



Fake News Prediction Model

Submitted
by: Vivek Rai

ACKNOWLEDGMENT

I would like to express my very great appreciation to my SME for her valuable and constructive suggestions during the planning and development of this research work. Her willingness to give his time so generously has been very much appreciated.

Separately, I would like to thank:

- FlipRobo Technologies team
- Data Trained Team

Research papers that helped me in this project was as follows:

- [\(PDF\) A Machine Learning Approach to Fake News Detection Using Knowledge Verification and Natural Language Processing \(researchgate.net\)](#)
- <https://theconversation.com/global/topics/fake-news-33438>

Articles that helped me in this project was as follows:

- [TF-IDF Vectorizerscikit-learn. Deep understanding TfIdfVectorizer by... | by Mukesh Chaudhary | Medium](#)

INTRODUCTION

Business Problem Framing

Fake news is a form of news consisting of deliberate disinformation or hoaxes spread via traditional news media or online social media. In this project, I have used different natural language processing (NLP) based machine learning and deep learning approaches including BERT to detect fake news from news headlines. Generally, a fake headline is a news headline which may read one way or state something as fact, but then the body of the article says something different. The Internet term for this type of misleading fake news is “**clickbait**” —headlines that catch a reader’s attention to make them click on the fake news. This type of fake news is misleading at best and untrue at worst.

Conceptual Background of the Domain Problem

The idea of fake news is often referred to as click-bait in social trends and is defined as a “made up story with an intention to deceive, geared towards getting clicks”, Tavernise (2016). Some news articles have titles which grab a reader’s interest. Yet, the author only emphasizes a specific part of the article in the title. If the article itself does not focus on or give much truth to what the title had written, the news may be misleading. The goal of this project is to use natural language processing techniques to automate stance detection, since it is not practical for humans to fact check every piece of information produced by the media

Review of Literature

If we look at some scholar work shows the issue that the fake news has been major concerned amongst scholar from various background. For instance, some authors have observed that fake news is no longer a preserve of the marketing and public relations departments. Instead there is a increasing risk of IT security, therefore, IT department is premised on the idea that it would help avert the various risks associated with the problem. So, if we good deeply into it we could find that the hackers use click bait with the help of fake news and make some professional of the organization downloads their malicious exploits in their system or leak sensitive information, albeit in an indirect manner. The user may, for instance, be tricked into believing that they are helping to disseminate the news further when, in the actual sense, they are providing the perpetrators with access to their emails, and we can also see that the fake news are worked extensively as they are using videos with original message and uses their facial structure to replace the message with false message they want us to believe, these fake news issues is bigger day by day and we need to implement more our research and extensive knowledge to solve the problem.

✚ Motivation for the Problem Undertaken

This project was highly motivated project as it includes the real time problem of fake news which if we see are getting bigger, as there various concern as people do good things work hard to build a reputation, and only one false news is enough to ruin it all, it also have inverse effect on the financial market as if we observe there will a good amount of fluctuation on stock markets based on news.

Analytical Problem Framing

✚ Data Sources and their formats

There are 6 columns in the dataset provided:

The description of each of the column is given below:

- ✚ “id”: Unique id of each news article
- ✚ “headline”: It is the title of the news.
- ✚ “news”: It contains the full text of the news article
- ✚ “Unnamed:0”: It is a serial number
- ✚ “written_by”: It represents the author of the news article
- ✚ “label”: It tells whether the news is fake (1) or not fake (0).

Features Present in the Dataset:

```
Index(['Unnamed: 0', 'id', 'headline', 'written_by', 'news', 'label'], dtype='object')
```

Total Number of Rows : 20800

Total Number of Features : 6

Data Types of Features :

```
Unnamed: 0    int64
id            int64
headline      object
written_by    object
news          object
label         int64
dtype: object
```

Dataset contains any NaN/Empty cells : True

Total number of empty rows in each feature:

```
Unnamed: 0    0
id            0
headline      558
written_by    1957
news          39
label         0
dtype: int64
```

Total number of unique values in each feature:

```
Number of unique values of Unnamed: 0 : 20800
Number of unique values of id : 20800
Number of unique values of headline : 19803
Number of unique values of written_by : 4201
Number of unique values of news : 20386
Number of unique values of label : 2
```

ham and spam counts

```
1    10413
0    10387
```

Name: label, dtype: int64

✚ Data Pre-processing Done

```

1  # Here I have made a function in which all the Data cleaning steps like removing data which is not useful like
2  # email address, mobile numbers, removing punctuations, converting all the documents into lowercase,
3  # using Lemmatization technique, filtering documents using Stopwords, using POS tagging,
4  # all these type of data preprocessing steps are being performed with the help of the function defined below.
5
6  # function to filter using POS tagging..
7  def get_pos(pos_tag):
8      if pos_tag.startswith('J'):
9          return wordnet.ADJ
10     elif pos_tag.startswith('N'):
11         return wordnet.NOUN
12     elif pos_tag.startswith('R'):
13         return wordnet.ADV
14     else:
15         return wordnet.NOUN
16
17 # Function for data cleaning...
18 def Processed_data(reviews):
19     # Replace email addresses with 'email'
20     review=re.sub(r'^.+@[^\s\.]+\.[a-z]{2,}$',' ', reviews)
21
22     # Replace 10 digit phone numbers (formats include paranthesis, spaces, no spaces, dashes) with 'phonenummer'
23     review=re.sub(r'^\d{3}\d{3}\d{4}$',' ',review)
24
25     # getting only words(i.e removing all the special characters)
26     review = re.sub(r'^\w+', ' ', review)
27
28     # getting only words(i.e removing all the " _ ")
29     review = re.sub(r'[_ ]+', ' ', review)
30
31     # getting rid of unwanted characters(i.e remove all the single characters left)
32     review=re.sub(r'^\s+[a-zA-Z]\s+',' ', review)
33
34     # Removing extra whitespaces
35     review=re.sub(r'\s+', ' ', review, flags=re.I)
36
37     #converting all the letters of the review into Lowercase
38     review = review.lower()
39
40     # splitting every words from the sentences
41     review = review.split()
42
43     # iterating through each words and checking if they are stopwords or not,
44     review=[word for word in review if not word in set(STOPWORDS)]
45
46     # remove empty tokens
47     review = [text for text in review if len(text) > 0]
48
49     # getting pos tag text
50     pos_tags = pos_tag(review)
51
52     # considering words having Length more than 3only
53     review = [text for text in review if len(text) > 3]
54
55     # performing Lemmatization operation and passing the word in get_pos function to get filtered using POS ...
56     review = [(WordNetLemmatizer().lemmatize(text[0], get_pos(text[1])) for text in pos_tags]
57
58     # considering words having Length more than 3 only
59     review = [text for text in review if len(text) > 3]
60     review = ' '.join(review)
61     return review

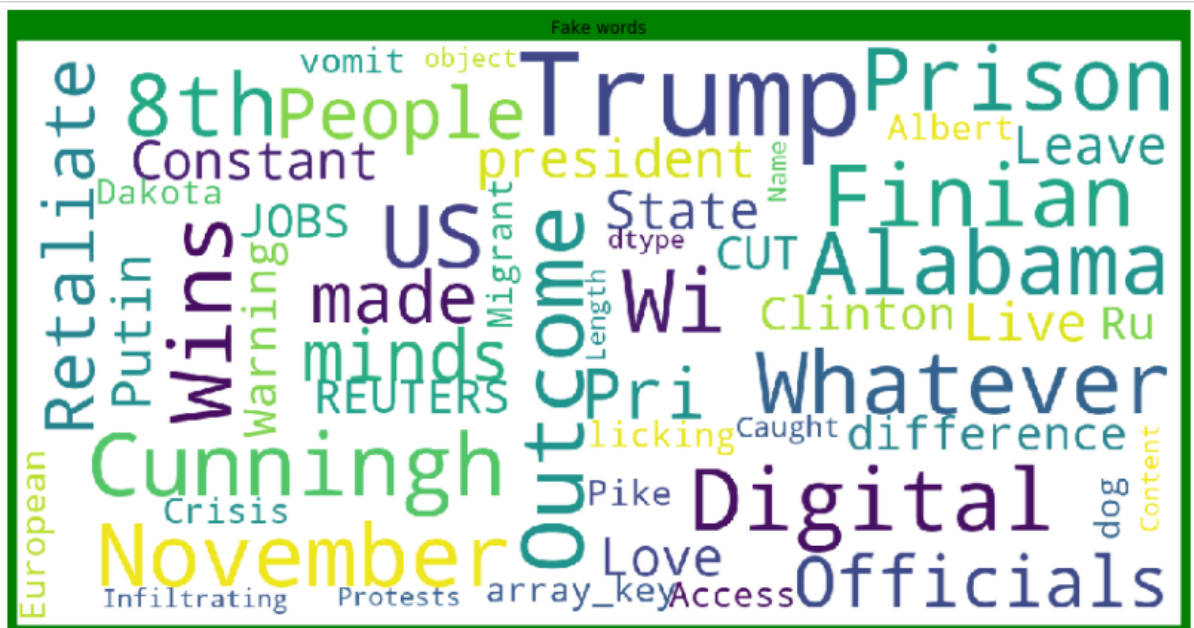
```

For Data pre-processing we did some data cleaning, where we used wordNetlemmatizer and porterStemmer to clean the words and removed special characters using Regexp Tokenizer and filter the words by removing stop words and then used lemmatizers and joined and return the filtered words.

Used TFIDF vectorizer to convert those text into vectors, and split the data and into test and train and trained various Machine learning algorithms.

Data Inputs- Logic- Output Relationships

Fake Words:



Not Fake Words:



From the above we can see that most frequent words on both labels and we can observe the words which are leading to fake news are trump, Clinton, prison, november, etc and words which are leading to real news are said, agriculture, police, questions etc, so we can clearly see that above dataset extensively deals with news around US presidential elections between Trump and Clinton.

🚧 Hardware and Software Requirements and Tools Used

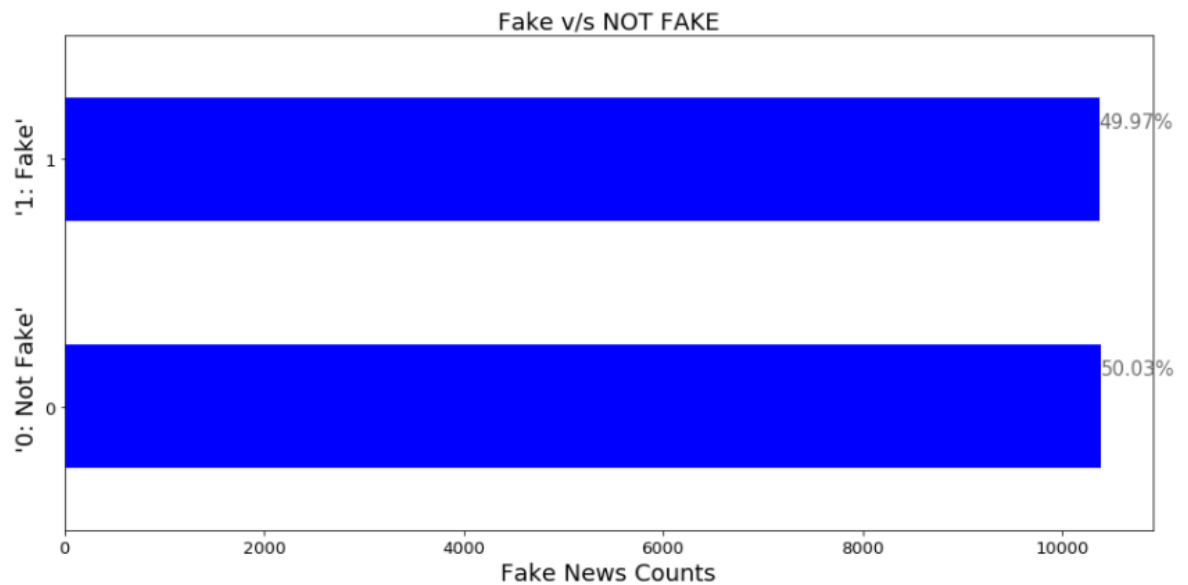
- 🚧 Hardware: 8GB RAM, 64-bit, 9th gen i7 processor.
- 🚧 Software: MS-Excel, Jupyter Notebook, python 3.6.

Libraries used:-

```
: 1 # importing useful Librarian fro the project
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import string
7 import re
8 # packages from gensim
9 from gensim import corpora
10 from gensim.utils import simple_preprocess
11 from gensim.parsing.preprocessing import STOPWORDS
12 from sklearn.feature_extraction.text import TfidfVectorizer
13
14 #packages from nltk
15 from nltk.corpus import wordnet
16 from nltk.stem import WordNetLemmatizer, SnowballStemmer
17 from nltk import pos_tag
18 from collections import Counter
19 import warnings
20 warnings.filterwarnings('ignore')
21
22
23 1 # Importing useful libraries for model training
24 2
25 3 from sklearn.linear_model import LogisticRegression
26 4 from sklearn.naive_bayes import MultinomialNB
27 5 from sklearn.tree import DecisionTreeClassifier
28 6
29 7 # Ensemble Techniques...
30 8
31 9 from sklearn.ensemble import RandomForestClassifier
32 10 from sklearn.linear_model import LogisticRegression
33 11 from sklearn.tree import DecisionTreeClassifier
34 12 from sklearn.neighbors import KNeighborsClassifier
35 13 from sklearn.svm import SVC
36 14 from sklearn.naive_bayes import MultinomialNB
37 15 from sklearn.naive_bayes import GaussianNB
38 16 from xgboost import XGBClassifier
39 17 from sklearn.ensemble import GradientBoostingClassifier
40 18 from sklearn.ensemble import AdaBoostClassifier
41 19 from sklearn.ensemble import RandomForestClassifier
42 20
43 21 # Model selection Libraries...
44 22 from sklearn.model_selection import cross_val_score, cross_val_predict, train_test_split
45 23 from sklearn.model_selection import GridSearchCV
46 24
47 25 # Importing some metrics we can use to evaluate our model performance...
48 26 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
49 27 from sklearn.metrics import roc_auc_score, roc_curve, auc
50 28 from sklearn.metrics import precision_score, recall_score, f1_score
51 29
52 30 # Creating instances for different Classifiers
53 31
54 32 RF=RandomForestClassifier()
55 33 LR=LogisticRegression()
56 34 MNB=MultinomialNB()
57 35 DT=DecisionTreeClassifier()
58 36 AD=AdaBoostClassifier()
59 37
60 38 XG=XGBClassifier()
61 39
62 40
```

Model/s Development and Evaluation

- ✚ Identification of possible problem-solving approaches (methods).



- ✚ From the above we can see that the dataset is balanced which is good as it will help our model to classify more accurately, so we should expect good accuracy score, and as the volume of data was also good.

✚ Testing of Identified Approaches (Algorithms)

- ✚ RF=RandomForestClassifier()
- ✚ LR=LogisticRegression()
- ✚ MNB=MultinomialNB()
- ✚ DT=DecisionTreeClassifier()
- ✚ AD=AdaBoostClassifier()
- ✚ XG=XGBClassifier()

Run and Evaluated selected models

```
1 # Putting Scikit-Learn machine Learning Models in a List so that it can be used for further evaluation in loop.
2 models=[]
3 models.append(('LogisticRegression',LR))
4 models.append(('MultinomialNB()',MNB))
5 models.append(('DecisionTreeClassifier',DT))
6 models.append(('RandomForestClassifier',RF))
7 models.append(('AdaBoostClassifier',AD))
8 models.append(('XGBClassifier',XG))
9
```



```

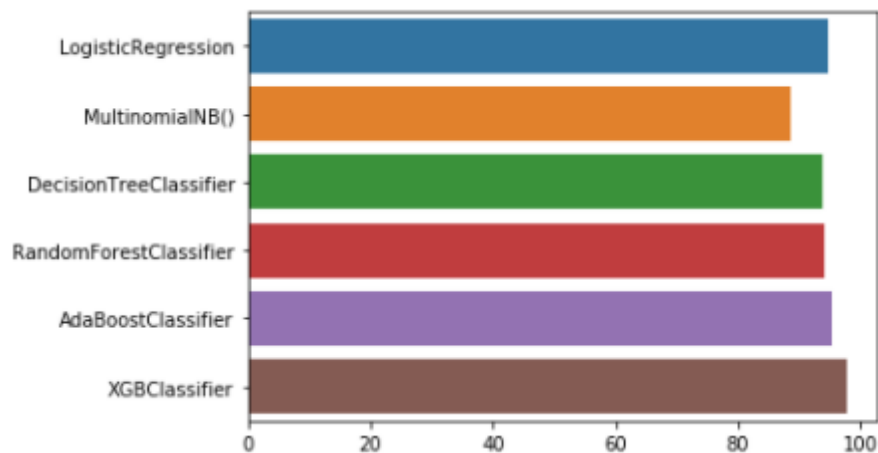
1 # Lists to store model name, Learning score, Accuracy score, cross_val_score, Auc Roc score .
2 Model=[]
3 Score=[]
4 Acc_score=[]
5 cvs=[]
6 rocscore=[]
7 # For Loop to Calculate Accuracy Score, Cross Val Score, Classification Report, Confusion Matrix
8
9 for name,model in models:
10     print('*****',name,'*****')
11     print('\n')
12     Model.append(name)
13     print(model)
14     print('\n')
15
16     # Now here I am calling a function which will calculate the max accuracy score for each model
17     # and return best random state.
18     r_state=max_acc_score(model,x,y)
19     x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.20,random_state=r_state,stratify=y)
20     model.fit(x_train,y_train)
21     #.....Learning Score.....
22     score=model.score(x_train,y_train)
23     print('Learning Score : ',score)
24     Score.append(score*100)
25     y_pred=model.predict(x_test)
26     acc_score=accuracy_score(y_test,y_pred)
27     print('Accuracy Score : ',acc_score)
28     Acc_score.append(acc_score*100)
29     #.....Finding Cross_val_score.....
30     cv_score=cross_val_score(model,x,y,cv=10,scoring='accuracy').mean()
31     print('Cross Val Score : ', cv_score)
32     cvs.append(cv_score*100)
33
34     #.....Roc auc score.....
35     false_positive_rate,true_positive_rate, thresholds=roc_curve(y_test,y_pred)
36     roc_auc=auc(false_positive_rate, true_positive_rate)
37     print('roc auc score : ', roc_auc)
38     rocscore.append(roc_auc*100)
39     print('\n')
40     print('Classification Report:\n',classification_report(y_test,y_pred))
41     print('\n')
42     print('Confusion Matrix:\n',confusion_matrix(y_test,y_pred))
43     print('\n')
44     plt.figure(figsize=(10,40))
45     plt.subplot(911)
46     plt.title(name)
47     plt.plot(false_positive_rate,true_positive_rate,label='AUC = %0.2f'% roc_auc)
48     plt.plot([0,1],[0,1],'r--')
49     plt.legend(loc='lower right')
50     plt.ylabel('True_positive_rate')
51     plt.xlabel('False_positive_rate')
52     print('\n\n')

```

Key Metrics for success in solving problem under consideration

	Model	Learning Score	Accuracy Score	Cross Val Score	Roc_Auc_curve
0	LogisticRegression	96.9722	94.7825	98.8231	94.7823
1	MultinomialNB()	90.6482	88.5375	97.0614	88.5331
2	DecisionTreeClassifier	100	93.755	93.4876	93.755
3	RandomForestClassifier	100	94.044	98.8315	94.0428
4	AdaBoostClassifier	95.0661	95.2801	98.6768	95.2801
5	XGBClassifier	99.9381	97.5919	99.6001	97.5918

Key Metrics used were the Accuracy Score, Crossvalidation Score and AUC & ROC Curve as this was binary classification. From the above we can see that there are various models out of which we few gave good accuracy score as more than 90%,



📊 Visualizations:

📊 Logistic regression:

***** LogisticRegression *****

LogisticRegression()

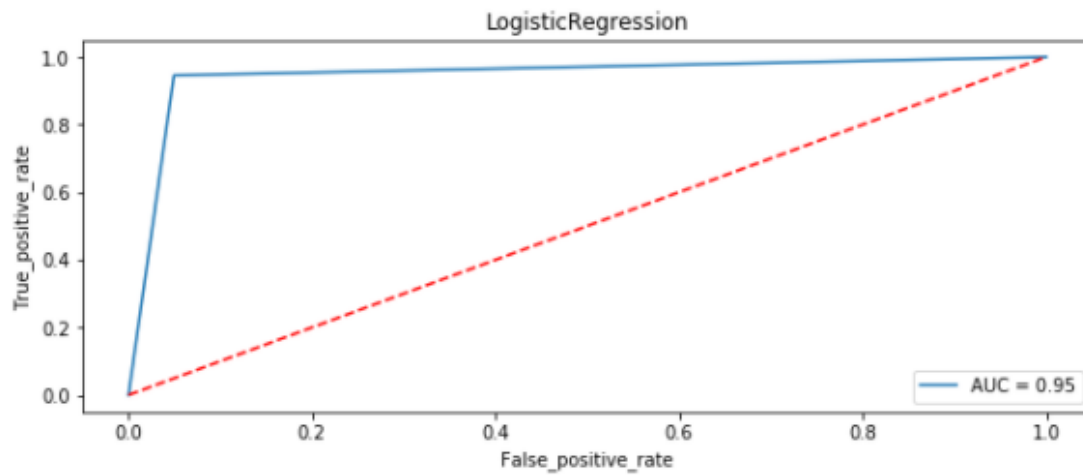
Max Accuracy Score corresponding to Random State 56 is: 0.9478246909616311

Learning Score : 0.9697219928433801
 Accuracy Score : 0.9478246909616311
 Cross Val Score : 0.9882306447402801
 roc auc score : 0.9478233644408907

Classification Report:

	precision	recall	f1-score	support
0	0.95	0.95	0.95	3116
1	0.95	0.95	0.95	3113
accuracy			0.95	6229
macro avg	0.95	0.95	0.95	6229
weighted avg	0.95	0.95	0.95	6229

Confusion Matrix:
 [[2962 154]
 [171 2942]]



Decision Tree Classifier:

***** DecisionTreeClassifier *****

DecisionTreeClassifier()

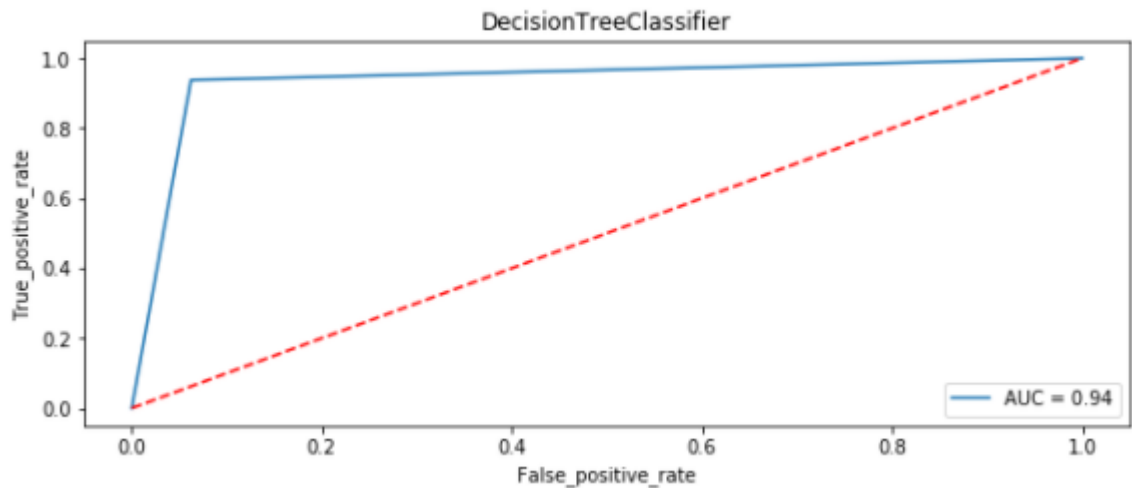
Max Accuracy Score corresponding to Random State 90 is: 0.9389950232782148

Learning Score : 1.0
 Accuracy Score : 0.9375501685663831
 Cross Val Score : 0.9348757400301622
 roc auc score : 0.9375500767620318

Classification Report:

	precision	recall	f1-score	support
0	0.94	0.94	0.94	3116
1	0.94	0.94	0.94	3113
accuracy			0.94	6229
macro avg	0.94	0.94	0.94	6229
weighted avg	0.94	0.94	0.94	6229

Confusion Matrix:
 [[2922 194]
 [195 2918]]



🚦 MultiNomial NB:

```
***** MultinomialNB() *****
```

```
MultinomialNB()
```

```
Max Accuracy Score corresponding to Random State 69 is: 0.8853748595280141
```

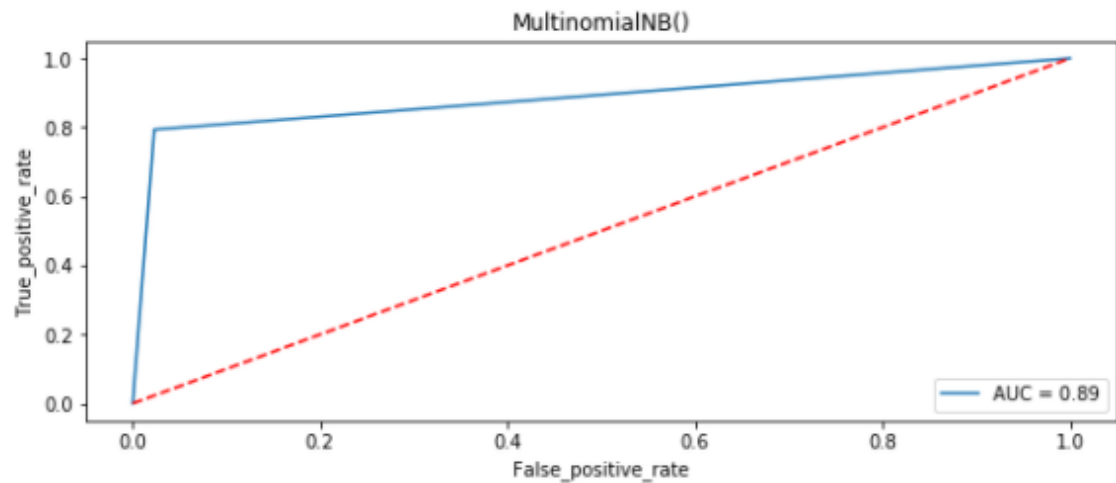
```
Learning Score : 0.9064822460776218
Accuracy Score : 0.8853748595280141
Cross Val Score : 0.9706136271516275
roc auc score : 0.8853306066282973
```

```
Classification Report:
              precision    recall  f1-score   support

     0           0.83       0.98       0.90       3116
     1           0.97       0.79       0.87       3113

   accuracy                   0.89       6229
  macro avg           0.90       0.89       0.88       6229
 weighted avg           0.90       0.89       0.88       6229
```

```
Confusion Matrix:
[[3045  71]
 [ 643 2470]]
```



✚ Random Forest Classifier:

***** RandomForestClassifier *****

RandomForestClassifier()

Max Accuracy Score corresponding to Random State 53 is: 0.9450955209503933

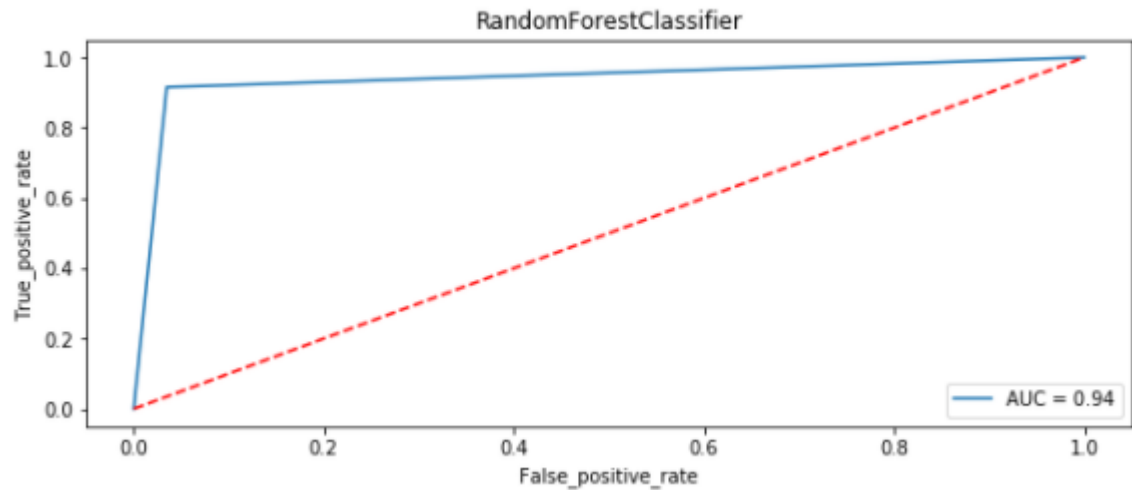
Learning Score : 1.0
 Accuracy Score : 0.9404398779900466
 Cross Val Score : 0.9883152982232619
 roc auc score : 0.9404278797720603

Classification Report:

	precision	recall	f1-score	support
0	0.92	0.97	0.94	3116
1	0.96	0.92	0.94	3113
accuracy			0.94	6229
macro avg	0.94	0.94	0.94	6229
weighted avg	0.94	0.94	0.94	6229

Confusion Matrix:

```
[[3008 108]
 [ 263 2850]]
```



Ada Boost Classifier:

***** AdaBoostClassifier *****

AdaBoostClassifier()

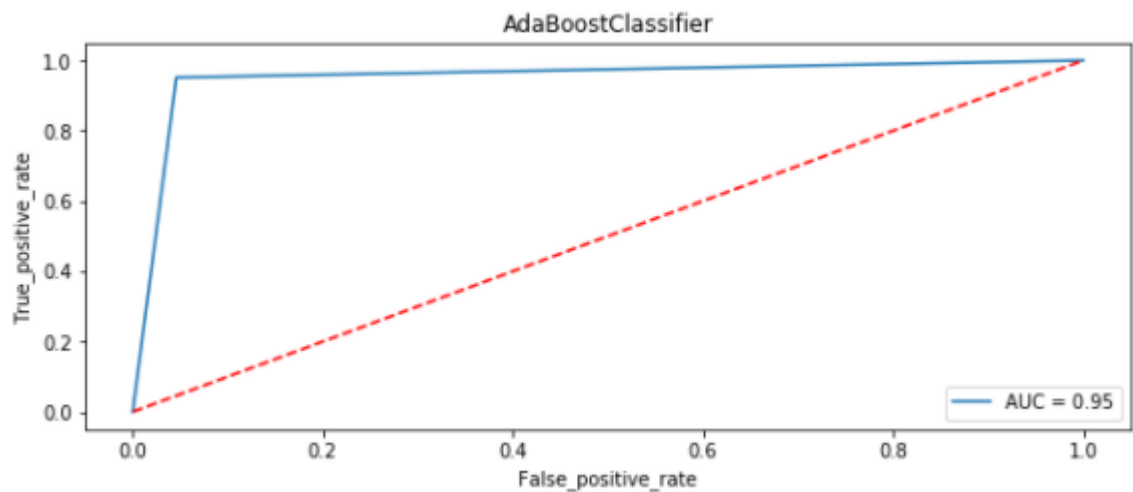
Max Accuracy Score corresponding to Random State 43 is: 0.9528014127468294

Learning Score : 0.9506606110652354
 Accuracy Score : 0.9528014127468294
 Cross Val Score : 0.9867683282736982
 roc auc score : 0.9528007832490113

Classification Report:

	precision	recall	f1-score	support
0	0.95	0.95	0.95	3116
1	0.95	0.95	0.95	3113
accuracy			0.95	6229
macro avg	0.95	0.95	0.95	6229
weighted avg	0.95	0.95	0.95	6229

Confusion Matrix:
 [[2973 143]
 [151 2962]]



✚ XGB Classifier:

***** XGBClassifier *****

```
XGBClassifier(base_score=None, booster=None, colsample_bylevel=None,
               colsample_bynode=None, colsample_bytree=None, gamma=None,
               gpu_id=None, importance_type='gain', interaction_constraints=None,
               learning_rate=None, max_delta_step=None, max_depth=None,
               min_child_weight=None, missing=nan, monotone_constraints=None,
               n_estimators=100, n_jobs=None, num_parallel_tree=None,
               random_state=None, reg_alpha=None, reg_lambda=None,
               scale_pos_weight=None, subsample=None, tree_method=None,
               validate_parameters=False, verbosity=None)
```

Max Accuracy Score corresponding to Random State 83 is: 0.9759190881361374

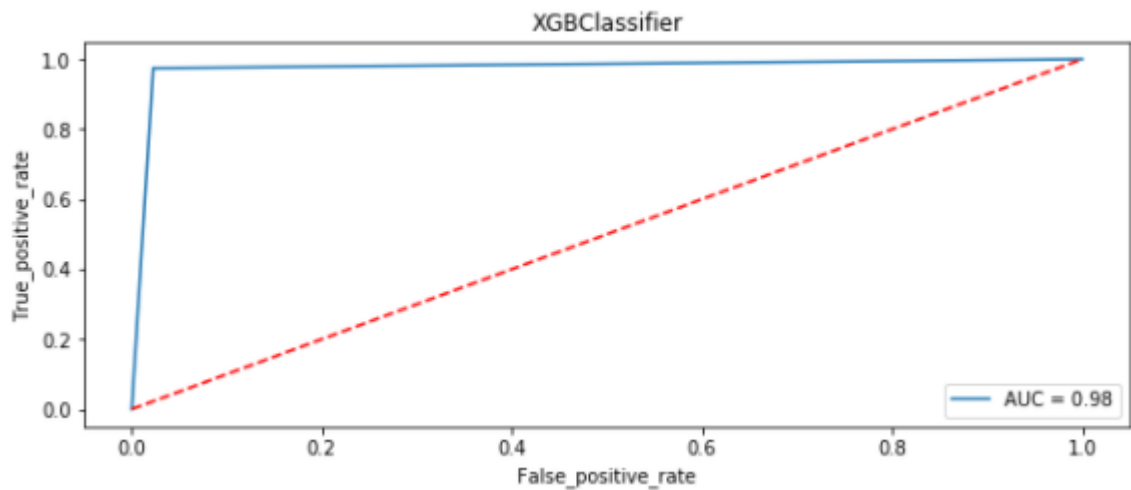
```
Learning Score : 0.9993806771263418
Accuracy Score : 0.9759190881361374
Cross Val Score : 0.9960013805958077
roc auc score : 0.9759183093631535
```

Classification Report:

	precision	recall	f1-score	support
0	0.97	0.98	0.98	3116
1	0.98	0.97	0.98	3113
accuracy			0.98	6229
macro avg	0.98	0.98	0.98	6229
weighted avg	0.98	0.98	0.98	6229

Confusion Matrix:

```
[[3046  70]
 [ 80 3033]]
```



After all this process conclusion is that XGB Classifier and Adaboost Classifier and Random Forest Classifier are performing well in terms of Accuracy score, Cross val score and Roc_Auc score as compared to other models.

Hyper Parameter Tuning Results:

```
1 # checking accuracy score(using max_acc_score function defined earlier) using best parameters which calculated from gridsearch
2 clf_lr = LogisticRegression(C=10,solver='liblinear',penalty='l2')
3 max_acc_score(clf_lr,x,y)
```

Max Accuracy Score corresponding to Random State 75 is: 0.9613100016053941

```
1 # checking accuracy score(using max_acc_score function defined earlier) using best parameters which calculated from gridsearch
2 clf_rf = RandomForestClassifier(n_estimators=500,max_features='auto')
3 max_acc_score(clf_rf,x,y)
```

Max Accuracy Score corresponding to Random State 53 is: 0.9510354792101461

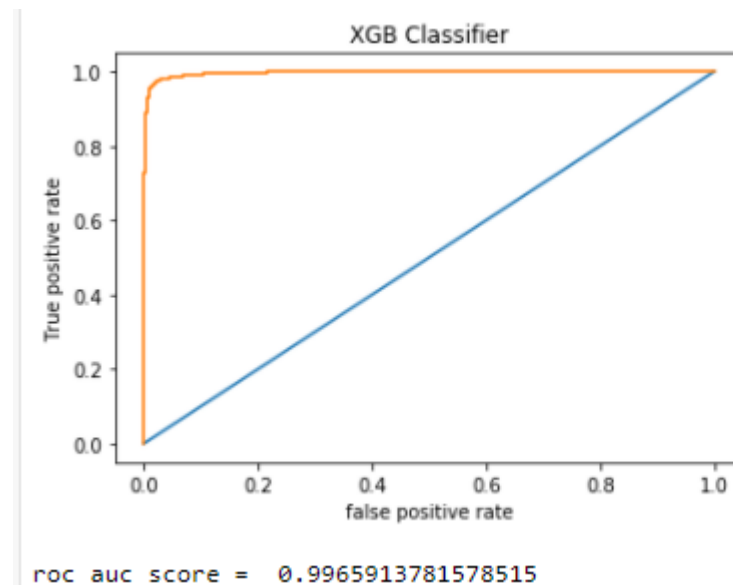
After all this process conclusion of Hyper Parameter is that Random Forest Classifier is giving accuracy of 95.10%, but XGB Classifier is giving an accuracy of 97% without tuning. So now I am making a final model using XGB Classifier.

Final Model:

```
1 # Using RandomForestClassifier for final model...
2 x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=83,test_size=.30,stratify=y)
3 XG=XGBClassifier()
4 XG.fit(x_train,y_train)
5 XG.score(x_train,y_train)
6 XGpred=XG.predict(x_test)
7 print('Accuracy Score:',accuracy_score(y_test,XGpred))
8 print('Confusion Matrix:',confusion_matrix(y_test,XGpred))
9 print('Classification Report:','\n',classification_report(y_test,XGpred))
```

Accuracy Score: 0.9759190881361374
Confusion Matrix: [[3046 70]
[80 3033]]
Classification Report:

	precision	recall	f1-score	support
0	0.97	0.98	0.98	3116
1	0.98	0.97	0.98	3113
accuracy			0.98	6229
macro avg	0.98	0.98	0.98	6229
weighted avg	0.98	0.98	0.98	6229



From the above visualization and matrices found that the XGBClassifier performed the best 99.6% AOC_ROC_SCORE, with precision accuracy score of 97% and recall 98%.

Interpretation of the Results

- From the above visualization and matrices found that the XGBClassifier performed the best AUC_ROC_SCORE i.e. **99.6%**.

CONCLUSION

Key Findings and Conclusions of the Study

From the whole evaluation we can see that the maximum number of words in fake news were regarding Trump, and Clinton and we can interpret that it was due to election campaign which was held during US presidential election and we know these adverse effects of the voters which were influenced by the fake news and most of the real news had said, trump and president, and fake news which was cleared by trump's campaign, but can hardly see any clarity or real news from the side of Clinton, and due to which the impact we already saw on election results and regarding the election advertisement and news Facebook's CEO Mark Zuckerberg also got extensively question by congress.

Learning Outcomes of the Study in respect of Data Science

So, from the words frequency chart we can clearly see that most of the news were related to US presidential election between Trump and Clinton, and by implementing passive aggressive algorithms we can see that the we have achieved a good score as it calculates the errors and updates its own learning rate which makes our model more reliable.

Limitations of this work and Scope for Future Work

- Machine Learning Algorithms like Gradient Boosting Classifier took enormous amount of time to build the model.
- Using Hyper-parameter tuning for XGB would have resulted in some more accuracy.
- Using Deep Learning for detection fake news may get some good results.