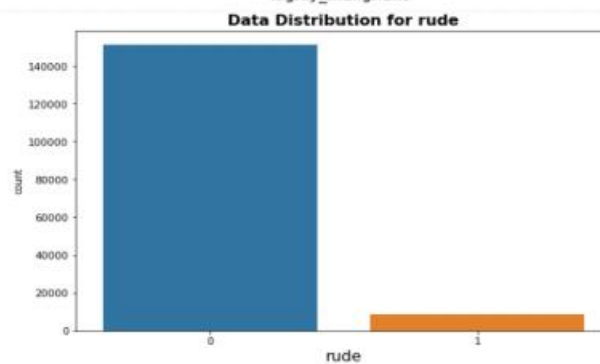
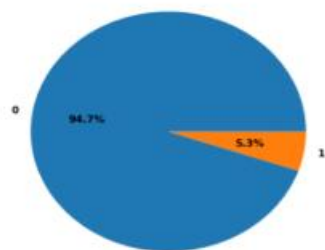
Data Distribution for malignant



Data Distribution for highly\_malignant

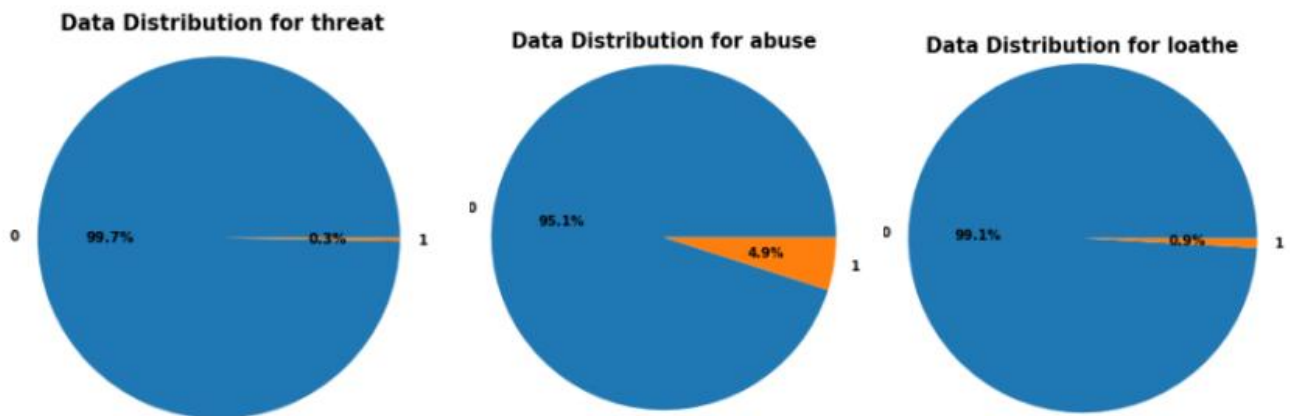


Data Distribution for rude



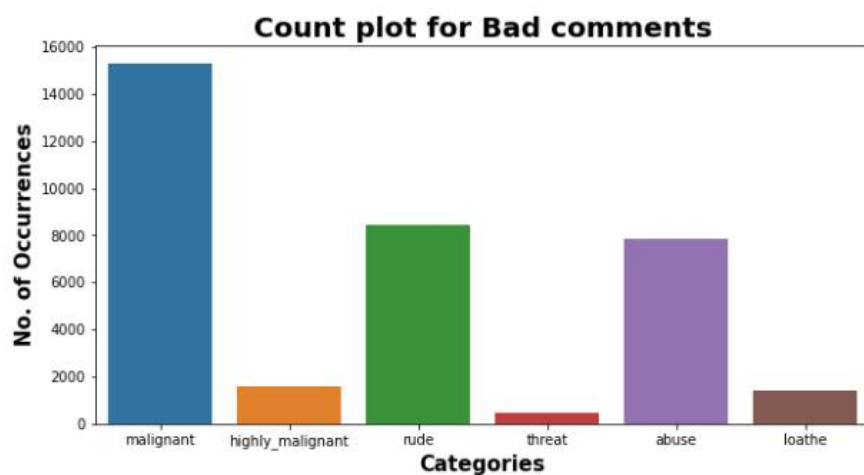
## Observations:

- For malignant comment distribution around 10% comment is malignant while 90% are good comments.
- For highly\_malignant comment distribution around 1% comment is highly\_malignant while 99% are good comments.
- For rude comment distribution around 5% comment is rude while 95% are good comments.



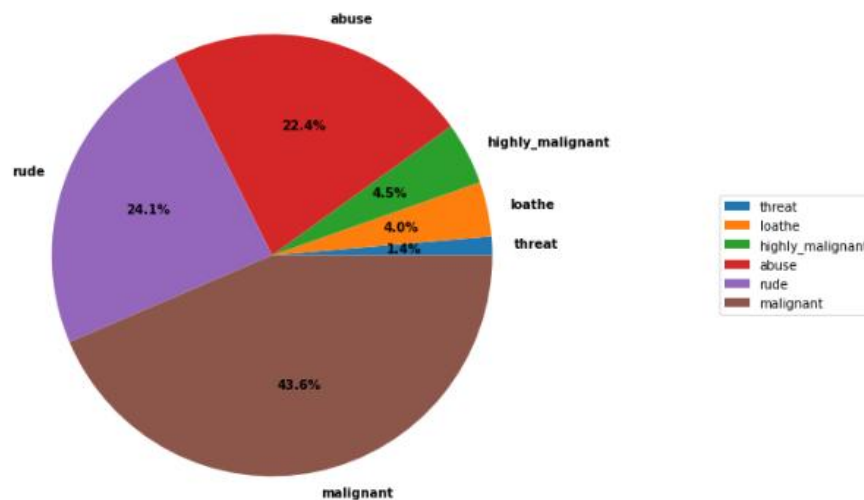
#### Observations:

- For threat comment distribution only 0.3% comment is highly\_malignant while 99.7% are not threatening comments.
- For abuse comment distribution around 5% comment is rude while 95% are good comments.
- For loathe comment distribution around 1% comment is rude while 99% are not loathe comments.





**Pie plot for Bad comments**

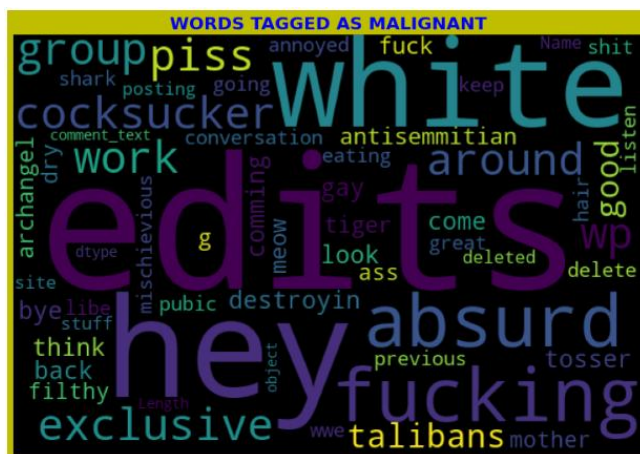


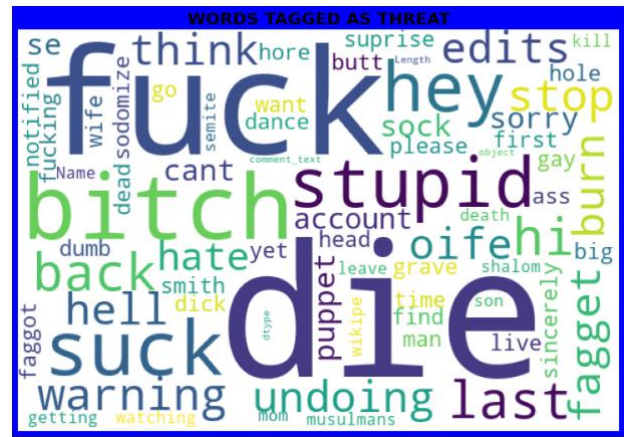
#### Observations:

- The maximum negative comments comes with Malignant in nature followed by rude categories.
- Very few comments comes with threatening nature.
- Total percentage of negative comment is 10.2% while good comment is 89.8%.
- Around 90% comments are good while rest 10% comments are Negative in nature.
- Out of total negative comments around 43.58% are malignant in nature followed by 24.07% are rude comments.

#### Word Cloud for different Feature:

We can see from the word clouds above that small texts are given less weight in their respective comment types than large texts are.





**Observations:**

- From word cloud of malignant comments, it is clear that it mostly consists of words like edits, hey, white, fucking, absurd, piss, cocksucker, Taliban etc.
- From word cloud of highly malignant comments, it is clear that it mostly consists of words like fuck, stupid, fucking, stupid, cocksucker, crow, piss, bitch, around, asshole etc.
- From word cloud of rude comments, it is clear that it mostly consists of words like shit, fucking, stuff, fucked, white, absurd, piece etc.
- From word cloud of threat comments, it is clear that it mostly consists of words like die, bitch, fuck, suck, stupid, back, hey, hi, back, last etc.



**Observations:**

- From word cloud of abuse comments, it is clear that it mostly consists of words like edits, white, ass, stuff, shit, piss, fucking, cocksucker, antisemmitian, gay etc.
- From word cloud of abuse comments, it is clear that it mostly consists of words like fuck, gay, jew, kill, antisemmitian, think etc.

## Interpretation of the Results

After all the pre-processing steps, the dataset is ready to train machine learning models. All unnecessary words from comment text are deleted as they might give overfitting problem as well as it also could increase the time complexity. Now apply this dataset on different ML Classification Model (as discussed on part 3.4 - 'Run and Evaluate selected models') and check the best model for this particular dataset.

# CONCLUSION

## **Key Findings and Conclusions of the Study**

Here, we observed the various detrimental effects that toxic or damaging social media remarks have on society. The ability to quickly and effectively identify remarks as hazardous could have a wide range of positive effects while also reducing the negative ones. We have also seen how easily accessible algorithms can be used in this way to handle this difficulty. It was shown in our particular investigation that a logistic regression solution offers a significant improvement in classification compared to any other approach.

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

Data cleansing is one of the most crucial phases; I attempted to make comments shorter and included all the relevant keywords in it. The power of visualisation is beneficial for converting data into a graphical representation; it helps me to comprehend what the data is trying to communicate. In this dataset, I utilised a variety of methods to find the best result and preserve that model. Logistic regression is the most effective algorithm.

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

Additionally, the following studies are examples that might be taken into account for future work in this field:

- We offer the following strategy to enhance NLP classifiers: Convolutional neural networks (CNN) and Support vector clustering (SVC) are two additional algorithms that can be used to enhance the performance of existing classifiers. In the present study, the issue was reduced to two classes, although it is worthwhile to pursue the primary objective of six classes of remarks.
- For text processing and text classification, we also advocate the use of SVM. To achieve the best results, a grid search is necessary for hyper-parameter optimization.