



NAME OF THE PROJECT

Malignant Comment Classifier

Submitted by:

Nomaan Sayed

## **ACKNOWLEDGMENT**

This includes mentioning of all the references, research papers, data sources, professionals and other resources that helped you and guided you in completion of the project.

# INTRODUCTION

- **Business Problem Framing**

The proliferation of social media enables people to express their opinions widely online. However, at the same time, this has resulted in the emergence of conflict and hate, making online environments uninviting for users. Although researchers have found that hate is a problem across multiple platforms, there is a lack of models for online hate detection. Online hate, described as abusive language, aggression, cyberbullying, hatefulness and many others has been identified as a major threat on online social media platforms. Social media platforms are the most prominent grounds for such toxic behaviour.

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

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

- **Conceptual Background of the Domain Problem**

Nowadays it is common that people show their hatred towards a particular person on social media by commenting in abusive language and other things.

- **Review of Literature**

It is common to abuse someone on social media, E-commerce site, twitter etc. We all have seen some or the other time, Machine learning helps us to get through this problem and detect these comments on real time.

## **Analytical Problem Framing**

- Data Sources and their formats  
Training and Testing data is provides separately.

The image displays two screenshots of a Jupyter Notebook interface, likely running on a local host (localhost:8888). The notebook is titled "Malignant comment classifier" and shows the process of loading and preprocessing data for a classification task.

**Top Screenshot:**

- In [2]:** `df = pd.read_csv('Malignant train.csv')`
- Out[2]:** A DataFrame with 159571 rows and 8 columns. The columns are: `id`, `comment_text`, `malignant`, `highly_malignant`, `rude`, `threat`, `abuse`, and `loathe`. The data is displayed in a table format with alternating light and dark rows.
- In [3]:** `df1 = pd.read_csv('Malignant test.csv')`
- Out[3]:** A DataFrame with 2 rows and 2 columns: `id` and `comment_text`.

**Bottom Screenshot:**

- In [3]:** `df1 = pd.read_csv('Malignant test.csv')`
- Out[3]:** A DataFrame with 153164 rows and 2 columns: `id` and `comment_text`. The data is displayed in a table format with alternating light and dark rows.
- In [4]:** `df.isnull().sum()`
- Out[4]:** `id 0`

- Data Preprocessing Done
  - 1- Load the data
  - 2- Checking null values
  - 3- Encoding dataset
  - 4- Data Visualization using seaborn

- 5- Describing dataset
- 6- Correlation of the dataset
- 7- Heatmap
- 8- Checking outliers and removing it
- 9- Transforming dataset
- 10- Scaling dataset for better understanding

- State the set of assumptions (if any) related to the problem under consideration  
No such assumptions are done.
- Hardware and Software Requirements and Tools Used

**Hardware- core i5 9<sup>th</sup> gen with 8gb ram**  
**Software-Jupyter Notebook by using Python**

### **Model/s Development and Evaluation**

- Testing of Identified Approaches (Algorithms)  
Logistic Regression is used
- Run and Evaluate selected models

```
import warnings
warnings.filterwarnings('ignore')

In [55]: for i in range(0,1000):
x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=i)
lr.fit(x_train,y_train)
pred_train=lr.predict(x_train)
pred_test=lr.predict(x_test)
if round(accuracy_score(y_train,pred_train)*100,1)==round(accuracy_score(y_test,pred_test)*100,1):
    print("At random state ",i,"The Model perform very well")
    print("At random state :-",i)
    print("Training r2_score is :-",accuracy_score(y_train,pred_train)*100)
    print("Testing r2_score is :-",accuracy_score(y_test,pred_test)*100)
```

At random state 4 The Model perform very well  
At random state :- 4  
Training r2\_score is :- 95.83966978057788  
Testing r2\_score is :- 95.79374827664002  
At random state 5 The Model perform very well  
At random state :- 5  
Training r2\_score is :- 95.82295827136149  
Testing r2\_score is :- 95.83134885819568  
At random state 11 The Model perform very well  
At random state :- 11  
Training r2\_score is :- 95.82713614866559  
Testing r2\_score is :- 95.8438823853809  
At random state 12 The Model perform very well  
At random state :- 12  
Training r2\_score is :- 95.84969668610772  
Testing r2\_score is :- 95.75113428421027  
At random state 17 The Model perform very well  
At random state :- 17  
Training r2\_score is :- 95.81961596951821

- Key Metrics for success in solving problem under consideration

```
In [56]: from sklearn.metrics import classification_report
print(classification_report(y_test,pred_test))
```

	precision	recall	f1-score	support
0	0.96	0.99	0.98	36007
1	0.92	0.62	0.74	3886
accuracy			0.96	39893
macro avg	0.94	0.81	0.86	39893
weighted avg	0.96	0.96	0.95	39893

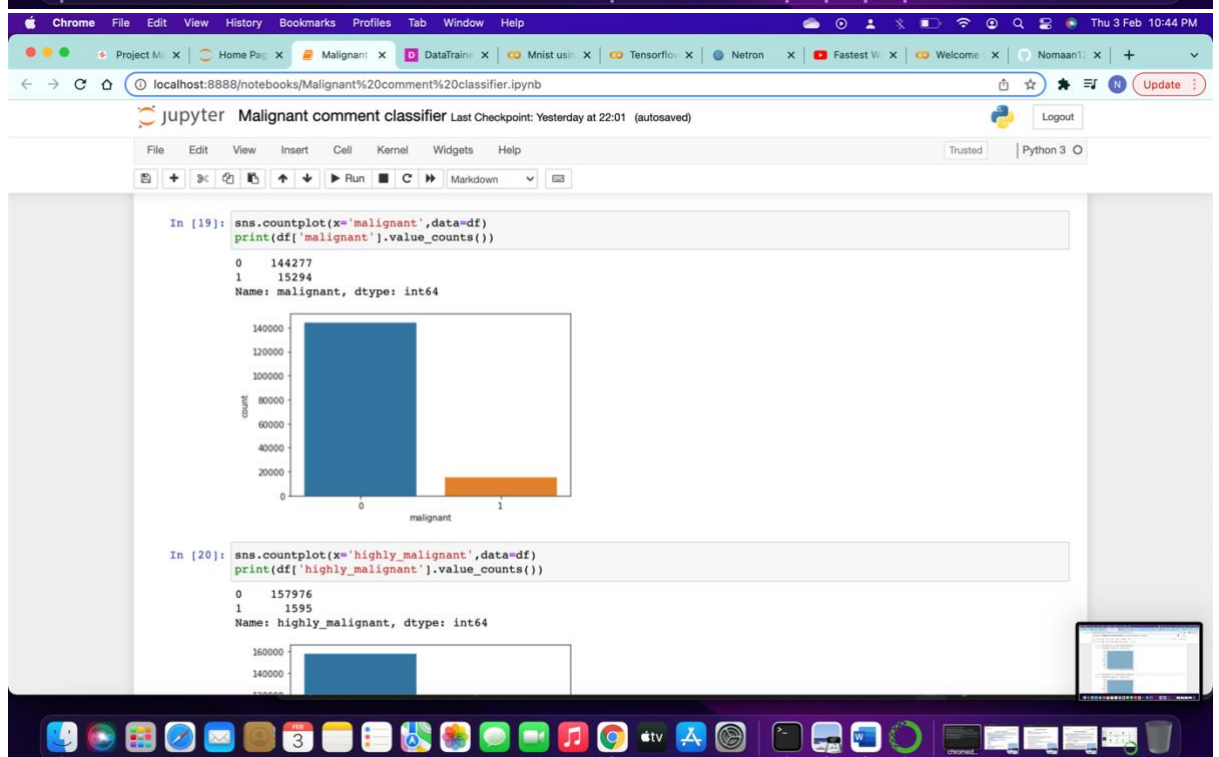
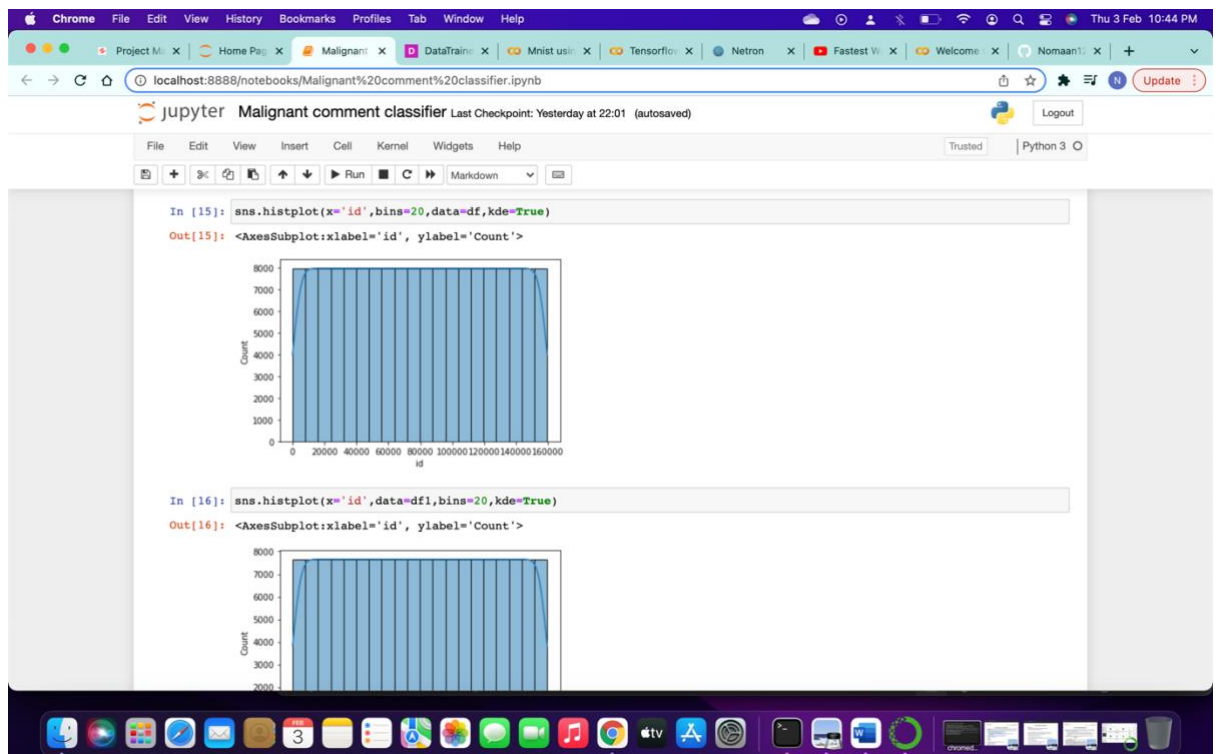
**Cross Validation of Dataset**

```
In [57]: pred_lr=lr.predict(x_test)
from sklearn.model_selection import cross_val_score
lss=accuracy_score(y_test,pred_lr)
for j in range(2,10):
    lsscore=cross_val_score(lr,x,y,cv=j)
    lsc=lsscore.mean()
    print("At cv :- ",j)
    print("Cross validation score is :-",lsc*100)
    print("accuracy_score is :-",lss*100)
    print("\n")
```

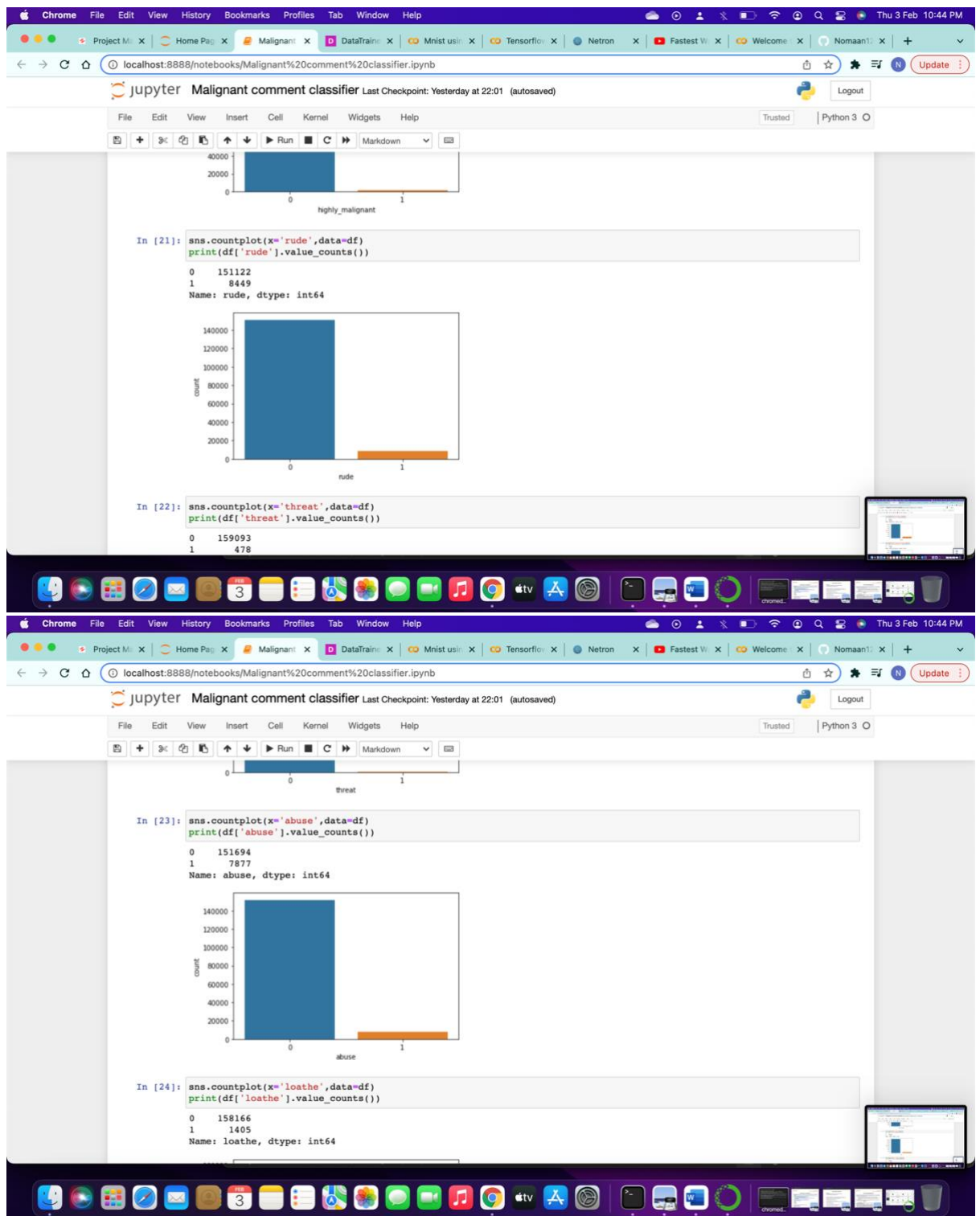
At cv :- 2  
Cross validation score is :- 95.82568263097457  
accuracy\_score is :- 95.73108064071391

At cv :- 3  
Cross validation score is :- 95.82568263097457  
accuracy\_score is :- 95.73108064071391

- Visualizations
- Interpretation of the Results







CONCLUSION

- Learning Outcomes of the Study in respect of Data Science

By using Data Science we can detect the real time malignant comments done which is very good. We can let our model thinks by using Datascience.

- Limitations of this work and Scope for Future Work

We can only detect it once it is commented, we cannot stop any person to comment from scratch.