



Malignant-Comments-Classifer

Submitted by:

ARGHA KR MAITI

ACKNOWLEDGMENT

I would like to express my thanks of gratitude to mentors [Swathi Mahasewth](#) as well as Flip Robo who gave me the golden opportunity to do this wonderful project on the topic [Malignant_Comments_Classifier](#), which also helped me in doing a lot of research and I came to know about so many new things. I am really very thankful to them.

I am making this project to increase my knowledge.

INTRODUCTION

Problem Statement:

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.

Objective:

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

Firstly, we will start by importing required libraries and databases.

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.neighbors import KNeighborsClassifier
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.ensemble import AdaBoostClassifier, RandomForestClassifier
10
11 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
12 from sklearn.model_selection import train_test_split
13 from sklearn.model_selection import cross_val_score
14
15 from nltk.stem import WordNetLemmatizer
16 import nltk
17 from nltk.corpus import stopwords
18 import string
19
20 import joblib
21 import warnings
22 warnings.filterwarnings('ignore')
```

```
train=pd.read_csv(r"C:\Users\User\Downloads\Malignant Comments Classifier Project\malignant_train.csv")
train=pd.DataFrame(data=train)
train
```

	id	comment_text	malignant	highly_malignant	rude	threat	abuse	loathe
0	0000997932d777bf	Explanation\nWhy the edits made under my usern...	0	0	0	0	0	0
1	000103f0d9cfb60f	D'aww! He matches this background colour I'm s...	0	0	0	0	0	0
2	000113f07ec002fd	Hey man, I'm really not trying to edit war. It...	0	0	0	0	0	0
3	0001b41b1c6bb37e	"\nMore\nI can't make any real suggestions on ...	0	0	0	0	0	0
4	0001d958c54c6e35	You, sir, are my hero. Any chance you remember...	0	0	0	0	0	0
...
159566	ffe987279560d7ff	"::::And for the second time of asking, when ...	0	0	0	0	0	0
159567	fea4adeee384e90	You should be ashamed of yourself \n\nThat is ...	0	0	0	0	0	0
159568	fee36eab5c267c9	Spitzer \n\nUmm, theres no actual article for ...	0	0	0	0	0	0
159569	fff125370e4aaaf3	And it looks like it was actually you who put ...	0	0	0	0	0	0
159570	fff46fc426af1f9a	"\nAnd ... I really don't think you understand...	0	0	0	0	0	0

159571 rows × 8 columns

Here is the list of all columns and size of dataset:

```
1 train.columns
```

```
Index(['id', 'comment_text', 'malignant', 'highly_malignant', 'rude', 'threat',
      'abuse', 'loathe'],
      dtype='object')
```

```
1 train.shape
```

```
(159571, 8)
```

Let's check the datatype of all columns:

```
1 train.dtypes
```

```
id                object
comment_text      object
malignant         int64
highly_malignant  int64
rude              int64
threat            int64
abuse             int64
loathe            int64
dtype: object
```

We can see both type of columns numerical and object type.

Then with the help of describe we will take a glimpse of data:

```
1 train.describe(include='all').T
```

	count	unique	top	freq	mean	std	min	25%	50%	75%	max
id	159571	159571	a2ae72dce45c1d4b	1	NaN	NaN	NaN	NaN	NaN	NaN	NaN
comment_text	159571	159571	"destroy me"" but will"	1	NaN	NaN	NaN	NaN	NaN	NaN	NaN
malignant	159571	NaN	NaN	NaN	0.0958445	0.294379	0	0	0	0	1
highly_malignant	159571	NaN	NaN	NaN	0.00999555	0.0994771	0	0	0	0	1
rude	159571	NaN	NaN	NaN	0.0529482	0.223931	0	0	0	0	1
threat	159571	NaN	NaN	NaN	0.00299553	0.0546496	0	0	0	0	1
abuse	159571	NaN	NaN	NaN	0.0493636	0.216827	0	0	0	0	1
loathe	159571	NaN	NaN	NaN	0.00880486	0.0934205	0	0	0	0	1

We can see in feature Id and comment_text all values are unique. Andd for all other features named malignant, highly_malignant, rude, threat, abuse and loathe we can see only 2 unique values 0 and 1 hence we getting 0 as min and 1 as max for these features.

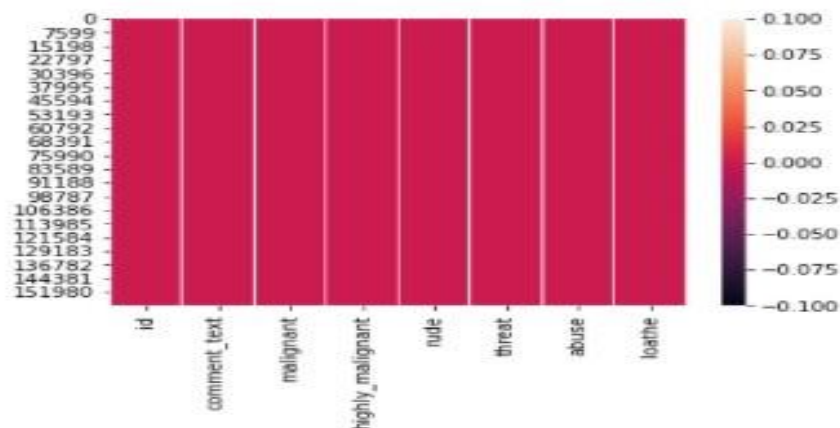
We can see count is all same for all variable, let's check for null values now:

```
1 train.isnull().sum()
```

```
id                0
comment_text      0
malignant         0
highly_malignant  0
rude              0
threat            0
abuse             0
loathe            0
dtype: int64
```

```
1 #visual repretation with the help of heatmap
2 sns.heatmap(train.isnull())
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1d91126c6d0>
```

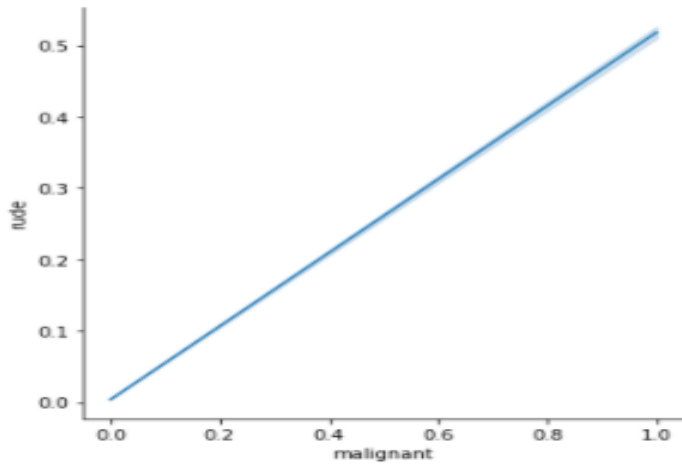


We can see id feature also has all unique values like serial number hence we dont require feature for future analysis

```
train.drop('id', axis=1, inplace= True)|
```

```
1 sns.relplot(x='malignant',y='rude',data=train,kind='line')
```

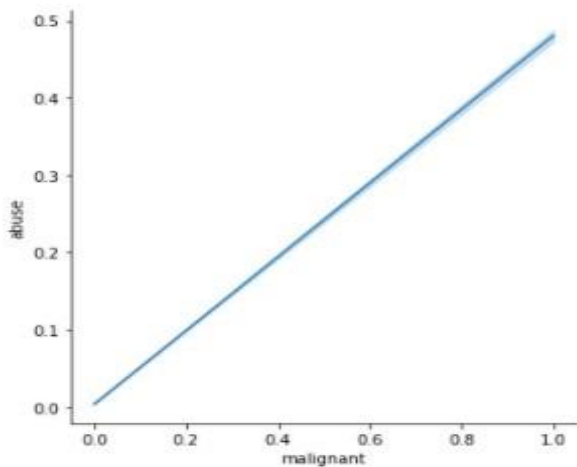
```
<seaborn.axisgrid.FacetGrid at 0x1d9112b16d0>
```



We can see positive relation in features malignant and rude.

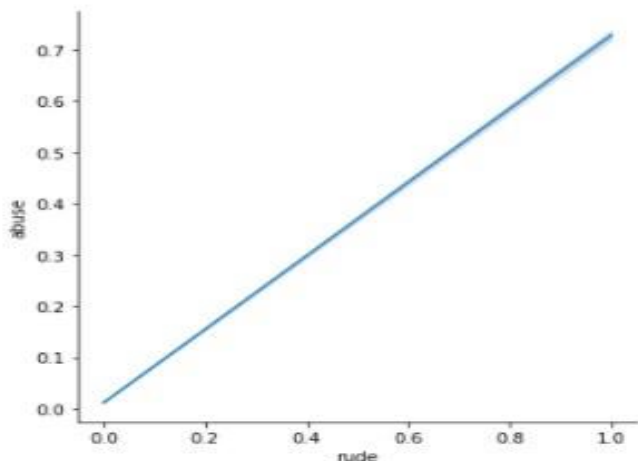
```
1 sns.relplot(x='malignant',y='abuse',data=train,kind='line')
```

```
<seaborn.axisgrid.FacetGrid at 0x1d9116a5220>
```



We can see positive relation in features malignant and abuse.

```
1 sns.relplot(x='rude',y='abuse',data=train,kind='line')
<seaborn.axisgrid.FacetGrid at 0x1d917354940>
```



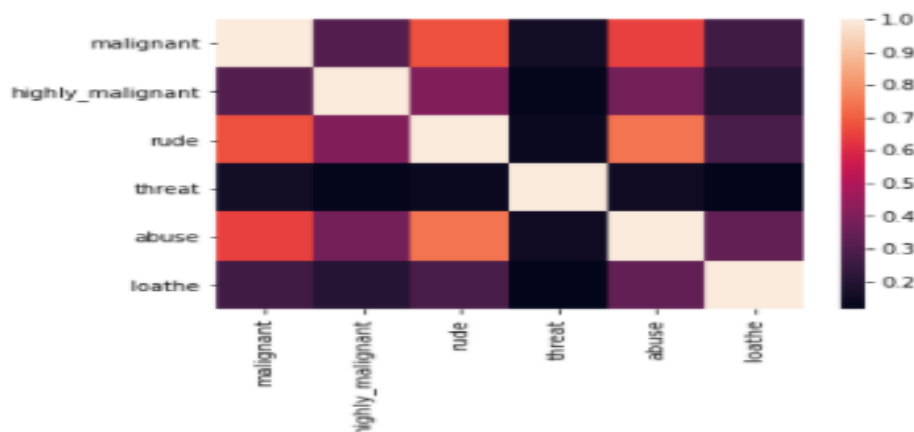
We can see positive relation in features rude and abuse.

Checking correlation in dataset:

```
2 print(train.corr())
3 print(sns.heatmap(train.corr()))
```

	malignant	highly_malignant	rude	threat	abuse	loathe
malignant	1.000000	0.308619	0.676515	0.157058	0.647518	0.266009
highly_malignant	0.308619	1.000000	0.403014	0.123601	0.375807	0.201600
rude	0.676515	0.403014	1.000000	0.141179	0.741272	0.286867
threat	0.157058	0.123601	0.141179	1.000000	0.150022	0.115128
abuse	0.647518	0.375807	0.741272	0.150022	1.000000	0.337736
loathe	0.266009	0.201600	0.286867	0.115128	0.337736	1.000000

```
loathe
malignant 0.266009
highly_malignant 0.201600
rude 0.286867
threat 0.115128
abuse 0.337736
loathe 1.000000
AxesSubplot(0.125,0.125;0.62x0.755)
```



In the above representation we can see dark or purple color which says low positive relationship of features and where we have orange going towards white is symbol of good positive relationship.

```
1 cat_col=[]
2 for i in train:
3     if train[i].nunique() <= 2:
4         cat_col.append(i)
5
6 print(cat_col)
```

```
['malignant', 'highly_malignant', 'rude', 'threat', 'abuse', 'loathe']
```

Though all the above features is on numeric type but it has only unique values in it. 0 represents 'No' and 1 represents 'Yes'.

```

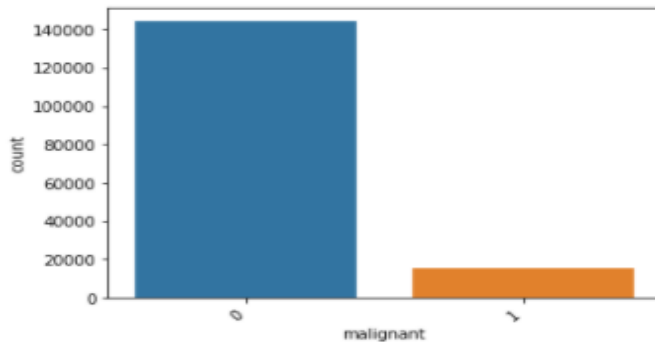
1 for i in cat_col:
2     plt.figure()
3     a=sns.countplot(train[i])
4     print(train[i].value_counts())
5     a.set_xticklabels(a.get_xticklabels(), rotation=40, ha="right")
6     plt.show()

```

```

0    144277
1     15294
Name: malignant, dtype: int64

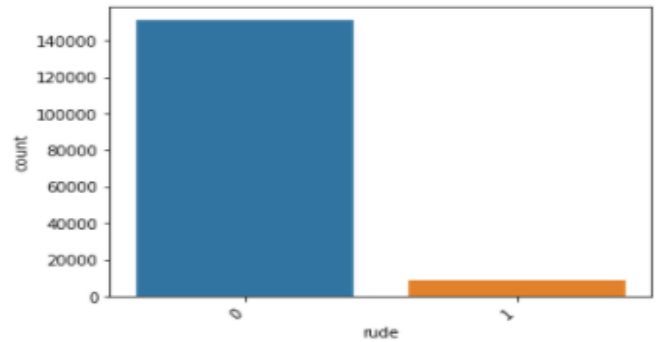
```



```

0    151122
1     8449
Name: rude, dtype: int64

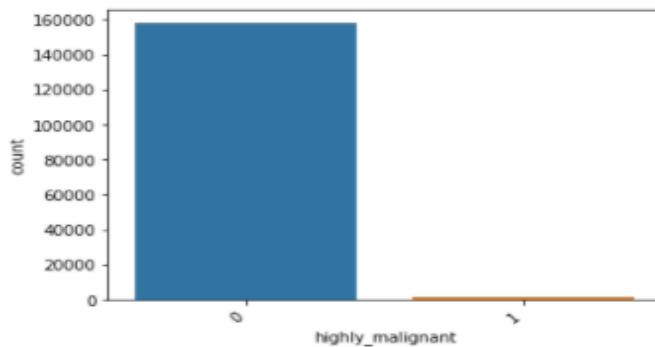
```



```

0    157976
1     1595
Name: highly_malignant, dtype: int64

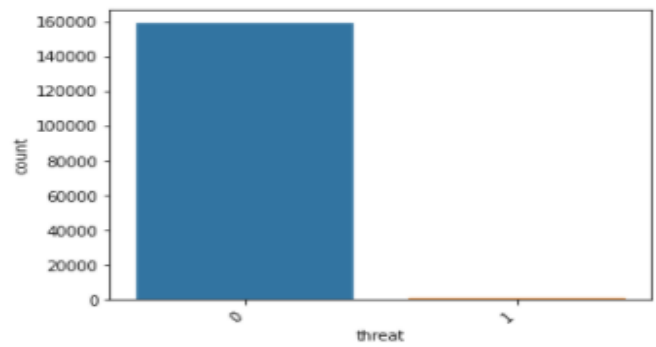
```



```

0    159093
1       478
Name: threat, dtype: int64

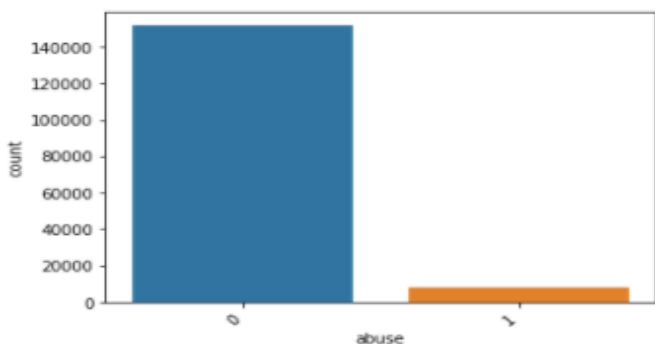
```



```

0    151694
1     7877
Name: abuse, dtype: int64

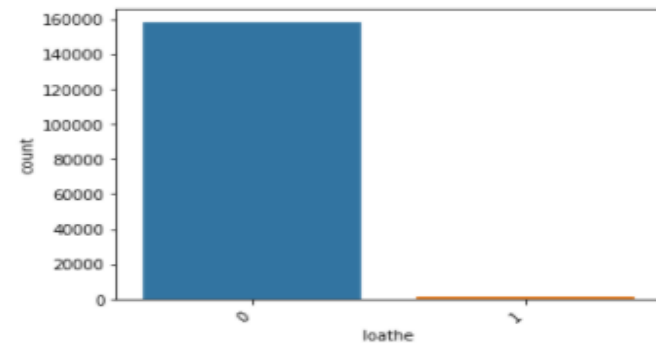
```



```

0    158166
1     1405
Name: loathe, dtype: int64

```



Below are the observations we can find from the above graph:

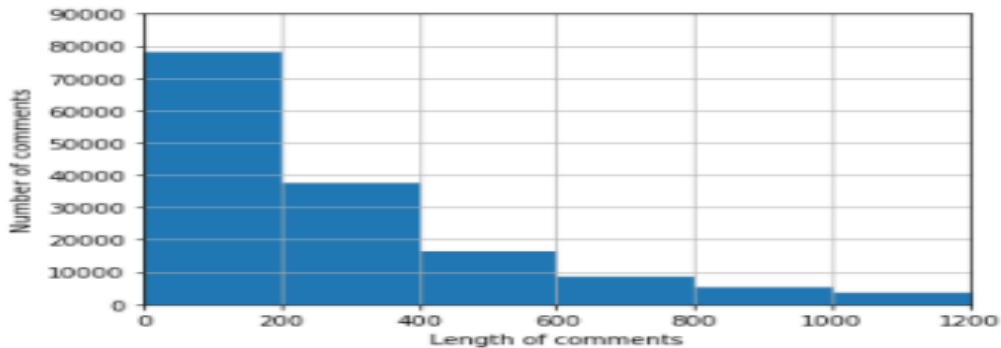
1. For feature malignant, we have 0 value 144277 times and 1 for 15294 times.
2. For feature highly_malignant, we have 0 value 157976 times and 1 for 1595 times.
3. For feature rude, we have 0 value 151122 times and 1 for 8449 times.
4. For feature threat, we have 0 value 159093 times and 1 for 478 times.
5. For feature abuse, we have 0 value 151694 times and 1 for 7877 times.
6. For feature loathe, we have 0 value 158166 times and 1 for 1405 times.


```

1 x = train['comment_text'].str.len()
2
3 print('average length of comment: {:.3f}'.format(sum(x)/len(x)) )
4 bins = [1,200,400,600,800,1000,1200]
5 plt.hist(x, bins=bins)
6 plt.xlabel('Length of comments')
7 plt.ylabel('Number of comments')
8 plt.axis([0, 1200, 0, 90000])
9 plt.grid(True)
10 plt.show()

```

average length of comment: 394.139



Above graph shows us the length of comments vs the number of comments. Maximum comments are with length 200 and very less comments are with length 1200.

```

1 train['length'] = train['comment_text'].str.len()
2 train.head(2)

```

	comment_text	malignant	highly_malignant	rude	threat	abuse	loathe	length
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

Added a new feature where we can see length of the particular comment.

```

# Convert all messages to lower case
train['comment_text'] = train['comment_text'].str.lower()

# Replace email addresses with 'email'
train['comment_text'] = train['comment_text'].str.replace(r'^.+@[^\.\.]*\.[a-z]{2,}$',
    'emailaddress')

# Replace URLs with 'webaddress'
train['comment_text'] = train['comment_text'].str.replace(r'^http://[a-zA-Z0-9\-\.\.]+\.[a-zA-Z]{2,3}(/\S*)?$',
    'webaddress')

# Replace money symbols with 'moneysymb' (£ can be typed with ALT key + 156)
train['comment_text'] = train['comment_text'].str.replace(r'£|\$', 'dollars')

# Replace 10 digit phone numbers (formats include paranthesis, spaces, no spaces, dashes) with 'phonenumber'
train['comment_text'] = train['comment_text'].str.replace(r'^\((?\d{3})\)?[\s-]?(\d{3})[\s-]?(\d{4})$',
    'phonenumber')

# Replace numbers with 'numbr'
train['comment_text'] = train['comment_text'].str.replace(r'\d+(\.\d+)?', 'numbr')

# Remove punctuations
train['comment_text'] = train['comment_text'].apply(lambda x: ' '.join(
    term for term in x.split() if term not in string.punctuation))

# Remove stopwords
stop_words = set(stopwords.words('english') + ['u', 'ü', 'ur', '4', '2', 'im', 'dont', 'doin', 'ure'])
train['comment_text'] = train['comment_text'].apply(lambda x: ' '.join(
    term for term in x.split() if term not in stop_words))

# Lemmetize words
lem=WordNetLemmatizer()
train['comment_text'] = train['comment_text'].apply(lambda x: ' '.join(
    lem.lemmatize(t) for t in x.split()))

```

After cleaning the comment text, let's have a look at comment length now.

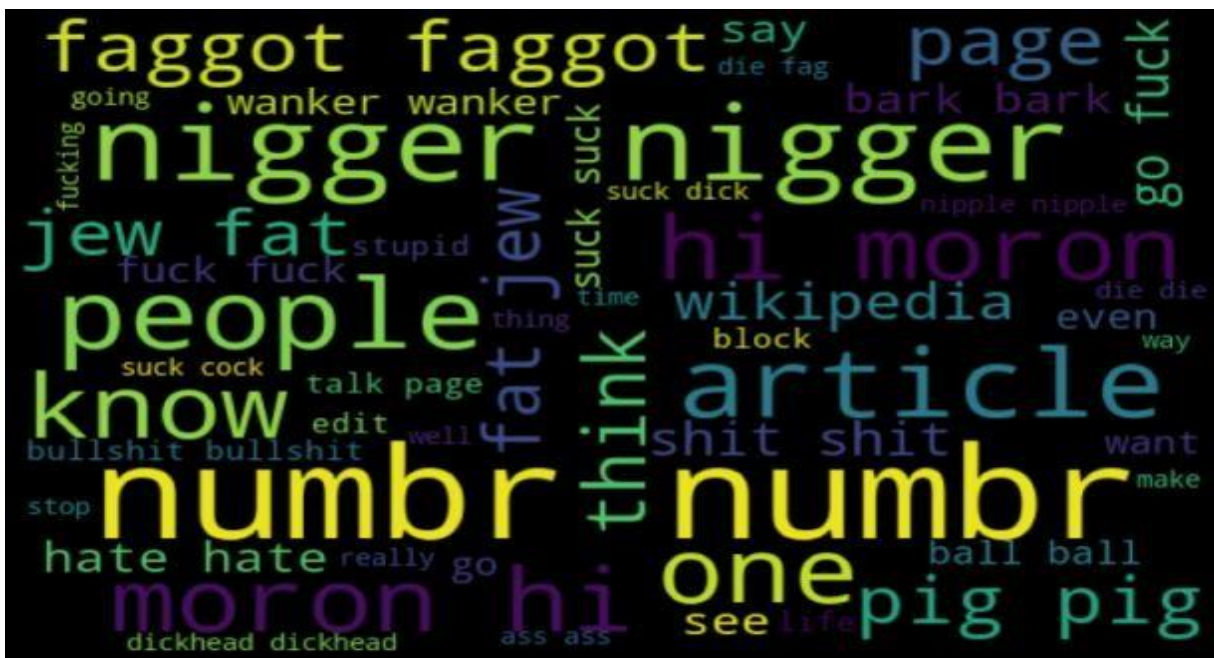
In the above pie chart, we can see very high number for malignant then rude and abuse are almost with same numbers. Then highly_malignant and loathe are almost with same numbers. Lastly, we can see very less numbers for feature threat.

```
1 target_data = train[cols_target]
2
3 train['bad'] = train[cols_target].sum(axis = 1)
4 print(train['bad'].value_counts())
5 train['bad'] = train['bad'] > 0
6 train['bad'] = train['bad'].astype(int)
7 print(train['bad'].value_counts())
```

```
0      143346
1       6360
3      4209
2      3480
4      1760
5       385
6       31
Name: bad, dtype: int64
0      143346
1      16225
Name: bad, dtype: int64
```

I have created a new feature name `bad` and it is the combination of features `malignant`, `highly_malignant`, `rude`, `threat`, `abuse`, `loathe` where the where value is 1.

```
1 #Getting sense of Loud words which are Loathe
2 from wordcloud import WordCloud
3 hams = train['comment_text'][train['bad']==1]
4 spam_cloud = WordCloud(width=600,height=400,background_color='black',max_words=50).generate(' '.join(hams))
5 plt.figure(figsize=(10,8),facecolor='k')
6 plt.imshow(spam_cloud)
7 plt.axis('off')
8 plt.tight_layout(pad=0)
9 plt.show()
```



Now with respect to bad feature the loud words are nigger, numbr, artivle, peoplr, hi moron one.

```

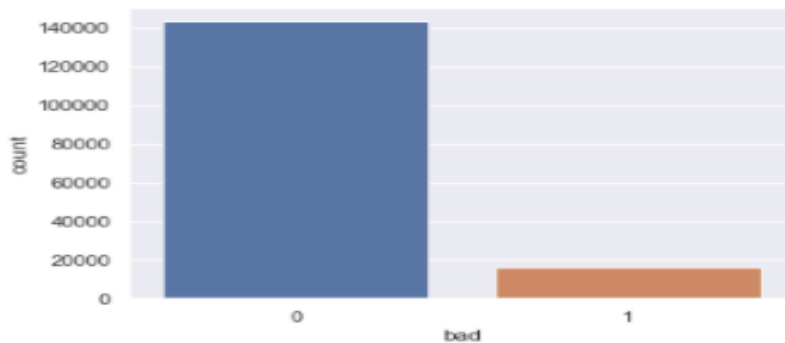
1 sns.set()
2 print(train['bad'].value_counts())
3 sns.countplot(x="bad" , data = train)
4 plt.show()

```

```

0    143346
1     16225
Name: bad, dtype: int64

```



We can see value count of target column is not same. We will work on this now.

Let's use TFidf method to transform data into vectors:

```

1 # Convert text into vectors using TF-IDF
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 tf_vec = TfidfVectorizer(max_features = 30000, stop_words='english')
4 features = tf_vec.fit_transform(train['comment_text'])
5 x = features

```

```
1 train.shape
```

```
(159571, 10)
```

```
1 y=train['bad']
```

```
1 print(x.shape,y.shape)
```

```
(159571, 30000) (159571,)
```

We will use under sampling method to balance the target column value counts:

```

1 from imblearn import under_sampling
2 from collections import Counter

```

```

1 from imblearn.under_sampling import RandomUnderSampler
2
3 rus=RandomUnderSampler(random_state=0)
4 x_resample,y_resampled=rus.fit_resample(x,y)
5 print(sorted(Counter(y_resampled)),y_resampled.shape)

```

```
[0, 1] (32450,)
```

```
1 y_resampled.value_counts()
```

```

1    16225
0    16225
Name: bad, dtype: int64

```

```
1 x_resample.shape, y_resampled.shape
```

```
((32450, 30000), (32450,))
```

```

1 x=x_resample
2 y=y_resampled

```

Let's split data into train and test set:

```
1 x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=56,test_size=.30)
```

```
1 x_train.shape,x_test.shape, y_train.shape,y_test.shape
```

```
((22715, 30000), (9735, 30000), (22715,), (9735,))
```

Model/s Development and Evaluation

Finding best random state:

```
1 maxAccu=0
2 maxRS=0
3 for i in range(1,200):
4     x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.20,random_state=i)
5     LR=LogisticRegression()
6     LR.fit(x_train,y_train)
7     pred=LR.predict(x_test)
8     acc=accuracy_score(y_test,pred)
9     if acc>maxAccu:
10         maxAccu=acc
11         maxRS=i
12 print("Best accuracy is ",maxAccu, " on Random State ",maxRS)
```

Best accuracy is 0.9061633281972264 on Random State 55

```
1 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.20,random_state=maxRS)
```

```
1 model=[LogisticRegression(),DecisionTreeClassifier(),KNeighborsClassifier(),AdaBoostClassifier(),RandomForestClassifier()]
2
3 for m in model:
4     m.fit(x_train,y_train)
5     #m.score(x_train,y_train)
6     pred=m.predict(x_test)
7     acc=accuracy_score(y_test,pred)
8     print('Accuracy Score of',m,'is:', acc)
9     print(confusion_matrix(y_test,pred))
10    print(classification_report(y_test,pred))
11    print('\n')
```

Accuracy Score of LogisticRegression() is: 0.9061633281972264
[[2835 205]
[404 2774]]

	precision	recall	f1-score	support
0	0.88	0.94	0.91	3312
1	0.93	0.87	0.90	3178
accuracy			0.91	6490
macro avg	0.91	0.91	0.91	6490
weighted avg	0.91	0.91	0.91	6490

Accuracy Score of KNeighborsClassifier() is: 0.6063174114021571
[[1964 1348]
[1207 1971]]

	precision	recall	f1-score	support
0	0.62	0.59	0.61	3312
1	0.59	0.62	0.61	3178
accuracy			0.61	6490
macro avg	0.61	0.61	0.61	6490
weighted avg	0.61	0.61	0.61	6490

Accuracy Score of DecisionTreeClassifier() is: 0.8469953775038521
[[2835 477]
[516 2662]]

	precision	recall	f1-score	support
0	0.85	0.86	0.85	3312
1	0.85	0.84	0.84	3178
accuracy			0.85	6490
macro avg	0.85	0.85	0.85	6490
weighted avg	0.85	0.85	0.85	6490

Accuracy Score of AdaBoostClassifier() is: 0.8305084745762712
[[3174 138]
[962 2216]]

	precision	recall	f1-score	support
0	0.77	0.96	0.85	3312
1	0.94	0.70	0.80	3178
accuracy			0.83	6490
macro avg	0.85	0.83	0.83	6490
weighted avg	0.85	0.83	0.83	6490

Accuracy Score of RandomForestClassifier() is: 0.8790446841294299
[[3084 228]
[557 2621]]

	precision	recall	f1-score	support
0	0.85	0.93	0.89	3312
1	0.92	0.82	0.87	3178
accuracy			0.88	6490
macro avg	0.88	0.88	0.88	6490
weighted avg	0.88	0.88	0.88	6490

- Accuracy for LogisticRegression model is 90.62%
- Accuracy for DecisionTreeClassifier model is 84.70%
- Accuracy for SVC model is 60.63%
- Accuracy for AdaBoostClassifier is 83.05%
- Accuracy for RandomForestClassifier is 87.90%

Hyper parameter tuning:

```
1 from sklearn.model_selection import RandomizedSearchCV
2 #creating parameter list to pass in RandomizedSearchCV
```

1. LogisticRegression:

```
1 #LogisticRegression parameters
2 parameters1={'penalty':['l2','l1','elasticnet', 'none'], 'dual':[True,False], 'tol':[0.0001], 'C':[1.0],
3             'fit_intercept':[True,False], 'intercept_scaling':[1], 'class_weight':[None], 'random_state':range(0,20),
4             'solver':['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'], 'max_iter':[100],
5             'multi_class':['auto', 'ovr', 'multinomial'], 'verbose':[0], 'warm_start':[True,False],
6             'n_jobs':[None], 'l1_ratio':[None]}
```

```
1 RSV1=RandomizedSearchCV(LogisticRegression(),parameters1,cv=5)
```

```
1 RSV1.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=LogisticRegression(),
                  param_distributions={'C': [1.0], 'class_weight': [None],
                                      'dual': [True, False],
                                      'fit_intercept': [True, False],
                                      'intercept_scaling': [1],
                                      'l1_ratio': [None], 'max_iter': [100],
                                      'multi_class': ['auto', 'ovr',
                                                    'multinomial'],
                                      'n_jobs': [None],
                                      'penalty': ['l2', 'l1', 'elasticnet',
                                                'none'],
                                      'random_state': range(0, 20),
                                      'solver': ['newton-cg', 'lbfgs',
                                                'liblinear', 'sag', 'saga'],
                                      'tol': [0.0001], 'verbose': [0],
                                      'warm_start': [True, False]})
```

```
1 RSV1.best_params_
```

```
{'warm_start': False,
 'verbose': 0,
 'tol': 0.0001,
 'solver': 'lbfgs',
 'random_state': 11,
 'penalty': 'l2',
 'n_jobs': None,
 'multi_class': 'multinomial',
 'max_iter': 100,
 'l1_ratio': None,
 'intercept_scaling': 1,
 'fit_intercept': False,
 'dual': False,
 'class_weight': None,
 'C': 1.0}
```

```
1 RSV_pred1=RSV1.best_estimator_.predict(x_test)
```

```
1 RSV_pred1
```

```
array([1, 0, 0, ..., 0, 0, 0])
```

```
1 score1=RSV1.score(x_train,y_train)
```

2. DecisionTreeClassifier

```
1 #DecisionTreeClassifier parameters
2 parameters2={'criterion':['gini', 'entropy'],'splitter':['best', 'random'],
3             'max_depth':[None],'random_state':range(0,20),
4             'max_leaf_nodes':[None],'class_weight':[None], 'min_samples_split':[1,2,3],
5             'min_samples_leaf':[1,2,3],'max_features':[None],'min_impurity_decrease':[0.0,1.0,2.0],
6             'min_weight_fraction_leaf':[0.0,1.0,2.0]}
```

```
1 RSV2=RandomizedSearchCV(DecisionTreeClassifier(),parameters2,cv=5)
```

```
1 RSV2.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=DecisionTreeClassifier(),
                  param_distributions={'class_weight': [None],
                                      'criterion': ['gini', 'entropy'],
                                      'max_depth': [None],
                                      'max_features': [None],
                                      'max_leaf_nodes': [None],
                                      'min_impurity_decrease': [0.0, 1.0, 2.0],
                                      'min_samples_leaf': [1, 2, 3],
                                      'min_samples_split': [1, 2, 3],
                                      'min_weight_fraction_leaf': [0.0, 1.0, 2.0],
                                      'random_state': range(0, 20),
                                      'splitter': ['best', 'random']})
```

```
1 RSV2.best_params_
```

```
{'splitter': 'random',
 'random_state': 2,
 'min_weight_fraction_leaf': 0.0,
 'min_samples_split': 3,
 'min_samples_leaf': 2,
 'min_impurity_decrease': 0.0,
 'max_leaf_nodes': None,
 'max_features': None,
 'max_depth': None,
 'criterion': 'gini',
 'class_weight': None}
```

```
1 RSV_pred2=RSV2.best_estimator_.predict(x_test)
```

```
1 RSV_pred2
```

```
array([1, 0, 0, ..., 0, 0, 0])
```

```
1 score2=RSV2.score(x_train,y_train)
```

3. KNeighborsClassifier

```
1 #KNeighborsClassifier parameters
2 parameters3={'n_neighbors':[4,5,6],'weights':['uniform','distance'],'algorithm':['auto','ball_tree','kd_tree','brute'],
3             'n_jobs':range(0,20)}
4
```

```
1 RSV3=RandomizedSearchCV(KNeighborsClassifier(),parameters3,cv=5)
```

```
1 RSV3.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=KNeighborsClassifier(),
                  param_distributions={'algorithm': ['auto', 'ball_tree',
                                                    'kd_tree', 'brute'],
                                      'n_jobs': range(0, 20),
                                      'n_neighbors': [4, 5, 6],
                                      'weights': ['uniform', 'distance']})
```

```
1 RSV3.best_params_
```

```
{'weights': 'distance', 'n_neighbors': 4, 'n_jobs': 18, 'algorithm': 'kd_tree'}
```

```
1 RSV_pred3=RSV3.best_estimator_.predict(x_test)
```

```
1 RSV_pred3
```

```
array([1, 0, 0, ..., 0, 0, 0])
```

```
1 score3=RSV3.score(x_train,y_train)
```

4. AdaBoostClassifier

```
1 parameters4={'base_estimator':[None], 'n_estimators':[20,40, 50,75,100],
2             'learning_rate':[0.1,0.01,1.0,2.0], 'algorithm':['SAMME.R', 'SAMME'],
3             'random_state':range(0,20)}
```

```
1 RSV4=RandomizedSearchCV(AdaBoostClassifier(),parameters4,cv=5)
```

```
1 RSV4.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=AdaBoostClassifier(),
                  param_distributions={'algorithm': ['SAMME.R', 'SAMME'],
                                      'base_estimator': [None],
                                      'learning_rate': [0.1, 0.01, 1.0, 2.0],
                                      'n_estimators': [20, 40, 50, 75, 100],
                                      'random_state': range(0, 20)})
```

```
1 RSV4.best_params_
```

```
{'random_state': 10,
 'n_estimators': 20,
 'learning_rate': 1.0,
 'base_estimator': None,
 'algorithm': 'SAMME.R'}
```

```
1 RSV_pred4=RSV4.best_estimator_.predict(x_test)
```

```
1 RSV_pred4
```

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

```
1 score4=RSV4.score(x_train,y_train)
```

5. RandomForestClassifier

```
1 parameters5={'n_estimators':[100], 'criterion':['gini', 'entropy'], 'max_depth':[None],
2             'min_samples_split':[2], 'min_samples_leaf':[1], 'min_weight_fraction_leaf':[0.0],
3             'max_features':['auto'], 'max_leaf_nodes':[None], 'min_impurity_decrease':[0.0],
4             'min_impurity_split':[None], 'bootstrap':[True, False], 'oob_score':[True, False], 'n_jobs':[None],
5             'random_state':range(0,20), 'verbose':[0], 'warm_start':[True, False], 'class_weight':[None],
6             'ccp_alpha':[0.0], 'max_samples':[None]}
```

```
1 RSV5=RandomizedSearchCV(RandomForestClassifier(),parameters5,cv=5)
```

```
1 RSV5.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=RandomForestClassifier(),
                  param_distributions={'bootstrap': [True, False],
                                      'ccp_alpha': [0.0],
                                      'class_weight': [None],
                                      'criterion': ['gini', 'entropy'],
                                      'max_depth': [None],
                                      'max_features': ['auto'],
                                      'max_leaf_nodes': [None],
                                      'max_samples': [None],
                                      'min_impurity_decrease': [0.0],
                                      'min_impurity_split': [None],
                                      'min_samples_leaf': [1],
                                      'min_samples_split': [2],
                                      'min_weight_fraction_leaf': [0.0],
                                      'n_estimators': [100], 'n_jobs': [None],
                                      'oob_score': [True, False],
                                      'random_state': range(0, 20),
                                      'verbose': [0],
                                      'warm_start': [True, False]})
```

```
1 RSV5.best_params_
```

```
{'warm_start': True,
 'verbose': 0,
 'random_state': 6,
 'oob_score': False,
 'n_jobs': None,
 'n_estimators': 100,
 'min_weight_fraction_leaf': 0.0,
 'min_samples_split': 2,
 'min_samples_leaf': 1,
 'min_impurity_split': None,
 'min_impurity_decrease': 0.0,
 'max_samples': None,
 'max_leaf_nodes': None,
 'max_features': 'auto',
 'max_depth': None,
 'criterion': 'entropy',
 'class_weight': None,
 'ccp_alpha': 0.0,
 'bootstrap': False}
```

```
1 RSV_pred5=RSV5.best_estimator_.predict(x_test)

1 RSV_pred5

array([1, 0, 0, ..., 0, 0, 0])

1 score5=RSV5.score(x_train,y_train)
```

```
1 print("Accuracy for LogisticRegression is ",score1*100,"%\n")
2 print("Accuracy for DecisionTreeClassifier is ",score2*100,"%\n")
3 print("Accuracy for KNeighborsClassifier is ",score3*100,"%\n")
4 print("Accuracy for AdaBoostClassifier is ",score4*100,"%\n")
5 print("Accuracy for RandomForestClassifier is ",score5*100,"%\n")
```

Accuracy for LogisticRegression is 94.79583975346686 %

Accuracy for DecisionTreeClassifier is 94.59167950693374 %

Accuracy for KNeighborsClassifier is 99.89214175654854 %

Accuracy for AdaBoostClassifier is 75.38135593220339 %

Accuracy for RandomForestClassifier is 99.8959938366718 %

AUC ROC:

```
1 from sklearn.metrics import roc_curve
2 from sklearn.metrics import roc_auc_score
```

```
1 #RandomizedSearchCV - LogisticRegression
2 fpr, tpr, thresholds=roc_curve(y_test, RSV_pred1)
```

```
1 plt.plot([0.1],[0.1], 'k--')
2 plt.plot(fpr, tpr, label=' LogisticRegression - RandomizedSearchCV')
3 plt.xlabel('False positive rate')
4 plt.ylabel('True positive rate')
5 plt.title('RandomizedSearchCV')
6 plt.show()
```



```
1 #RandomizedSearchCV - DecisionTreeClassifier
2 fpr, tpr, thresholds=roc_curve(y_test, RSV_pred2)
```

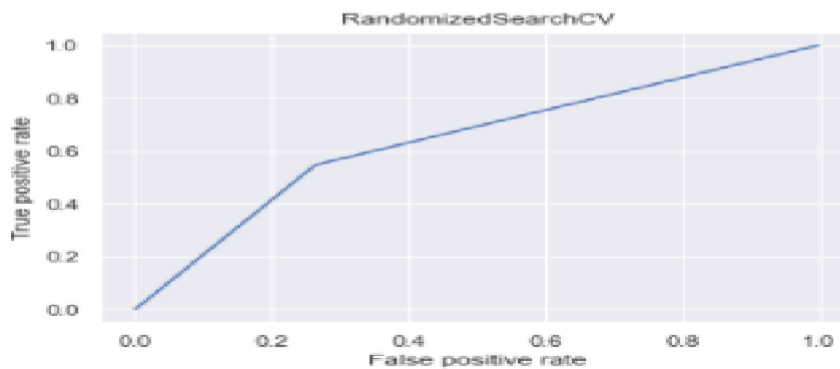
```
1 #RandomizedSearchCV - DecisionTreeClassifier
2 fpr, tpr, thresholds=roc_curve(y_test, RSV_pred2)
```

```
1 plt.plot([0.1],[0.1], 'k--')
2 plt.plot(fpr, tpr, label=' DecisionTreeClassifier - RandomizedSearchCV')
3 plt.xlabel('False positive rate')
4 plt.ylabel('True positive rate')
5 plt.title('RandomizedSearchCV')
6 plt.show()
```



```
1 #RandomizedSearchCV - KNeighborsClassifier
2 fpr,tpr,thresholds=roc_curve(y_test,RSV_pred3)
```

```
1 plt.plot([0.1],[0.1],'k--')
2 plt.plot(fpr,tpr,label=' KNeighborsClassifier - RandomizedSearchCV')
3 plt.xlabel('False positive rate')
4 plt.ylabel('True positive rate')
5 plt.title('RandomizedSearchCV')
6 plt.show()
```

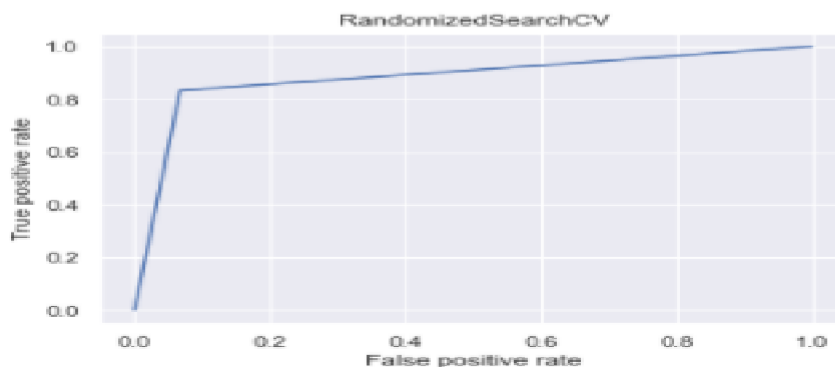


```
1 plt.plot([0.1],[0.1],'k--')
2 plt.plot(fpr,tpr,label=' AdaBoostClassifier - RandomizedSearchCV')
3 plt.xlabel('False positive rate')
4 plt.ylabel('True positive rate')
5 plt.title('RandomizedSearchCV')
6 plt.show()
```



```
1 #RandomizedSearchCV - RandomForestClassifier
2 fpr,tpr,thresholds=roc_curve(y_test,RSV_pred5)
```

```
1 plt.plot([0.1],[0.1],'k--')
2 plt.plot(fpr,tpr,label=' RandomForestClassifier - RandomizedSearchCV')
3 plt.xlabel('False positive rate')
4 plt.ylabel('True positive rate')
5 plt.title('RandomizedSearchCV')
6 plt.show()
```



Considering LogisticRegression with RandomizedSearchCV as final model for saving model as we have got best accuracy score with RandomForestClassifier - logistic regression and AUC_roc score is graph also is better than other models.

Saving the model:

```
1 #saving the model
2 import pickle
3 filename='RSVMalignant.pkl'
4 pickle.dump(RSV1,open(filename,'wb'))
```

We will import the test file now:

```
test=pd.read_csv(r"C:\Users\User\Downloads\Malignant Comments Classifier Project\malignant_test.csv")
test=pd.DataFrame(data=test)
test
```

	id	comment_text
	0	00001cee341fdb12 Yo bitch Ja Rule is more succesful then you'll...
	1	0000247867823ef7 == From RfC == \n\n The title is fine as it is...
	2	00013b17ad220c46 " \n\n == Sources == \n\n * Zawe Ashton on Lap...
	3	00017563c3f7919a :If you have a look back at the source, the in...
	4	00017695ad8997eb I don't anonymously edit articles at all.
...
153159	ffcd0960ee309b5	. \n i totally agree, this stuff is nothing bu...
153160	fffd7a9a6eb32c16	== Throw from out field to home plate. == \n\n...
153161	ffda9e8d6fafa9e	" \n\n == Okinotorishima categories == \n\n I ...
153162	ffe8f1340a79fc2	" \n\n == ""One of the founding nations of the...
153163	ffffce3fb183ee80	" \n :::Stop already. Your bullshit is not wel...

153164 rows × 2 columns

```
1 test.columns #list of columns
```

```
Index(['id', 'comment_text'], dtype='object')
```

```
1 test.shape
```

```
(153165, 2)
```

Test dataset consists of 153164 rows and 2 columns.

```
1 test.dtypes
```

```
id          object
comment_text object
dtype: object
```

We can see both type of columns object type.

```
1 test.describe(include='all').T
```

	count	unique	top	freq
id	153164	153164	c3efd845b70c67d3	1
comment_text	153165	153165	" \n ***So images of ass and shit are okay, bu...	1

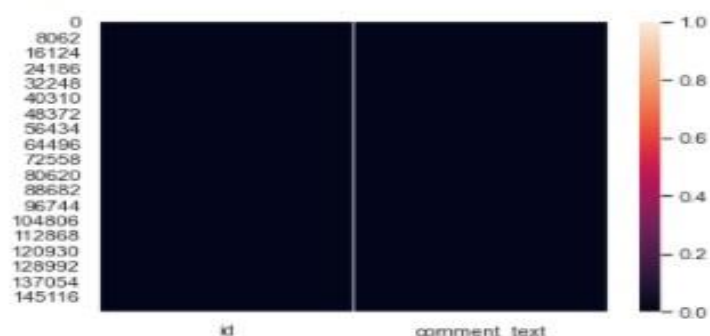
We can see count is all same for all variable, let's check for null values now:

```
1 test.isnull().sum()
```

```
id          1
comment_text 0
dtype: int64
```

```
1 #visual representation with the help of heatmap
2 sns.heatmap(test.isnull())
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1d911ac8e20>
```

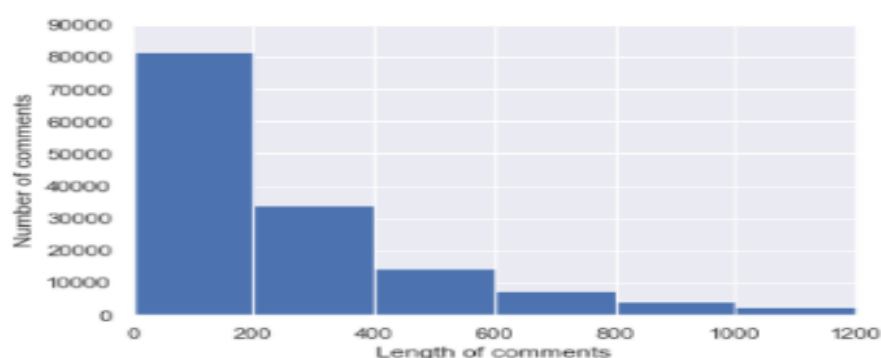


Id feature also has all unique values like serial number hence we dont require feature for future analysis

```
2 test.drop('id', axis=1, inplace= True)
```

```
1 x = test['comment_text'].str.len()
2
3 print('average length of comment: {:.3f}'.format(sum(x)/len(x)) )
4 bins = [1,200,400,600,800,1000,1200]
5 plt.hist(x, bins=bins)
6 plt.xlabel('Length of comments')
7 plt.ylabel('Number of comments')
8 plt.axis([0, 1200, 0, 90000])
9 plt.grid(True)
10 plt.show()
```

average length of comment: 364.873



Above graph shows us the length of comments vs the number of comments. Maximum comments are with length 200 and very less comments are with length 1200.

```
1 test['length'] = test['comment_text'].str.len()
2 test.head(2)
```

	comment_text	length
0	Yo bitch Ja Rule is more succesful then you'll...	367
1	== From RfC == \n\n The title is fine as it is...	50

Added a new feature where we can see length of the particular comment.

```

1 # Convert all messages to lower case
2 test['comment_text'] = test['comment_text'].str.lower()
3
4 # Replace email addresses with 'email'
5 test['comment_text'] = test['comment_text'].str.replace(r'^.+@[^\.\.]*\.[a-z]{2,}$',
6 'emailaddress')
7
8 # Replace URLs with 'webaddress'
9 test['comment_text'] = test['comment_text'].str.replace(r'^http://[a-zA-Z0-9\-\.\.]+\.[a-zA-Z]{2,3}(/[^*]*)?$',
10 'webaddress')
11
12 # Replace money symbols with 'moneysymb' (£ can be typed with ALT key + 156)
13 test['comment_text'] = test['comment_text'].str.replace(r'£|$', 'dollars')
14
15 # Replace 10 digit phone numbers (formats include paranthesis, spaces, no spaces, dashes) with 'phonenumber'
16 test['comment_text'] = test['comment_text'].str.replace(r'^\(?[\d]{3}\)?[\s-]?[\d]{3}[\s-]?[\d]{4}$',
17 'phonenumber')
18
19
20 # Replace numbers with 'numbr'
21 test['comment_text'] = test['comment_text'].str.replace(r'\d+(\.\d+)?', 'numbr')
22
23 # Remove punctuations
24 test['comment_text'] = test['comment_text'].apply(lambda x: ' '.join(
25 term for term in x.split() if term not in string.punctuation))
26
27 # Remove stopwords
28 stop_words = set(stopwords.words('english') + ['u', 'ü', 'ur', '4', '2', 'im', 'dont', 'doin', 'ure'])
29 test['comment_text'] = test['comment_text'].apply(lambda x: ' '.join(
30 term for term in x.split() if term not in stop_words))
31
32 # Lemmetize words
33 lem=WordNetLemmatizer()
34 test['comment_text'] = test['comment_text'].apply(lambda x: ' '.join(
35 lem.lemmatize(t) for t in x.split()))

```

After cleaning the comment text, let's have a look at comment length now.

```

1 test['clean_length'] = test.comment_text.str.len()
2 test.head()

```

	comment_text	length	clean_length
0	yo bitch ja rule succesful ever whats hating s...	367	249
1	== rfc == title fine is. imo.	50	29
2	== source == zawe ashton lapland —	54	34
3	:if look back source, information updated corr...	205	117
4	anonymously edit article all.	41	29

```

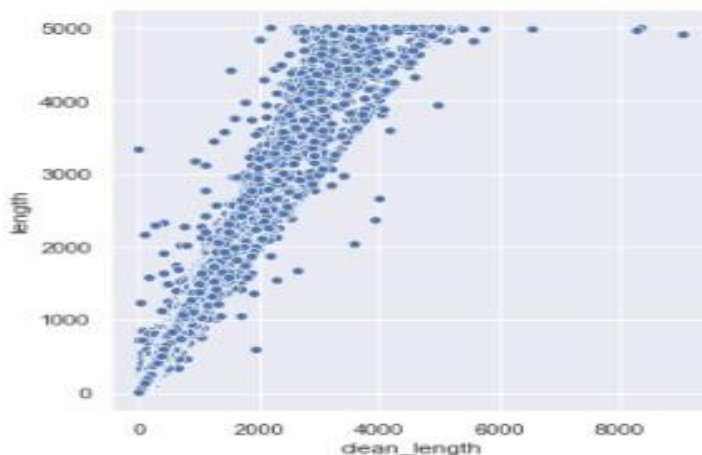
1 # Total Length removal
2 print ('Origion Length', test.length.sum())
3 print ('Clean Length', test.clean_length.sum())

```

Origion Length 55885736
Clean Length 39400455

```
1 sns.relplot(x='clean_length',y='length',data=test,kind='scatter')
```

<seaborn.axisgrid.FacetGrid at 0x1d919ecd040>



Dropping newly created columns because it's not required for further analysis:


```
1 test.drop('length', axis=1, inplace= True)
```

```
1 test.drop('clean_length', axis=1, inplace= True)
```

Let's use TFidf method to transform data into vectors:

```
1 # Convert text into vectors using TF-IDF
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 tf_vec = TfidfVectorizer(max_features = 30000, stop_words='english')
4 features = tf_vec.fit_transform(test['comment_text'])
5 x = features
```

```
1 x.shape
```

```
(153165, 30000)
```

load the model

```
1 fitted_model=pickle.load(open("RSVMalignant.pkl","rb"))
```

```
1 fitted_model
```

```
RandomizedSearchCV(cv=5, estimator=LogisticRegression(),
    param_distributions={'C': [1.0], 'class_weight': [None],
                        'dual': [True, False],
                        'fit_intercept': [True, False],
                        'intercept_scaling': [1],
                        'l1_ratio': [None], 'max_iter': [100],
                        'multi_class': ['auto', 'ovr',
                                        'multinomial'],
                        'n_jobs': [None],
                        'penalty': ['l2', 'l1', 'elasticnet',
                                    'none'],
                        'random_state': range(0, 20),
                        'solver': ['newton-cg', 'lbfgs',
                                  'liblinear', 'sag', 'saga'],
                        'tol': [0.0001], 'verbose': [0],
                        'warm_start': [True, False]})
```

prediction over test dataset

```
1 prediction=fitted_model.predict(x)
```

```
1 prediction=fitted_model.best_estimator_.predict(x)
```

```
1 prediction
```

```
array([1, 0, 0, ..., 0, 0, 1])
```

```
1 test['bad'] = prediction
2 test.to_csv('test.csv')
```

CONCLUSION

In Malignant Comments Classifier project our goal was to build a prototype of online hate and abusive 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. Our database consists of 159571 rows and 8 columns in which we have 2 object type columns and all other columns are numeric. First, I have dropped some columns which are not required also I have dropped id column as it has all unique values like serial number hence we don't require it as a feature for future analysis.

I have performed some visualization showing relations between some features. Then I have checked correlation of data. With the help of countplot I tried to show numbers with all features. And with the help of NLP methods, I cleaned comments text feature. Wordcloud helped in showing allowed words with respect to all features. Used pie chart to show contribution of all features. Then created a new feature combining all features.

With the help of TFidf method to transform data into vectors. Our target column does not contain equal number of data hence I performed under sampling method to equalize. Used 5 methods for model building. With the help of RandomizedSearchCV I have tried to improve accuracy. I decided to go ahead with Logistic Regression and saved model.

Performed similar operations on test dataset as well according to the requirement. After that loaded the model and performed prediction over test dataset and then saved the predictions in the test file.