



A project on Spam Email classification

Submitted by:

Sankalp Mahapatra

Internship-29

ACKNOWLEDGMENT

I would like to express my sincere thanks of gratitude to my SME as well as “Flip Robo Technologies” team for letting me work on “Email spam detection classifier” project. Their suggestions and directions have helped me in the completion of this project successfully. This project also helped me in doing lots of research wherein I came to know about so many new things.

Finally, I would like to thank my family and friends who have helped me with their valuable suggestions and guidance and have been very helpful in various stages of project completion.

TABLE OF CONTENTS:

1. Introduction

- 1.1 Business Problem Framing
- 1.2 Conceptual Background of the Domain Problem
- 1.3 Review of Literature
- 1.4 Motivation for the Problem Undertaken

2. Analytical Problem Framing

- 2.1 Mathematical/ Analytical Modelling of the Problem
- 2.2 Data Sources and their formats
- 2.3 Data Pre-processing Done
- 2.4 Data Inputs- Logic- Output Relationships
- 2.5 Hardware & Software Requirements & Tools Used

3. Model/s Development and Evaluation

- 3.1 Identification of possible Problem-solving approaches (Methods)
- 3.2 Visualizations
- 3.3 Testing of Identified Approaches (Algorithms)
- 3.4 Run and Evaluate Selected Models
- 3.5 Key Metrics for success in solving problem under consideration
- 3.6 Interpretation of the Results

4. Conclusion

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

1. INTRODUCTION

1.1 Business Problem Framing

You were recently hired in a Start-up Company and was asked to build a system to identify spam emails. We will explore and understand the process of classifying Emails as Spam or Not Spam by build Machine Learning and NPL model to detect the HAM and SPAM mails. The model will detect the unsolicited and unwanted emails and thus we can prevent them from creeping into user's inbox and therefore, increase the user Experience.

1.2 CONCEPTUAL BACKGROUND OF THE DOMAIN PROBLEM

As we know how a machine translates language, or how voice assistants respond to questions, or how mail gets automatically classified into spam or not spam, all these tasks are done through Natural Language Processing (NLP), which processes text into useful insights that can be applied to future data. In the field of artificial intelligence, NLP is one of the most complex areas of research due to the fact that text data is contextual. It needs modification to make it machine-interpretable and requires multiple stages of processing for feature extraction.

Classification problems can be broadly split into two categories: binary classification problems, and multi-class classification problems. Binary classification means there are only two possible label classes, e.g. a patient's condition is cancerous or it isn't, or a financial transaction is fraudulent or it is not. Multi-class classification refers to cases where there are more than two label classes. An example of this is classifying the sentiment of a movie review into positive, negative, or neutral.

There are many types of NLP problems, and one of the most common types is the classification of strings. Examples of this include the classification of movies/news articles into different genres and the automated classification of emails into a spam or not spam. We shall be looking into this last example in more detail for this project.

1.3 REVIEW OF LITERATURE

In recent times, unwanted commercial / promotional bulk emails also known as spam has become a huge problem on the internet and for our mail inbox. An individual / organization sending the spam messages are referred to as the spammers. Such a person gathers email addresses from different websites, chatrooms, and other sources to send the mail to bulk audience. Spam prevents the user from making full and good use of time, storage capacity and network bandwidth. The huge volume of spam mails flowing through the computer networks have destructive effects on the memory space of email servers, communication bandwidth, CPU power and user time. The menace of spam email is on the increase on yearly basis and is responsible for over 80% of the whole global email traffic (Source google).

Users who receive spam emails that they did not request find it very irritating. It is also resulted to untold financial loss to many users who have fallen victim of internet scams and other fraudulent practices of spammers who send emails pretending to be from reputable companies with the intention to persuade individuals to disclose sensitive personal information like passwords, Bank Verification Number (BVN) and credit card numbers.

1.4 MOTIVATION FOR THE PROBLEM UNDERTAKEN

Motivation for this project has been undertaken because it is a project which is assigned to me during my internship at Flip Robo Technologies. This project will help Start-up companies to detect and filter the SPAM mails in their Email inbox and therefore, increase the user experience and save their server from unwanted mails, phishing mails or other viruses.

2. ANALYTICAL PROBLEM FRAMING

2.1 MATHEMATICAL/ ANALYTICAL MODELING OF THE PROBLEM

Throughout the project multiple mathematical and analytical models have been used, first we have checked the ratio of spam and ham emails in our dataset. The shape of our data set is 2893 rows and 3 columns. Then we have used regular expressions to clean the message column which contained body of the email. Then we have used TfidfVectorizer, to transforms text to feature vectors that can be used as input to estimator.

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   v1               5572 non-null   object
1   v2               5572 non-null   object
2   Unnamed: 2       50 non-null     object
3   Unnamed: 3       12 non-null     object
4   Unnamed: 4       6 non-null      object
dtypes: object(5)
memory usage: 217.8+ KB
```

2.2 DATA SOURCES AND THEIR FORMATS

The data was provided to us from the FlipRobo Technologies as a part of our Internship assignment. The data was provided in CSV format and there were 3 attributes and 5572 rows in the data set.

```

1
2 #Let's Load the CSV file
3
4 df=pd.read_csv("https://raw.githubusercontent.com/Sankalpmahapatra10/-Email-Spam-Classifer-project/main/spam.csv", encoding
5 df.head()

```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only ...	NaN	NaN	NaN
1	ham	Ok lar... Joking wif u oni...	NaN	NaN	NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	NaN	NaN	NaN
3	ham	U dun say so early hor... U c already then say...	NaN	NaN	NaN
4	ham	Nah I don't think he goes to usf, he lives aro...	NaN	NaN	NaN

```
1 df.tail()
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
5567	spam	This is the 2nd time we have tried 2 contact u...	NaN	NaN	NaN
5568	ham	Will I_b going to esplanade fr home?	NaN	NaN	NaN
5569	ham	Pity, * was in mood for that. So...any other s...	NaN	NaN	NaN
5570	ham	The guy did some bitching but I acted like i'd...	NaN	NaN	NaN
5571	ham	Rofl. Its true to its name	NaN	NaN	NaN

2.3 DATA PREPROCESSING DONE

After loading all the data we will proceed with the data pre-processing.

Following Steps were followed during data pre-processing:

Removing unwanted attribute from Dataset :

It's quite hard to find whether a mail is a spam or not just by looking at the subject. So we started by replacing the null values.

```

1
2 #Let's check if there are null values in our dataset
3
4 df.isnull().sum()

```

```

v1          0
v2          0
Unnamed: 2   5522
Unnamed: 3   5560
Unnamed: 4   5566
dtype: int64

```

```

1 #Lets delete thee null values in the Unnamed:2,Unnamed: 3 and Unnamed: 4 columns
2 df.drop(columns=['Unnamed: 2','Unnamed: 3', 'Unnamed: 4'], inplace = True)

```

```

1 #Let's check the columns in our dataset
2
3 df.columns

```

```
Index(['v1', 'v2'], dtype='object')
```

```

1 #Let's re check that we have sucessfully imputed all the null values
2
3 df.isnull().sum()

```

```

v1    0
v2    0
dtype: int64

```

Adding additional attribute :

In order to analyse the data in a better way while doing pre-processing, we have added an attribute 'Length' which shows length of the message against it. This was done just to compare the length of text before and after preprocessing and to get idea about the memory optimization.

```
1 df['num_words'] = df['text'].apply(lambda x : len(nltk.word_tokenize(x)))
```

```
1 df.head()
```

	target	text	num_characters	num_words
0	0	Go until jurong point, crazy.. Available only ...	111	24
1	0	Ok lar... Joking wif u oni...	29	8
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37
3	0	U dun say so early hor... U c already then say...	49	13
4	0	Nah I don't think he goes to usf, he lives aro...	61	15

```
1 df['num_sentences'] = df['text'].apply(lambda x : len(nltk.sent_tokenize(x)))
```

```
1 df.head()
```

	target	text	num_characters	num_words	num_sentences
0	0	Go until jurong point, crazy.. Available only ...	111	24	2
1	0	Ok lar... Joking wif u oni...	29	8	2
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37	2
3	0	U dun say so early hor... U c already then say...	49	13	1
4	0	Nah I don't think he goes to usf, he lives aro...	61	15	1

Converting all the messages to lower case:

All messages in the 'message' attribute was converted to small case since keeping words in large case does not make sense as same word with small and large case conveys same meaning.

```
1 def transform_text(text):
2     text = text.lower()
3     text = nltk.word_tokenize(text)
4
5     y = []
6     for i in text:
7         if i.isalnum():
8             y.append(i)
9
10    text = y[:] # Copying List we have to clone it
11    y.clear()
12
13    for i in text:
14        if i not in stopwords.words('english') and i not in string.punctuation:
15            y.append(i)
16    text = y[:]
17    y.clear()
18
19    for i in text:
20        y.append(ps.stem(i))
21
22    return " ".join(y)
```

```
1 transform_text('Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat')
```

```
'go jurong point crazi avail bugi n great world la e buffet cine got amor wat'
```

```
1 df['transformed_text'] = df['text'].apply(transform_text)
```

```
1 df.head()
```

	target	text	num_characters	num_words	num_sentences	transformed_text
0	0	Go until jurong point, crazy.. Available only ...	111	24	2	go jurong point crazi avail bugi n great world...
1	0	Ok lar... Joking wif u oni...	29	8	2	ok lar joke wif u oni
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37	2	free entri 2 wkli comp win fa cup final tkt 21...
3	0	U dun say so early hor... U c already then say...	49	13	1	u dun say earli hor u c already say
4	0	Nah I don't think he goes to usf, he lives aro...	61	15	1	nah think goe usf live around though

DATA INPUTS- LOGIC- OUTPUT RELATIONSHIPS

We have analysed the words that were present in the spam and ham mails, based on the words present and the data we already have which says if the mail is ham or spam, we are going to train the model to predict the same.

HARDWARE AND SOFTWARE REQUIREMENTS AND TOOLS USED

sankalp
TUF Gaming FX505GT_FX505GT

Rename this PC

① Device specifications

Copy ^

Device name	sankalp
Processor	Intel(R) Core(TM) i5-9300H CPU @ 2.40GHz 2.40 GHz
Installed RAM	8.00 GB (7.85 GB usable)
Device ID	58BCD2FF-BCE9-40CA-A5CA-FBDAB9E81B18
Product ID	00327-35855-14973-AAOEM
System type	64-bit operating system, x64-based processor
Pen and touch	No pen or touch input is available for this display

Related links

[Domain or workgroup](#)

[System protection](#)

[Advanced system settings](#)

Windows specifications

Copy ^

Edition	Windows 11 Home Single Language
Version	21H2
Installed on	28-12-2021
OS build	22000.1219
Experience	Windows Feature Experience Pack 1000.22000.1219.0

[Microsoft Services Agreement](#)
[Microsoft Software License Terms](#)

LIBRARIES:

```
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 %matplotlib inline
7 import seaborn as sns
8
9
10
11 import nltk, re
12 import string
13 from nltk.corpus import stopwords
14 from collections import Counter
15 from nltk.corpus import stopwords
16 from wordcloud import WordCloud
17 from nltk.stem.porter import PorterStemmer
18 from sklearn.feature_extraction.text import CountVectorizer
19 from nltk.stem import WordNetLemmatizer, SnowballStemmer
20 from sklearn.feature_extraction.text import TfidfVectorizer
21
22 from sklearn.model_selection import train_test_split
23 from sklearn.linear_model import LogisticRegression
24 from sklearn.naive_bayes import MultinomialNB
25 from sklearn.tree import DecisionTreeClassifier
26 from sklearn.neighbors import KNeighborsClassifier
27 from sklearn.ensemble import RandomForestClassifier
28 from sklearn.ensemble import AdaBoostClassifier
29 from sklearn.ensemble import BaggingClassifier
```

MODEL/S DEVELOPMENT AND EVALUATION

IDENTIFICATION OF POSSIBLE PROBLEM-SOLVING APPROACHES (METHODS)

As the target column was Bivariant data and the algorithm that we choose depends on this target variable. So, we have chosen classification analysis for this project.

TESTING OF IDENTIFIED APPROACHES (ALGORITHMS)

We have used the following algorithms

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.svm import SVC
3 from sklearn.naive_bayes import MultinomialNB
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.neighbors import KNeighborsClassifier
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.ensemble import AdaBoostClassifier
8 from sklearn.ensemble import BaggingClassifier
9 from sklearn.ensemble import ExtraTreesClassifier
10 from sklearn.ensemble import GradientBoostingClassifier
11 from xgboost import XGBClassifier
```

```
1 svc = SVC(kernel='sigmoid', gamma=1.0)
2 knn = KNeighborsClassifier()
3 mnb = MultinomialNB()
4 dtc = DecisionTreeClassifier(max_depth=5)
5 lrc = LogisticRegression(solver='liblinear', penalty = 'l1')
6 rfc = RandomForestClassifier(n_estimators=50, random_state=2)
7 abc = AdaBoostClassifier(n_estimators=50, random_state=2)
8 bc = BaggingClassifier(n_estimators=50, random_state=2)
9 etc = ExtraTreesClassifier(n_estimators=50, random_state=2)
10 gbd = GradientBoostingClassifier(n_estimators=50, random_state=2)
11 xgb = XGBClassifier(n_estimators=50, random_state=2)
```

RUN AND EVALUATE SELECTED MODELS

```
1 clfs = {
2     'SVC': svc,
3     'KN': knn,
4     'NB': nb,
5     'DT': dtc,
6     'LR': lrc,
7     'RF': rfc,
8     'AdaBoost': abc,
9     'BgC': bc,
10    'ETC': etc,
11    'GBDT': gbd,
12    'xgb': xgb
13 }
```

```
1 def train_classifier(clf, X_train, y_train):
2     clf.fit(X_train, y_train)
3     y_pred = clf.predict(X_test)
4     accuracy = accuracy_score(y_test, y_pred)
5     precision = precision_score(y_test, y_pred)
6
7     return accuracy, precision
```

```
1
2 train_classifier(svc, X_train, y_train)
```

(0.9758220502901354, 0.9747899159663865)

```
1 accuracy_scores = []
2 precision_scores = []
3
4 for name, clf in clfs.items():
5     current_accuracy, current_precision = train_classifier(clf, X_train, y_train)
6     print("For", name)
7     print("Accuracy", current_accuracy)
8     print("Precision", current_precision)
9
10
11     accuracy_scores.append(current_accuracy)
12     precision_scores.append(current_precision)
```

For SVC
Accuracy 0.9758220502901354
Precision 0.9747899159663865
For KN
Accuracy 0.9052224371373307
Precision 1.0
For NB
Accuracy 0.9709864603481625
Precision 1.0
For DT
Accuracy 0.9274661508704062
Precision 0.8118811881188119
For LR
Accuracy 0.9584139264990329
Precision 0.9702970297029703
For RF
Accuracy 0.9758220502901354
Precision 0.9829059829059829
For AdaBoost
Accuracy 0.960348162475822
Precision 0.9292035398230089
For BgC
Accuracy 0.9584139264990329
Precision 0.8682170542635659
For ETC
Accuracy 0.9748549323017408

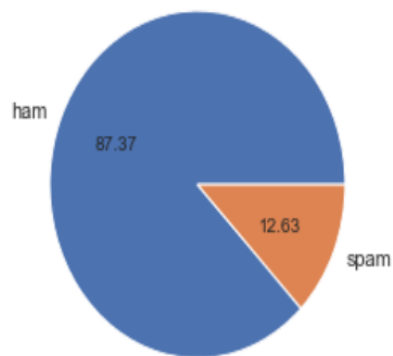
For ETC
Accuracy 0.9748549323017408
Precision 0.9745762711864406
For GBDT
Accuracy 0.9468085106382979
Precision 0.9191919191919192
For xgb
Accuracy 0.9671179883945842
Precision 0.9333333333333333


```
1 y_pred = voting.predict(X_test)
2 print("Accuracy", accuracy_score(y_test, y_pred))
3 print("Precision", precision_score(y_test,y_pred))
```

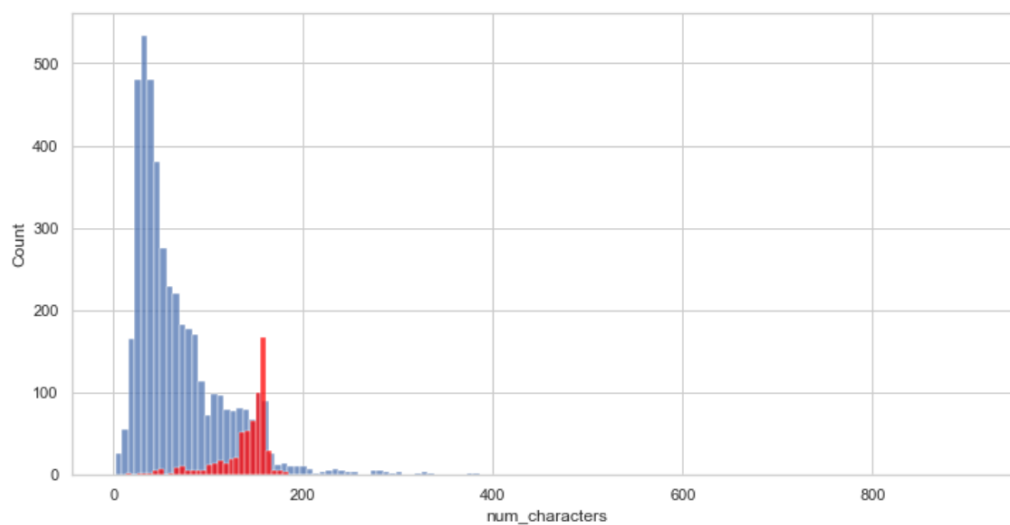
Accuracy 0.9816247582205029
Precision 0.9917355371900827

Visualizations

```
1 plt.pie(df['target'].value_counts(), labels= ['ham', 'spam'], autopct='%0.2f')
2 plt.show()
```



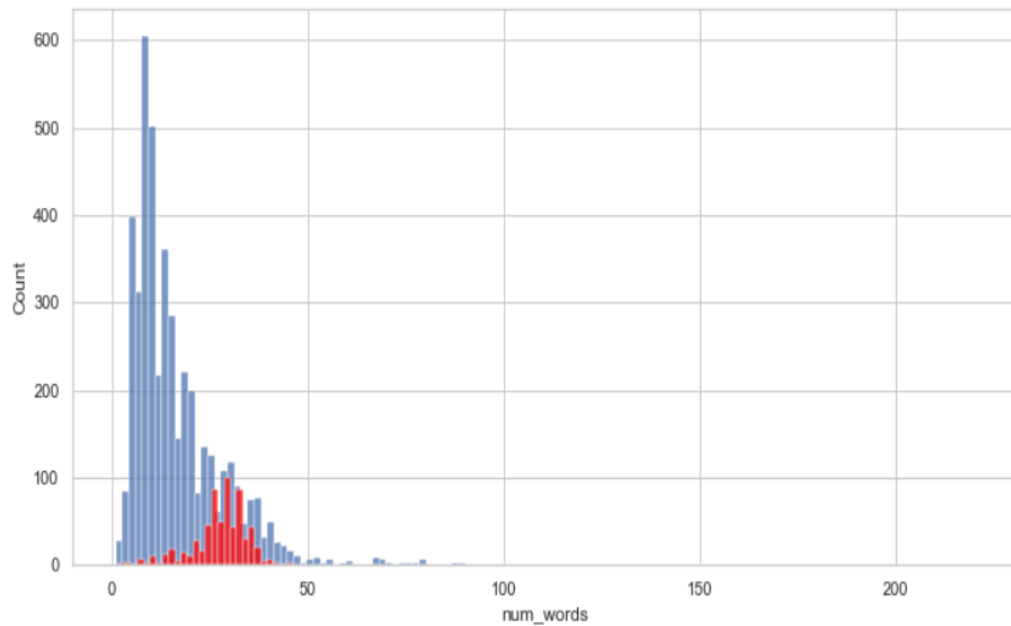
```
1 plt.figure(figsize =(12,6))
2 sns.histplot(df[df['target'] == 0]['num_characters']);
3 sns.histplot(df[df['target'] == 1]['num_characters'], color = 'red');
```



```

1 plt.figure(figsize =(12,6))
2 sns.histplot(df[df['target'] == 0]['num_words']);
3 sns.histplot(df[df['target'] == 1]['num_words'], color = 'red');

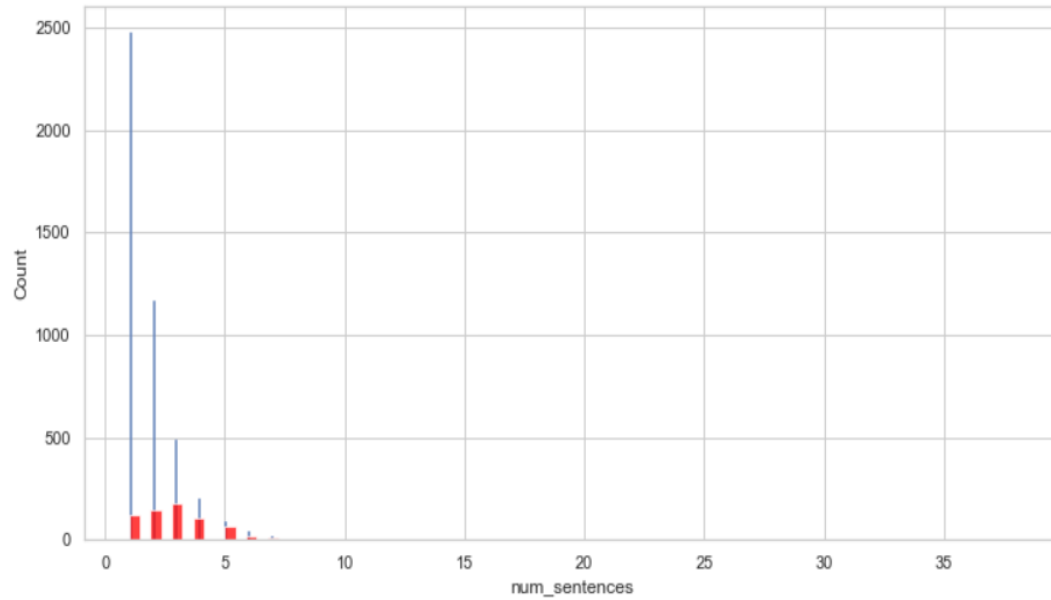
```



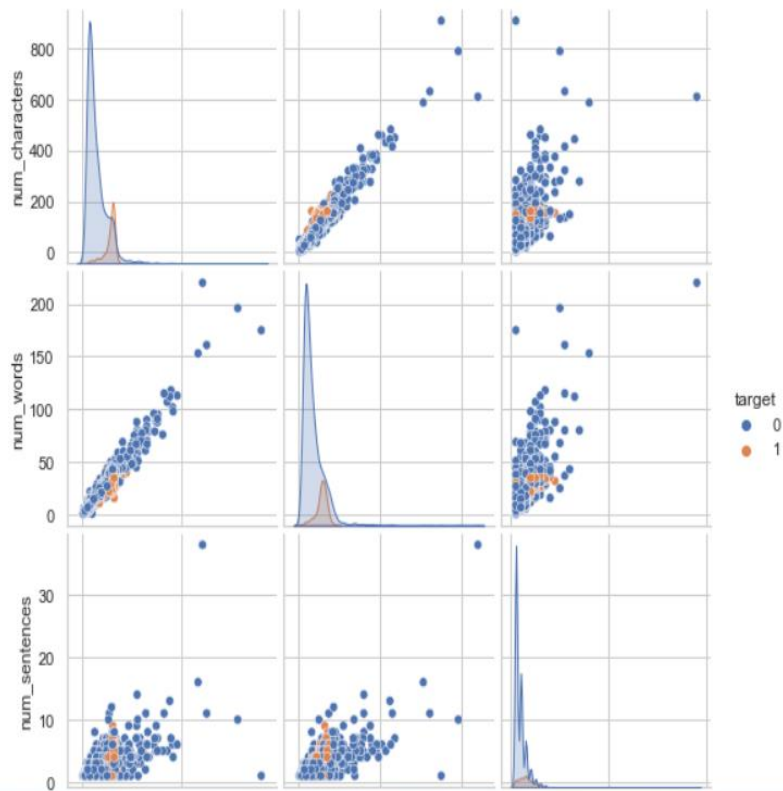
```

1 plt.figure(figsize =(12,6))
2 sns.histplot(df[df['target'] == 0]['num_sentences']);
3 sns.histplot(df[df['target'] == 1]['num_sentences'], color = 'red');

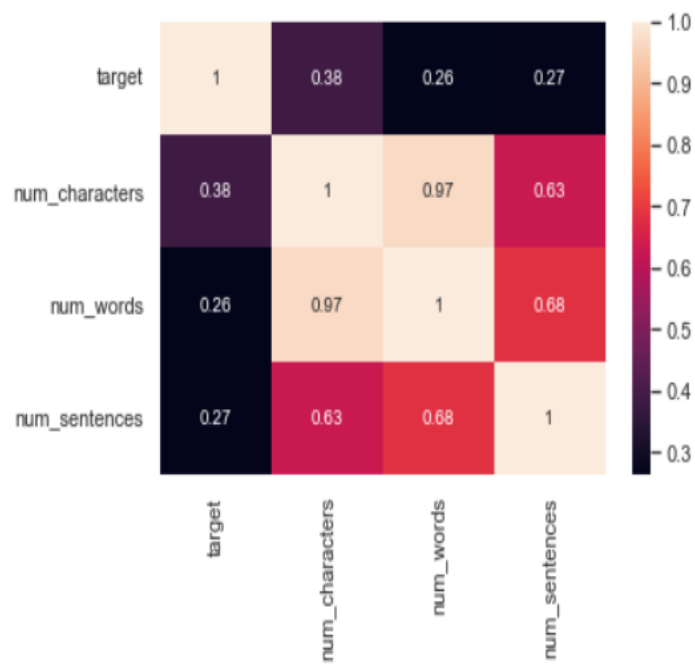
```




```
1 sns.pairplot(df,hue='target');
```

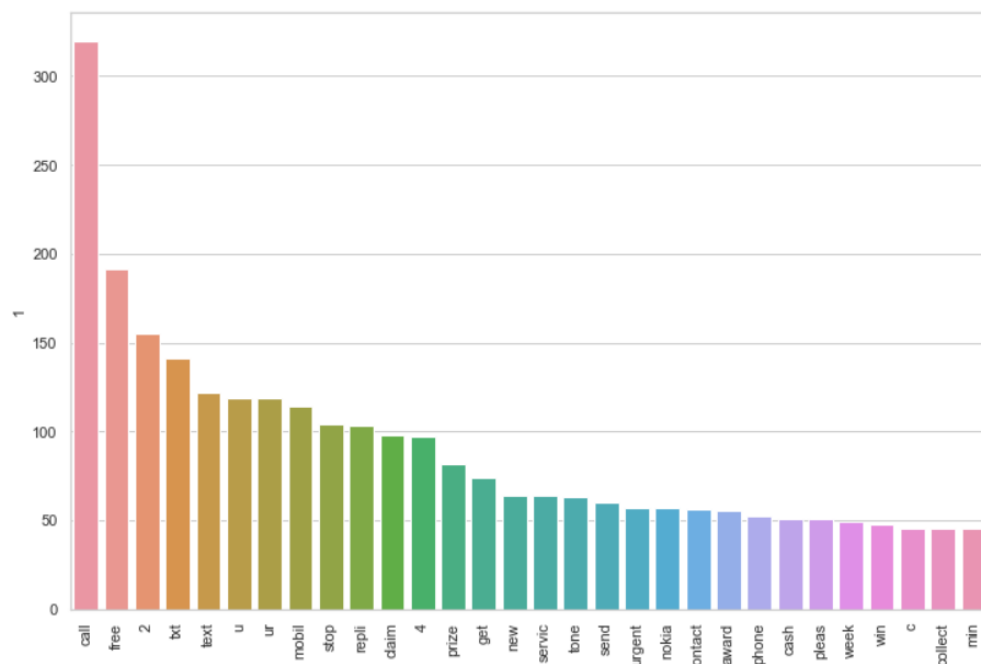


```
1 sns.heatmap(df.corr(), annot = True);
```





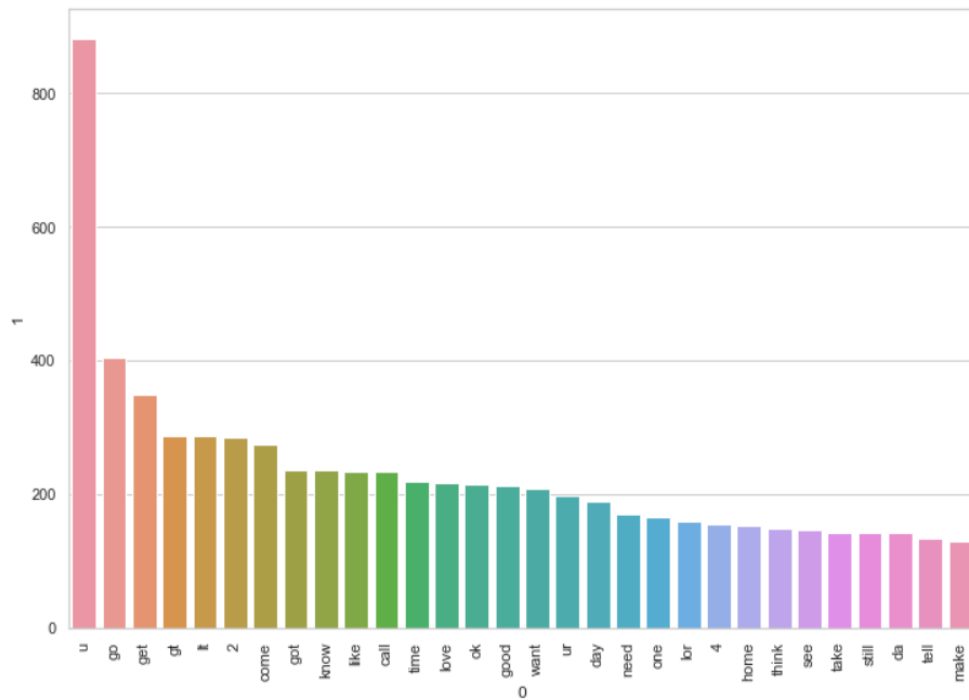
```
1 from collections import Counter
2 plt.figure(figsize=(12,8))
3 sns.barplot(pd.DataFrame(Counter(spam_corpus).most_common(30))[0], pd.DataFrame(Counter(spam_corpus).most_common(30))[1])
4 plt.xticks(rotation = 'vertical');
```



```

1 plt.figure(figsize=(12,8))
2 sns.barplot(pd.DataFrame(Counter(ham_corpus).most_common(30))[0], pd.DataFrame(Counter(ham_corpus).most_common(30))[1])
3 plt.xticks(rotation = 'vertical');

```



CONCLUSION

KEY FINDINGS AND CONCLUSIONS OF THE STUDY

From the whole evaluation we found out that the spam emails can be classified and can be stopped doing harm to the users.

LEARNING OUTCOMES OF THE STUDY IN RESPECT OF DATA SCIENCE

I found visualisation a very useful technique to infer insights from dataset. The ROC AUC plot gives large info about the false positive rate and True positive rate at various thresholds.

We are able to classify the emails as spam or non-spam. With high number of emails lots of people using the system it will be difficult to handle all possible mails as our project deals with only limited amount of corpus

LIMITATIONS OF THIS WORK AND SCOPE FOR FUTURE WORK

Since the data contained less number of '1' target labels. The trained model will be limited in scope for this label. More data of spam can definitely improve the model's performance on identification of Spam mails.