# Extreme Gradient boosting on Amazon fine food dataset

Data Source: https://www.kaggle.com/snap/amazon-fine-food-reviews (https://www.kaggle.com/snap/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:

To perform Extreme gradient boosting on different vectors like BOW, Tf-idf, Avg-W2vec & Tf-idf_W2vec.

```
1   %matplotlib inline
2   import warnings
3   warnings.filterwarnings("ignore")
4
5   import sqlite3
6   import pandas as pd
7   import numpy as np
8   import nltk
9   import string
10  import matplotlib.pyplot as plt
11  import seaborn as sns
12  from sklearn.feature_extraction.text import TfidfTransformer
13  from sklearn.feature_extraction.text import TfidfVectorizer
14
15  from sklearn.feature_extraction.text import CountVectorizer
16  from sklearn.metrics import confusion_matrix
17  from sklearn import metrics
18  from sklearn.metrics import roc_curve, auc
19  from nltk.stem.porter import PorterStemmer
20
21  import re
22
23  import string
24  from nltk.corpus import stopwords
25  from nltk.stem import PorterStemmer
26  from nltk.stem.wordnet import WordNetLemmatizer
27
28  from gensim.models import Word2Vec
29  from gensim.models import KeyedVectors
30  import pickle
```

```python
#Importing Train and test dataset
train_data=pd.read_csv("E:/Applied AI assignments/Amazon_fine_train_data.csv")
test_data=pd.read_csv("E:/Applied AI assignments/Amazon_fine_test_data.csv")
```

```python
train_data=train_data.astype(str)
test_data=test_data.astype(str)
```

```python
train_data.shape
```

```
(80000, 13)
```

```python
train_data['Score'].value_counts()
```

```
positive    70407
negative     9593
Name: Score, dtype: int64
```

```python
test_data.shape
```

```python
test_data['Score'].value_counts()
```

```
positive    17322
negative     2678
Name: Score, dtype: int64
```

```python
#Train data
y_train = train_data['Score']
x_train = train_data['CleanedText']

#Test data
y_test = test_data['Score']
x_test = test_data['CleanedText']
```

```python
1   #Replacing Positive score with 0 and negative score with 1
2   y_train.replace('negative',1,inplace=True)
3   y_train.replace('positive',0,inplace=True)
4
5   y_test.replace('negative',1,inplace=True)
6   y_test.replace('positive',0,inplace=True)
```

```python
1   from xgboost import XGBClassifier
2   from sklearn.model_selection import RandomizedSearchCV
3   from sklearn.model_selection import TimeSeriesSplit
4   from sklearn.metrics import accuracy_score
5   from sklearn.metrics import recall_score
6   from sklearn.metrics import precision_score
7   from sklearn.metrics import f1_score
8   from sklearn.metrics import make_scorer
9   from sklearn.metrics import confusion_matrix
10  from sklearn.cross_validation import cross_val_score
11  from collections import Counter
12  from sklearn import cross_validation
13  from wordcloud import WordCloud
14  import matplotlib.pyplot as plt
```

Randomisedsearch CV

```
 1  n_estimators=list(range(100,1100,100))
 2  max_depth=list(range(1,50))
 3  learning_rate=[0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]
 4
 5  T= TimeSeriesSplit(n_splits=5)
 6
 7  param_distributions  = dict(n_estimators=n_estimators,max_depth=max_depth,learning_rate=learning_rate)
 8  print(param_distributions)
 9
10  # instantiate and fit the grid
11  grid = RandomizedSearchCV(XGBClassifier(), param_distributions, cv=T, scoring='f1', return_train_score=I
```

```
{'n_estimators': [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000], 'max_depth': [1, 2, 3, 4, 5, 6, 7, 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], 'learning_rate': [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0]}
```

## Binary Bow

```
 1  count_vect = CountVectorizer(binary=True)
 2
 3  #Train data
 4  vocabulary = count_vect.fit(x_train) #in scikit-learn
 5  Bow_x_train= count_vect.transform(x_train)
 6  print("the type of count vectorizer ",type(Bow_x_train))
 7  print("the shape of out text BOW vectorizer ",Bow_x_train.get_shape())
 8  print("the number of unique words ", Bow_x_train.get_shape()[1])
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (80000, 33433)
the number of unique words  33433
```

```
1  #Test data
2  Bow_x_test = count_vect.transform(x_test)
3  print("the type of count vectorizer ",type(Bow_x_test))
4  print("the shape of out text BOW vectorizer ",Bow_x_test.get_shape())
5  print("the number of unique words ", Bow_x_test.get_shape()[1])
```
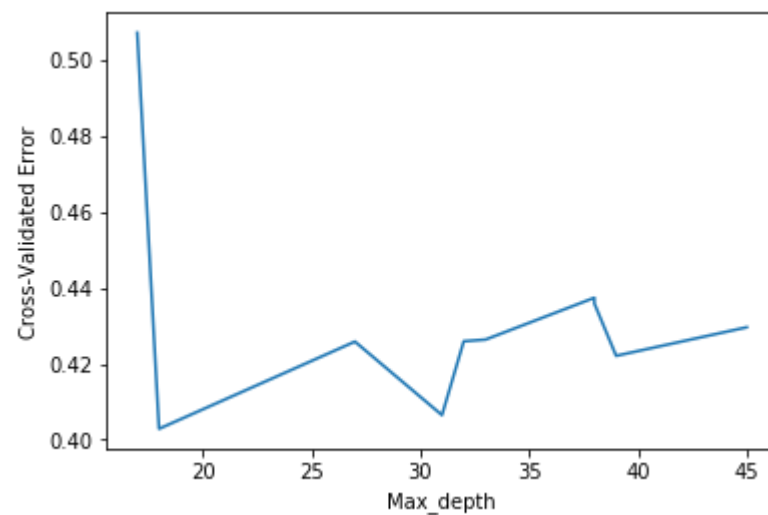
```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (20000, 33433)
the number of unique words  33433
```

```
1  grid.fit(Bow_x_train, y_train)
2
3  # examine the best model
4  print(grid.best_score_)
5  print(grid.best_params_)
```

```
0.601488197668569
{'n_estimators': 900, 'max_depth': 15, 'learning_rate': 0.5}
```

```python
1   #Plotting Max_depth v/s CV_error
2   a=pd.DataFrame(grid.cv_results_)[['mean_test_score', 'std_test_score', 'params']]
3   a['max_depth'] = [d.get('max_depth') for d in a['params']]
4   b=a.sort_values(['max_depth'])
5   CV_Error=1-b['mean_test_score']
6   max_depth =b['max_depth']
7
8
9   plt.plot(max_depth,CV_Error)
10  plt.xlabel('Max_depth')
11  plt.ylabel('Cross-Validated Error')
```

Text(0,0.5,'Cross-Validated Error')

```
1   from sklearn.metrics import f1_score
2   #{'n_estimators': 900, 'max_depth': 15, 'learning_rate': 0.5}
3   xgb_optimal=XGBClassifier(n_estimators=900,max_depth=15,learning_rate=0.5,gamma=25,n_jobs=-1)
4
5   # fitting the model
6   xgb_optimal.fit(Bow_x_train, y_train)
7
8   # predict the response
9   pred_bow = xgb_optimal.predict(Bow_x_test)
10
11  # evaluate f1_score
12  f1_score = f1_score(y_test, pred_bow)
13
14  # Train & Test Error
15  print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,xgb_optimal.predict(Bow_
16  print("The overall f1_score for the Test Data is : ", metrics.f1_score(y_test,pred_bow))
```

```
The overall f1_score for the Train Data is :  0.5730344780313935
The overall f1_score for the Test Data is :  0.5529241179313679
```
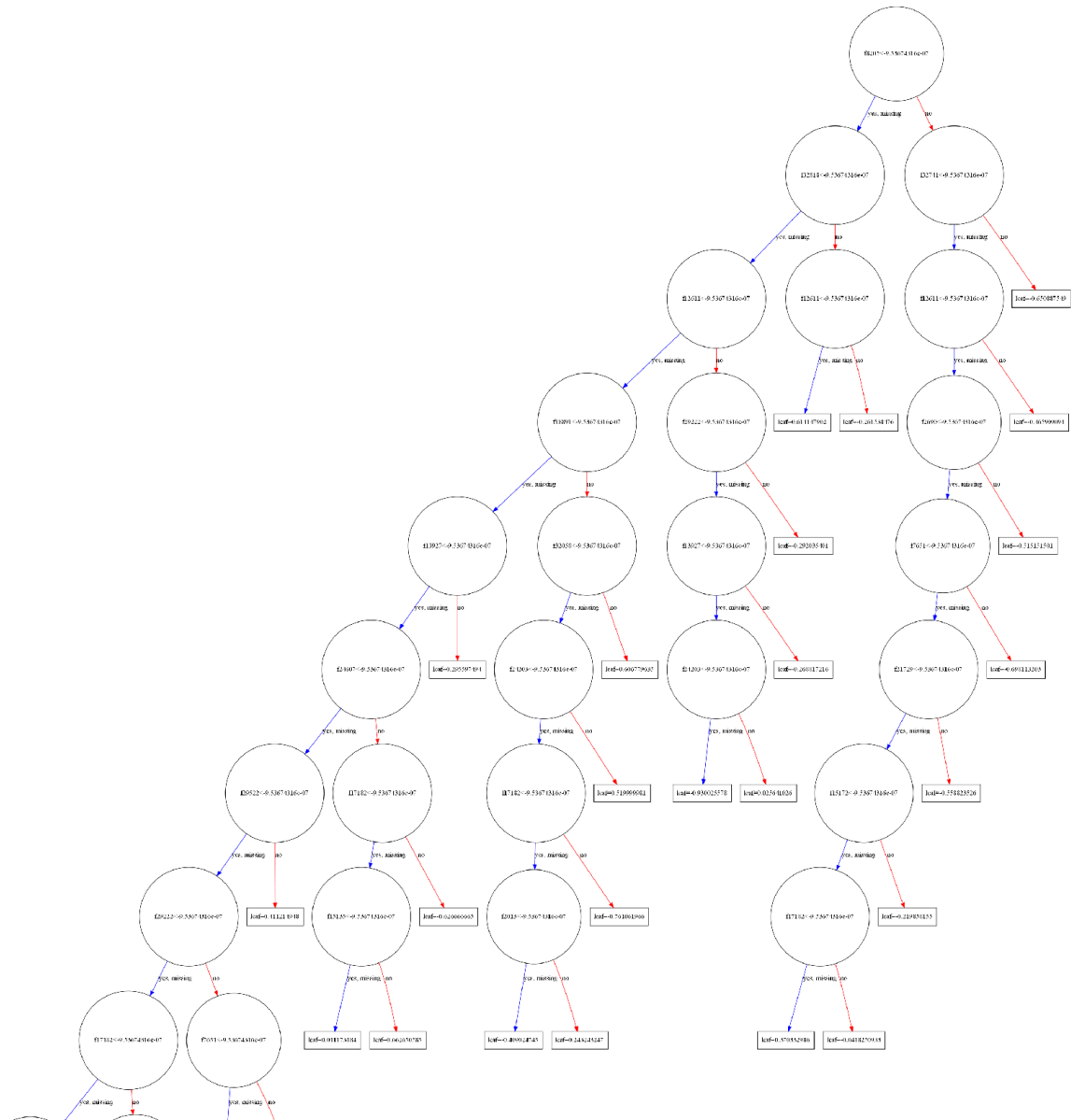
```
1  from xgboost import plot_tree
2  import xgboost as xgb
3  xgb.plot_tree(xgb_optimal,num_trees=0)
4  plt.rcParams['figure.figsize'] = [50, 50]
5  plt.show()
```

```
1   #Feature importance
2   feature_names = np.array(vocabulary.get_feature_names())
3   indices = np.argsort(xgb_optimal.feature_importances_)[::-1]
4
5   #Top 20 features
6   f=feature_names[indices[:20]]
7
8   sf = ""
9   for i in f:
10      sf += str(i)+","
11
12  print("************ Top 20 Features *******************")
13  wordcloud = WordCloud(width = 800, height = 800,
14                  background_color ='black',
15                  min_font_size = 10).generate(sf)
16
17  # plot the WordCloud image
18  plt.figure(figsize = (5,5), facecolor = None)
19  plt.imshow(wordcloud)
20  plt.axis("off")
21  plt.tight_layout(pad = 0)
22  plt.show()
23
```
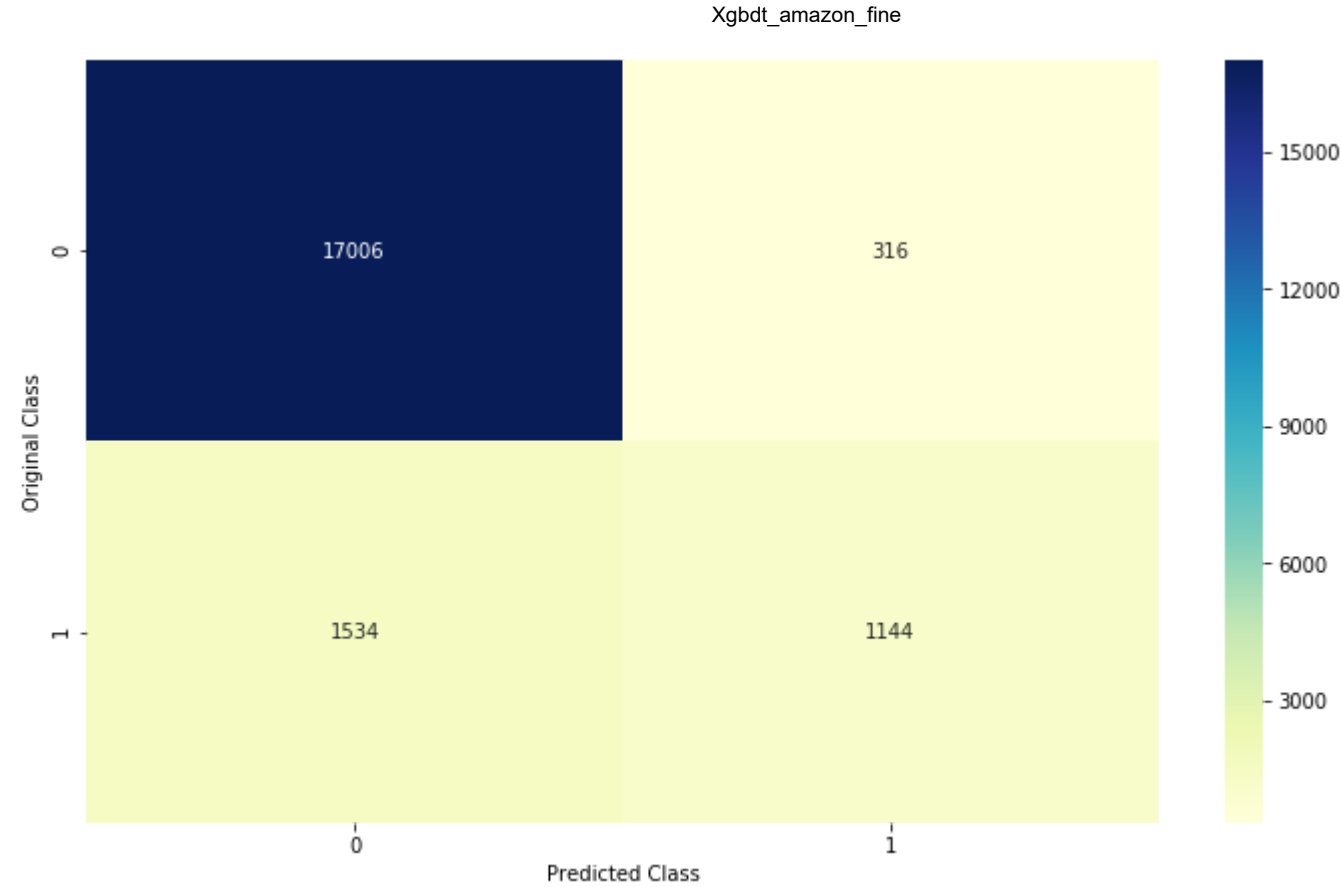
```
************ Top 20 Features *******************
```
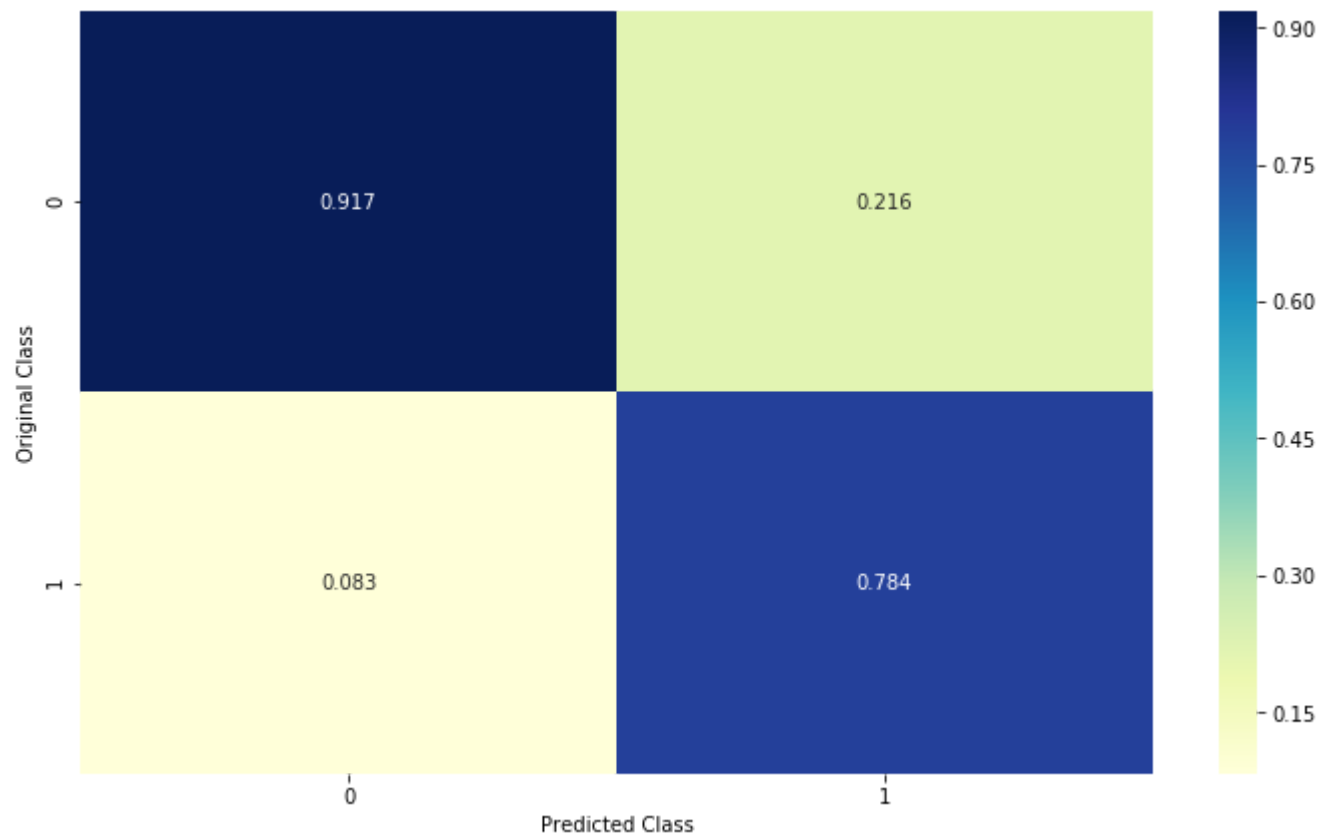
```
1   #Confusion matrix
2   C = confusion_matrix(y_test, pred_bow)
3   A =(((C.T)/(C.sum(axis=1))).T)
4   B =(C/C.sum(axis=0))
5   labels = [0,1]
```

```python
1    print("-"*20, "Confusion matrix", "-"*20)
2   plt.figure(figsize=(12,7))
3   sns.heatmap(C, annot=True, cmap="YlGnBu", fmt="g", xticklabels=labels, yticklabels=labels)
4   plt.xlabel('Predicted Class')
5   plt.ylabel('Original Class')
6   plt.show()
7
8   print("-"*20, "Precision matrix (Columm Sum=1)", "-"*20)
9   plt.figure(figsize=(12,7))
10  sns.heatmap(B, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
11   plt.xlabel('Predicted Class')
12  plt.ylabel('Original Class')
13  plt.show()
14
15      # representing B in heatmap format
16  print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
17  plt.figure(figsize=(12,7))
18  sns.heatmap(A, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
19  plt.xlabel('Predicted Class')
20  plt.ylabel('Original Class')
21  plt.show()
```
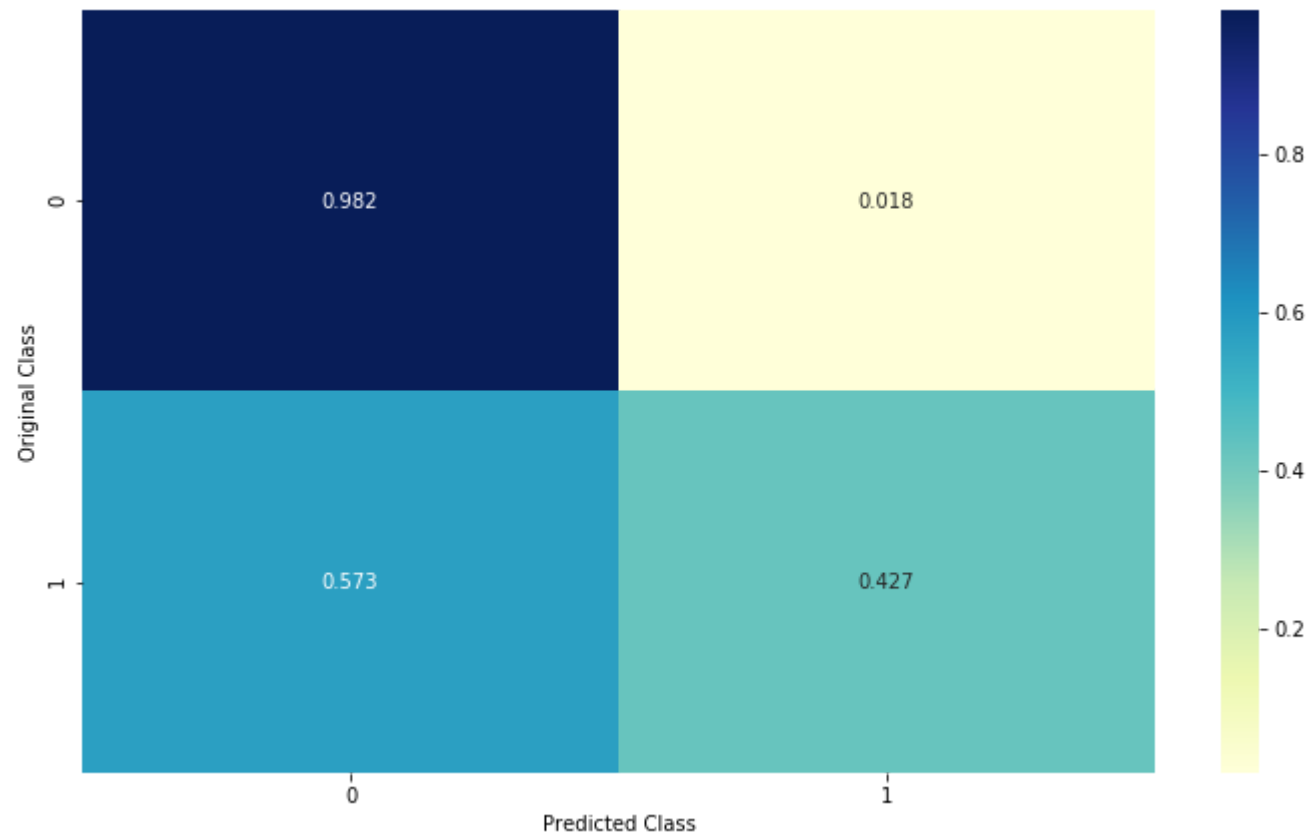
```
------------------- Confusion matrix --------------------
```

-------------------- Precision matrix (Columm Sum=1) --------------------

-------------------- Recall matrix (Row sum=1) --------------------

Tf-idf

```
1  #Initiating Vectorizer
2  count_vect = TfidfVectorizer(ngram_range=(1,2))
3
4  #Train data
5  vocabulary = count_vect.fit(x_train)
6  Tfidf_x_train= count_vect.transform(x_train)
7  print("the type of count vectorizer ",type(Tfidf_x_train))
8  print("the shape of out text BOW vectorizer ",Tfidf_x_train.get_shape())
9  print("the number of unique words ", Tfidf_x_train.get_shape()[1])
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (80000, 1013943)
the number of unique words  1013943
```

```
1  #Test data
2  Tfidf_x_test= count_vect.transform(x_test)
3  print("the type of count vectorizer ",type(Tfidf_x_test))
4  print("the shape of out text BOW vectorizer ",Tfidf_x_test.get_shape())
5  print("the number of unique words ", Tfidf_x_test.get_shape()[1])
```
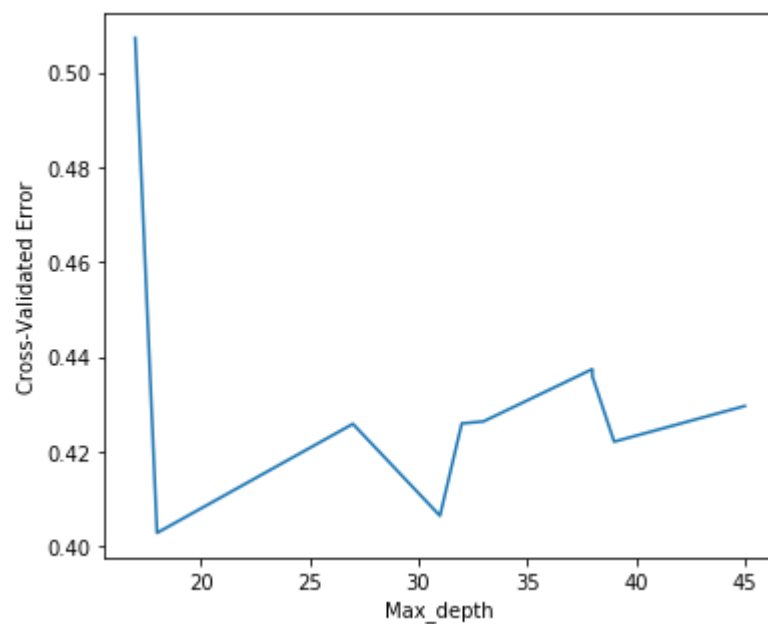
```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (20000, 1013943)
the number of unique words  1013943
```

```
1  grid_tfidf=grid.fit(Tfidf_x_train, y_train)
2
3  # examine the best model
4  print(grid_tfidf.best_score_)
5  print(grid_tfidf.best_params_)
```

```
0.5971658522910812
{'n_estimators': 600, 'max_depth': 18, 'learning_rate': 0.7}
```

```
 1  #Plotting Max_depth v/s CV_error
 2  a=pd.DataFrame(grid_tfidf.cv_results_)[['mean_test_score', 'std_test_score', 'params']]
 3  a['max_depth'] = [d.get('max_depth') for d in a['params']]
 4  b=a.sort_values(['max_depth'])
 5  CV_Error=1-b['mean_test_score']
 6  max_depth =b['max_depth']
 7
 8  plt.figure(figsize=(6,5))
 9  plt.plot(max_depth,CV_Error)
10  plt.xlabel('Max_depth')
11  plt.ylabel('Cross-Validated Error')
12
```

Text(0,0.5,'Cross-Validated Error')

```python
1
2   #{'n_estimators': 600, 'max_depth': 18, 'learning_rate': 0.7}
3   xgb_optimal=XGBClassifier(n_estimators=600,max_depth=18,learning_rate=0.7,gamma=30,n_jobs=-1)
4
5   # fitting the model
6   xgb_optimal.fit(Tfidf_x_train, y_train)
7
8   # predict the response
9   pred_tfidf = xgb_optimal.predict(Tfidf_x_test)
10
11  # evaluate accuracy
12  f1_score = f1_score(y_test, pred_tfidf)
```

```python
1   # Train & Test Error
2   print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,xgb_optimal.predict(Tfid
3   print("The overall f1_score for the Test Data is : ", metrics.f1_score(y_test,pred_tfidf))
```

```
The overall f1_score for the Train Data is :  0.5975351179432813
The overall f1_score for the Test Data is :  0.5673825820155771
```
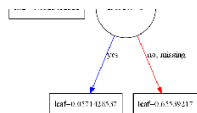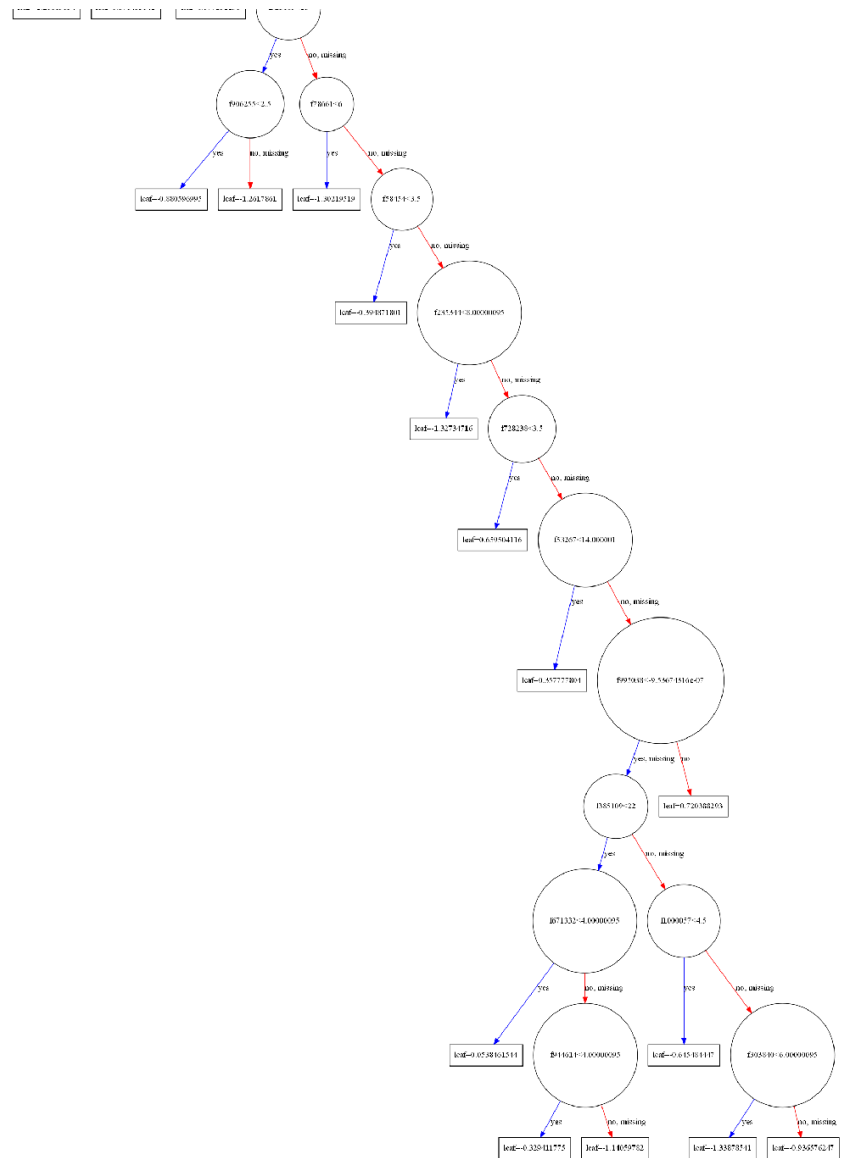
```
1  from xgboost import plot_tree
2  import xgboost as xgb
3
4  xgb.plot_tree(xgb_optimal,num_trees=0)
5  plt.rcParams['figure.figsize'] = [50, 50]
6  plt.show()
```

```
 1   #Feature importance
 2   feature_names = np.array(vocabulary.get_feature_names())
 3   indices = np.argsort(xgb_optimal.feature_importances_)[::-1]
 4
 5   #Top 20 features
 6   f=feature_names[indices[:20]]
 7
 8   sf = ""
 9   for i in f:
10       sf += str(i)+","
11
12   print("*********** Top 20 Features ******************")
13   wordcloud = WordCloud(width = 800, height = 800,
14                   background_color ='black',
15                   min_font_size = 10).generate(sf)
16
17   # plot the WordCloud image
18   plt.figure(figsize = (5,5), facecolor = None)
19   plt.imshow(wordcloud)
20   plt.axis("off")
21   plt.tight_layout(pad = 0)
22   plt.show()
23
```
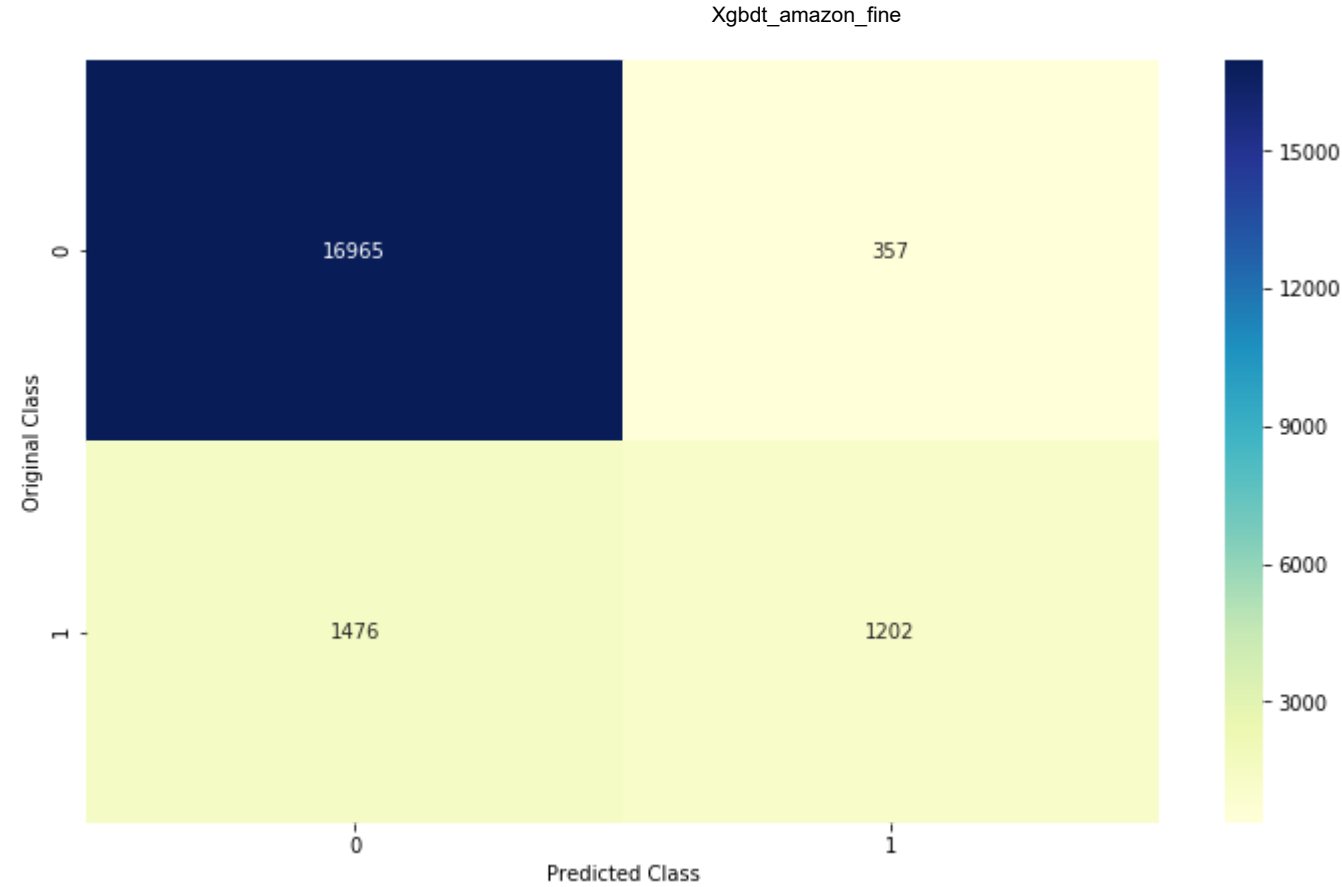
*********** Top 20 Features ******************

```
1   #Confusion matrix
2   C = confusion_matrix(y_test, pred_tfidf)
3   A =(((C.T)/(C.sum(axis=1))).T)
4   B =(C/C.sum(axis=0))
5   labels = [0,1]
```
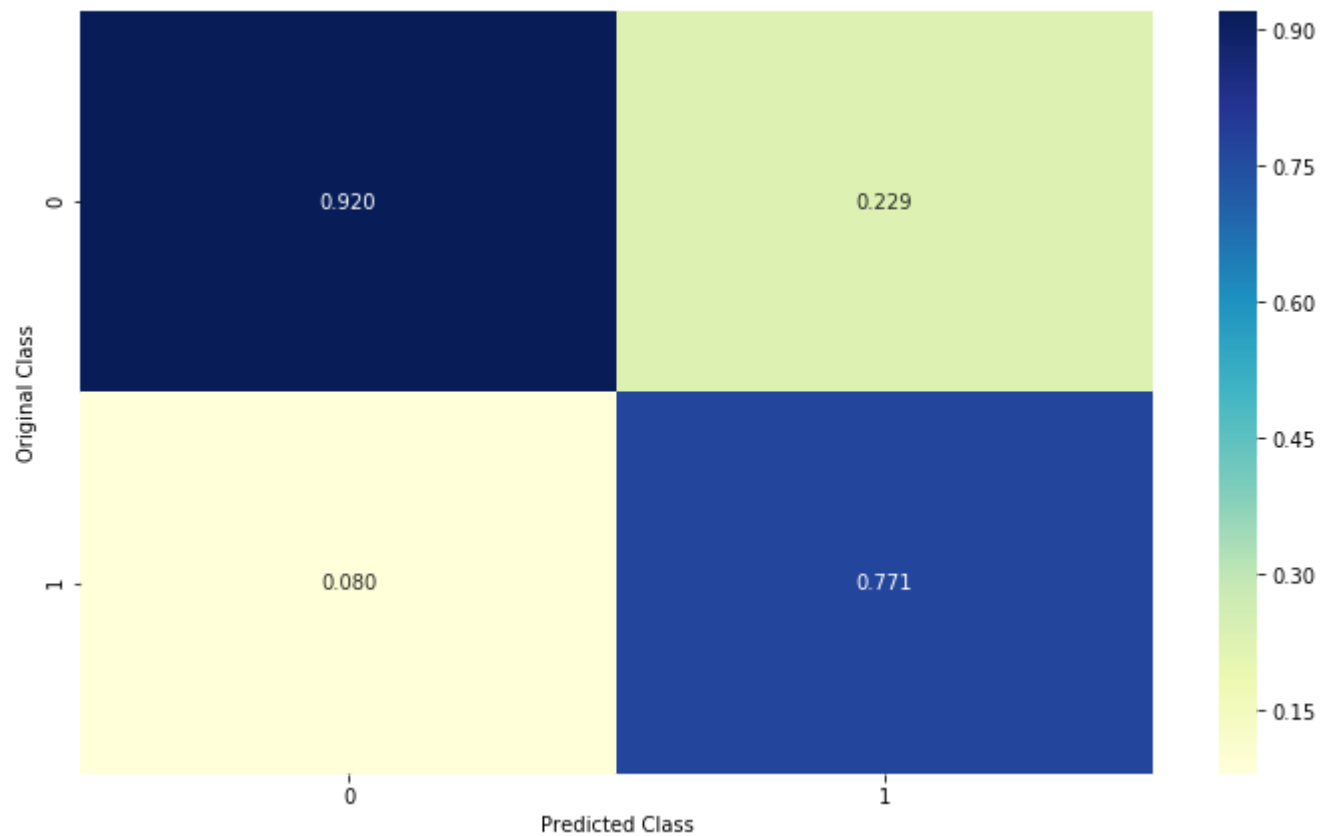
```python
 1    print("-"*20, "Confusion matrix", "-"*20)
 2   plt.figure(figsize=(12,7))
 3   sns.heatmap(C, annot=True, cmap="YlGnBu", fmt="g", xticklabels=labels, yticklabels=labels)
 4   plt.xlabel('Predicted Class')
 5   plt.ylabel('Original Class')
 6   plt.show()
 7
 8   print("-"*20, "Precision matrix (Columm Sum=1)", "-"*20)
 9   plt.figure(figsize=(12,7))
10   sns.heatmap(B, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
11    plt.xlabel('Predicted Class')
12   plt.ylabel('Original Class')
13   plt.show()
14
15       # representing B in heatmap format
16   print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
17   plt.figure(figsize=(12,7))
18   sns.heatmap(A, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
19   plt.xlabel('Predicted Class')
20   plt.ylabel('Original Class')
21   plt.show()
```
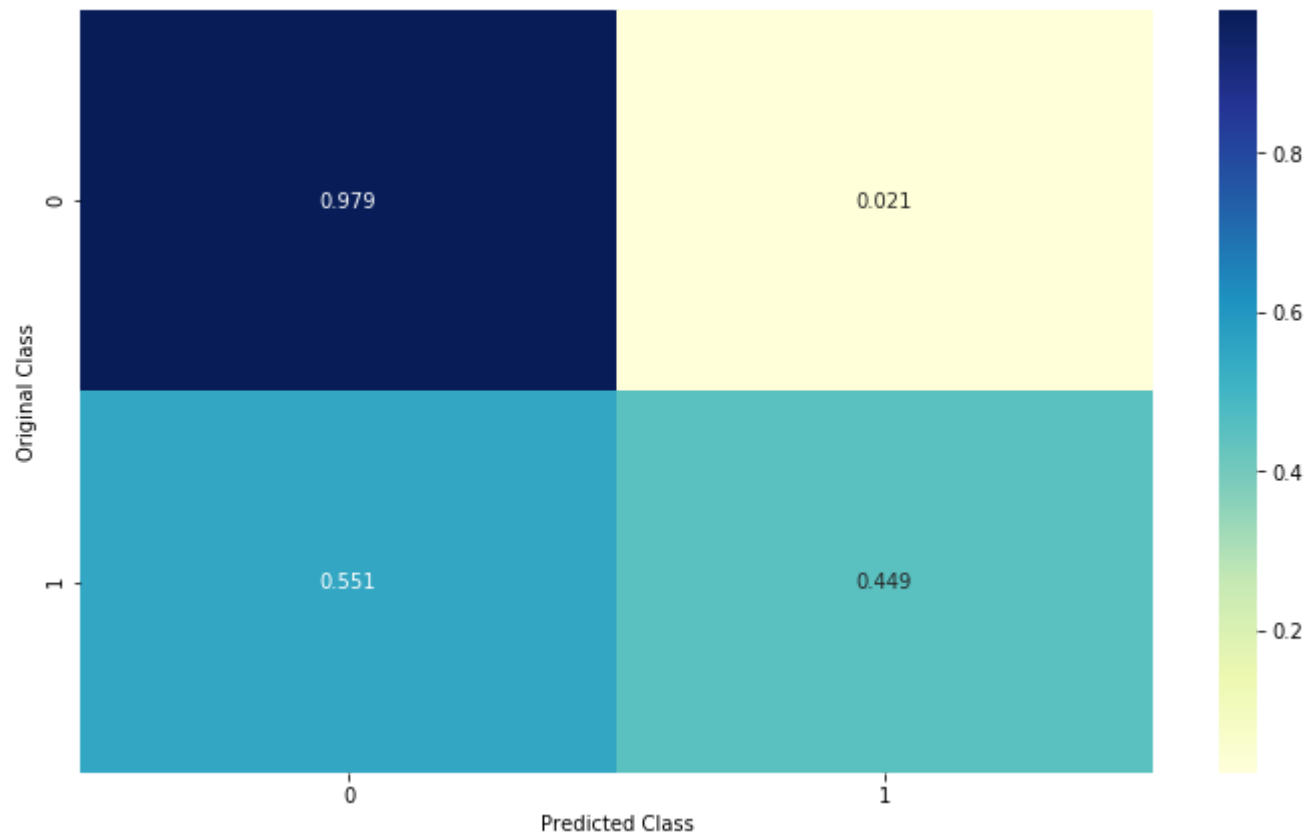
```
-------------------- Confusion matrix --------------------
```

-------------------- Precision matrix (Columm Sum=1) --------------------

-------------------- Recall matrix (Row sum=1) --------------------

## Avg-W2Vec

```
1  #W2V list of Training data
2  i=0
3  list_of_sent_train=[]
4  for sent in train_data['CleanedText'].values:
5      list_of_sent_train.append(sent.split())
```

```
1   #W2V List of Test data
2   i=0
3   list_of_sent_test=[]
4   for sent in test_data['CleanedText'].values:
5       list_of_sent_test.append(sent.split())
```

```
1   #Training W2V train model
2   # min_count = 5 considers only words that occured atleast 5 times
3   w2v_model_train=Word2Vec(list_of_sent_train,min_count=5,size=200, workers=6)
```

```
1   w2v_words_train = list(w2v_model_train.wv.vocab)
2   print("number of words that occured minimum 5 times ",len(w2v_words_train))
3   print("sample words ", w2v_words_train[0:50])
```

number of words that occured minimum 5 times  11361
sample words  ['witti', 'littl', 'book', 'make', 'son', 'laugh', 'loud', 'car', 'drive', 'along', 'alway', 'sing', 'refrai
n', 'hes', 'learn', 'whale', 'india', 'droop', 'love', 'new', 'word', 'introduc', 'silli', 'classic', 'will', 'bet', 'stil
l', 'abl', 'memori', 'colleg', 'rememb', 'see', 'show', 'air', 'televis', 'year', 'ago', 'child', 'sister', 'later', 'bough
t', 'day', 'thirti', 'someth', 'use', 'seri', 'song', 'student', 'teach', 'preschool']

```
1   #Train data
2   # average Word2Vec
3   # compute average word2vec for each review.
4   sent_vectors_train_avgw2v = []; # the avg-w2v for each sentence/review is stored in this list
5   for sent in list_of_sent_train: # for each review/sentence
6       sent_vec = np.zeros(200) # as word vectors are of zero length
7       cnt_words =0; # num of words with a valid vector in the sentence/review
8       for word in sent: # for each word in a review/sentence
9           if word in w2v_words_train:
10              vec = w2v_model_train.wv[word]
11              sent_vec += vec
12              cnt_words += 1
13      if cnt_words != 0:
14          sent_vec /= cnt_words
15      sent_vectors_train_avgw2v.append(sent_vec)
16  print(len(sent_vectors_train_avgw2v))
17  print(len(sent_vectors_train_avgw2v[0]))
```

```
80000
200
```

```
1   #Test data
2   # average Word2Vec
3   # compute average word2vec for each review.
4   sent_vectors_test_avgw2v = []; # the avg-w2v for each sentence/review is stored in this list
5   for sent in list_of_sent_test: # for each review/sentence
6       sent_vec = np.zeros(200) # as word vectors are of zero length
7       cnt_words =0; # num of words with a valid vector in the sentence/review
8       for word in sent: # for each word in a review/sentence
9           if word in w2v_words_train:
10              vec = w2v_model_train.wv[word]
11              sent_vec += vec
12              cnt_words += 1
13      if cnt_words != 0:
14          sent_vec /= cnt_words
15      sent_vectors_test_avgw2v.append(sent_vec)
16  print(len(sent_vectors_test_avgw2v))
17  print(len(sent_vectors_test_avgw2v[0]))
```

```
20000
200
```

```
1   sent_vectors_train_avgw2v=np.asarray(sent_vectors_train_avgw2v)
2   sent_vectors_test_avgw2v=np.asarray(sent_vectors_test_avgw2v)
```
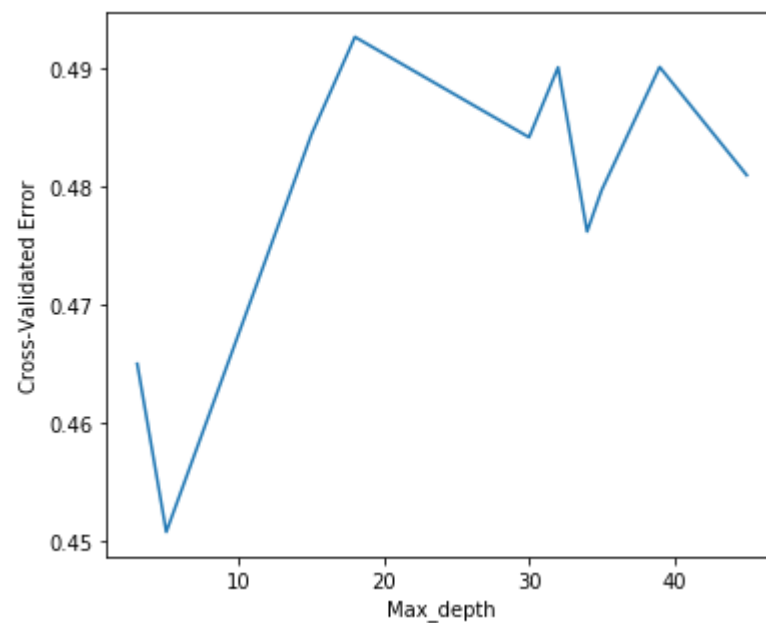
```
1   grid.fit(sent_vectors_train_avgw2v, y_train)
2
3   # examine the best model
4   print(grid.best_score_)
5   print(grid.best_params_)
```

```
0.5492236670295039
{'n_estimators': 600, 'max_depth': 5, 'learning_rate': 0.6}
```

```python
1   #Plotting Max_depth v/s CV_error
2   a=pd.DataFrame(grid.cv_results_)[['mean_test_score', 'std_test_score', 'params']]
3   a['max_depth'] = [d.get('max_depth') for d in a['params']]
4   b=a.sort_values(['max_depth'])
5   CV_Error=1-b['mean_test_score']
6   max_depth =b['max_depth']
7
8   plt.figure(figsize=(6,5))
9   plt.plot(max_depth,CV_Error)
10  plt.xlabel('Max_depth')
11  plt.ylabel('Cross-Validated Error')
```

Text(0,0.5,'Cross-Validated Error')

```python
1
2   #{'n_estimators': 600, 'max_depth': 5, 'learning_rate': 0.6}
3   xgb_optimal=XGBClassifier(n_estimators=600,max_depth=5,learning_rate=0.6,gamma=20,n_jobs=-1)
4
5   # fitting the model
6   xgb_optimal.fit(sent_vectors_train_avgw2v, y_train)
7
8   # predict the response
9   pred_avg_w2v = xgb_optimal.predict(sent_vectors_test_avgw2v)
10
11  # evaluate f1_score
12  f1_score = f1_score(y_test, pred_avg_w2v)
13
14  # Train & Test Error
15  print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,xgb_optimal.predict(sen
16  print("The overall f1_score for the Test Data is : ", metrics.f1_score(y_test,pred_avg_w2v))
```

```
The overall f1_score for the Train Data is :  0.6463444553483808
The overall f1_score for the Test Data is :  0.5497440462942356
```
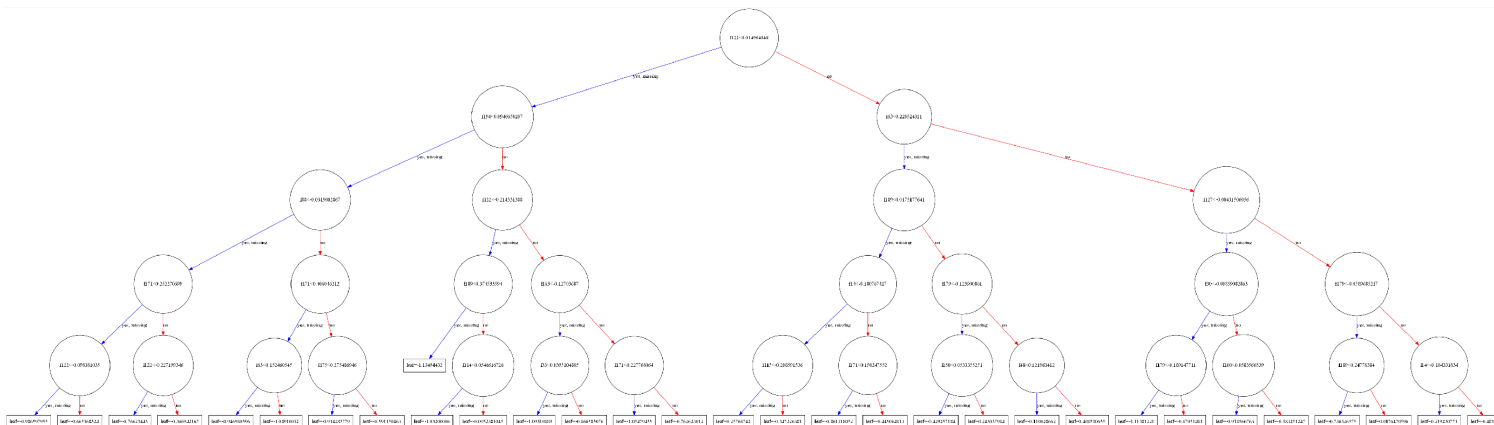
```
1   from xgboost import plot_tree
2   import xgboost as xgb
3   xgb.plot_tree(xgb_optimal,num_trees=0)
4   plt.figure(figsize=(50,50))
5   plt.show()
```
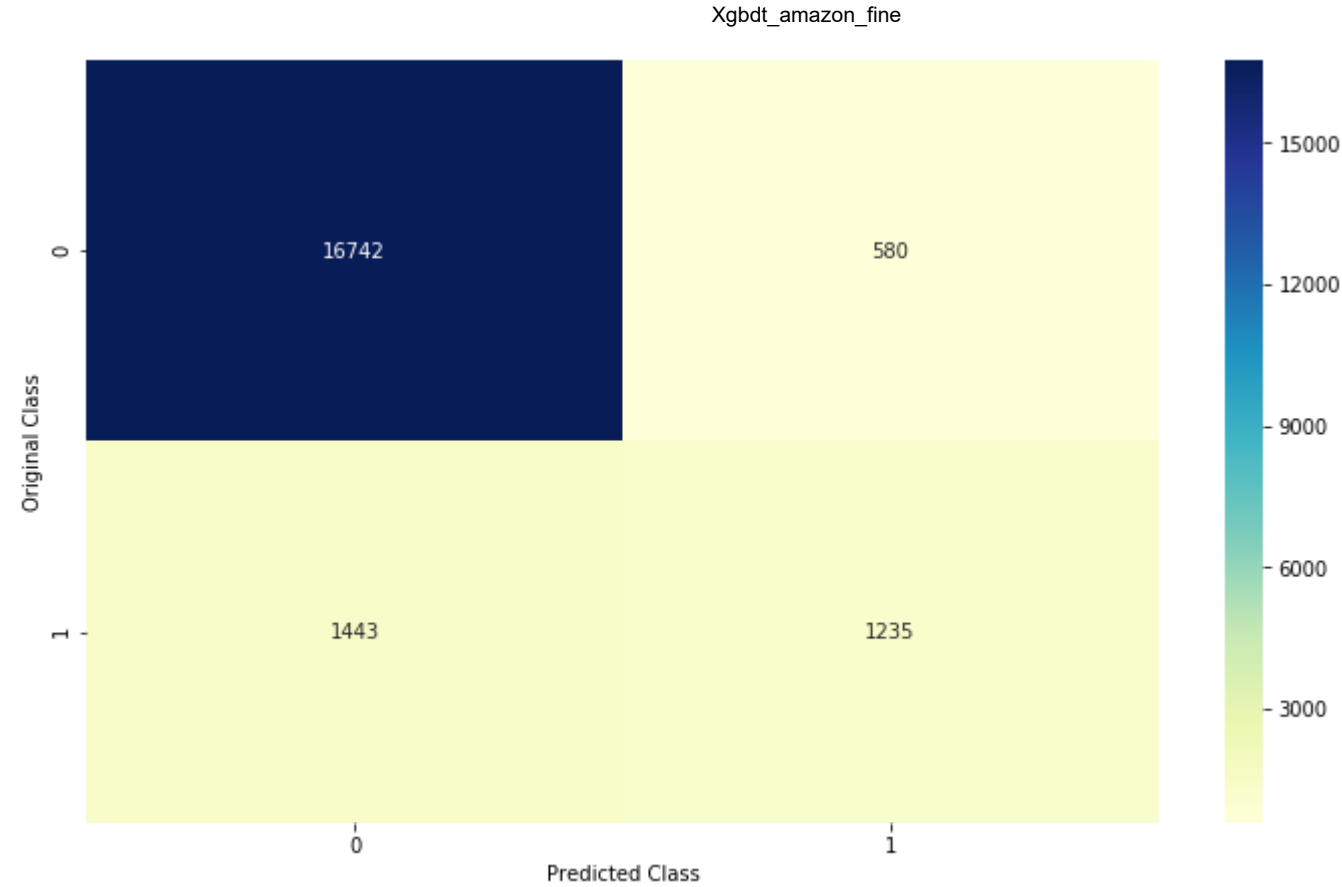


<Figure size 3600x3600 with 0 Axes>

```
1   #Confusion matrix
2   C = confusion_matrix(y_test, pred_avg_w2v)
3   A =(((C.T)/(C.sum(axis=1))).T)
4   B =(C/C.sum(axis=0))
5   labels = [0,1]
```
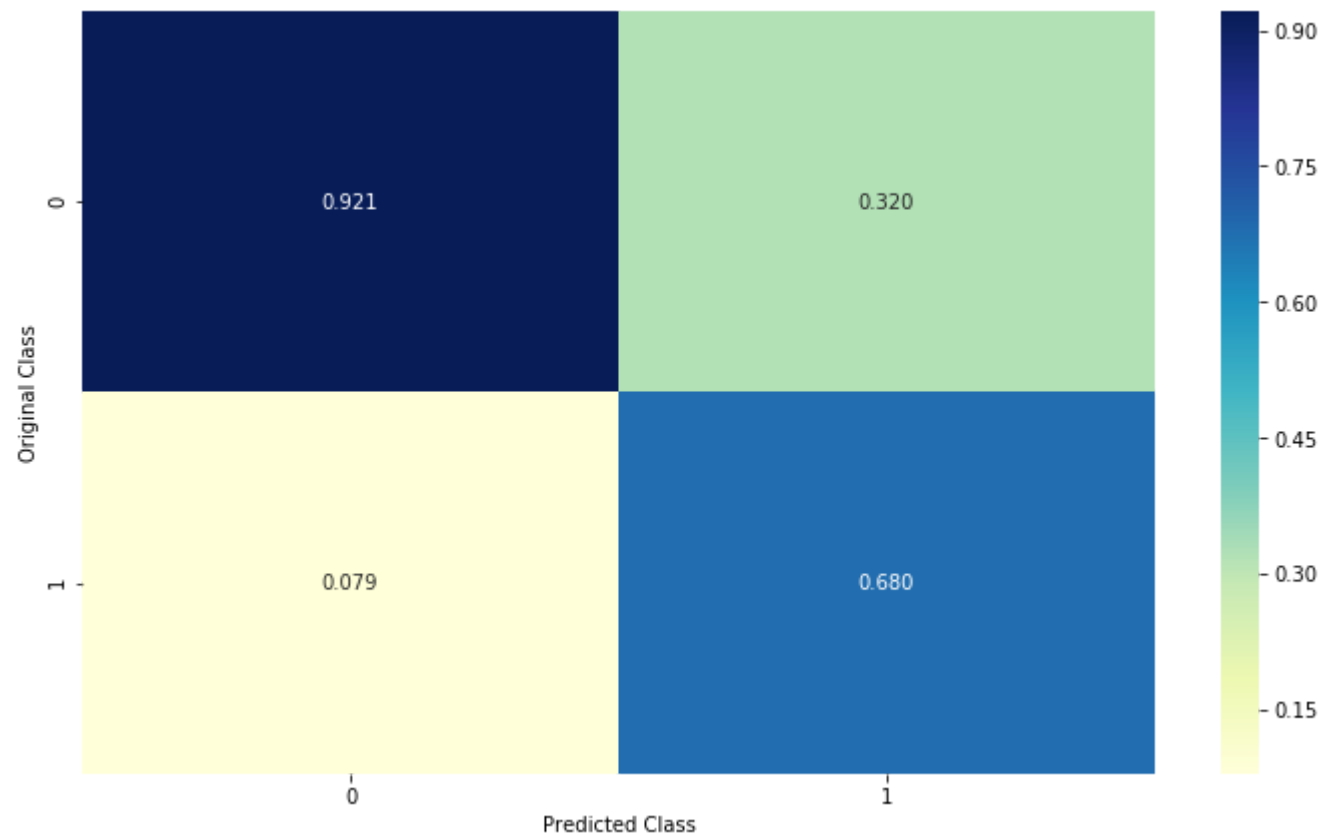
```
1    print("-"*20, "Confusion matrix", "-"*20)
2   plt.figure(figsize=(12,7))
3   sns.heatmap(C, annot=True, cmap="YlGnBu", fmt="g", xticklabels=labels, yticklabels=labels)
4   plt.xlabel('Predicted Class')
5   plt.ylabel('Original Class')
6   plt.show()
7
8   print("-"*20, "Precision matrix (Columm Sum=1)", "-"*20)
9   plt.figure(figsize=(12,7))
10  sns.heatmap(B, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
11   plt.xlabel('Predicted Class')
12  plt.ylabel('Original Class')
13  plt.show()
14
15      # representing B in heatmap format
16  print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
17  plt.figure(figsize=(12,7))
18  sns.heatmap(A, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
19  plt.xlabel('Predicted Class')
20  plt.ylabel('Original Class')
21  plt.show()
```
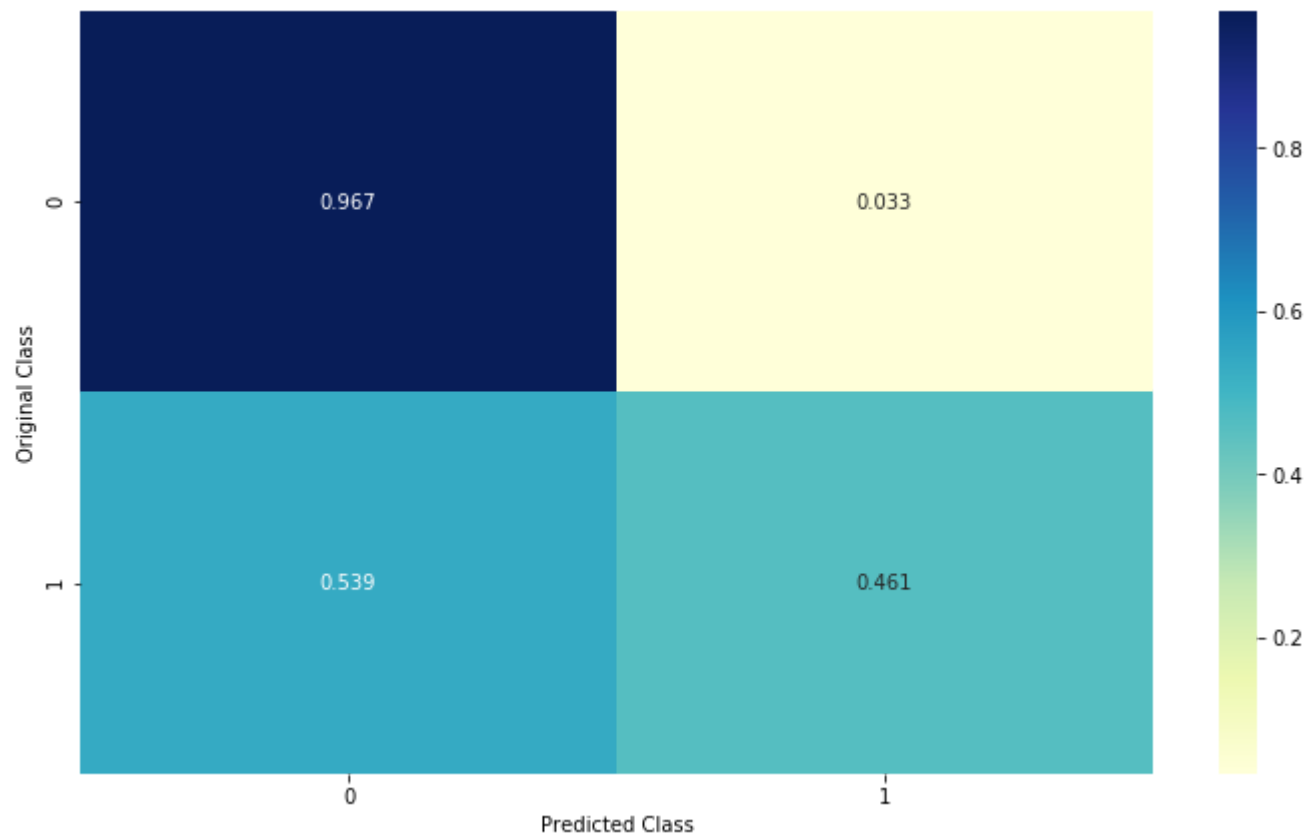
```
-------------------- Confusion matrix --------------------
```

-------------------- Precision matrix (Columm Sum=1) --------------------

-------------------- Recall matrix (Row sum=1) --------------------

## Tfidf-W2V

Importing Pre-featurized pickle file

```
1  #Train
2  with open('train.pickle', 'rb') as f:
3      tfidf_w2v_sent_vectors_train = pickle.load(f)
```

```
1  #Test
2  with open('test.pickle', 'rb') as f:
3      tfidf_w2v_sent_vectors_test = pickle.load(f)
```

```
1  tfidf_w2v_sent_vectors_train=np.asarray(tfidf_w2v_sent_vectors_train)
2  tfidf_w2v_sent_vectors_test=np.asarray(tfidf_w2v_sent_vectors_test)
```
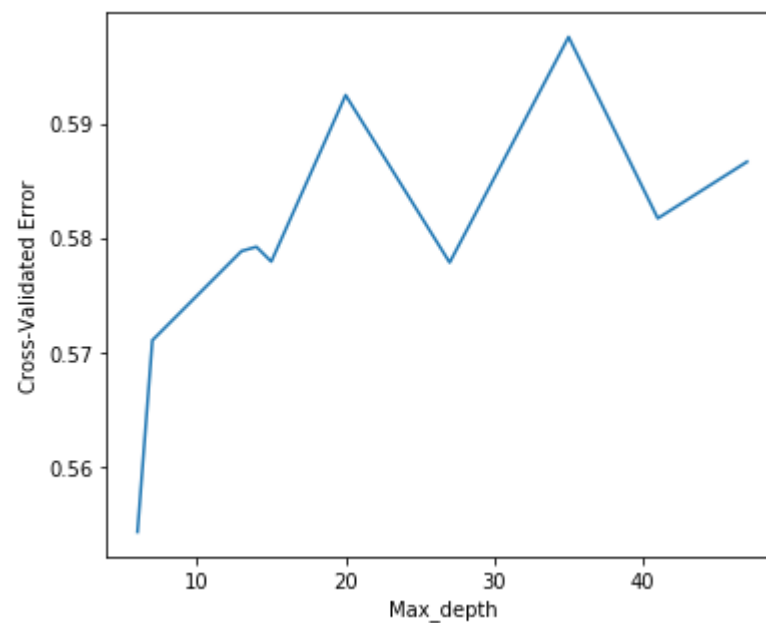
```
1  grid.fit(tfidf_w2v_sent_vectors_train, y_train)
2
3  # examine the best model
4  print(grid.best_score_)
5  print(grid.best_params_)
```

```
0.44570093765091245
{'n_estimators': 900, 'max_depth': 6, 'learning_rate': 1.0}
```

```
1   #Plotting Max_depth v/s CV_error
2   a=pd.DataFrame(grid.cv_results_)[['mean_test_score', 'std_test_score', 'params']]
3   a['max_depth'] = [d.get('max_depth') for d in a['params']]
4   b=a.sort_values(['max_depth'])
5   CV_Error=1-b['mean_test_score']
6   max_depth =b['max_depth']
7
8   plt.figure(figsize=(6,5))
9   plt.plot(max_depth,CV_Error)
10  plt.xlabel('Max_depth')
11  plt.ylabel('Cross-Validated Error')
```

Text(0,0.5,'Cross-Validated Error')

```python
1  from sklearn.metrics import f1_score
2  #{'n_estimators': 900, 'max_depth': 6, 'learning_rate': 1.0}
3  xgb_optimal=XGBClassifier(n_estimators=900,max_depth=6,learning_rate=1.0,gamma=20 ,n_jobs=-1)
4
5  # fitting the model
6  xgb_optimal.fit(tfidf_w2v_sent_vectors_train, y_train)
7
8  # predict the response
9  pred_tfidf_w2v = xgb_optimal.predict(tfidf_w2v_sent_vectors_test)
10
11  # evaluate f1_score
12  f1_score = f1_score(y_test, pred_tfidf_w2v)
13
14  # Train & Test Error
15  print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,xgb_optimal.predict(tfi(
16  print("The overall f1_score for the Test Data is : ", metrics.f1_score(y_test,pred_tfidf_w2v))
```

```
The overall f1_score for the Train Data is :  0.565834065520271
The overall f1_score for the Test Data is :  0.44323835368611486
```
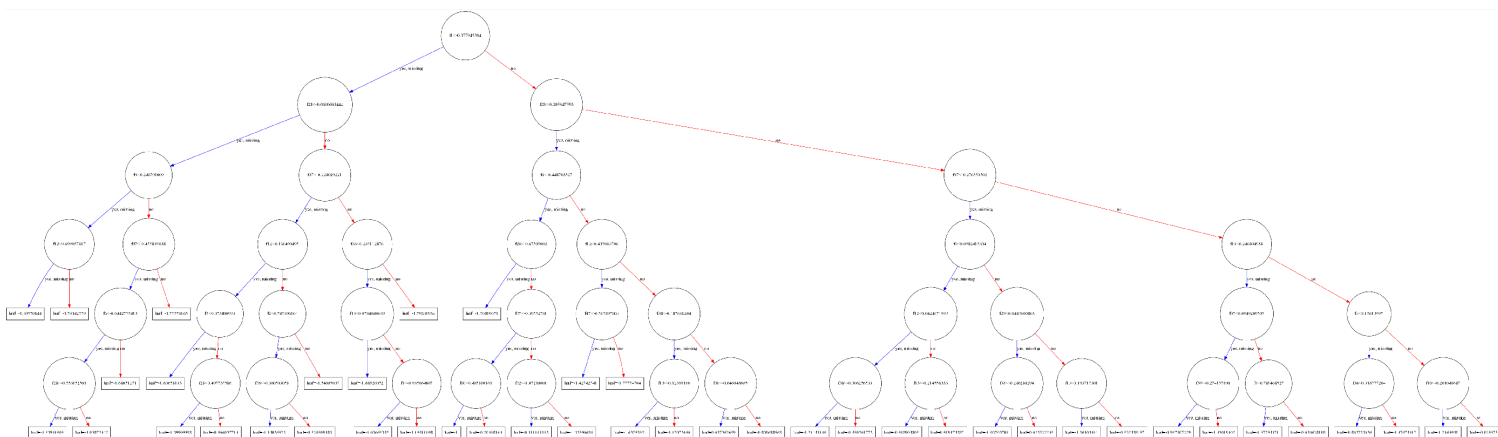
```
1   from xgboost import plot_tree
2   import xgboost as xgb
3   xgb.plot_tree(xgb_optimal,num_trees=0)
4   plt.figure(figsize=(50,50))
5   plt.show()
```
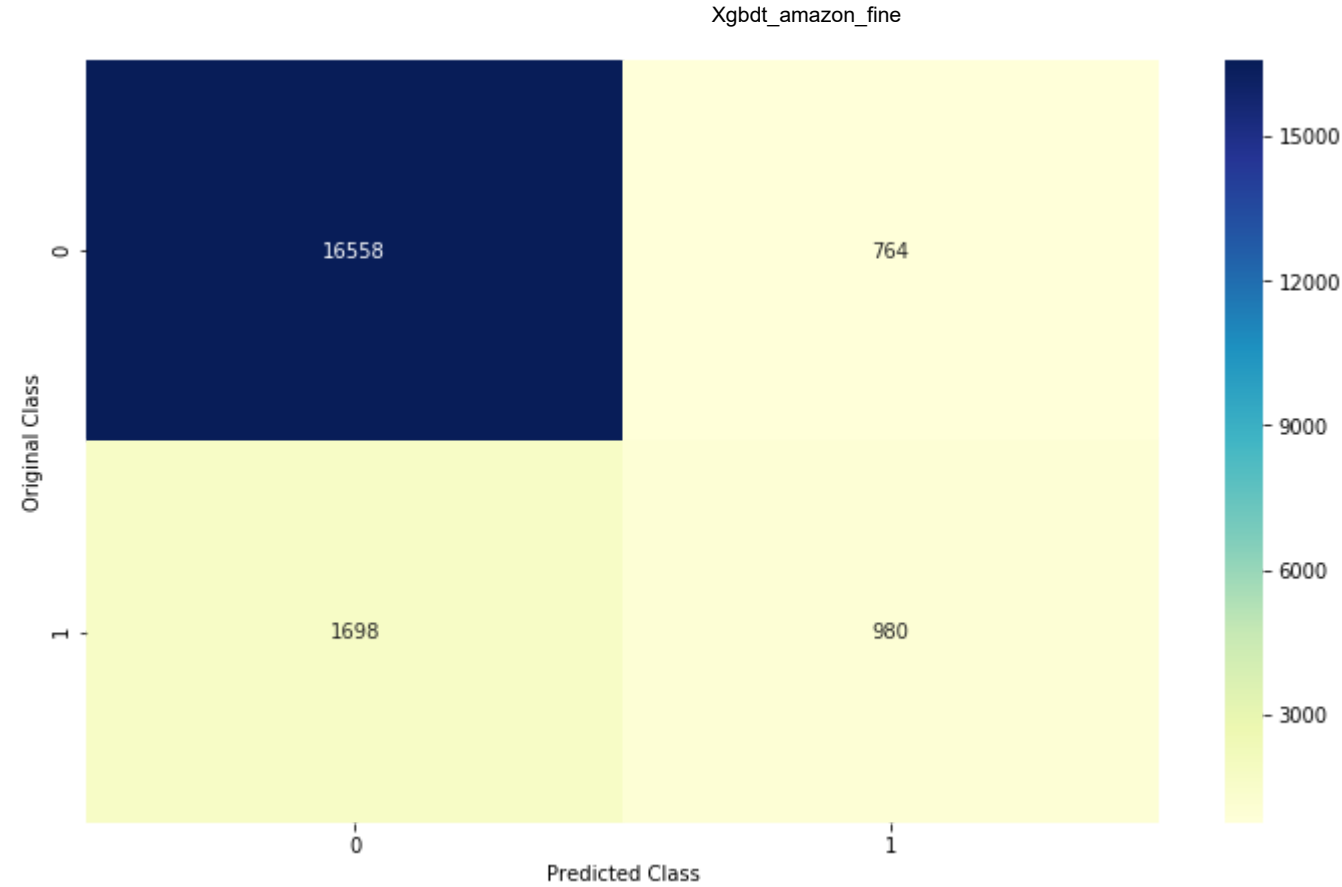


<Figure size 3600x3600 with 0 Axes>

```
1   #Confusion matrix
2   C = confusion_matrix(y_test, pred_tfidf_w2v)
3   A =(((C.T)/(C.sum(axis=1))).T)
4   B =(C/C.sum(axis=0))
5   labels = [0,1]
```
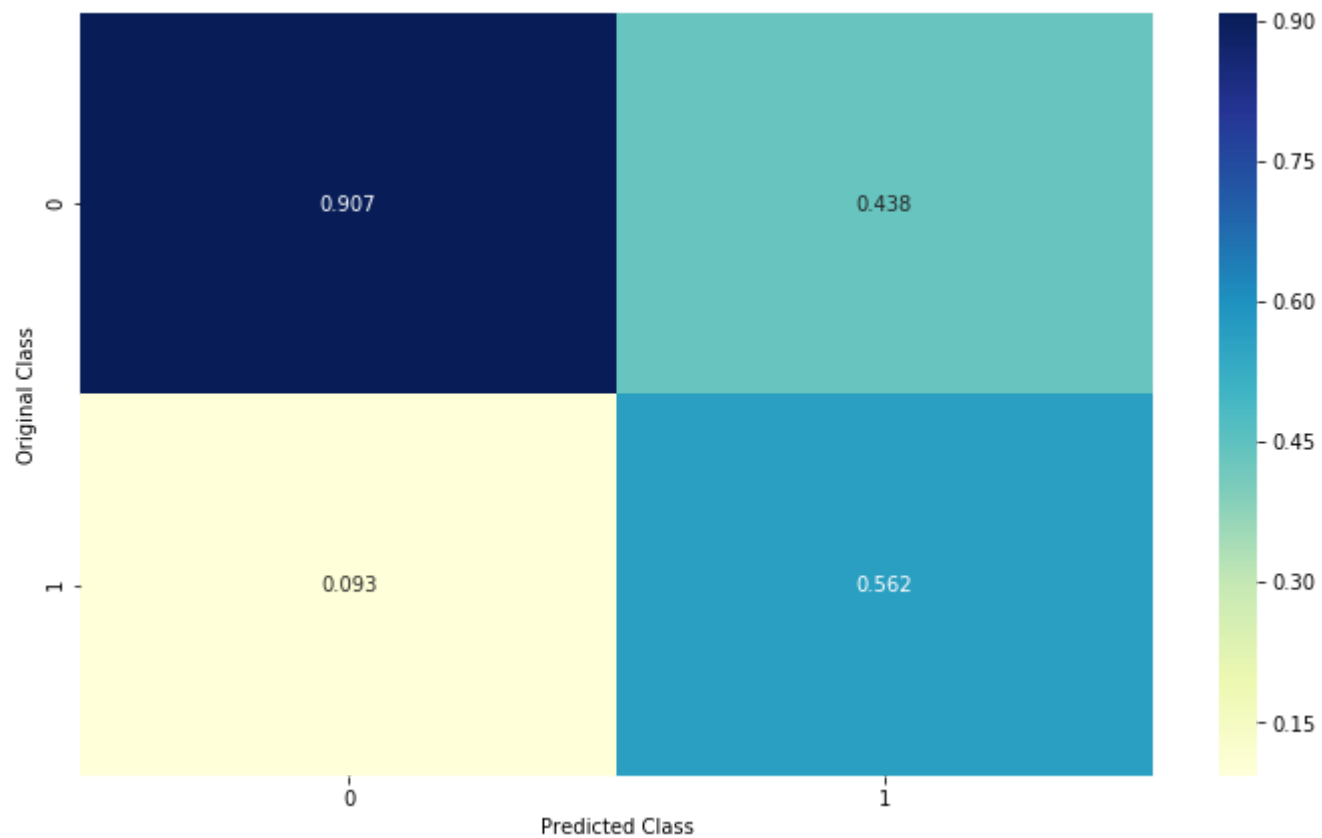
```
 1    print("-"*20, "Confusion matrix", "-"*20)
 2   plt.figure(figsize=(12,7))
 3   sns.heatmap(C, annot=True, cmap="YlGnBu", fmt="g", xticklabels=labels, yticklabels=labels)
 4   plt.xlabel('Predicted Class')
 5   plt.ylabel('Original Class')
 6   plt.show()
 7
 8   print("-"*20, "Precision matrix (Columm Sum=1)", "-"*20)
 9   plt.figure(figsize=(12,7))
10   sns.heatmap(B, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
11    plt.xlabel('Predicted Class')
12   plt.ylabel('Original Class')
13   plt.show()
14
15       # representing B in heatmap format
16   print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
17   plt.figure(figsize=(12,7))
18   sns.heatmap(A, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
19   plt.xlabel('Predicted Class')
20   plt.ylabel('Original Class')
21   plt.show()
```

```
-------------------- Confusion matrix --------------------
```

-------------------- Precision matrix (Columm Sum=1) --------------------

-------------------- Recall matrix (Row sum=1) --------------------