



SUBMITTED BY
V Samyukta

ACKNOWLEDGMENT

It gives me immense pleasure to deliver this report. Working on this project was a great learning experience that helped me attain in-depth knowledge on data analysis process.

Flip Robo Technologies (Bangalore) provided all of the necessary information and datasets, required for the completion of the project.

I express my gratitude to my SME, Ms. Khushboo Garg, for providing the dataset and directions for carrying out the project report procedure.

My heartfelt gratitude to DataTrained institute and FlipRobo company for providing me this internship opportunity. Last but not least to my sincere thanks to my family and all those who helped me directly or indirectly in completion this project.

CONTENTS

1. Introduction

- Business Problem Framing:
- Conceptual Background of the Domain Problem
- Review of Literature
- Motivation for the Problem Undertaken

2. Analytical Problem Framing

- Mathematical/ Analytical Modelling of the Problem
- Data Sources and their formats
- Data Preprocessing Done
- Data Inputs-Logic-Output Relationships
- Hardware and Software Requirements and Tools Used

3. Data Analysis and Visualization

- Identification of possible problem-solving approaches (methods)
- Testing of Identified Approaches (Algorithms)
- Key Metrics for success in solving problem under consideration
- Visualization
- Run and Evaluate selected models
- Interpretation of the Results

4. Conclusion

- Key Findings and Conclusions of the Study
- Learning Outcomes of the Study in respect of Data Science
- Limitations of this work and Scope for Future Work

5. Reference

1.INTRODUCTION

- **Business Problem Framing:**

The proliferation of social media enables people to express their opinions widely online. However, at the same time, this has resulted in the emergence of conflict and hate, making online environments uninviting for users. Although researchers have found that hate is a problem across multiple platforms, there is a lack of models for online hate detection. Online hate, described as abusive language, aggression, cyberbullying, hatefulness and many others has been identified as a major threat on online social media platforms. Social media platforms are the most prominent grounds for such toxic behaviour. There has been a remarkable increase in the cases of cyberbullying and trolls on various social media platforms. Many celebrities and influences are facing backlashes from people and have to come across hateful and offensive comments. This can take a toll on anyone and affect them mentally leading to depression, mental illness, self-hatred and suicidal thoughts. Internet comments are bastions of hatred and vitriol. While online anonymity has provided a new outlet for aggression and hate speech, machine learning can be used to fight it. The problem we sought to solve was the tagging of internet comments that are aggressive towards other users. This means that insults to third parties such as celebrities will be tagged as unoffensive, but “u are an idiot” is clearly offensive. Our goal is to build a prototype of online hate and abuse comment classifier which can 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:**

In the past few years its seen that the cases related to social media hatred have increased exponentially. The social media is turning into a dark venomous pit for people now a days. Online hate is the result of difference in opinion, race, religion, occupation, nationality etc. In social media the people spreading or involved in such kind of activities uses filthy languages, aggression, images etc. to offend and gravely hurt the person on the other side. This is one of the major concerns now. The result of such activities can be dangerous. It gives mental trauma to the victims making their lives miserable. People who are not well aware of mental health online hate or cyber bullying 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**

Nowadays users leave numerous comments on different social networks, news portals, and forums. Some of the comments are toxic or abusive. Due to numbers of comments, it is unfeasible to manually moderate them, so most of the systems use some kind of automatic discovery of toxicity using machine learning models. In this work, we performed

a systematic review of the state-of-the-art in toxic comment classification using machine learning methods. First, we have investigated when and where the papers were published and their maturity level. In our analysis of every primary study, we investigated: data set used, evaluation metric, used machine learning methods, classes of toxicity, and comment language.

- **Motivation for the Problem Undertaken:**

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

2.ANALYTICAL PROBLEM FRAMING

- **Mathematical / Analytical Modelling of the Problem:**

In this given problem the label can be either 0 or 1, where 0 denotes a NO while 1 denotes a YES. So clearly it is a binary classification problem and I have to use all classification algorithms while building the model. We would perform one type of supervised learning algorithms: Classification. While it seems more reasonable to perform Classification since we have 5-6 class to predict. Here, we will only perform classification. Since there only 1 feature in the dataset, filtering the words is needed to prevent overfit. In order to determine the regularization parameter, throughout the project in classification part, we would first remove email, phone number, web address, spaces and stops words etc. In order to further improve our models, we also performed TFIDF in order to convert the tokens from the train documents into vectors so that machine can do further processing. I have used all the classification algorithms while building model then tuned the best model and saved the best model. At last, I have predicted the Malignance using saved model.

- **Data Sources and their formats:**

The data set contains the training set, which has approximately 1,59,000 samples and the test set which contains nearly 1,53,000 samples. All the data samples contain 8 fields which includes 'Id', 'Comments', '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.

The data set includes:

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

This project is more about exploration, feature engineering and classification that can be done on this data. Since the data set is huge and includes many categories of comments, we can do good amount of data exploration and derive some interesting features using the comments text column available.

We need to build a model that can differentiate between comments and its categories.

- **Data Pre-processing Done:**

- ❖ As a first step I have imported required libraries and I have imported the dataset which was in csv format.
- ❖ Cleaned the data from junk values. Replace multiple spaces with single space So that it will be easy to classify it.
- ❖ I am creating a function for feature engineering and making three different columns using comment_text column Length: indicating the length of the text. Exclamation: indicates whether '!' is present in the text or not. Question: indicates whether '?' is present in the text or not.
- ❖ By observing these comments, we can say that we need to do lot of text processing as there are many words which are not important for prediction, as well as numbers and other stuff.

- **Hardware and Software Requirements and Tools Used:**

Hardware technology being Used:-

- ❖ CPU: HP Pavilion
- ❖ Chip: intel core13 8th Gen
- ❖ RAM: 8 GB

Software Technology being Used:-

- ❖ Programming language: Python
- ❖ Distribution: Anaconda Navigator
- ❖ Browser based language shell: Jupyter Notebook

Libraries/Packages Used:-

- ❖ **import pandas as pd:** **pandas** is a popular Python-based data analysis toolkit which can be imported using **import pandas as pd**. It presents a diverse range of utilities, ranging from parsing multiple file formats to converting an entire data table into a numpy matrix array. This makes pandas a trusted ally in data science and machine learning.
- ❖ **import numpy as np:** NumPy is the fundamental package for scientific computing in Python. It is a Python library that provides a multidimensional array object, various derived objects (such as masked arrays and matrices), and an assortment of routines for fast operations on arrays, including mathematical, logical, shape manipulation, sorting, selecting, I/O, discrete Fourier transforms, basic linear algebra, basic statistical operations, random simulation and much more.
- ❖ **import seaborn as sns:** Seaborn is a data visualization library built on top of matplotlib and closely integrated with pandas data structures in Python. Visualization is the central part of Seaborn which helps in exploration and understanding of data.
- ❖ **Import matplotlib.pyplot as plt:** matplotlib.pyplot is a collection of functions that make matplotlib work like MATLAB. Each pyplot function makes some change to a figure: e.g., creates a figure, creates a plotting area in a figure, plots some lines in a plotting area, decorates the plot with labels, etc.
- ❖ `from sklearn.svm import LinearSVC`
- ❖ `from sklearn.naive_bayes import MultinomialNB`
- ❖ `from sklearn.linear_model import LogisticRegression`
- ❖ `from lightgbm import LGBMClassifier`
- ❖ `from sklearn.linear_model import SGDClassifier`
- ❖ `from sklearn.metrics import classification_report`
- ❖ `from sklearn.metrics import accuracy_score`
- ❖ `from sklearn.model_selection import cross_val_score`

With this sufficient libraries we can go ahead with our model building.

3.DATA ANALYSIS AND VISUALIZATION

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

Just make the comments more appropriate so that we'll get less word to process and get more accuracy. Removed extra spaces, converted email address into email keyword, likely wise phone number etc. Tried to make Comments small and more appropriate as much as it was possible.

- **Testing of Identified Approaches (Algorithms):**

In this nlp based project we need to predict multiple targets which are binary. I have converted the text into vectors using TFIDF vectorizer and separated our feature and labels then build the model using One Vs Rest Classifier. Among all the algorithms which I have used for this purpose I have chosen LinearSVC as best suitable algorithm for our final model as it is performing well compared to other algorithms while evaluating with different metrics I have used following algorithms and evaluated them

- ❖ LinearSVC
- ❖ LogisticRegression
- ❖ MultinomialNB
- ❖ LightGBMClassifier
- ❖ SGDClassifier

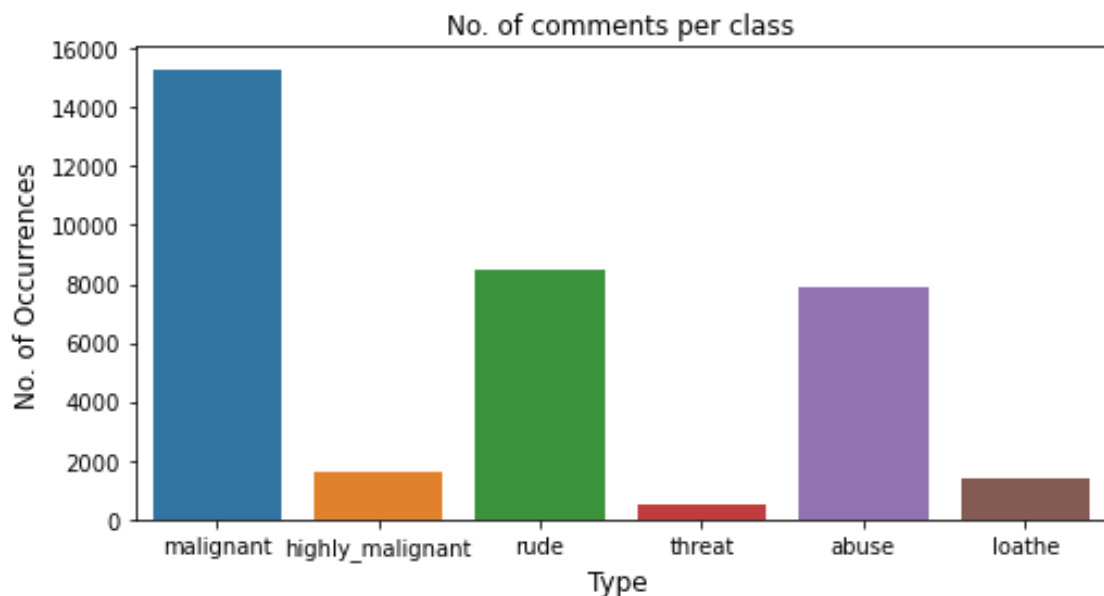
From all of these above models LinearSVC was giving me good performance.

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

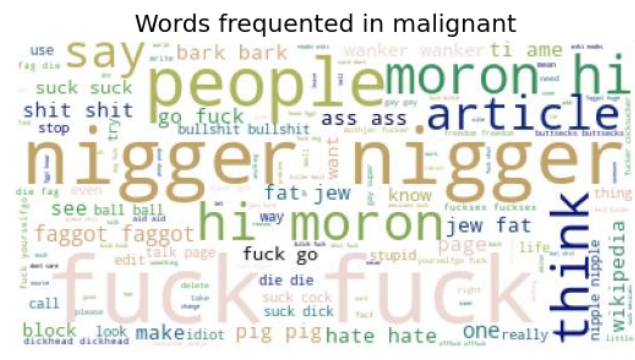
I have used the following metrics for evaluation:

- I have used `f1_score`, `precision_score`, `recall_score`, `multilabel_confusion_matrix` and `hamming loss` all these evaluation metrics to select best suitable algorithm for our final model.
- **Precision** can be seen as a measure of quality, higher precision means that an algorithm returns more relevant results than irrelevant ones.
- **Recall** is used as a measure of quantity and high recall means that an algorithm returns most of the relevant results.
- **Accuracy score** is used when the True Positives and True negatives are more important. Accuracy can be used when the class distribution is similar.
- **F1-score** is used when the False Negatives and False Positives are crucial. While F1-score is a better metric when there are imbalanced classes.

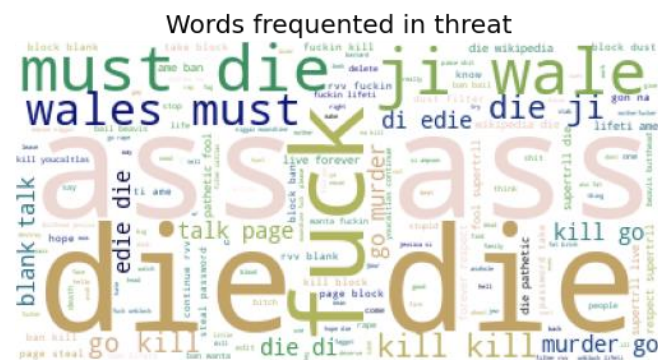
- **Visualization:**



The above figure represents count plot for all our labels. Looking at this plot we can conclude that more number of comments has been labelled as malignant compared to others. Very less number of comments has been labelled as threat.



- ❖ The above both figures are representing the word occurrence in case of malignant and highly malignant comments respectively.



- ❖ The above both figures are representing the word occurrence in case of threat and highly rude comments respectively.

- **Run and evaluate selected models**

1. Model building:

```

: #Lets define different algorithms
    svc = LinearSVC()
    lr = LogisticRegression(solver='lbfgs')
    mnb = MultinomialNB()
    lgb = LGBMClassifier()
    sgd = SGDClassifier()

: #function for printing score
    def print_score(y_pred,clf):
        print('classifier:',clf.__class__.__name__)
        print("Jaccard score: {}".format(jaccard_score(y_test,y_pred,average='micro')))
        print("Accuracy score: {}".format(accuracy_score(y_test,y_pred)))
        print("f1_score: {}".format(f1_score(y_test,y_pred,average='micro')))
        print("Precision : ", precision_score(y_test,y_pred,average='micro'))
        print("Recall: {}".format(recall_score(y_test,y_pred,average='micro')))
        print("Hamming loss: ", hamming_loss(y_test,y_pred))
        print("Confusion matrix:\n ", multilabel_confusion_matrix(y_test,y_pred))
        print('=====\\n')

: #models with evaluation using OneVsRestClassifier
    for classifier in [svc,lr,mnb,sgd,lgb]:
        clf = OneVsRestClassifier(classifier)
        clf.fit(x_train,y_train)
        y_pred = clf.predict(x_test)
        print score(y_pred, classifier)

```

classifier: LinearSVC
Jaccard score: 0.545426673479816
Accuracy score: 0.9194846213621437
f1_score: 0.7058590133580215
Precision : 0.8491646778042959
Recall: 0.6039379880049791
Hamming loss: 0.018583042973286876

Confusion matrix:

```
[[[35695  383]
 [ 1267 2548]]
```

```
[[39421  66]
 [ 306 100]]
```

```
[[37589 161]
 [ 694 1449]]
```

```
[[39769 19]
 [ 82 23]]
```

```
[[37610 272]
 [ 873 1138]]
```

```
[[39489 47]
 [ 278 79]]]
```

=====

classifier: LogisticRegression
Jaccard score: 0.5223303001778057
Accuracy score: 0.9198355601233299
f1_score: 0.6862246650635521
Precision : 0.8733823015040224
Recall: 0.565123910829467
Hamming loss: 0.01908020621830062

Confusion matrix:

```
[[[35821 257]
 [ 1430 2385]]
```

```
[[39416 71]
 [ 290 116]]
```

```
[[37620 130]
 [ 774 1369]]
```

```
[[39781 7]
 [ 92 13]]
```

```
[[37650 232]
 [ 960 1051]]
```

```
[[39509 27]
 [ 297 60]]]
```

=====

classifier: MultinomialNB
Jaccard score: 0.43909710391822826

Accuracy score: 0.9130925224976812
f1_score: 0.6102397158922758
Precision : 0.8813849113058346
Recall: 0.4666742107049904
Hamming loss: 0.022008873737247137
Confusion matrix:

```
[[[35901 177]
 [ 1827 1988]]]
```

```
[[39440 47]
 [ 328 78]]]
```

```
[[37624 126]
 [ 1002 1141]]]
```

```
[[39788 0]
 [ 105 0]]]
```

```
[[37696 186]
 [ 1123 888]]]
```

```
[[39517 19]
 [ 328 29]]]
```

=====

classifier: SGDClassifier
Jaccard score: 0.4520326952032695
Accuracy score: 0.9145714787055373
f1_score: 0.6226205466261758
Precision : 0.9011578044596913
Recall: 0.4756138961185923
Hamming loss: 0.021286107002899422
Confusion matrix:

```
[[[35967 111]
 [ 1856 1959]]]
```

```
[[39487 0]
 [ 406 0]]]
```

```
[[37618 132]
 [ 833 1310]]]
```

```
[[39788 0]
 [ 105 0]]]
```

```
[[37668 214]
 [ 1099 912]]]
```

```
[[39532 4]
 [ 335 22]]]
```

=====

classifier: LGBMClassifier
Jaccard score: 0.5359082679541339
Accuracy score: 0.9174291229037675
f1_score: 0.6978388998035363

Precision : 0.8282294419399969
Recall: 0.6029195428312776
Hamming loss: 0.019276564810869076
Confusion matrix:
[[[35781 297]
[1412 2403]]

[[39390 97]
[308 98]]

[[37534 216]
[617 1526]]

[[39735 53]
[82 23]]

[[37498 384]
[828 1183]]

[[39478 58]
[262 95]]]

2. Hyper Parameter Tunning:

I have did hyperparameter tuning for LinearSVC for the parameters like
'estimator__penalty', 'estimator__loss', 'estimator__multi_class', 'estimator__dual',
'estimator__intercept_scaling', 'estimator__C'

Hyperparameter Tuning:

```
In [82]: 1 param = {  
2         'estimator__penalty': ['l1'],  
3         'estimator__loss': ['hinge', 'squared_hinge'],  
4         'estimator__multi_class': ['ovr', 'crammer_singer'],  
5         'estimator__dual': [False],  
6         'estimator__intercept_scaling': [2,4,5],  
7         'estimator__C': [2]  
8     }  
  
In [83]: 1 #train the model with given parameters using GridSearchCV  
2 svc = OneVsRestClassifier(LinearSVC())  
3 GCV = GridSearchCV(svc,param,cv = 3, verbose =0,n_jobs=-1)  
4 GCV.fit(x_train,y_train)  
  
Out[83]: GridSearchCV(cv=3, estimator=OneVsRestClassifier(estimator=LinearSVC()),  
n_jobs=-1,  
param_grid={'estimator__C': [2], 'estimator__dual': [False],  
'estimator__intercept_scaling': [2, 4, 5],  
'estimator__loss': ['hinge', 'squared_hinge'],  
'estimator__multi_class': ['ovr', 'crammer_singer'],  
'estimator__penalty': ['l1']})  
  
In [84]: 1 #printing the best parameters found by GridSearchCV  
2 GCV.best_params_  
  
Out[84]: {'estimator__C': 2,  
'estimator__dual': False,  
'estimator__intercept_scaling': 4,  
'estimator__loss': 'squared_hinge',  
'estimator__multi_class': 'ovr',  
'estimator__penalty': 'l1'}
```

- ❖ And after doing hyperparameter tuning I got above parameters as best suitable parameters for our final model.
- ❖ I have tested my final model using these parameters and got better results compared to earlier results for my final model.

Final Model:

```
1 model = OneVsRestClassifier(LinearSVC(C=2,dual = False, loss='hinge',multi_class='crammer_singer', penalty = 'l1',intercept_scaling=1))
2 model.fit(x_train,y_train)
3 y_pred = model.predict(x_test)
4
5 print("Jaccard score: {}".format(jaccard_score(y_test,y_pred,average='micro')))
6 print("Accuracy score: {}".format(accuracy_score(y_test,y_pred)))
7 print("f1_score: {}".format(f1_score(y_test,y_pred,average='micro')))
8 print("Precision : ", precision_score(y_test,y_pred,average='micro'))
9 print("Recall: {}".format(recall_score(y_test,y_pred,average='micro')))
10 print("Hamming loss: ", hamming_loss(y_test,y_pred))
11 print("\nConfusion matrix: \n", multilabel_confusion_matrix(y_test,y_pred))
```

Jaccard score: 0.5463542718642699
Accuracy score: 0.9179555310455468
f1_score: 0.7066353187042843
Precision : 0.8355212355212355
Recall: 0.6121987099694467
Hamming loss: 0.018766868038670108

Confusion matrix:
[[[35683 395]
[1249 2566]]

[[39426 61]
[327 79]]

[[37556 194]
[660 1483]]

[[39758 30]
[79 26]]

[[37552 330]
[845 1166]]

[[39481 55]
[267 90]]]

- ❖ After training and building our final model, I used this model to make predictions for test dataset. Before doing predictions, the test dataset has been cleaned and processed with the same functions which are used for train dataset. And then doing vectorization I have predicted the output labels with our final model.

2. Saving the model and Predicting Ticket Price for test data:

Model Saving:

```
1 import joblib
2 joblib.dump(model,"Malignant_comment.pkl")
```

```
87]: ['Malignant_comment.pkl']
```

```
1 #Loading the model
2 model = joblib.load('Malignant_comment.pkl')
```

- I have saved my best model using .pkl as follows.
- Now loading my saved model and predicting the test values.

- Now loading my saved model and predicting the values for test dataset.

```

1 pred=pd.DataFrame(predictions, columns = ['malignant','highly_malignant','rude','threat','abuse','loathe'])
2 pred

```

37]:

	malignant	highly_malignant	rude	threat	abuse	loathe
0	0	0	0	0	0	0
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0
...
153159	0	0	0	0	0	0
153160	0	0	0	0	0	0
153161	0	0	0	0	0	0
153162	0	0	0	0	0	0
153163	0	0	0	0	0	0

4.CONCLUSION

- **Key Findings and Conclusions of the Study:**

For this project we have provided with huge amount of comments with multiple targets which are binary in nature. I observe that there are many words with incorrect spellings. At first I have created three columns one is with the length of the text, another as 'question' whether the comment contains '?' mark or not and third as 'exclamation' whether the comment contains '!' mark. To clean the column comment_text I have gone through different text processing steps like lowercasing the text, removing unwanted elements like stopwords, '\n', Urls, numbers, punctuations etc. As the text column is with many miss-spelled words and the problem is multi-labelled so we are getting slightly lower accuracy for this task. However we have selected best model among all the algorithms. There are some comments which are from different language other than English we can try the same approach by removing those comments with other languages.

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

I found that the dataset was quite interesting to handle. Improvement in computing technology has made it possible to examine social information that cannot previously be captured, processed and analysed. New analytical techniques of machine learning can be used in property research. The power of visualization has helped us in understanding the data by graphical representation it has made me to understand what data is trying to say. Data cleaning is one of the most important steps to remove unrealistic values and stopwords. This study is an exploratory attempt to use four machine learning algorithms in estimating malignant comments, and then compare their results.

To conclude, the application of machine learning in malignant classification is still at an early stage. We hope this study has moved a small step ahead in providing some methodological and empirical contributions to crediting institutes, and presenting an alternative approach to the valuation of malignance.

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

Additionally, the followings are some suggested studies to be considered as future work in this area:

a) We suggest a plan to improve the NLP classifiers: first by using other algorithms which such as Support Vector Clustering (SVC) and Convolutional Neural Networks (CNN); secondly, extend the classifiers to the overall goal which is multi-label classifiers. in the current study, the problem simplified into two classes but it worth to pursue a main goal which is 6 classes of comments.

b) We also suggest using SVM for text processing and text classification. It requires a grid search for hyper-parameter tuning to get the best results.

5.REFERENCE

www.fliprobo.com