# ▾ RANDOM FOREST GBDT ALGORITHM on Amazon Fine Food Reviews

The Amazon Fine Food Reviews dataset consists of reviews of fine foods from Amazon.

Number of reviews: 568,454
Number of users: 256,059
Number of products: 74,258
Timespan: Oct 1999 - Oct 2012
Number of Attributes/Columns in data: 10

Attribute Information:

1. Id
2. ProductId - unique identifier for the product
3. UserId - unqiue identifier for the user
4. ProfileName
5. HelpfulnessNumerator - number of users who found the review helpful
6. HelpfulnessDenominator - number of users who indicated whether they found the review helpful or not
7. Score - rating between 1 and 5
8. Time - timestamp for the review
9. Summary - brief summary of the review
10. Text - text of the review

Objective:

Given a review, determine whether the review is positive (Rating of 4 or 5) or negative (rating of 1 or 2).

# ▾ Loading,Cleaning & Preprocessing the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

   In order to load the data, We have used the SQLITE dataset as it easier to query the data and visualise the data efficiently.

   Here as we only want to get the global sentiment of the recommendations (positive or negative), we will purposefully ignc
   equal to 3. If the score id above 3, then the recommendation wil be set to "positive". Otherwise, it will be set to "negative".

```
1  %matplotlib inline
2  import warnings
3  import sqlite3
4  import pandas as pd
5  import numpy as np
6  import nltk
7  import string
8  import matplotlib.pyplot as plt
9  import seaborn as sns
10 from sklearn.feature_extraction.text import TfidfTransformer
11 from sklearn.feature_extraction.text import TfidfVectorizer
12 from sklearn.feature_extraction.text import CountVectorizer
13 from sklearn.metrics import confusion_matrix
14 from sklearn import metrics
15 from sklearn.metrics import roc_curve, auc
16 from nltk.stem.porter import PorterStemmer
17
```

```
18 import re
19 import string
20 from nltk.corpus import stopwords
21 from nltk.stem import PorterStemmer
22 from nltk.stem.wordnet import WordNetLemmatizer
23
24 from gensim.models import Word2Vec
25 from gensim.models import KeyedVectors
26 import pickle
27
28 from tqdm import tqdm
29 import os
30
31 warnings.filterwarnings("ignore")
```

⟶  /usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning: pandas.util.testing is d
       import pandas.util.testing as tm

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

⟶  Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pf

    Enter your authorization code:
    ..........
    Mounted at /content/drive

```
1 con = sqlite3.connect("/content/drive/My Drive/Colab Notebooks/database.sqlite")
2
3 filtered_data=pd.read_sql_query("""SELECT * FROM Reviews WHERE Score != 3""",con);
4 filtered_data.head(3)
```

⟶

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator |
|---|---|---|---|---|---|---|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |
| 1 | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 0 |
| 2 | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | 1 |

```
1 filtered_data.shape
2
3 def partition(x):
4   if x < 3 :
5     return 'negative'
6   return 'positive'
7
8 actualScore=filtered_data['Score']
9 positive_negative=actualScore.map(partition)
10 filtered_data['Score']=positive_negative
11 print("Number of datapoints",filtered_data.shape)
12 filtered_data.head(3)
```

⟶

Number of datapoints (525814, 10)

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator |
|---|---|---|---|---|---|---|
| **0** | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |
| **1** | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 0 |
| **2** | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | 1 |

```
1 display = pd.read_sql_query("""
2 SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
3 FROM Reviews
4 GROUP BY UserId
5 HAVING COUNT(*)>1
6 """, con)
```

```
1 print(display.shape)
2 display.head(3)
```

⊳  (80668, 7)

| | UserId | ProductId | ProfileName | Time | Score | |
|---|---|---|---|---|---|---|
| **0** | #oc-R115TNMSPFT9I7 | B007Y59HVM | Breyton | 1331510400 | 2 | Overall its just OK when considering t |
| **1** | #oc-R11D9D7SHXIJB9 | B005HG9ET0 | Louis E. Emory "hoppy" | 1342396800 | 5 | My wife has recurring extreme muscle sp |
| **2** | #oc-R11DNU2NBKQ23Z | B007Y59HVM | Kim Cieszykowski | 1348531200 | 1 | This coffee is horrible and unfortuna |

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5 ORDER BY ProductID
6 """, con)
7 display.head()
```

⊳

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | |
|---|---|---|---|---|---|---|---|---|
| **0** | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |
| **1** | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |
| **2** | 138277 | B000HDOPYM | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |
| **3** | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |
| **4** | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |

```
1 sorted_data=filtered_data.sort_values('ProductId',axis=0,ascending=True,inplace=False,kind='quicksort',na_position
2
3 final_data=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"},keep='first',inplace=False)
4 final_data.shape
```

⊳  (364173, 10)

```
1 (final_data['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
69.25890143662969
```

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND Id=44737 OR Id=64422
5 ORDER BY ProductID
6 """, con)
7
8 display.head()
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score |
|---|---|---|---|---|---|---|---|
| 0 | 64422 | B000MIDROQ | A161DK06JJMCYF | J. E. Stephens "Jeanne" | 3 | 1 | 5 |
| 1 | 44737 | B001EQ55RW | A2V0I904FH7ABY | Ram | 3 | 2 | 4 |

```
1 final_data=final_data[final_data.HelpfulnessNumerator<=final_data.HelpfulnessDenominator]
```

```
1 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
True
```

```
1 stopping_words = set(stopwords.words('english'))
2 print(stopping_words)
```

```
{'y', "mustn't", 'you', 'other', 'ours', 'hadn', "wouldn't", 'during', 'through', "you've", "hadn't", 'she', 'an
```

```
1 def clean_html(text):
2     clean_r = re.compile('<,*?>')
3     clean_text = re.sub(clean_r,'',text)
4     return clean_text
5
6 def Clean_punc(text):
7     clean_sentence = re.sub(r'[?|!|\'|"|#]',r' ',text)
8     clean_data = re.sub(r'[.|,|)|(|\|/)]',r' ',clean_sentence)
9     return clean_data
```

```
1 from tqdm import tqdm
2 import os
3 import pdb
4 import pickle
5
6 from tqdm import tqdm
7 import os
8 import pdb
9 import pickle
10
11 stem_no = nltk.stem.SnowballStemmer('english')
12
13 if not os.path.isfile('final_data.sqlite'):
14     final_string=[]
15     all_positive_words=[]
16     all_negative_words=[]
17     for i,sentence in enumerate(tqdm(final_data['Text'].values)):
18         filtered_sentence=[]
```

```
19              sent_without_html_tags=clean_html(sentence)
20              #pdb.set_trace()
21              for w in sent_without_html_tags.split():
22                  for cleaned_words in Clean_punc(w).split():
23                      if ((cleaned_words.isalpha()) & (len(cleaned_words) > 2)):
24                          if(cleaned_words.lower() not in stopping_words) :
25                              stemming=(stem_no.stem(cleaned_words.lower())).encode('utf8')
26                              filtered_sentence.append(stemming)
27                              if(final_data['Score'].values)[i]=='positive':
28                                  all_positive_words.append(stemming)
29                              if(final_data['Score'].values)[i]=='negative':
30                                  all_negative_words.append(stemming)
31          str1 = b" ".join(filtered_sentence)
32          final_string.append(str1)
33
34      final_data['Cleaned_text']=final_string
35      final_data['Cleaned_text']=final_data['Cleaned_text'].str.decode("utf-8")
36
37      conn = sqlite3.connect('final_data.sqlite')
38      cursor=conn.cursor
39      conn.text_factory = str
40      final_data.to_sql('Reviews',conn,schema=None,if_exists='replace',index=True,index_label=None,chunksize=None,dt
41      conn.close()
42
43
44      with open('positive_words.pkl','wb') as f :
45          pickle.dump(all_positive_words,f)
46      with open('negative_words.pkl','wb') as f :
47          pickle.dump(all_negative_words,f)
```

```
100%|██████████| 364171/364171 [05:47<00:00, 1048.66it/s]
```

```
1 final_data['total_words'] = [len(x.split()) for x in final_data['Cleaned_text'].tolist()]
```

```
1 final_data.sort_values(by=['Time'], inplace=True, ascending=True)
```

```
1 final_data['Score'].value_counts()
```

```
positive    307061
negative     57110
Name: Score, dtype: int64
```

```
1 count_positive,count_negative=final_data['Score'].value_counts()
```

```
1 final_data_positive_class=final_data[final_data['Score']=='positive']
2 final_data_negative_class=final_data[final_data['Score']=='negative']
```

## ▾ RANDOM UP SAMPLING

Note : In Up Sampling , there will be duplicate random records from the minority class which can caus

```
1 #final_data_negative=final_data_negative_class.sample(count_positive,replace=True)
2 #final_data_after_Sampling=pd.concat([final_data_positive_class,final_data_negative], axis=0 )
3 final_data_positive=final_data_positive_class.sample(count_negative)
4
5 fina_data_after_Sampling=pd.concat([final_data_positive,final_data_negative_class], axis=0)
```

```
1 fina_data_after_Sampling['Score'].value_counts()
```

```
negative    57110
positive    57110
Name: Score, dtype: int64
```

```
1 fina_data_after_Sampling.shape
```

```
(114220, 12)
```

```
1 final_data_100K=fina_data_after_Sampling[0:100000]
2 amazon_polarity_labels=final_data_100K['Score'].values
3 final_data_100K.head(2)
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Sc |
|---|---|---|---|---|---|---|---|
| **134333** | 145807 | B003D4F1QS | A1XU4U4159OFXD | Liz | 2 | 2 | pos |
| **455421** | 492385 | B000F9ZDKI | AYCLRXW5VOIOQ | penny pincher "two cents" | 13 | 13 | pos |

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.metrics import accuracy_score
3 from sklearn.model_selection import cross_val_score
4 from collections import Counter
5 from sklearn.metrics import confusion_matrix
6 from sklearn.metrics import classification_report
7
8 X_1,X_Test,Y_1,Y_Test = train_test_split(final_data_100K,amazon_polarity_labels,test_size=0.2,random_state=0)
9 X_Train,X_CV,Y_Train,Y_CV = train_test_split(X_1,Y_1,test_size=0.2)
```

## ▾ RANDOM FOREST ALGORITHM

▾ APPLY BAG OF WORDS VECTORIZATION TECHNIQUE USING RANDOM FOREST TO FIND HYPERPARAMETER

```
1 print(X_Train.shape, Y_Train.shape)
2 print(X_CV.shape, Y_CV.shape)
3 print(X_Test.shape, Y_Test.shape)
4
5 print("="*100)
6
7
8 count_vector=CountVectorizer(min_df=1)
9 X_Train_data_bow=(count_vector.fit_transform(X_Train['Cleaned_text'].values))
10 X_Test_data_bow=(count_vector.transform(X_Test['Cleaned_text'].values))
11 X_CV_data_bow=(count_vector.transform(X_CV['Cleaned_text'].values))
12
13 print("After vectorizations")
14 print(X_Train_data_bow.shape, Y_Train.shape)
15 print(X_CV_data_bow.shape, Y_CV.shape)
16 print(X_Test_data_bow.shape, Y_Test.shape)
17 print("="*100)
```

```
(64000, 12) (64000,)
(16000, 12) (16000,)
(20000, 12) (20000,)
==========================================================================================
After vectorizations
(64000, 30318) (64000,)
(16000, 30318) (16000,)
(20000, 30318) (20000,)
==========================================================================================
```

```python
1 from sklearn.model_selection import GridSearchCV
2 from scipy.stats import randint as sp_randint
3 from sklearn.model_selection import cross_val_score
4 from sklearn.ensemble import RandomForestClassifier
5
6 max_depth = [1,5,10,50,100,500,1000]
7 n_estimators =  [20,40,60,80,100]
8 def Random_Forest_Classifier(x_training_data,y_training_data,x_cv_training,y_cv_training):
9   grid_params = { 'max_depth' : max_depth,
10                     'n_estimators' : n_estimators
11                 }
12   Classifier_RFT = RandomForestClassifier(random_state=None, class_weight ='balanced',oob_score=True)
13   clf=GridSearchCV(Classifier_RFT,grid_params,scoring='roc_auc',return_train_score=True,cv=3,n_jobs=-1)
14   clf.fit(x_training_data,y_training_data)
15   results = pd.DataFrame.from_dict(clf.cv_results_)
16   results = results.sort_values(['param_max_depth'])
17   results = results.sort_values(['param_n_estimators'])
18   train_auc= results['mean_train_score']
19   train_auc_std= results['std_train_score']
20   cv_auc = results['mean_test_score']
21   cv_auc_std= results['std_test_score']
22   best_depth =  results['param_max_depth']
23   best_estimators =  results['param_n_estimators']
24   #log_alpha=np.log10(list(results["param_alpha"]))
25   print(clf.best_score_)
26   print(clf.best_params_)
27   Build_HeatMap(max_depth,n_estimators,x_training_data,y_training_data,x_cv_training,y_cv_training)
28   return results,clf,Classifier_RFT
```
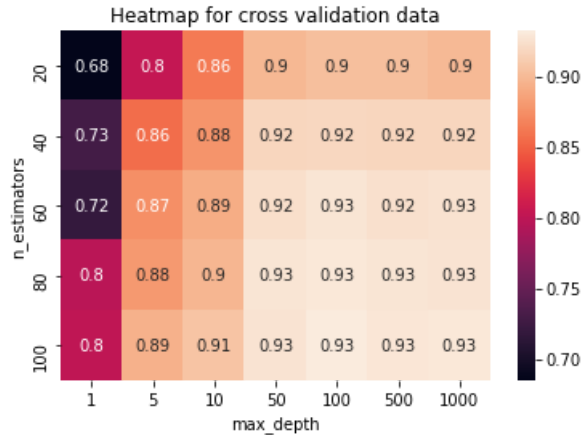
```python
1 import seaborn as sns
2 from sklearn.metrics import roc_auc_score
3
4 def Build_HeatMap(max_depth,n_estimators,X_Train,Y_Train,X_CV,Y_CV) :
5   base_learner= []
6   depths =[]
7   auc_score_cv =[]
8   for i in n_estimators :
9     for j in max_depth :
10       RFClassifier=RandomForestClassifier(random_state=None, class_weight ='balanced',max_depth = j,n_estimators =
11       RFClassifier.fit(X_Train,Y_Train)
12       predict_cv=RFClassifier.predict_proba(X_CV)[:,1]
13       predict_train=RFClassifier.predict_proba(X_Train)[:,1]
14       base_learner.append(i)
15       depths.append(j)
16       auc_score_cv.append(roc_auc_score(Y_CV,predict_cv))
17   data_frame = pd.DataFrame({'n_estimators': base_learner, 'max_depth': depths, 'AUC': auc_score_cv})
18   data_pivoted = data_frame.pivot("n_estimators", "max_depth", "AUC")
19   heat_map = sns.heatmap(data_pivoted,annot=True)
20   plt.title('Heatmap for cross validation data')
21   plt.show()
```

```python
1 print ('-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS')
2 result.best_depth.RF_Classiier=Random_Forest_Classifier(X_Train_data_bow,Y_Train,X_CV_data_bow,Y_CV)
```

```
result,best_depth,rf_classifier=random_forest_classifier(X_Train_data_bow,Y_Train,X_CV_data_bow,Y_CV)
```

```
------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS
0.9299056267764548
{'max_depth': 100, 'n_estimators': 100}
```

Heatmap for cross validation data



```
1 def Find_best_Depth_estimator(best_depth,best_estimators) :
2   best_depth = best_depth.best_params_
3   best_estimators = best_estimators.best_params_
4   best_depth=best_depth.get("max_depth")
5   best_estimators = best_estimators.get("n_estimators")
6   print(best_depth)
7   print(best_estimators)
8   return best_depth,best_estimators
```

```
 1 from sklearn.metrics import roc_curve, auc
 2
 3 RFClassifier= RandomForestClassifier(random_state=None, class_weight ='balanced',max_depth = 100,n_estimators = 10
 4 clf=RFClassifier.fit(X_Train_data_bow,Y_Train)
 5 pred_test_data=RFClassifier.predict(X_Test_data_bow)
 6 y_train_predicted_prob = RFClassifier.predict_proba(X_Train_data_bow)[:,1]
 7 y_test_predicted_prob=RFClassifier.predict_proba(X_Test_data_bow)[:,1]
 8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
 9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
17 plt.show()
```



```
1 pip install WordCloud
```

```
Requirement already satisfied: WordCloud in /usr/local/lib/python3.6/dist-packages (1.5.0)
Requirement already satisfied: numpy>=1.6.1 in /usr/local/lib/python3.6/dist-packages (from WordCloud) (1.18.3)
Requirement already satisfied: pillow in /usr/local/lib/python3.6/dist-packages (from WordCloud) (7.0.0)
```

```python
1  feature_names = count_vector.get_feature_names()
2  display_data=''
3  coef=RFClassifier.feature_importances_
4  features=np.argsort(coef)[::-1]
5
6  for i in features[0:20]:
7    display_data+=feature_names[i]
8    display_data+=' '
9
10 from wordcloud import WordCloud
11 wordcloud = WordCloud(background_color='white').generate(display_data)
12 plt.figure(figsize=(8,8),facecolor=None)
13 plt.imshow(wordcloud)
14 plt.axis("off")
15 plt.tight_layout(pad = 0)
16 plt.show()
```



## TF-IDF VECTORIZATION TECHNIQUE USING RANDOM FOREST TO FIND THE BEST HYPERPARAMETER

```python
1  from sklearn.feature_extraction.text import TfidfVectorizer
2
3  print(X_Train.shape, Y_Train.shape)
4  print(X_CV.shape, Y_CV.shape)
5  print(X_Test.shape, Y_Test.shape)
6
7  print("="*100)
8
9
10 tfidf_vector=TfidfVectorizer(min_df=10)
11 X_Train_data_tfidf=(tfidf_vector.fit_transform(X_Train['Cleaned_text'].values))
12 X_Test_data_tfidf=(tfidf_vector.transform(X_Test['Cleaned_text'].values))
13 X_CV_data_tfidf=(tfidf_vector.transform(X_CV['Cleaned_text'].values))
14
15 print("After vectorizations")
16 print(X_Train_data_tfidf.shape, Y_Train.shape)
17 print(X_CV_data_tfidf.shape, Y_CV.shape)
18 print(X_Test_data_tfidf.shape, Y_Test.shape)
19 print("="*100)
```

```
(64000, 12) (64000,)
(16000, 12) (16000,)
(20000, 12) (20000,)
==============================================================================================
After vectorizations
(64000, 6990) (64000,)
(16000, 6990) (16000,)
(20000, 6990) (20000,)
==============================================================================================
```

```
1 print ('-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS')
2 result,best_depth,RF_Classiier=Random_Forest_Classifier(X_Train_data_tfidf,Y_Train,X_CV_data_tfidf,Y_CV)
```

```
-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS
0.9265478624970153
{'max_depth': 100, 'n_estimators': 100}
```
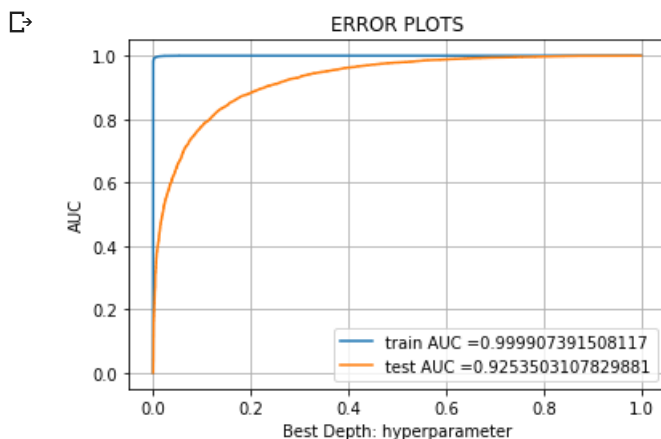

Heatmap for cross validation data

```
1 RFClassifier= RandomForestClassifier(random_state=None, class_weight ='balanced',max_depth = 100,n_estimators = 10
2 clf=RFClassifier.fit(X_Train_data_tfidf,Y_Train)
3 pred_test_data=RFClassifier.predict(X_Test_data_tfidf)
4 y_train_predicted_prob = RFClassifier.predict_proba(X_Train_data_tfidf)[:,1]
5 y_test_predicted_prob=RFClassifier.predict_proba(X_Test_data_tfidf)[:,1]
6 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
7 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
8 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
9 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
10 plt.legend()
11 plt.xlabel("Best Depth: hyperparameter")
12 plt.ylabel("AUC")
13 plt.title("ERROR PLOTS")
14 plt.grid()
15 plt.show()
```


ERROR PLOTS

```
1 feature_names = tfidf_vector.get_feature_names()
2 display_data=''
3 coef=RFClassifier.feature_importances_
4 features=np.argsort(coef)[::-1]
5
6 for i in features[0:20]:
7   display_data+=feature_names[i]
8   display_data+=' '
9
10 from wordcloud import WordCloud
11 wordcloud = WordCloud(background_color='white').generate(display_data)
12 plt.figure(figsize=(8,8),facecolor=None)
13 plt.imshow(wordcloud)
14 plt.axis("off")
15 plt.tight_layout(pad = 0)
16 plt.show()
```



## ▾ Avg Word2Vec Vectorization Technique using Random Forest

```
1 from gensim.models import Word2Vec
2 from gensim.models import KeyedVectors
3 import pickle
4
5 list_of_sent_train_avgw2v=[]
6 list_of_sent_test_avgw2v=[]
7 list_of_sent_cv_avgw2v=[]
8 for sent_train_avgw2v in tqdm(X_Train['Cleaned_text'].values):
9     list_of_sent_train_avgw2v.append(sent_train_avgw2v.split())
```

```
100%|████████████| 64000/64000 [00:01<00:00, 36198.00it/s]
```

```
1 for sent_test_avgw2v in tqdm(X_Test['Cleaned_text'].values):
2     list_of_sent_test_avgw2v.append(sent_test_avgw2v.split())
3
4 for sent_cv_avgw2v in tqdm(X_CV['Cleaned_text'].values):
5     list_of_sent_cv_avgw2v.append(sent_cv_avgw2v.split())
```

```
100%|████████████| 20000/20000 [00:00<00:00, 137861.46it/s]
100%|████████████| 16000/16000 [00:00<00:00, 139946.29it/s]
```

```
1 w2v_model_train = Word2Vec(list_of_sent_train_avgw2v,min_count=5,size=50,workers=4)
2 w2v_words_svm_train=list(w2v_model_train.wv.vocab)
```

```
1 w2v_model_test = Word2Vec(list_of_sent_test_avgw2v,min_count=5,size=50,workers=4)
```

```
2 w2v_words_svm_test=list(w2v_model_test.wv.vocab)
```

```
1 w2v_model_cv = Word2Vec(list_of_sent_cv_avgw2v,min_count=5,size=50,workers=4)
2 w2v_words_svm_cv=list(w2v_model_cv.wv.vocab)
```

```
1 train_vectors=[];
2 for sent in list_of_sent_train_avgw2v:
3     sent_vec=np.zeros(50)
4     cnt_words=0;
5     for word in sent:
6         if word in w2v_words_svm_train:
7             vec=w2v_model_train.wv[word]
8             sent_vec+=vec
9             cnt_words+=1
10    if cnt_words !=0:
11        sent_vec/=cnt_words
12    train_vectors.append(sent_vec)
13 print(len(train_vectors))
14 print(len(train_vectors[0]))
```

```
64000
50
```

```
1 test_vectors=[];
2 for sent in tqdm(list_of_sent_test_avgw2v):
3     sent_vec=np.zeros(50)
4     cnt_words=0;
5     for word in sent:
6         if word in w2v_words_svm_test:
7             vec=w2v_model_test.wv[word]
8             sent_vec+=vec
9             cnt_words+=1
10    if cnt_words !=0:
11        sent_vec/=cnt_words
12    test_vectors.append(sent_vec)
13 print(len(test_vectors))
14 print(len(test_vectors[0]))
```

```
100%|██████████| 20000/20000 [00:16<00:00, 1210.51it/s]20000
50
```

```
1 cv_vectors=[];
2 for sent in tqdm(list_of_sent_cv_avgw2v):
3     sent_vec=np.zeros(50)
4     cnt_words=0;
5     for word in sent:
6         if word in w2v_words_svm_cv:
7             vec=w2v_model_cv.wv[word]
8             sent_vec+=vec
9             cnt_words+=1
10    if cnt_words !=0:
11        sent_vec/=cnt_words
12    cv_vectors.append(sent_vec)
13 print(len(cv_vectors))
14 print(len(cv_vectors[0]))
```
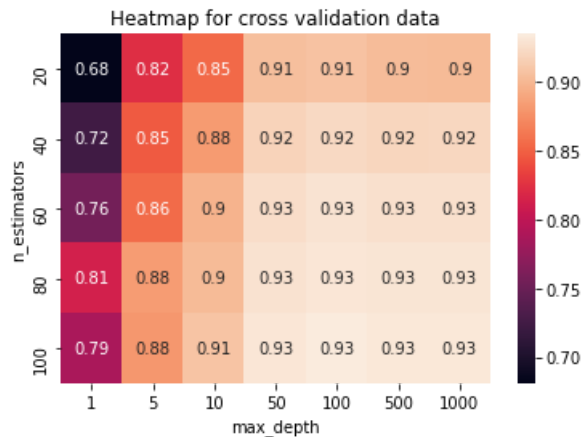
```
100%|██████████| 16000/16000 [00:12<00:00, 1241.77it/s]16000
50
```

```
1 print ('-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS')
2 result,best_depth,RF_Classiier=Random_Forest_Classifier(train_vectors,Y_Train,cv_vectors,Y_CV)
```

```
------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS
0.9072753950284099
{'max_depth': 50, 'n_estimators': 100}
```
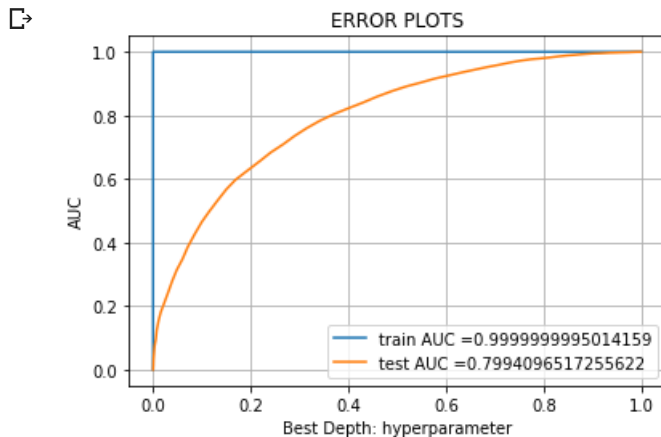
Heatmap for cross validation data



```
 1 RFClassifier= RandomForestClassifier(random_state=None, class_weight ='balanced',max_depth = 50,n_estimators = 100
 2 clf=RFClassifier.fit(train_vectors,Y_Train)
 3 pred_test_data=RFClassifier.predict(test_vectors)
 4 y_train_predicted_prob = RFClassifier.predict_proba(train_vectors)[:,1]
 5 y_test_predicted_prob=RFClassifier.predict_proba(test_vectors)[:,1]
 6 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
 7 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
 8 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
 9 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
10 plt.legend()
11 plt.xlabel("Best Depth: hyperparameter")
12 plt.ylabel("AUC")
13 plt.title("ERROR PLOTS")
14 plt.grid()
15 plt.show()
```



## ▾ TF-IDF Word2Vec Vectorization Technique using Random Forest

```
1 model_Avgw2v = TfidfVectorizer()
2 X_Train_Avgw2v=model_Avgw2v.fit_transform(X_Train['Cleaned_text'].values)
```

```
1 X_Test_Avgw2v=model_Avgw2v.transform(X_Test['Cleaned_text'].values)
2 X_CV_Avgw2v=model_Avgw2v.transform(X_CV['Cleaned_text'].values)
```

```
1 dictionary = dict(zip(model_Avgw2v.get_feature_names(), list(model_Avgw2v.idf_)))
```

```
1 tfidf_feature=model_Avgw2v.get_feature_names()
2
3 tfidf_sent_vectors_train=[];
4 #final_tf_idf = [];
5 row=0;
6
7 for sent in tqdm(list_of_sent_train_avgw2v):
8     sent_vec=np.zeros(50)
9     weight_sum=0;
10    for word in sent :
11        if word in w2v_words_svm_train and word in tfidf_feature :
12            vec=w2v_model_train.wv[word]
13            #tf_idf=final_tf_idf[row,tfidf_feature.index(word)]
14            tf_idf=dictionary[word]*(sent.count(word)/len(sent))
15            sent_vec+=(vec*tf_idf)
16            weight_sum+=tf_idf
17
18    if weight_sum!=0:
19        sent_vec/=weight_sum
20    tfidf_sent_vectors_train.append(sent_vec)
21    row+=1
```

100%|████████████| 64000/64000 [13:27<00:00, 79.30it/s]

```
1 tfidf_sent_vectors_test=[];
2 #final_tf_idf = [];
3 row=0;
4
5 for sent in tqdm(list_of_sent_test_avgw2v):
6     sent_vec=np.zeros(50)
7     weight_sum=0;
8     for word in sent :
9         if word in w2v_words_svm_test and word in tfidf_feature :
10            vec=w2v_model_test.wv[word]
11            #tf_idf=final_tf_idf[row,tfidf_feature.index(word)]
12            tf_idf=dictionary[word]*(sent.count(word)/len(sent))
13            sent_vec+=(vec*tf_idf)
14            weight_sum+=tf_idf
15
16    if weight_sum!=0:
17        sent_vec/=weight_sum
18    tfidf_sent_vectors_test.append(sent_vec)
19    row+=1
```

100%|████████████| 20000/20000 [04:08<00:00, 80.46it/s]

```
1 print ('-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS')
2 result,best_depth,RF_Classiier=Random_Forest_Classifier(tfidf_sent_vectors_train,Y_Train,tfidf_sent_vectors_test,Y
```
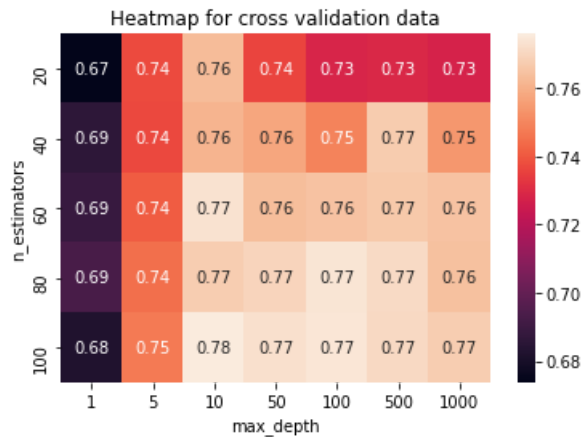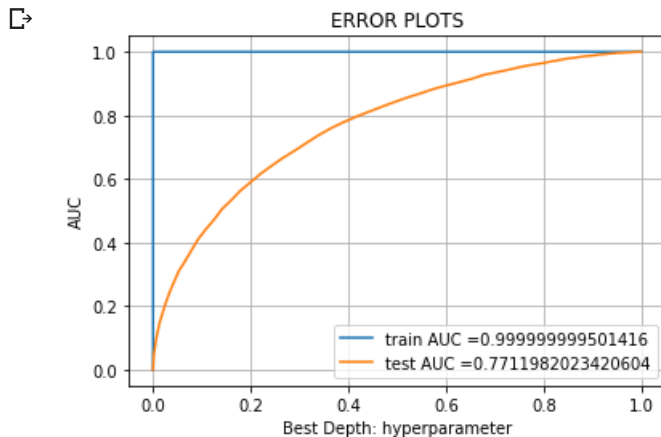
```
-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS
0.8817480033471647
{'max_depth': 100, 'n_estimators': 100}
```

Heatmap for cross validation data



```
 1 RFClassifier= RandomForestClassifier(random_state=None, class_weight ='balanced',max_depth = 100,n_estimators = 10
 2 clf=RFClassifier.fit(tfidf_sent_vectors_train,Y_Train)
 3 pred_test_data=RFClassifier.predict(tfidf_sent_vectors_test)
 4 y_train_predicted_prob = RFClassifier.predict_proba(tfidf_sent_vectors_train)[:,1]
 5 y_test_predicted_prob=RFClassifier.predict_proba(tfidf_sent_vectors_test)[:,1]
 6 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
 7 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
 8 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
 9 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
10 plt.legend()
11 plt.xlabel("Best Depth: hyperparameter")
12 plt.ylabel("AUC")
13 plt.title("ERROR PLOTS")
14 plt.grid()
15 plt.show()
```



## ▾ PRETTY TABLE FOR RANDOM FOREST

```
1 pip install -U PTable
```

```
Collecting PTable
  Downloading https://files.pythonhosted.org/packages/ab/b3/b54301811173ca94119eb474634f120a49cd370f257d1aae5a4a
Building wheels for collected packages: PTable
  Building wheel for PTable (setup.py) ... done
  Created wheel for PTable: filename=PTable-0.9.2-cp36-none-any.whl size=22908 sha256=92b2471788d31ed83df5850fc3
  Stored in directory: /root/.cache/pip/wheels/22/cc/2e/55980bfe86393df3e9896146a01f6802978d09d7ebcba5ea56
Successfully built PTable
Installing collected packages: PTable
Successfully installed PTable-0.9.2
```

```
1 from prettytable import PrettyTable
2
3 x= PrettyTable()
4 x.field_names = ["Vectorizer" , "Hyperparameter(Best Depth)", "n-estimators","AUC"]
5 x.add_row(["Bag Of Words",100,100,0.9299])
6 x.add_row(["Tf-Idf",100,100,0.9265])
7 x.add_row(["Avg Word2Vec",50,100,0.9072])
8 x.add_row(["Tf-Idf Word2Vec",100,100,0.8817])
9 print(x)
```

```
+-----------------+----------------------------+--------------+--------+
|    Vectorizer   | Hyperparameter(Best Depth) | n-estimators |  AUC   |
+-----------------+----------------------------+--------------+--------+
|   Bag Of Words  |            100             |     100      | 0.9299 |
|      Tf-Idf     |            100             |     100      | 0.9265 |
|   Avg Word2Vec  |             50             |     100      | 0.9072 |
| Tf-Idf Word2Vec |            100             |     100      | 0.8817 |
+-----------------+----------------------------+--------------+--------+
```

## ▾ *GRADIENT BOOSTING ALGORITHM - XGBOOST*

```
1 pip install xgboost
```

```
Requirement already satisfied: xgboost in /usr/local/lib/python3.6/dist-packages (0.90)
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from xgboost) (1.18.3)
Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packages (from xgboost) (1.4.1)
```

## ▾ APPLY BAG OF WORDS VECTORIZATION TECHNIQUE USING GBDT TO FIND THE BEST DEPTH(HYPERPARAMETER)

```
 1 from sklearn.model_selection import GridSearchCV
 2 from scipy.stats import randint as sp_randint
 3 from sklearn.model_selection import cross_val_score
 4 import xgboost as xgb
 5
 6
 7 max_depth = [1,5,10,50,100,500,1000]
 8 n_estimators =  [20,40,60,80,100]
 9 def XGBoost_Classifier(x_training_data,y_training_data,x_cv_training,y_cv_training):
10   grid_params = { 'max_depth' : max_depth,
11                   'n_estimators' :n_estimators
12                 }
13   Classifier_XGBoost = xgb.XGBClassifier(random_state=0, booster='gbtree')
14   clf=GridSearchCV(Classifier_XGBoost,grid_params,scoring='roc_auc',return_train_score=True,cv=3,n_jobs=-1)
15   clf.fit(x_training_data,y_training_data)
16   results = pd.DataFrame.from_dict(clf.cv_results_)
17   results = results.sort_values(['param_max_depth'])
18   results = results.sort_values(['param_n_estimators'])
19   train_auc= results['mean_train_score']
```

```
20    train_auc_std= results['std_train_score']
21    cv_auc = results['mean_test_score']
22    cv_auc_std= results['std_test_score']
23    best_depth =  results['param_max_depth']
24    best_estimators =  results['param_n_estimators']
25    print(clf.best_score_)
26    print(clf.best_params_)
27    Build_HeatMap(max_depth,n_estimators,x_training_data,y_training_data,x_cv_training,y_cv_training)
28    return results,clf,Classifier_XGBoost
```
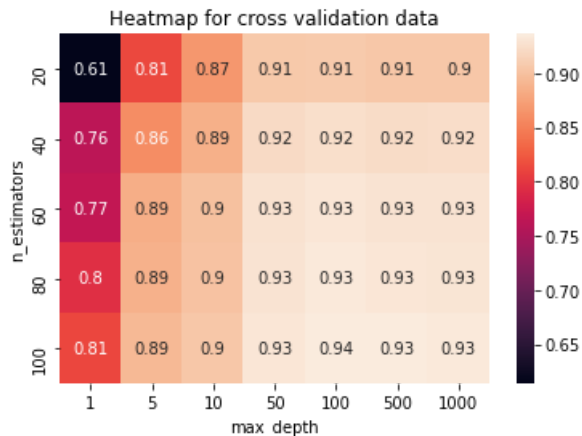
```
1 print ('-------------BEST DEPTH---------------')
2 result,best_depth,XGBoostClassiier=XGBoost_Classifier(X_Train_data_bow,Y_Train,X_CV_data_bow,Y_CV)
```

⤷      -------------BEST DEPTH---------------
       0.928610806402034
       {'max_depth': 100, 'n_estimators': 100}



Heatmap for cross validation data

```
1 from sklearn.metrics import roc_curve, auc
2
3 GBDTClassifier= xgb.XGBClassifier(random_state=0,max_depth = 100,n_estimators = 100,booster='gbtree')
4 clf=GBDTClassifier.fit(X_Train_data_bow,Y_Train)
5 pred_test_data=GBDTClassifier.predict(X_Test_data_bow)
6 y_train_predicted_prob = GBDTClassifier.predict_proba(X_Train_data_bow)[:,1]
7 y_test_predicted_prob=GBDTClassifier.predict_proba(X_Test_data_bow)[:,1]
8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
17 plt.show()
```

⤷

ERROR PLOTS

```
1 feature_names = count_vector.get_feature_names()
2 display_data=''
3 coef=GBDTClassifier.feature_importances_
4 features=np.argsort(coef)[::-1]
5
6 for i in features[0:20]:
7   display_data+=feature_names[i]
8   display_data+=' '
9
10 from wordcloud import WordCloud
11 wordcloud = WordCloud(background_color='white').generate(display_data)
12 plt.figure(figsize=(8,8),facecolor=None)
13 plt.imshow(wordcloud)
14 plt.axis("off")
15 plt.tight_layout(pad = 0)
16 plt.show()
```
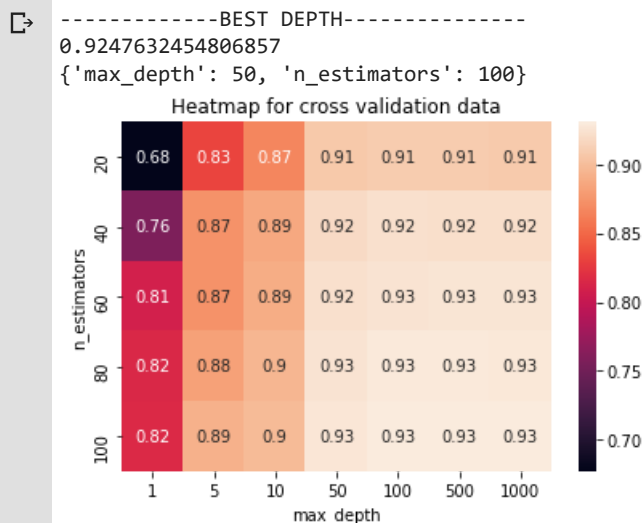


## APPLY TF-IDF VECTORIZATION TECHNIQUE USING XGBOOST TO FIND BEST HYPERPAF

```
1 print ('-------------BEST DEPTH---------------')
2 result,best_depth,XGBoostClassiier=XGBoost_Classifier(X_Train_data_tfidf,Y_Train,X_CV_data_tfidf,Y_CV)
```

```
-------------BEST DEPTH---------------
0.9247632454806857
{'max_depth': 50, 'n_estimators': 100}
```
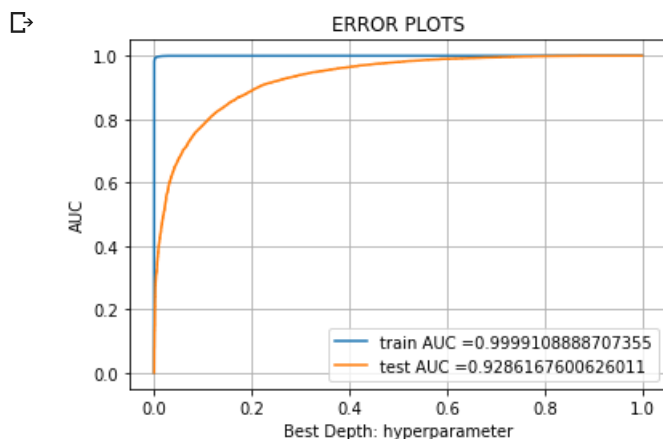


Heatmap for cross validation data

```
1 from sklearn.metrics import roc_curve, auc
2
3 GBDTClassifier= xgb.XGBClassifier(random_state=0,max_depth = 50,n_estimators = 100,booster='gbtree')
```

```
 4 clf=GBDTClassifier.fit(X_Train_data_tfidf,Y_Train)
 5 pred_test_data=GBDTClassifier.predict(X_Test_data_tfidf)
 6 y_train_predicted_prob = GBDTClassifier.predict_proba(X_Train_data_tfidf)[:,1]
 7 y_test_predicted_prob=GBDTClassifier.predict_proba(X_Test_data_tfidf)[:,1]
 8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
 9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
17 plt.show()
```
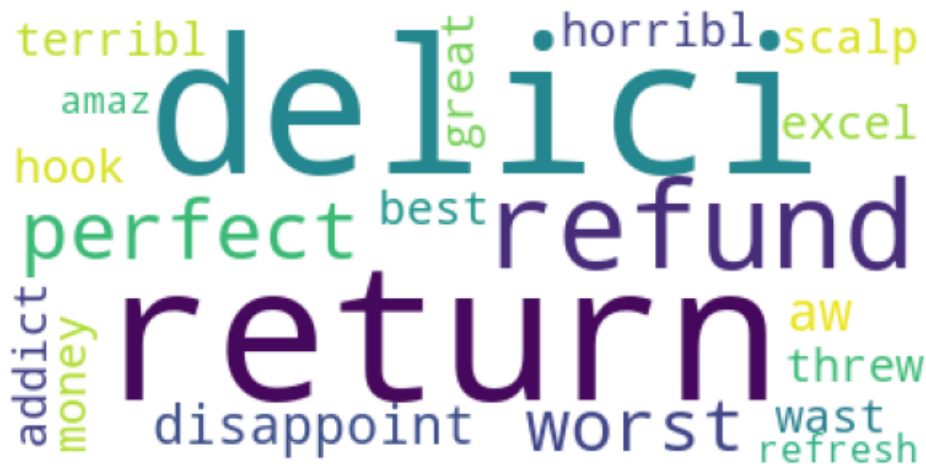


### FIND THE IMPORTANT FEATURES AND PRINT THEM USING WORD CLOUD

```
 1 feature_names = tfidf_vector.get_feature_names()
 2 display_data=''
 3 coef=GBDTClassifier.feature_importances_
 4 features=np.argsort(coef)[::-1]
 5
 6 for i in features[0:20]:
 7   display_data+=feature_names[i]
 8   display_data+=' '
 9
10 from wordcloud import WordCloud
11 wordcloud = WordCloud(background_color='white').generate(display_data)
12 plt.figure(figsize=(8,8),facecolor=None)
13 plt.imshow(wordcloud)
14 plt.axis("off")
15 plt.tight_layout(pad = 0)
16 plt.show()
```
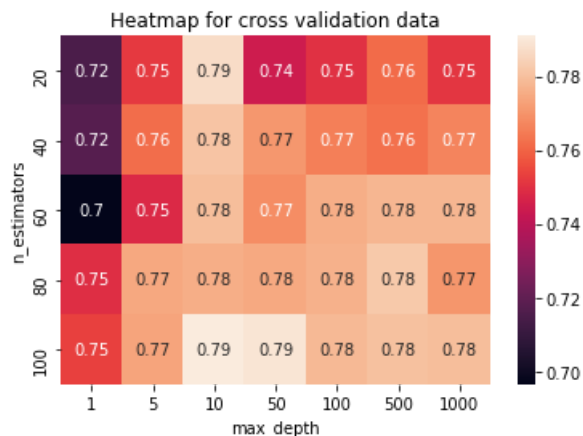
## Avg Word2Vec Vectorization Technique using XGBoost Algorithm

```
1 XGBoost_Train_AvgWord2Vec=np.array(train_vectors)
2 XGBoost_Test_AvgWord2Vec = np.array(test_vectors)
3 XGBoost_CV_AvgWord2Vec = np.array(cv_vectors)
```
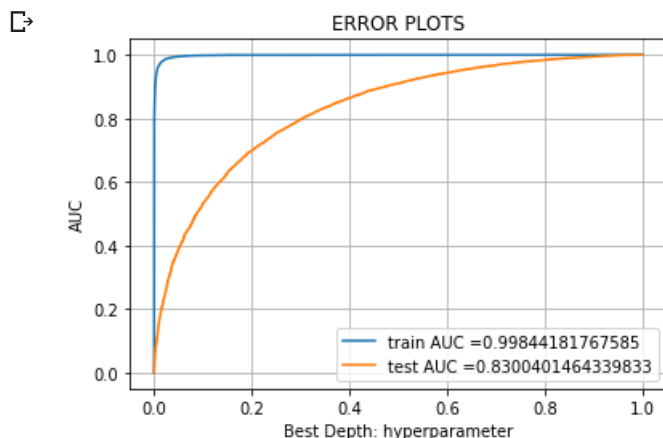
```
1 print ('-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS')
2 result,best_depth,XGBoostClassiier=XGBoost_Classifier(XGBoost_Train_AvgWord2Vec,Y_Train,XGBoost_CV_AvgWord2Vec,Y_C
```

```
-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS
0.9186298659126494
{'max_depth': 10, 'n_estimators': 100}
```



Heatmap for cross validation data

```
 1 from sklearn.metrics import roc_curve, auc
 2
 3 GBDTClassifier= xgb.XGBClassifier(random_state=0,max_depth = 10,n_estimators = 100,booster='gbtree')
 4 clf=GBDTClassifier.fit(XGBoost_Train_AvgWord2Vec,Y_Train)
 5 pred_test_data=GBDTClassifier.predict(XGBoost_Test_AvgWord2Vec)
 6 y_train_predicted_prob = GBDTClassifier.predict_proba(XGBoost_Train_AvgWord2Vec)[:,1]
 7 y_test_predicted_prob=GBDTClassifier.predict_proba(XGBoost_Test_AvgWord2Vec)[:,1]
 8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
 9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
```
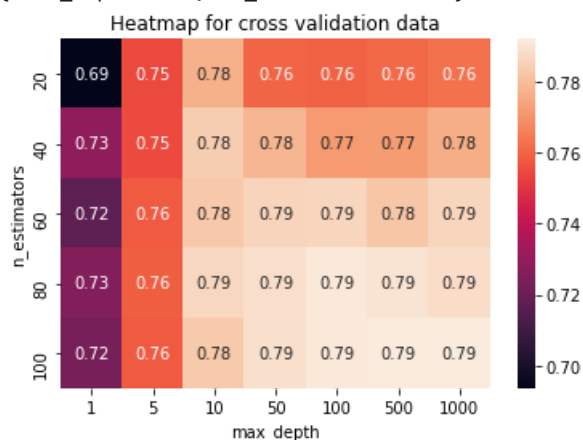
```
17 plt.show()
```



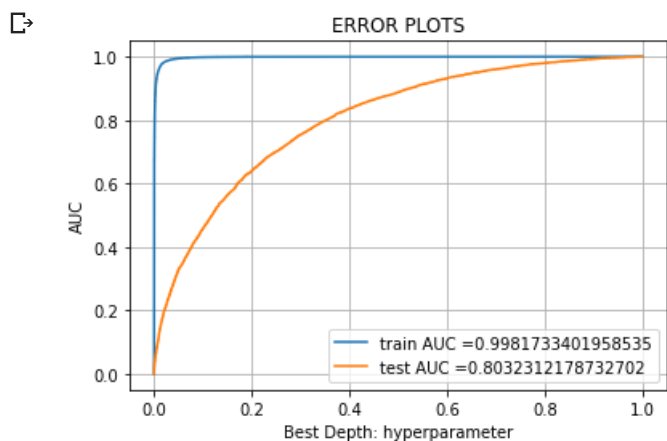## TF-IDF Word2Vec Vectorization Technique using XGBoost Algorithm

```
1 XGBoost_Train_TfidfWord2Vec=np.array(tfidf_sent_vectors_train)
2 XGBoost_Test_TfidfWord2Vec = np.array(tfidf_sent_vectors_test)
```

```
1 print ('-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS')
2 result,best_depth,XGBoostClassiier=XGBoost_Classifier(XGBoost_Train_TfidfWord2Vec,Y_Train,XGBoost_Test_TfidfWord2V
```

```
-------------BEST DEPTH USING DIFFERENT RANGE OF ESTIMATORS
0.8986430183498193
{'max_depth': 10, 'n_estimators': 100}
```



```
 1 from sklearn.metrics import roc_curve, auc
 2
 3 GBDTClassifier= xgb.XGBClassifier(random_state=0,max_depth = 10,n_estimators = 100,booster='gbtree')
 4 clf=GBDTClassifier.fit(XGBoost_Train_TfidfWord2Vec,Y_Train)
 5 pred_test_data=GBDTClassifier.predict(XGBoost_Test_TfidfWord2Vec)
 6 y_train_predicted_prob = GBDTClassifier.predict_proba(XGBoost_Train_TfidfWord2Vec)[:,1]
 7 y_test_predicted_prob=GBDTClassifier.predict_proba(XGBoost_Test_TfidfWord2Vec)[:,1]
 8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
 9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
17 plt.show()
```

## ▾ PRETTY TABLE FOR XGBOOST

```
1 from prettytable import PrettyTable
2
3 x= PrettyTable()
4 x.field_names = ["Vectorizer" , "Hyperparameter(Best Depth)", "n-estimators","AUC"]
5 x.add_row(["Bag Of Words",100,100,0.9286])
6 x.add_row(["Tf-Idf",50,100,0.9247])
7 x.add_row(["Avg Word2Vec",10,100,0.9186])
8 x.add_row(["Tf-Idf Word2Vec",10,100,0.8986])
9 print(x)
```

```
+-----------------+---------------------------+--------------+--------+
|   Vectorizer    | Hyperparameter(Best Depth) | n-estimators |  AUC   |
+-----------------+---------------------------+--------------+--------+
|   Bag Of Words  |            100            |     100      | 0.9286 |
|     Tf-Idf      |             50            |     100      | 0.9247 |
|   Avg Word2Vec  |             10            |     100      | 0.9186 |
| Tf-Idf Word2Vec |             10            |     100      | 0.8986 |
+-----------------+---------------------------+--------------+--------+
```

CONCLUSION

BY LOOKING AT THE TABLE FOR RANDOM FOREST AND XGBOOST , TF-IDF VECTORIZATION PERFOR
WITH AUC OF 0.9247 HAVING BEST DEPTH AS 50

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.