



Malignant Comment Classifier Project

Submitted by:
Megha Singh

ACKNOWLEDGMENT

The internship opportunity I have with Flip Robo Technologies is a great chance for learning and professional development. I perceive this opportunity as a big milestone in my career development. I will strive to use gained skills acknowledge in the best possible way. I would like to extend my appreciation and thanks for the mentors from Data Trained and professionals from FlipRobo Technologies who had extended their help and support.

INTRODUCTION

- **Business Problem Framing**

With the proliferation of social media there has been an emergence of conflict and hate, making online environments uninviting for users. 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. 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**

Predictive modelling, Classification algorithms are some of the machine learning techniques used along with the various libraries of the NLTK suite for Classification of comments.

Using NLTK tools, the frequencies of malignant words occurring in textual data were estimated and given appropriate weightage, whilst filtering out words, and other noise which do not have any impact on the semantics of the comments and reducing the words to their base lemmas for efficient processing and accurate classification of the comments.

- **Review of Literature**

Two research papers titled: "Toxic Comment Classification" by Sara Zaheri and "Machine learning methods for toxic comment classification: a systematic review" by Darko Androcec were reviewed and studied to gain insights into the nature of malignant comments, their impact on social media platforms and the various

methods that are employed for training models to detect, identify and classify them.

- **Motivation for the Problem Undertaken**

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. Automatic recognition of malignant comments on online forums, and social media serves as a useful provision for moderators of public platforms as well as users who could receive warnings and filter unwanted contents. The need of advanced methods and techniques to improve identification of different types of comments posted online motivated the current project.

Analytical Problem Framing

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

Various Classification analysis techniques were used to build Classification models to determine whether an input Message content is benign or malignant. Machine Learning Algorithms such as Multinomial Naïve Bayes and Complement Naïve Bayes were employed which are based on the Bayes Theorem:

$$P(\text{message is malignant} \mid \text{message content}) = \frac{P(\text{message content} \mid \text{malignant}) \cdot P(\text{malignant})}{P(\text{message content})}$$

The probability of message being Malignant, knowing that Message Content has occurred could be calculated. Event of “Message Content” represents the evidence and “Message is Malignant”, the hypothesis to be approved. The theorem runs on the assumption that all predictors/features are independent and the presence of one would not affect the other.

The approach to classify a comment as malignant would depend on training data labelled as various categories of malignant messages and benign messages.

- **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.

	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

Figure 1 Train Dataset

	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.

Figure 2 Test Dataset

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.

- Data Preprocessing Done

The dataset was checked to see if there were any null values or random characters present. None were found.

Column: **ID** was dropped since they don't contribute to building a good model for predicting the target variable values.

```

1 # Convert all messages to lower case
2 trainDF['comment_text'] = trainDF['comment_text'].str.lower()
3
4
5 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'^(?!\.)*[a-z]{2,}$', 'emailaddress') # Replace email add
6
7 # Replace URLs with 'webaddress'
8 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'^http://[a-zA-Z0-9\-\.\+]\.[a-zA-Z]{2,3}(/[^/]*)?$', 'webaddress')
9
10
11 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'£|$', 'dollars') # Replace money symbols with 'moneysymb'
12
13 # Replacing 10 digit phone numbers with 'phonenumber'
14 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'^(?[\d]{3})?[\s-]?[\d]{3}[\s-]?[\d]{4}$', 'phonenumber')
15
16 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'\d+(\.\d+)?', 'num') # Replace numbers with 'num'
17
18
19 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'[^\w\d\s]', ' ') #removing punctuations
20
21 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'[_]', ' ') #removing underscore characters
22
23 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'[a-zA-Z]\s+', ' ') #removing single characters
24
25 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'\s+', ' ') #removing whitespace between terms with a single
26
27 trainDF['comment_text'] = trainDF['comment_text'].str.replace(r'^\s+|\s+$', ' ') #removing leading and trailing whitespace
28

```

```

1 trainDF.head()

```

	comment_text	malignant	highly_malignant	rude	threat	abuse	loathe	Stringlength
0	explanation why the edits made under my usern...	0	0	0	0	0	0	264
1	d aww he matches this background colour m seem...	0	0	0	0	0	0	112
2	hey man m really not trying to edit war it jus...	0	0	0	0	0	0	233
3	more can make any real suggestions on improve...	0	0	0	0	0	0	622
4	you sir are my hero any chance you remember wh...	0	0	0	0	0	0	87

```

1 import nltk
2 from nltk.corpus import stopwords,wordnet

```

```

1 from nltk.stem import WordNetLemmatizer

```

```

1 stop_words = set(stopwords.words('english') + ['u','m', 'u', 'ur', 'd', '2', 'im', 'dont', 'doin','u're", 'ure'])
2 trainDF['comment_text'] = trainDF['comment_text'].apply(lambda x: ' '.join(term for term in x.split() if term not in stop_wo

```

```

1 lem=WordNetLemmatizer()
2 trainDF['comment_text'] = trainDF['comment_text'].apply(lambda x: ' '.join(lem.lemmatize(t) for t in x.split()))

```

```

1 trainDF['Clesned_Stringlength'] = trainDF['comment_text'].str.len()
2 trainDF.head()

```

	comment_text	malignant	highly_malignant	rude	threat	abuse	loathe	Stringlength	Cleaned_Stringlength
0	explanation edits made username hardcore metal...	0	0	0	0	0	0	264	164
1	awww match background colour seemingly stuck th...	0	0	0	0	0	0	112	83
2	hey man really trying edit war guy constantly ...	0	0	0	0	0	0	233	141
3	make real suggestion improvement wondered sect...	0	0	0	0	0	0	622	364
4	sir hero chance remember page	0	0	0	0	0	0	87	29

The train and test dataset contents were then converted into lowercase. Punctuations, unnecessary characters etc were removed, currency symbols, phone numbers, web urls, email addresses etc were replaced with single words. Tokens that contributed nothing to semantics of the messages were removed

as Stop words. Finally retained tokens were lemmatized using WordNetLemmatizer().

The string lengths of original comments and the cleaned comments were then compared.

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

The comment tokens so vectorised using TfidfVectorizer are input and classified as benign(0) or malignant(1) as output by classification models.

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

The comment content made available in Train and Test Dataset is assumed to be written in English Language in the standard Greco-Roman script. This is so that the Stopword package and WordNetLemmatizer can be effectively used.

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

Hardware Used:

- Processor: Intel core i3-2348M, 2.3GHz
- Physical Memory: 4.0GB
- GPU: NVIDIA GeForce 710M, 2GB .

Software Used:

- Windows 10 Operating System
- Anaconda Package and Environment Manager:
Anaconda is a distribution of the Python and R programming languages for scientific computing, that aims to simplify package management and deployment. The distribution includes data science packages suitable for Windows and provides a host of tools and

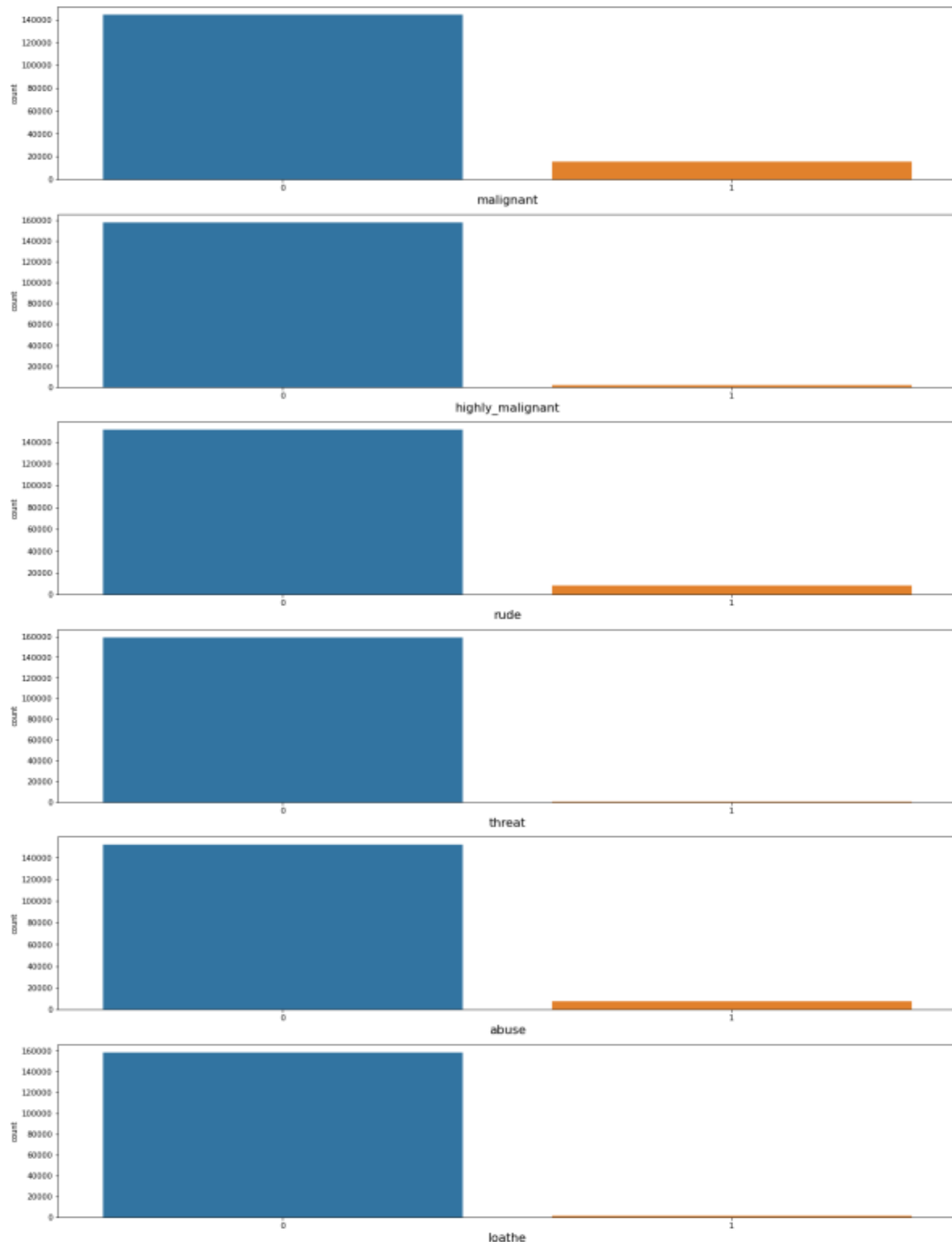
environment for conducting Data Analytical and Scientific works. Anaconda provides all the necessary Python packages and libraries for Machine learning projects.

- Jupyter Notebook: The Jupyter Notebook is an open-source web application that allows data scientists to create and share documents that integrate live code, equations, computational output, visualizations, and other multimedia resources, along with explanatory text in a single document.
- Python3: It is open source, interpreted, high level language and provides great approach for object-oriented programming. It is one of the best languages used for Data Analytics And Data science projects/application. Python provides numerous libraries to deal with mathematics, statistics and scientific function.
- Python Libraries used:
 - Pandas: For carrying out Data Analysis, Data Manipulation, Data Cleaning etc
 - Numpy: For performing a variety of operations on the datasets.
 - matplotlib.pyplot, Seaborn: For visualizing Data and various relationships between Feature and Label Columns
 - sklearn for Modelling Machine learning algorithms, Evaluation metrics, Data Transformation etc
 - imblearn.over_sampling: To employ SMOTE technique for balancing out the classes.
 - re, string: To perform regex operations
 - Wordcloud: For Data Visualization
 - NLTK: To use various Natural Language Processing Tools.

Exploratory Data Analysis Visualizations

Barplots, Countplots, Distplots, WordClouds were used to visualise the data of all the columns and their relationships with Target variable.

Analyzing the Feature Columns



From the graphs about it is observed that majority of the comments are benign.

Unprocessed vs Cleaned string lengths

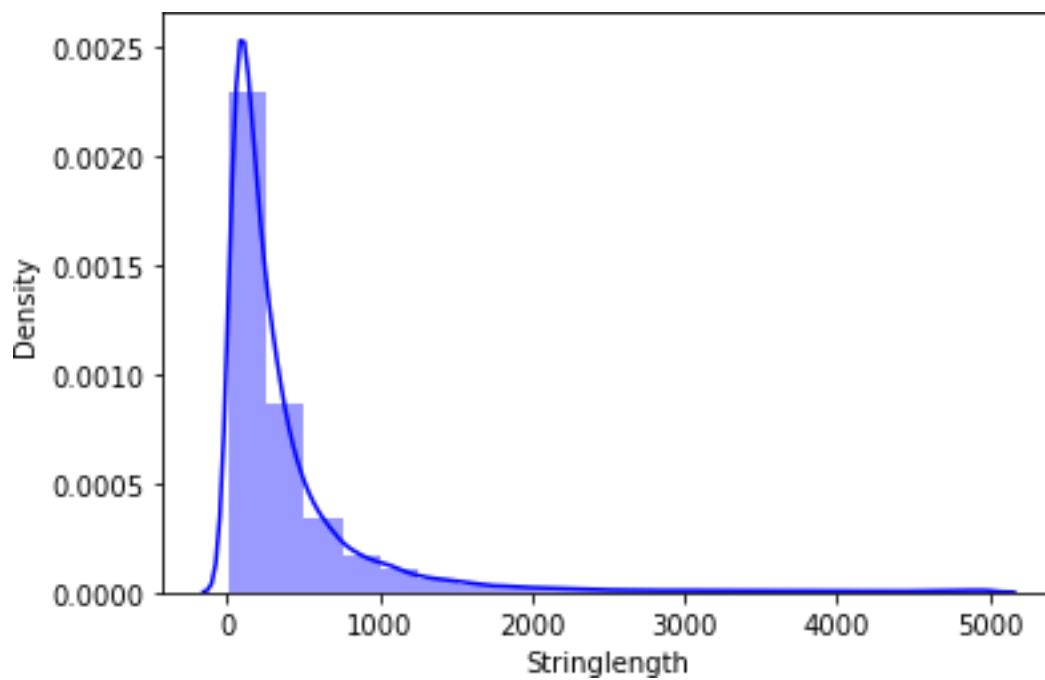


Figure 3 String Length of unprocessed comments

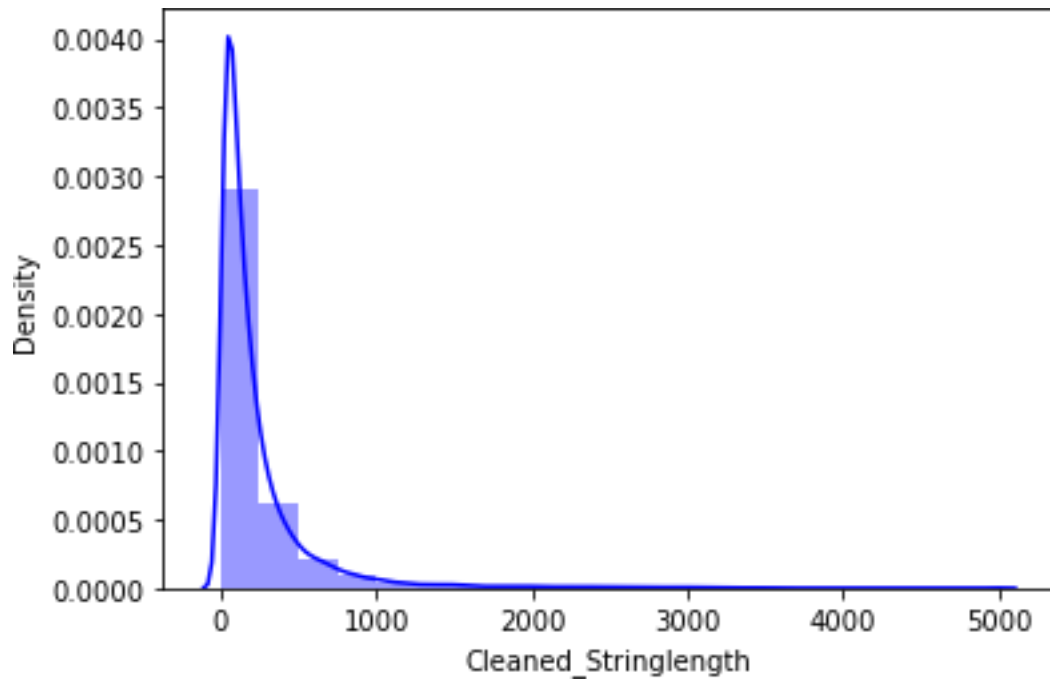
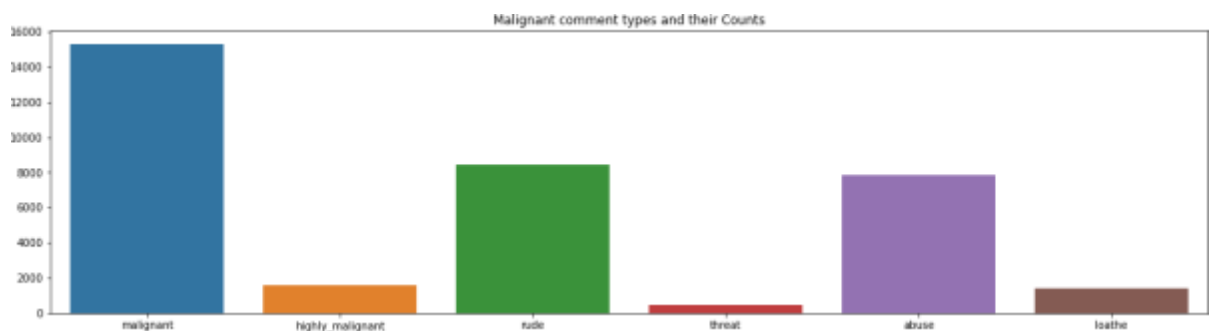


Figure 4 Cleaned Comments String Length

Above graphs show that the string length of comments was drastically brought down after processing.



The above graph shows the composition of toxic comments, of which majority are malignant followed by rude comments, abusive comments, highly malignant comments, hateful comments and threats.

Feature Engineering

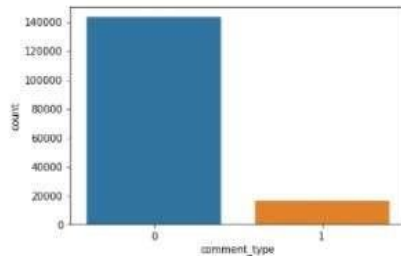
The comments data could belong to more than one label simultaneously (rude comments are at the same time malignant and in some cases can also be deemed hateful, abusive comments are hateful and can be highly malignant at the same time, threats are highly malignant too etc.)

Since each of the categories had very small data available to work with, a new column: „comment_type“ was created which only had binary classes: 0 which represented all the benign comments and 1 which represented all the comments which fell under malignant, highly malignant, abusive, hateful, rude, threat features. This column acted as Target Label column for malignant comment classification.

Visualising data in Target column

```
1 sns.countplot(trainDF['comment_type'])
```

<AxesSubplot:xlabel='comment_type', ylabel='count'>



```
1 print('Benign comment ratio = ',round(len(trainDF[trainDF['comment_type']==0])/len(trainDF.comment_type),2)*100,'%')
```

Benign comment ratio = 90.0 %

```
1 print('Malignant comment ratio = ',round(len(trainDF[trainDF['comment_type']==1])/len(trainDF.comment_type),2)*100,'%')
```

Malignant comment ratio = 10.0 %

Classes are imbalanced

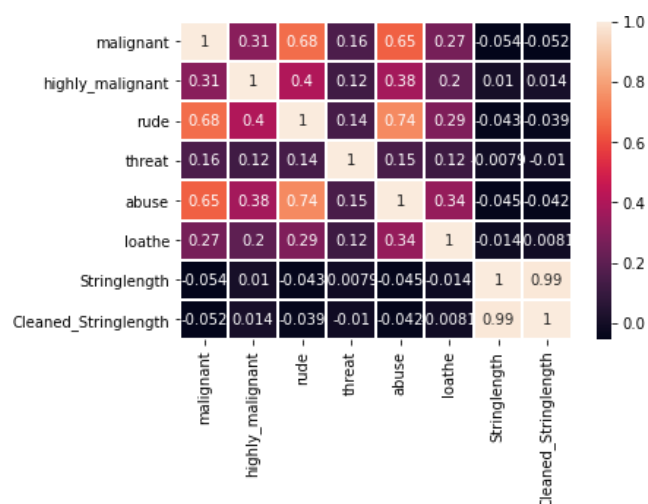
The classes appear to be imbalanced with 90% of comments being benign (0) and only 10% being malignant (1).

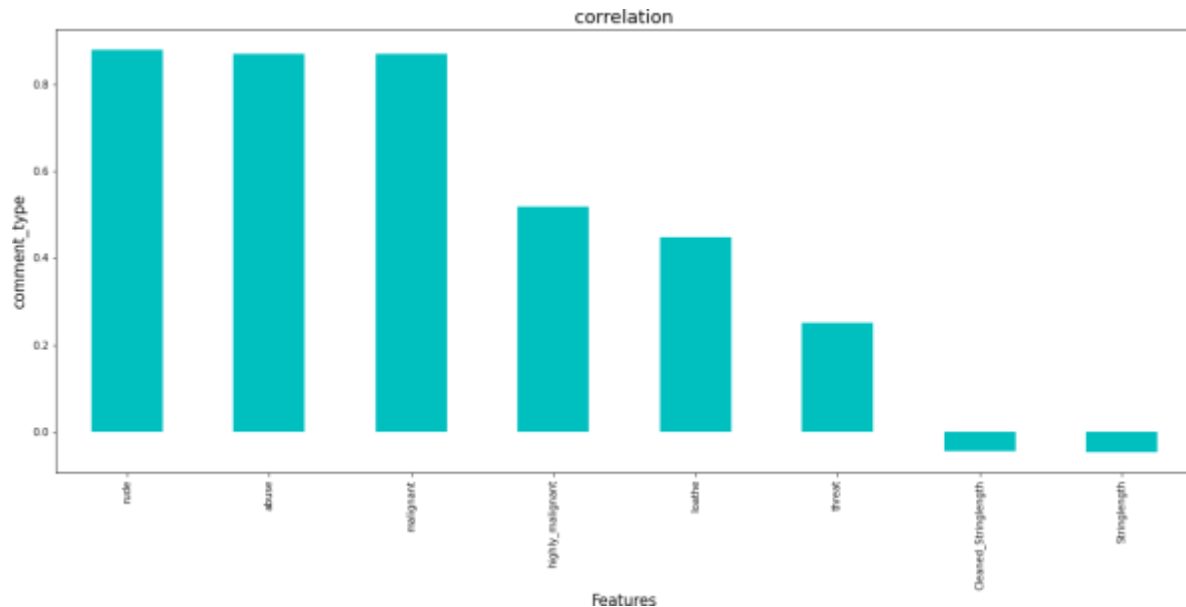
Smote Technique was used to balance out the classes

Balancing out classes in Label column using SMOTE technique.

```
1 from imblearn.over_sampling import SMOTE as sm
2
3 smt_x,smt_y = sm().fit_resample(X,y)
```

Finding Correlation





From the graphs above it is observed that columns: Rude, Abuse, Malignant have highest positive correlation with comment_type.

Model/s Development and Evaluation

- Identification of possible problem-solving approaches (methods)

The model algorithms used were as follows:

- Logistic Regression: It is a classification algorithm used to find the probability of event success and event failure. It is used when the dependent variable is binary(0/1, True/False, Yes/No) in nature. It supports categorizing data into discrete classes by studying the relationship from a given set of labelled data. It learns a linear relationship from the given dataset and then introduces a non-linearity in the form of the Sigmoid function. It not only provides a measure of how appropriate a predictor(coefficient size)is, but also its direction of association (positive or negative).
- Multinomial Naïve Bayes Classifier: Multinomial Naive Bayes algorithm is a probabilistic learning method that is mostly used in Natural Language Processing (NLP). The algorithm is based on the Bayes theorem. It calculates the probability of each tag for a

given sample and then gives the tag with the highest probability as output.

- **RandomForestClassifier:** A random forest is a meta estimator that fits a number of classifying decision trees on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. A random forest produces good predictions that can be understood easily. It reduces overfitting and can handle large datasets efficiently. The random forest algorithm provides a higher level of accuracy in predicting outcomes over the decision tree algorithm.
- **Complement Naïve Bayes Classifier:** Complement Naive Bayes is somewhat an adaptation of the standard Multinomial Naive Bayes algorithm. Complement Naive Bayes is particularly suited to work with imbalanced datasets. In complement Naive Bayes, instead of calculating the probability of an item belonging to a certain class, we calculate the probability of the item belonging to all the classes.
- **Passive Aggressive Classifier:** Passive-Aggressive algorithms do not require a learning rate and are called so because if the prediction is correct, keep the model and do not make any changes. i.e., the data in the example is not enough to cause any changes in the model. If the prediction is incorrect, make changes to the model. i.e., some change to the model may correct it.
- **AdaBoost Classifier:** The basis of this algorithm is the [Boosting](#) main core: give more weight to the misclassified observations. the meta-learner adapts based upon the results of the weak classifiers, giving more weight to the misclassified observations of the last weak learner. The individual learners can be weak, but as long as the performance of each weak learner is better than random guessing, the final model can converge to a strong learner (a learner not influenced by outliers and with a great generalization power, in order to have strong performances on unknown data).

Best Random state was found to be 56

```
1 from sklearn.naive_bayes import MultinomialNB
2 maxAcc = 0
3 maxRS=0
4 for i in range(0,100):
5     x_train,x_test,y_train,y_test = train_test_split(smt_x,smt_y,test_size = .30, random_state = i)
6     modRF = MultinomialNB()
7     modRF.fit(x_train,y_train)
8     pred = modRF.predict(x_test)
9     acc = accuracy_score(y_test,pred)
10    if acc>maxAcc:
11        maxAcc=acc
12        maxRS=i
13 print(f"Best Accuracy is: {maxAcc} on random_state: {maxRS}")
```

Best Accuracy is: 0.909566551948656 on random_state: 56

Training the Models

```
1 RFC.fit(x_train,y_train)
2 XGBC.fit(x_train,y_train)
3 adbc.fit(x_train,y_train)
4 LOGR.fit(x_train,y_train)
5 MNB.fit(x_train,y_train)
6 CNB.fit(x_train,y_train)
```

```
1 pc.fit(x_train,y_train)
PassiveAggressiveClassifier()
```

All Models have been trained.

Model Cross Validation

Cross validation is a technique for assessing how the statistical analysis generalises to an independent data set. It is a technique for evaluating machine learning models by training several models on subsets of the available input data and evaluating them on the complementary subset of the data. Using cross-validation, there are high chances that we can detect over-fitting with ease. Model Cross Validation scores were then obtained for assessing how the statistical analysis generalises to an independent data set. The models were evaluated by training several models on subsets of the available input data and evaluating them on the complementary subset of the data.

Hyper Parameter Tuning

GridSearchCV was used for Hyper Parameter Tuning of the Random Forest Classifier model.

```

1 parameter = {'n_estimators': [50, 100, 300], 'max_depth': [10, 20], 'min_samples_leaf': [2, 5, 30], 'min_samples_split': [2, 5, 10], 'criterion': ['gini', 'entropy']}
2
3 GridCV = GridSearchCV(RandomForestClassifier(), parameter, cv=5, n_jobs = -1, verbose = 1)
4
5 GridCV.fit(x_train, y_train)
6
7 GridCV.best_params_
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1
```

After Tuning the hyper parameters and based on the input parameter values and after fitting the train datasets it is found that Logistic Regression model performs the best.

The model was saved and the Test Dataset was then prepared for final classification work by the model. This model was then tested using the Test Dataset. The model performed with good amount of accuracy.

Comment Classification 1 means Malignant and 0 means Benign

```
1 Prediction_accuracy = pd.DataFrame({'Classification': mod.predict(test), 'Comments': testDF['comment_text']})  
2 Prediction_accuracy.head()
```

	Classification	Comments
0	1	yo bitch ja rule succesful ever whats hating s..
1	0	rfo title fine imo
2	0	source zawe ashton lapland
3	0	look back source information updated correct f..
4	0	anonymously edit article

CONCLUSION

- Key Findings and Conclusions of the Study

The final model offered 1.03% performance boost over the benchmark logistic regression model.

The Model has 95.72% accuracy. But since the dataset was highly imbalanced that is not the best metric for measuring its efficiency. Recall score of 0.93 for Benign (0) and 0.98 for Malignant(1), on the other hand, means that the model is optimized better to detect actual malignant comments. However, there is a need to strike a balance between precision and recall and have low false positives, which unnecessarily consume time and low false negatives which means only very few toxic comments deceive the model. F1 score of 0.96 provides a nuanced way to catch positive results without harming the usefulness of the model.

- Learning Outcomes of the Study in respect of Data Science

The various data pre-processing and feature engineering steps in the project lent cognizance to various efficient methods for processing textual data. The NLTK suite is very useful in pre-processing text-based data and building classification models.

- Limitations of this work and Scope for Future Work

The models were trained on a highly imbalanced dataset where the total malignant comments formed only 10% of the entire available data, which seriously affected the training and accuracy of the models. By training the models on more diverse data sets, longer comments, and a more balanced dataset, more accurate and efficient classification models can be built.

Thank You