

Support Vector Machine 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 - unique 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 Support Vector Machine 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
31
32 from tqdm import tqdm
```

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

```
1 train_data=train_data.astype(str)
2 test_data=test_data.astype(str)
```

```
1 train_data.shape

(80000, 13)
```

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

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

```
1 test_data.shape

(20000, 13)
```

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

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

```
1 #Train data
2 y_train = train_data['Score']
3 x_train = train_data['CleanedText']
4
5 #Test data
6 y_test = test_data['Score']
7 x_test = test_data['CleanedText']
```

```
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)
```

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

Since SVM is computationally expensive therefore running SGD classifier on different featurization like BOW, Tfidf(bigram), Avg-W2V & Tfidf-W2V. For the featurization where we will get best performance metric(f1_score), we will run SVC classifier

RandomisedSearch CV

```

1 alpha_range = [0.000001,0.00001,0.0001,0.001,0.01,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0,1,2,3,4,5,6,7,
2 epsilon_range=[0.0000001,0.000001,0.00001,0.0001,0.001,0.01,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]
3 T= TimeSeriesSplit(n_splits=5)
4 weight=[None,'balanced']
5
6 param_distributions = dict(alpha=alpha_range,epsilon=epsilon_range,class_weight=weight)
7 print(param_distributions)
8
9 # instantiate and fit the grid
10 grid = RandomizedSearchCV(SGDClassifier(), param_distributions, cv=T, scoring='f1',n_iter=30, return_tr:

```

```

{'alpha': [1e-06, 1e-05, 0.0001, 0.001, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1, 2, 3, 4, 5, 6, 7, 8, 9,
10, 15, 20, 30, 40, 50], 'epsilon': [1e-07, 1e-06, 1e-05, 0.0001, 0.001, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9,
1.0], 'class_weight': [None, 'balanced']}

```

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 #Standardizing Bow_x_train and Bow_x_test
2 from sklearn.preprocessing import StandardScaler
3 Standard=StandardScaler(with_mean=False)
4 Bow_x_train = Standard.fit_transform(Bow_x_train)
5 Bow_x_test = Standard.transform(Bow_x_test)
6
7 print(Bow_x_train.shape)
8 print(Bow_x_test.shape)
```

```
(80000, 33433)
(20000, 33433)
```

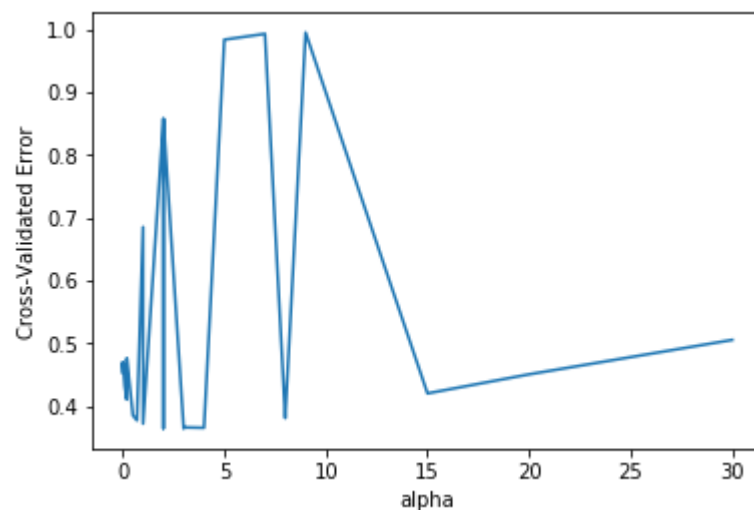
Fitting Randomised search CV on BOW

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

```
0.6363102574032368
{'epsilon': 0.4, 'class_weight': 'balanced', 'alpha': 3}
```

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

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



```
1  #{'epsilon': 0.4, 'class_weight': 'balanced', 'alpha': 3}
2  svm_optimal=SGDClassifier(alpha=3,epsilon=0.4,class_weight='balanced')
3
4  # fitting the model
5  svm_optimal.fit(Bow_x_train, y_train)
6
7  # predict the response
8  pred_bow = svm_optimal.predict(Bow_x_test)
9
10 # evaluate f1_score
11 f1_score = f1_score(y_test, pred_bow)
12
13 # Train & Test Error
14 print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,svm_optimal.predict(Bow_
15 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 : 0.7488488356795158

The overall f1_score for the Test Data is : 0.6735686274509805

Pertubation test

```
1  # Re-training the model after adding noise
2  Epsilon = np.random.normal(loc=0,scale =0.01)
3  Noise_Bow_x_train=Bow_x_train
4  Noise_Bow_x_train.data+=Epsilon
```

```
1  Noise_Bow_x_train.shape
```

(80000, 33433)


```
1  #{'epsilon': 0.4, 'class_weight': 'balanced', 'alpha': 3}
2  svm_optimal_noise=SGDClassifier(alpha=3,epsilon=0.4,class_weight='balanced')
3
4  # fitting the model
5  svm_optimal_noise.fit(Noise_Bow_x_train, y_train)
6
7  # predict the response
8  pred_bow = svm_optimal_noise.predict(Bow_x_test)
9
10 # evaluate f1_score
11 f1_score = f1_score(y_test, pred_bow)
12
13 # Train & Test Error
14 print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,svm_optimal_noise.predict(Bow_x_train)))
15 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 : 0.7488608482299334

The overall f1_score for the Test Data is : 0.6734597899357265

```
1  #Features
2  feature_names = np.array(vocabulary.get_feature_names())
3  feature_names.shape
```

(33433,)

```
1  #Weights before adding noise
2  svm_optimal.coef_.shape
```

(1, 33433)

```
1  #Weights after adding noise
2  svm_optimal_noise.coef_.shape
```

(1, 33433)

```

1 merge_arr = np.concatenate([svm_optimal.coef_, svm_optimal_noise.coef_], axis=0)
2 merge=pd.DataFrame(data=merge_arr,columns=feature_names).transpose()
3 merge[2]=((merge[1]-merge[0])/merge[0])*100
4 merge

```

	0	1	2
aaa	-0.000568	-0.000568	-0.006081
aaaaaaaaagghh	-0.000670	-0.000670	-0.004300
aaaaah	-0.000758	-0.000757	-0.006081
aaaaahhhhhhhhhhhhhhh	-0.000268	-0.000268	-0.004300
aaaah	-0.000134	-0.000268	99.991400
aaah	-0.000387	-0.000464	19.991062
aachen	0.001966	0.001966	-0.004300
aad	0.000000	0.000000	NaN
aadp	-0.000536	-0.000536	-0.004300
aafco	0.000438	0.000438	-0.010533
aagh	-0.000670	-0.000670	-0.004300
aah	-0.000852	-0.000852	-0.006081
aahh	-0.000670	-0.000670	-0.004300
aand	-0.000670	-0.000670	-0.004300
aardvark	-0.000670	-0.000670	-0.004300
aarrgh	0.003931	0.003931	-0.004300
ab	-0.000506	-0.000506	-0.011377
aback	-0.000867	-0.000915	5.445340
abandon	-0.000282	-0.000119	-57.928600
abaolut	-0.000670	-0.000670	-0.004300
abattoir	-0.000134	-0.000134	-0.004300
abba	-0.000773	-0.000773	-0.007448
abbey	-0.000568	-0.000663	16.659572
abbi	-0.001377	-0.001317	-4.357023

	0	1	2
abbott	-0.000095	-0.000095	-0.006081
abbrevi	-0.000077	-0.000155	99.985104
abc	-0.000460	-0.000391	-15.048432
abcstor	-0.000670	-0.000670	-0.004300
abd	-0.000473	-0.000473	-0.006081
abdomen	-0.000387	-0.000387	-0.007448
...
zot	-0.001004	-0.001004	-0.008600
zotz	-0.000599	-0.000599	-0.012162
zour	0.001788	0.001787	-0.008600
zout	-0.000947	-0.000947	-0.006081
zowi	-0.000402	-0.000402	-0.004300
zreport	-0.000670	-0.000670	-0.004300
zsweet	-0.000670	-0.000670	-0.004300
zuc	-0.000670	-0.000670	-0.004300
zucchini	-0.000297	-0.000363	22.550771
zuccini	-0.001620	-0.001569	-3.136021
zuccnini	0.000000	0.000000	NaN
zuchinni	-0.000670	-0.000670	-0.004300
zuke	0.000402	0.000512	27.418087
zulu	-0.000670	-0.000670	-0.004300
zum	0.000000	-0.000134	-inf
zummi	0.000000	-0.000134	-inf
zune	-0.000670	-0.000670	-0.004300
zupreem	-0.000268	-0.000268	-0.004300
zurich	-0.000670	-0.000670	-0.004300
zwar	-0.000134	-0.000134	-0.004300
zwieback	0.002064	0.002219	7.483452
zwiebeck	-0.000670	-0.000670	-0.004300

	0	1	2
zydeco	-0.000670	-0.000670	-0.004300
zzzzzs	-0.000663	-0.000568	-14.290927
zzzzzz	0.000000	0.000000	NaN
zzzzzzzz	0.002948	0.002948	-0.004300
zzzzzzzzzz	-0.000268	-0.000268	-0.004300
zzzzzzzzzzzz	0.000000	0.000000	NaN
zzzzzzzzzzzzzz	-0.000670	-0.000670	-0.004300
çay	-0.000268	-0.000402	49.993550

33433 rows x 3 columns

```
1 merge[merge[2]>30].shape
```

```
(1319, 3)
```

1319 features out of 33433 shows percentage change > 30 post pertubation test i.e 3.94%

We can say that our data isn't much affected by multicollinearity

```
1 #Features
2 feature_names = np.array(vocabulary.get_feature_names())
3 sorted_coef_index = svm_optimal.coef_[0].argsort()
```

```
1 #Top 20 positive features
2 p=feature_names[sorted_coef_index[:20]]
3
4 sp = ""
5 for i in p:
6     sp += str(i)+", "
7 print(sp)
```

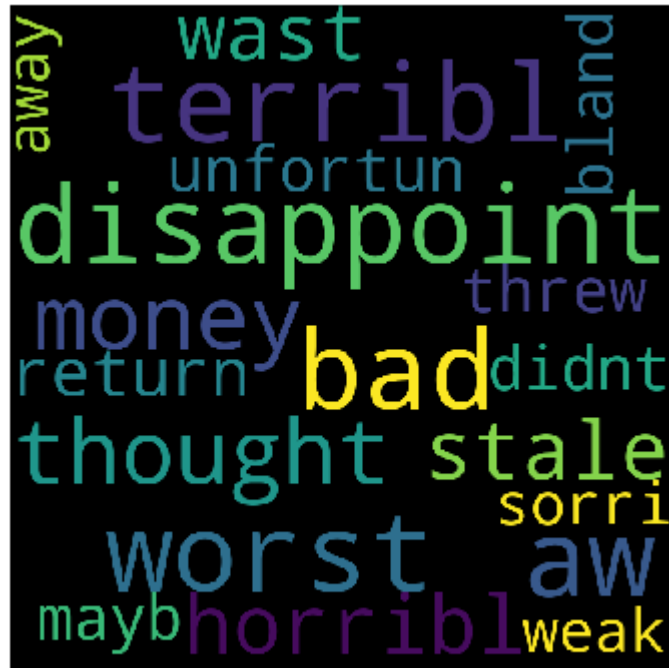
great,love,best,delici,perfect,favorit,excel,good,find,easi,nice,wonder,snack,make,thank,alway,tasti,keep,year,high,

```
1  #Top 20 negative features
2  n=feature_names[sorted_coef_index[:-21:-1]]
3
4  sn = ""
5  for i in n:
6      sn += str(i)+", "
7  print(sn)
```

disappoint,worst,terribl,bad,aw,thought,money,horribl,stale,would,wast,return,unfortun,threw,bland,didnt,mayb,sorri,weak,aw
ay,

```
1 print("***** Top 20 Negative words *****")
2 wordcloud = WordCloud(width = 800, height = 800,
3                       background_color = 'black',
4                       min_font_size = 10).generate(sn)
5
6 # plot the WordCloud image
7 plt.figure(figsize = (5,5), facecolor = None)
8 plt.imshow(wordcloud)
9 plt.axis("off")
10 plt.tight_layout(pad = 0)
11 plt.show()
12
13
14 print("***** Top 20 Positive words *****")
15 wordcloud = WordCloud(width = 800, height = 800,
16                       background_color = 'black',
17                       min_font_size = 10).generate(sp)
18
19 # plot the WordCloud image
20 plt.figure(figsize = (5,5), facecolor = None)
21 plt.imshow(wordcloud)
22 plt.axis("off")
23 plt.tight_layout(pad = 0)
24 plt.show()
```

***** Top 20 Negative words *****



***** Top 20 Positive words *****

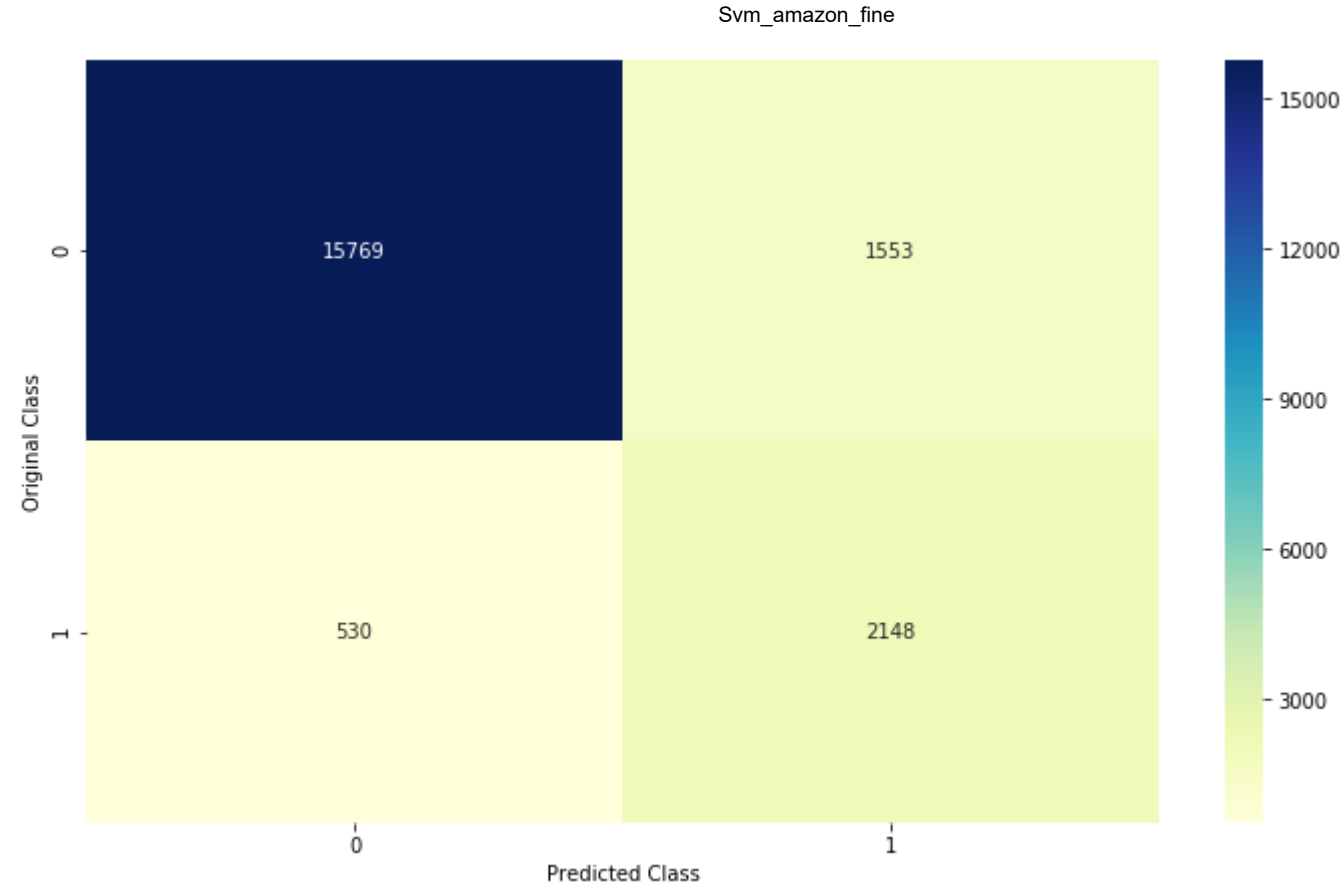


```
1 #Confusion matrix
2 C = confusion_matrix(y_test, pred_bow)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

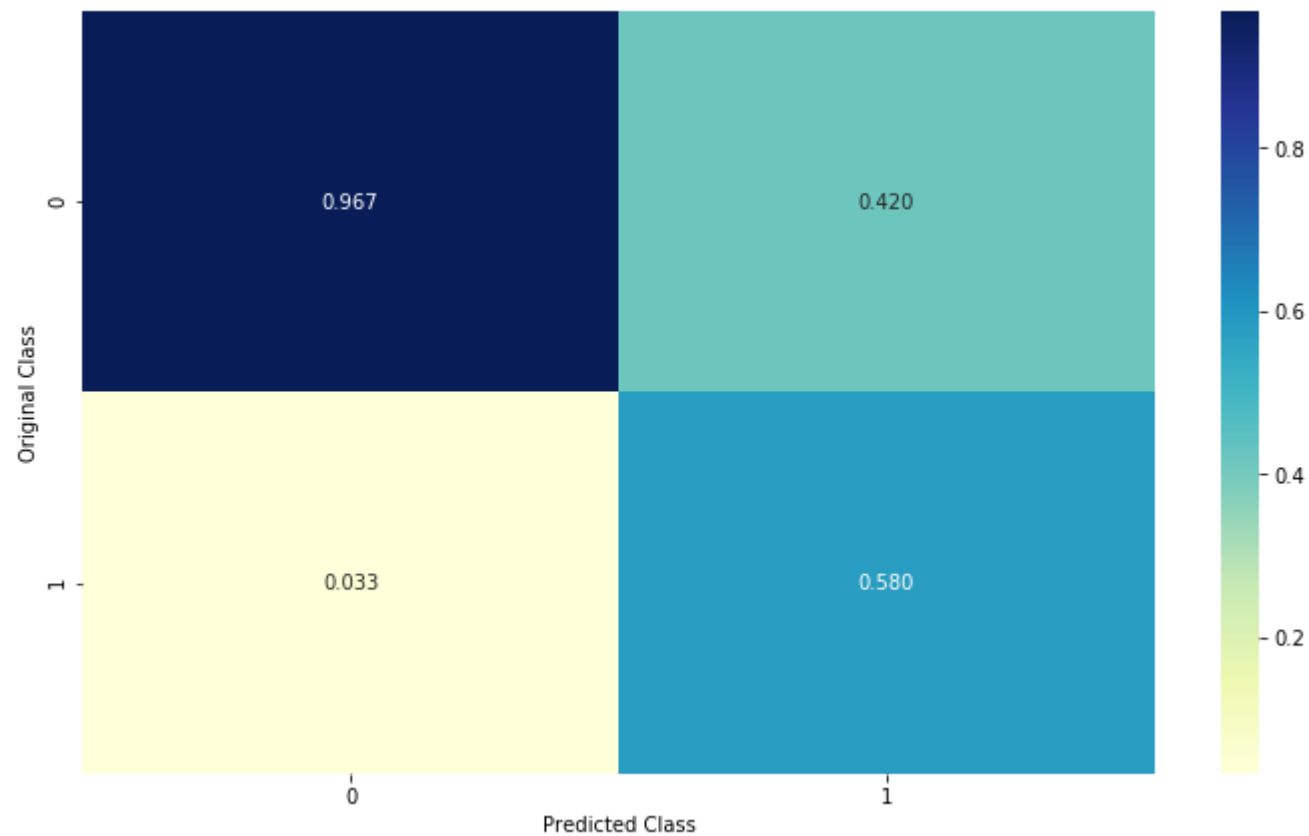


```
1  print("-"*20, "Confusion matrix", "-"*20)
2  plt.figure(figsize=(12,7))
3  sns.heatmap(C, annot=True, cmap="YlGnBu", fmt="g", xticklabels=labels, yticklabels=labels)
4  plt.xlabel('Predicted Class')
5  plt.ylabel('Original Class')
6  plt.show()
7
8  print("-"*20, "Precision matrix (Column Sum=1)", "-"*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()
```

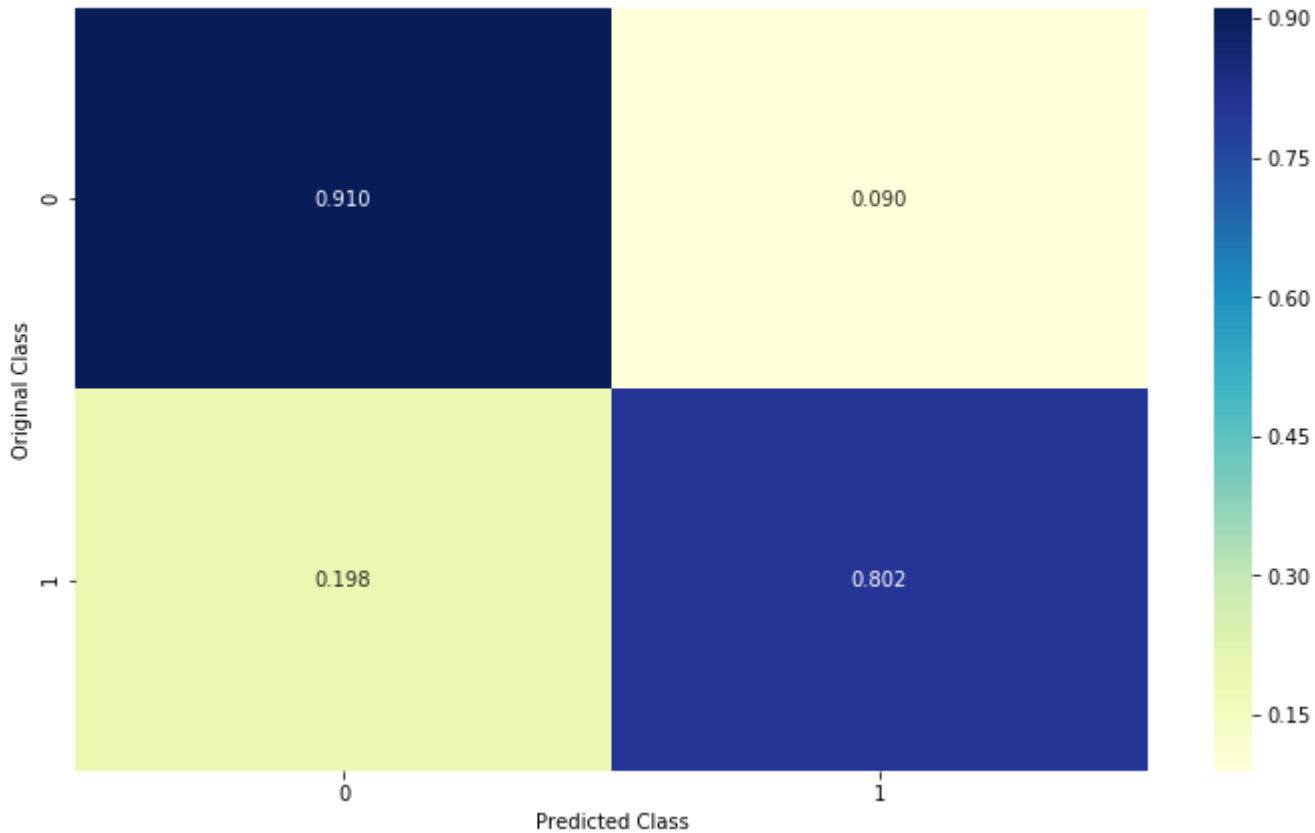
----- Confusion matrix -----



----- Precision matrix (Column 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 #Standardizing
2 from sklearn.preprocessing import StandardScaler
3 Standard=StandardScaler(with_mean=False)
4 Tfidf_x_train = Standard.fit_transform(Tfidf_x_train)
5 Tfidf_x_test = Standard.transform(Tfidf_x_test)
6
7 print(Tfidf_x_train.shape)
8 print(Tfidf_x_test.shape)
```

(80000, 1013943)
(20000, 1013943)

Fitting Grid Search on Tf-Idf

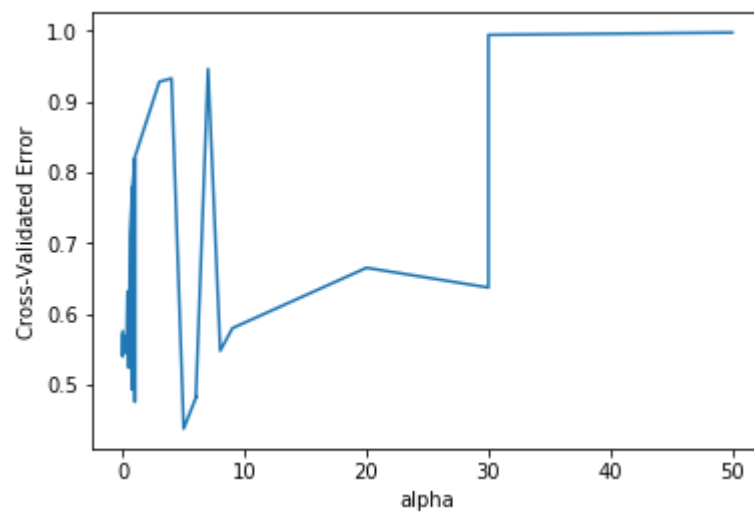
```
1 grid.fit(Tfidf_x_train, y_train)
2
3 # examine the best model
4 print(grid.best_score_)
5 print(grid.best_params_)
```

0.5625001303946727

{'epsilon': 0.2, 'class_weight': 'balanced', 'alpha': 5}

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

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



```
1
2  #{'epsilon': 0.2, 'class_weight': 'balanced', 'alpha': 5}
3  svm_optimal=SGDClassifier(alpha=5,epsilon=0.2,class_weight='balanced')
4
5  # fitting the model
6  svm_optimal.fit(Tfidf_x_train, y_train)
7
8  # predict the response
9  pred_tfidf = svm_optimal.predict(Tfidf_x_test)
10
11 # evaluate accuracy
12 f1_score = f1_score(y_test, pred_tfidf)
13
14 # Train & Test Error
15 print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,svm_optimal.predict(Tfidf_x_train)))
16 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 : 0.9989060790748555

The overall f1_score for the Test Data is : 0.5788245279015488

Pertubation test

```
1  # Re-training the model after adding noise
2  Epsilon = np.random.normal(loc=0,scale =0.01)
3  Noise_Tfidf_x_train=Tfidf_x_train
4  Noise_Tfidf_x_train.data+=Epsilon
```



```
1
2  #{'epsilon': 0.2, 'class_weight': 'balanced', 'alpha': 5}
3  svm_optimal_noise=SGDClassifier(alpha=5,epsilon=0.2,class_weight='balanced')
4
5  # fitting the model
6  svm_optimal_noise.fit(Noise_Tfidf_x_train, y_train)
7
8  # predict the response
9  pred_tfidf = svm_optimal_noise.predict(Tfidf_x_test)
10
11 # evaluate accuracy
12 f1_score = f1_score(y_test, pred_tfidf)
13
14 # Train & Test Error
15 print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,svm_optimal.predict(Noi:
16 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 : 0.9989060790748555

The overall f1_score for the Test Data is : 0.5773239736226335

```
1 #Features
2 feature_names = np.array(vocabulary.get_feature_names())
3 feature_names.shape
```

(1013943,)

```
1 svm_optimal.coef_.shape
```

(1, 1013943)

```
1 svm_optimal_noise.coef_.shape
```

(1, 1013943)

```

1 merge_arr = np.concatenate([svm_optimal.coef_, svm_optimal_noise.coef_], axis=0)
2 merge=pd.DataFrame(data=merge_arr,columns=feature_names).transpose()
3 merge[2]=((merge[1]-merge[0])/merge[0])*100
4 merge

```

	0	1	2
aaa	-0.000105	-0.000179	71.285216
aaa condit	-0.000080	-0.000161	99.996881
aaa perfect	-0.000080	-0.000080	-0.001559
aaaaaaaaagghh	-0.000161	-0.000161	-0.001559
aaaaah	-0.000113	-0.000113	-0.002219
aaaaah awak	-0.000080	-0.000080	-0.001559
aaaaah satisfi	-0.000080	-0.000080	-0.001559
aaaaahhhhhhhhhhhhhhhhh	-0.000080	-0.000080	-0.001559
aaaaahhhhhhhhhhhhhhhhh angel	-0.000080	-0.000080	-0.001559
aaaah	-0.000080	-0.000080	-0.001559
aaaah snob	-0.000080	-0.000080	-0.001559
aaah	-0.000134	-0.000134	-0.002801
aaah inhal	-0.000080	-0.000080	-0.001559
aaah miss	-0.000080	-0.000080	-0.001559
aaah sip	-0.000080	-0.000080	-0.001559
aachen	0.000590	0.000590	-0.001559
aachen munich	0.000590	0.000590	-0.001559
aad	0.000000	0.000000	NaN
aad sausag	0.000000	0.000000	NaN
aadp	-0.000080	-0.000080	-0.001559
aafco	0.000142	0.000184	29.714666
aafco also	-0.000080	-0.000080	-0.001559
aafco certifi	0.000590	0.000590	-0.001559
aafco countri	-0.000080	-0.000080	-0.001559

	0	1	2
aafco definit	0.000590	0.000590	-0.001559
aafco dog	-0.000080	0.000000	-100.000000
aafco guidelin	0.000000	0.000000	NaN
aafco requir	-0.000080	-0.000080	-0.001559
aagh	-0.000080	-0.000080	-0.001559
aagh yelp	-0.000080	-0.000080	-0.001559
...
zum heal	-0.000080	-0.000080	-0.001559
zummi	-0.000080	-0.000080	-0.001559
zummi love	-0.000080	-0.000080	-0.001559
zummi tast	-0.000080	-0.000080	-0.001559
zummi tri	-0.000080	-0.000080	-0.001559
zune	-0.000241	-0.000241	-0.001559
zune video	-0.000241	-0.000241	-0.001559
zupreem	-0.000080	-0.000080	-0.001559
zupreem ferret	-0.000080	-0.000080	-0.001559
zurich	-0.000080	-0.000080	-0.001559
zurich schnatzlet	-0.000080	-0.000080	-0.001559
zwar	-0.000080	-0.000080	-0.001559
zwar billig	-0.000080	-0.000080	-0.001559
zwieback	0.000288	0.000288	-0.002321
zwieback toast	0.000288	0.000288	-0.002321
zwiebeck	-0.000161	-0.000161	-0.001559
zwiebeck toast	-0.000161	-0.000161	-0.001559
zydeco	-0.000161	-0.000161	-0.001559
zydeco saturday	-0.000161	-0.000161	-0.001559
zzzzzs	-0.000170	-0.000170	-0.002208
zzzzzs larg	-0.000161	-0.000161	-0.001559
zzzzzz	0.000000	0.000000	NaN

	0	1	2
zzzzzz say	0.000000	0.000000	NaN
zzzzzzzz	0.000590	0.000590	-0.001559
zzzzzzzz high	0.000590	0.000590	-0.001559
zzzzzzzzzz	-0.000080	-0.000080	-0.001559
zzzzzzzzzzzz	-0.000080	-0.000080	-0.001559
zzzzzzzzzzzz final	-0.000080	-0.000080	-0.001559
zzzzzzzzzzzzzz	-0.000080	-0.000080	-0.001559
çay	-0.000080	-0.000080	-0.001559

1013943 rows x 3 columns

```
1 merge[merge[2]>30].shape
```

(30883, 3)

30883 features out of 1013943 shows percentage change > 30 post pertubation test i.e 3.04%

We can say that our data isn't much affected by multicollinearity

```
1 sorted_coef_index = svm_optimal.coef_[0].argsort()
```

```
1 #Top 20 positive features
2 p=feature_names[sorted_coef_index[:20]]
3
4 sp = ""
5 for i in p:
6     sp += str(i)+", "
7 print(sp)
```

great,love,good,best,delici,excel,favorit,find,use,make,perfect,wonder,tasti,nice,flavor,price,easi,enjoy,high recommend,thank,

```
1 #Top 20 negative features
2 n=feature_names[sorted_coef_index[:-21:-1]]
3
4 sn = ""
5 for i in n:
6     sn += str(i)+", "
7 print(sn)
```

disappoint,worst,wast money,horribl,terribl,aw,wont buy,wast,threw,return,stale,refund,bland,money,wors,two star,disappoint
product,bad,never buy,disgust,

```
1 print("***** Top 20 Negative words *****")
2 wordcloud = WordCloud(width = 800, height = 800,
3                       background_color = 'black',
4                       min_font_size = 10).generate(sn)
5
6 # plot the WordCloud image
7 plt.figure(figsize = (5,5), facecolor = None)
8 plt.imshow(wordcloud)
9 plt.axis("off")
10 plt.tight_layout(pad = 0)
11 plt.show()
12
13
14 print("***** Top 20 Positive words *****")
15 wordcloud = WordCloud(width = 800, height = 800,
16                       background_color = 'black',
17                       min_font_size = 10).generate(sp)
18
19 # plot the WordCloud image
20 plt.figure(figsize = (5,5), facecolor = None)
21 plt.imshow(wordcloud)
22 plt.axis("off")
23 plt.tight_layout(pad = 0)
24 plt.show()
```

***** Top 20 Negative words *****



***** Top 20 Positive words *****

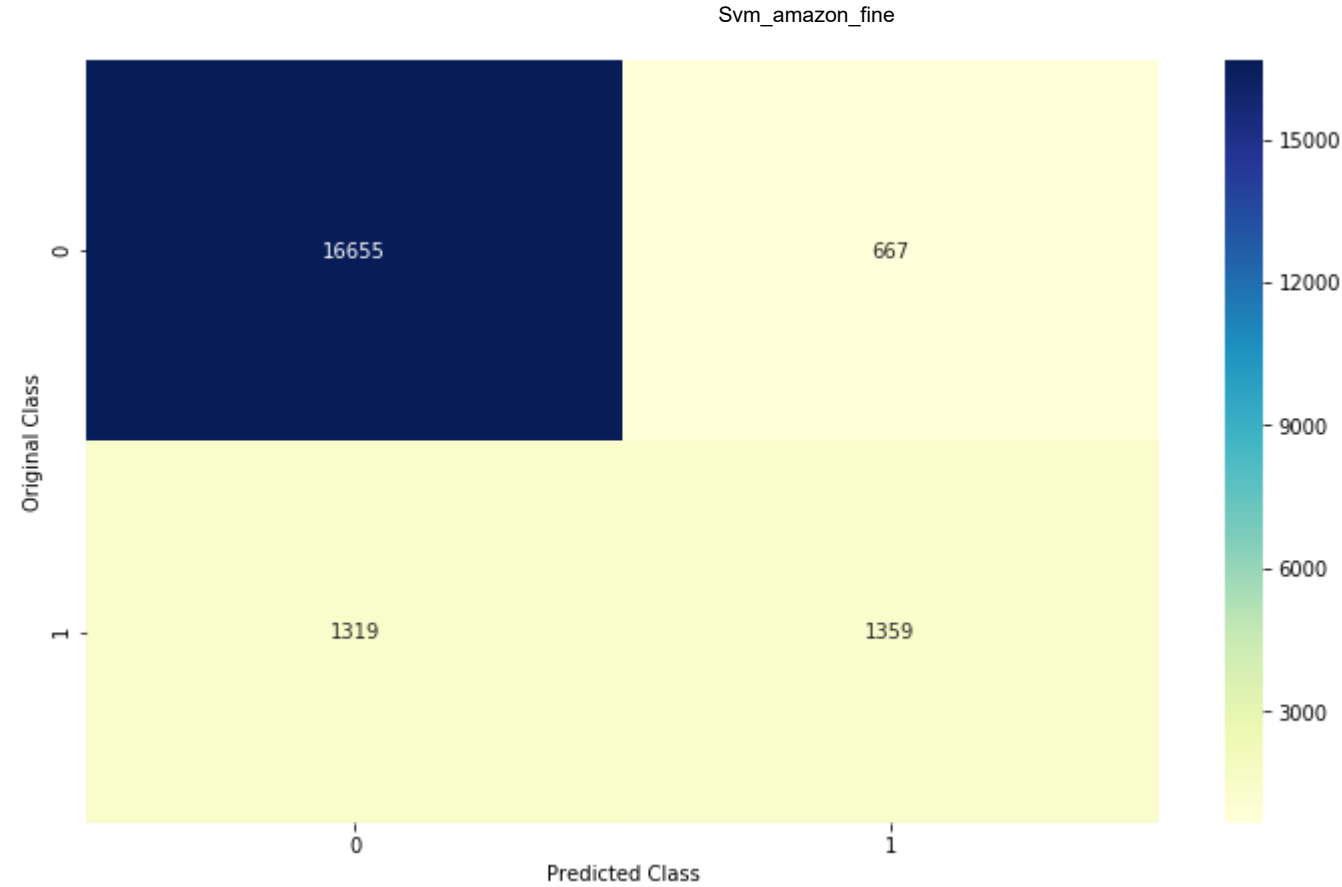


```
1 #Confusion matrix
2 C = confusion_matrix(y_test, pred_tfidf)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

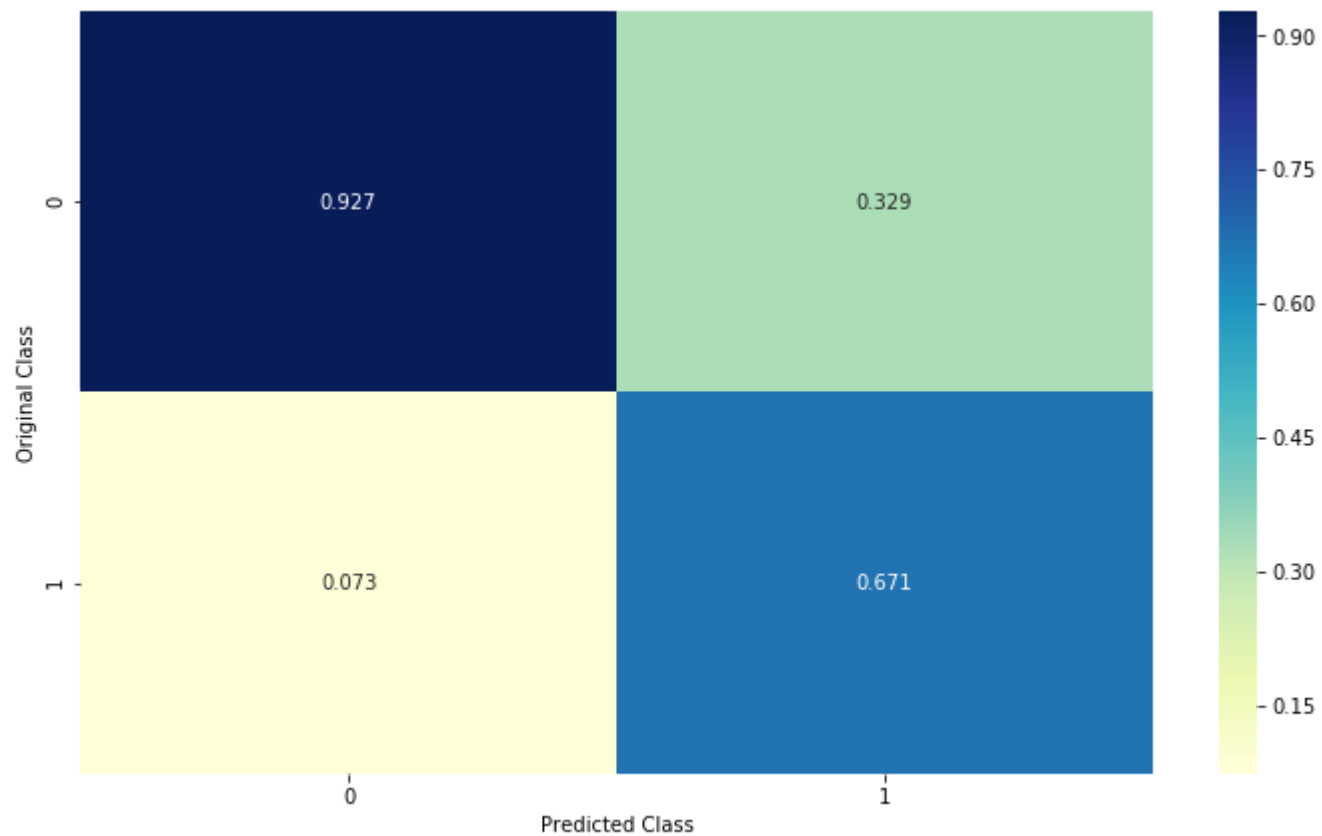


```
1  print("-"*20, "Confusion matrix", "-"*20)
2  plt.figure(figsize=(12,7))
3  sns.heatmap(C, annot=True, cmap="YlGnBu", fmt="g", xticklabels=labels, yticklabels=labels)
4  plt.xlabel('Predicted Class')
5  plt.ylabel('Original Class')
6  plt.show()
7
8  print("-"*20, "Precision matrix (Column Sum=1)", "-"*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()
```

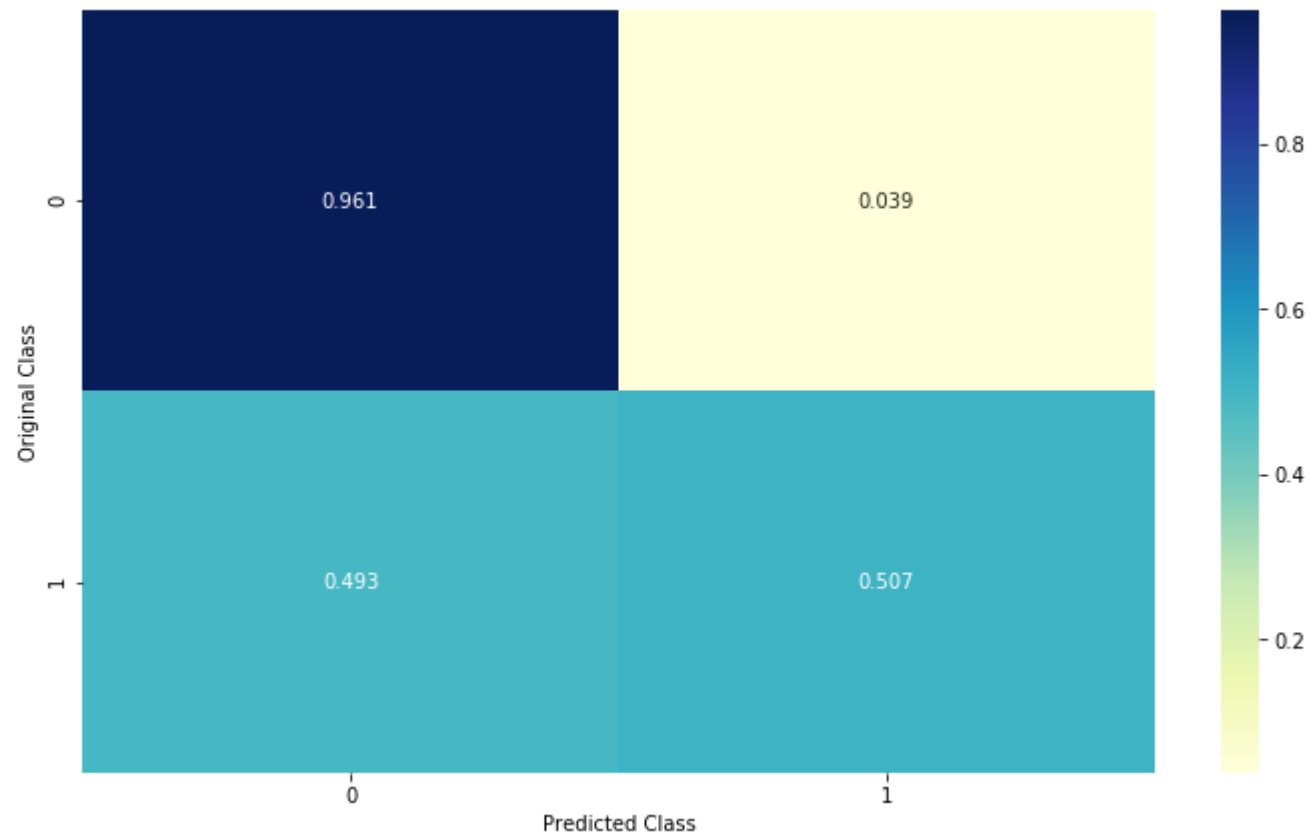
----- Confusion matrix -----



----- Precision matrix (Column 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 occurred atleast 5 times
3 w2v_model_train=Word2Vec(list_of_sent_train,min_count=5,size=50, workers=6)
```

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

number of words that occurred minimum 5 times 11361

sample words ['witti', 'littl', 'book', 'make', 'son', 'laugh', 'loud', 'car', 'drive', 'along', 'alway', 'sing', 'refrain', 'hes', 'learn', 'whale', 'india', 'droop', 'love', 'new', 'word', 'introduc', 'silli', 'classic', 'will', 'bet', 'still', 'abl', 'memori', 'colleg', 'rememb', 'see', 'show', 'air', 'televis', 'year', 'ago', 'child', 'sister', 'later', 'bought', '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(50) # 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

50

```

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(50) # 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

50

```

1  #Standardizing Avg-W2v
2  from sklearn.preprocessing import StandardScaler
3  Standard=StandardScaler()
4  sent_vectors_train_avgw2v = Standard.fit_transform(sent_vectors_train_avgw2v)
5  sent_vectors_test_avgw2v = Standard.transform(sent_vectors_test_avgw2v)
6
7  print(sent_vectors_train_avgw2v.shape)
8  print(sent_vectors_test_avgw2v.shape)

```

(80000, 50)

(20000, 50)

Fitting grid search on Avg-W2V

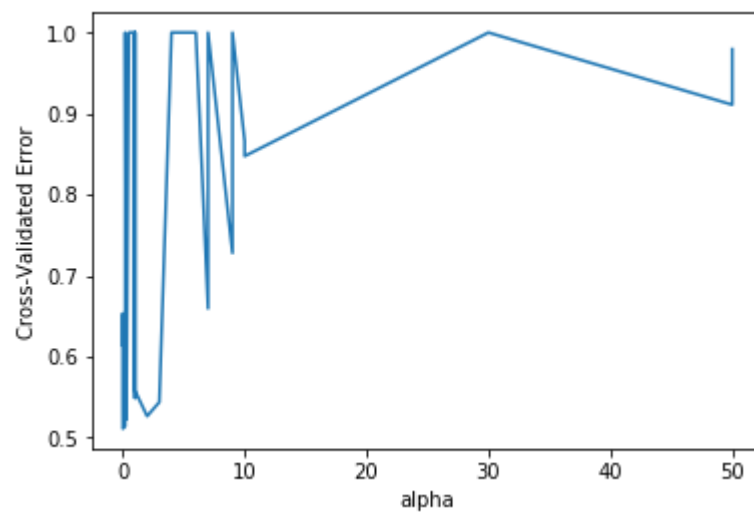
```
1 grid.fit(sent_vectors_train_avgw2v, y_train)
2
3 # examine the best model
4 print(grid.best_score_)
5 print(grid.best_params_)
```

0.48896063270159695

{'epsilon': 0.7, 'class_weight': 'balanced', 'alpha': 0.01}


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

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



```
1  #{'epsilon': 0.7, 'class_weight': 'balanced', 'alpha': 0.01}
2  svm_optimal=SGDClassifier(alpha=0.01,epsilon=0.7,class_weight='balanced')
3
4  # fitting the model
5  svm_optimal.fit(sent_vectors_train_avgw2v, y_train)
6
7  # predict the response
8  pred_avg_w2v = svm_optimal.predict(sent_vectors_test_avgw2v)
9
10 # evaluate f1_score
11 f1_score = f1_score(y_test, pred_avg_w2v)
12
13 # Train & Test Error
14 print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,svm_optimal.predict(sent_vectors_train_avgw2v)))
15 print("The overall f1_score for the Test Data is : ", metrics.f1_score(y_test,pred_avg_w2v))
```

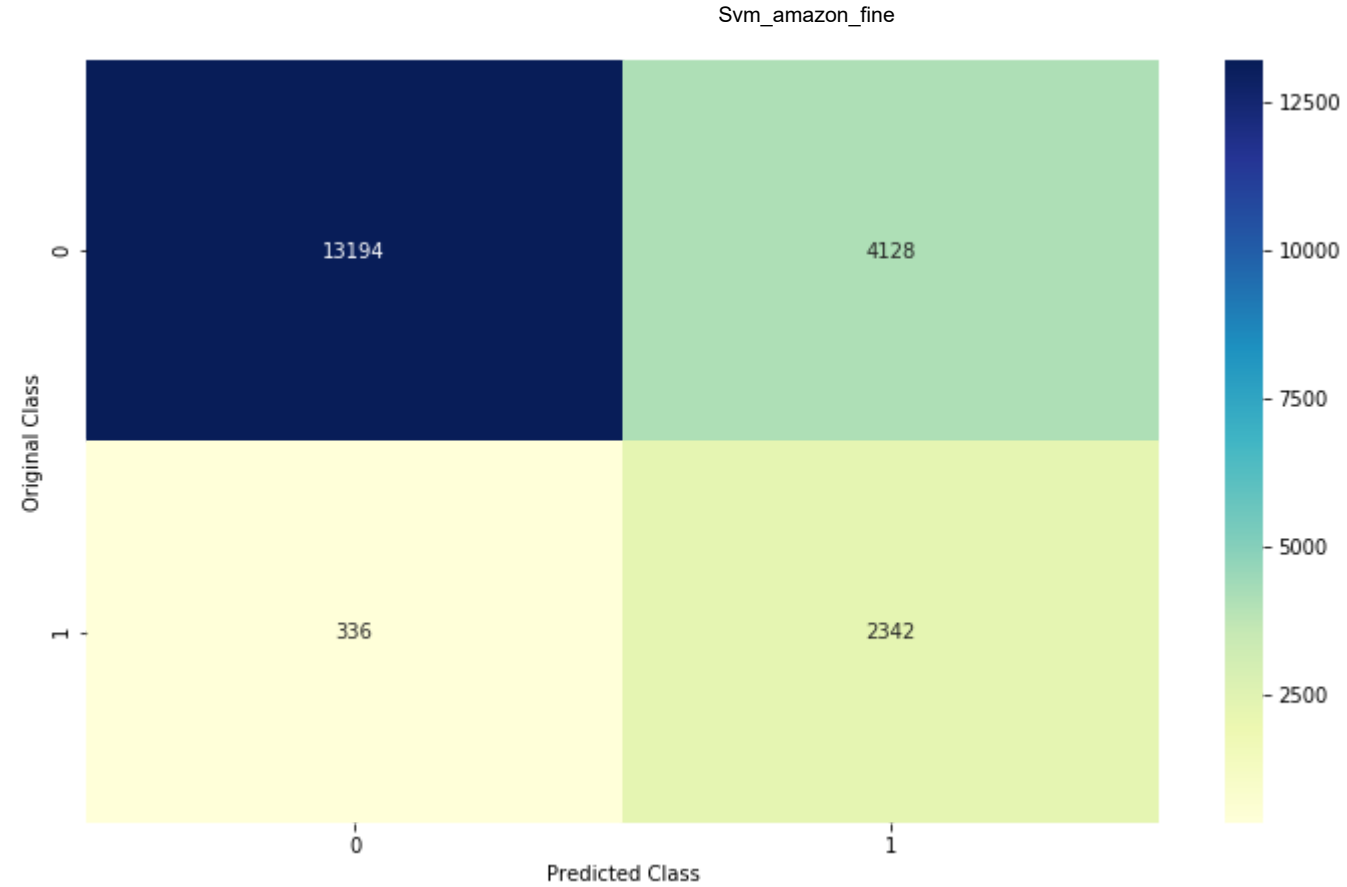
The overall f1_score for the Train Data is : 0.4781548524054046

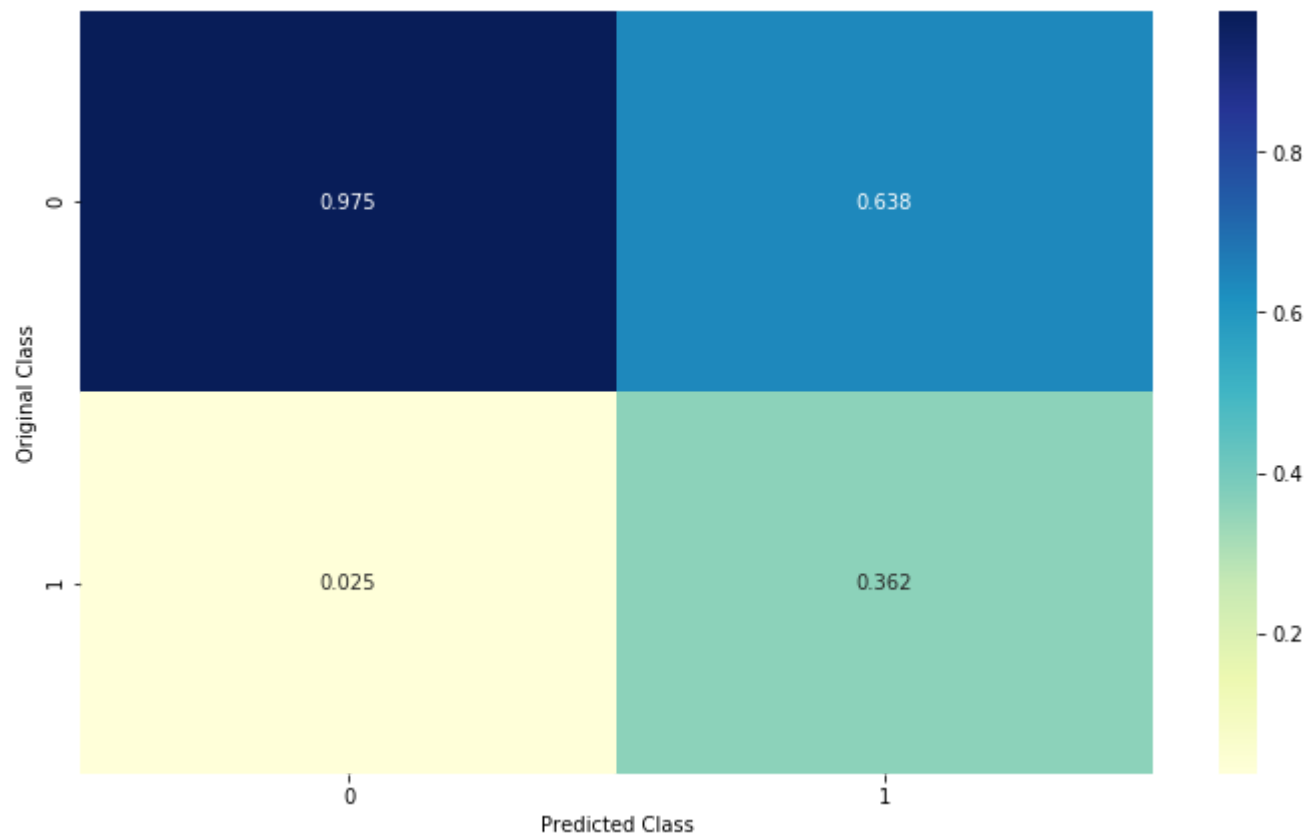
The overall f1_score for the Test Data is : 0.5021829411138323

```
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]
```

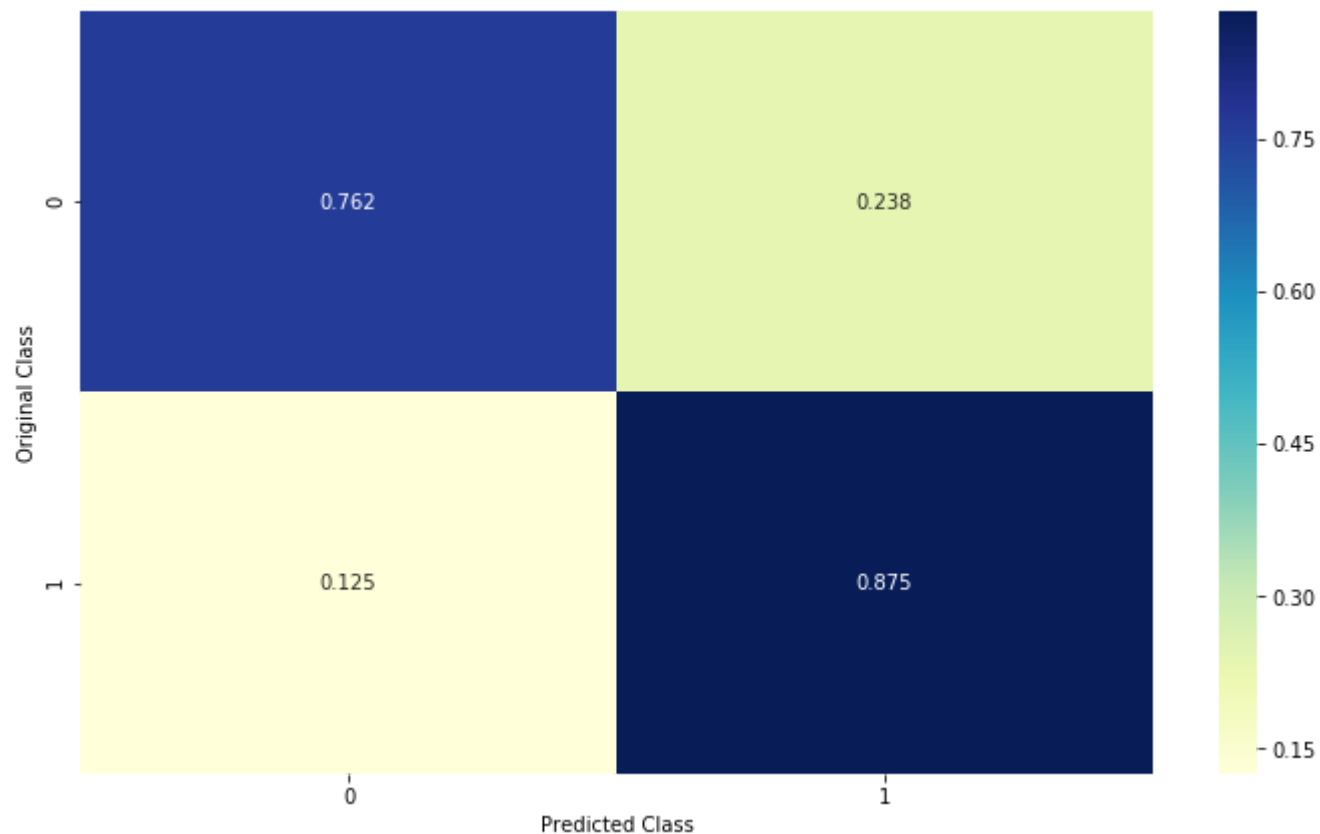
```
1  print("-"*20, "Confusion matrix", "-"*20)
2  plt.figure(figsize=(12,7))
3  sns.heatmap(C, annot=True, cmap="YlGnBu", fmt="g", xticklabels=labels, yticklabels=labels)
4  plt.xlabel('Predicted Class')
5  plt.ylabel('Original Class')
6  plt.show()
7
8  print("-"*20, "Precision matrix (Column Sum=1)", "-"*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 -----





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



TF-Idf W2Vec

```
1 tf_idf_vect = TfidfVectorizer()
2 vocabulary = tf_idf_vect.fit(train_data['CleanedText'])
3 final_tf_idf= tf_idf_vect.transform(train_data['CleanedText'])
4
5 # we are converting a dictionary with word as a key, and the idf as a value
6 dictionary = dict(zip(vocabulary.get_feature_names(), list(tf_idf_vect.idf_)))
```

```

1  # TF-IDF weighted Word2Vec
2  tfidf_feat = tf_idf_vect.get_feature_names() # tfidf words/col-names
3  # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf
4
5  tfidf_w2v_sent_vectors_train = []; # the tfidf-w2v for each sentence/review is stored in this list
6  row=0;
7  for sent in tqdm(list_of_sent_train): # for each review/sentence
8      sent_vec = np.zeros(50) # as word vectors are of zero length
9      weight_sum =0; # num of words with a valid vector in the sentence/review
10     for word in sent: # for each word in a review/sentence
11         if word in w2v_words_train:
12             vec = w2v_model_train.wv[word]
13             # tfidf = tf_idf_matrix[row, tfidf_feat.index(word)]
14             # to reduce the computation we are
15             # dictionary[word] = idf value of word in whole corpus
16             # sent.count(word) = tf value of word in this review
17             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
18             sent_vec += (vec * tf_idf)
19             weight_sum += tf_idf
20     if weight_sum != 0:
21         sent_vec /= weight_sum
22     tfidf_w2v_sent_vectors_train.append(sent_vec)
23     row += 1

```



```
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_)
```

0.46358605223865923

{'epsilon': 0.1, 'class_weight': 'balanced', 'alpha': 0.01}

```
1 from sklearn.metrics import f1_score
2 {'epsilon': 0.1, 'class_weight': 'balanced', 'alpha': 0.01}
3 svm_optimal=SGDClassifier(alpha=0.01,epsilon=0.1,class_weight='balanced')
4
5 # fitting the model
6 svm_optimal.fit(tfidf_w2v_sent_vectors_train, y_train)
7
8 # predict the response
9 pred_tfidf_w2v_sent_vectors_test = svm_optimal.predict(tfidf_w2v_sent_vectors_test)
10
11 # evaluate f1_score
12 f1_score = f1_score(y_test, pred_tfidf_w2v_sent_vectors_test)
13
14 # Train & Test Error
15 print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,svm_optimal.predict(tfidf_w2v_sent_vectors_train)))
16 print("The overall f1_score for the Test Data is : ", metrics.f1_score(y_test,pred_tfidf_w2v_sent_vectors_test))
```

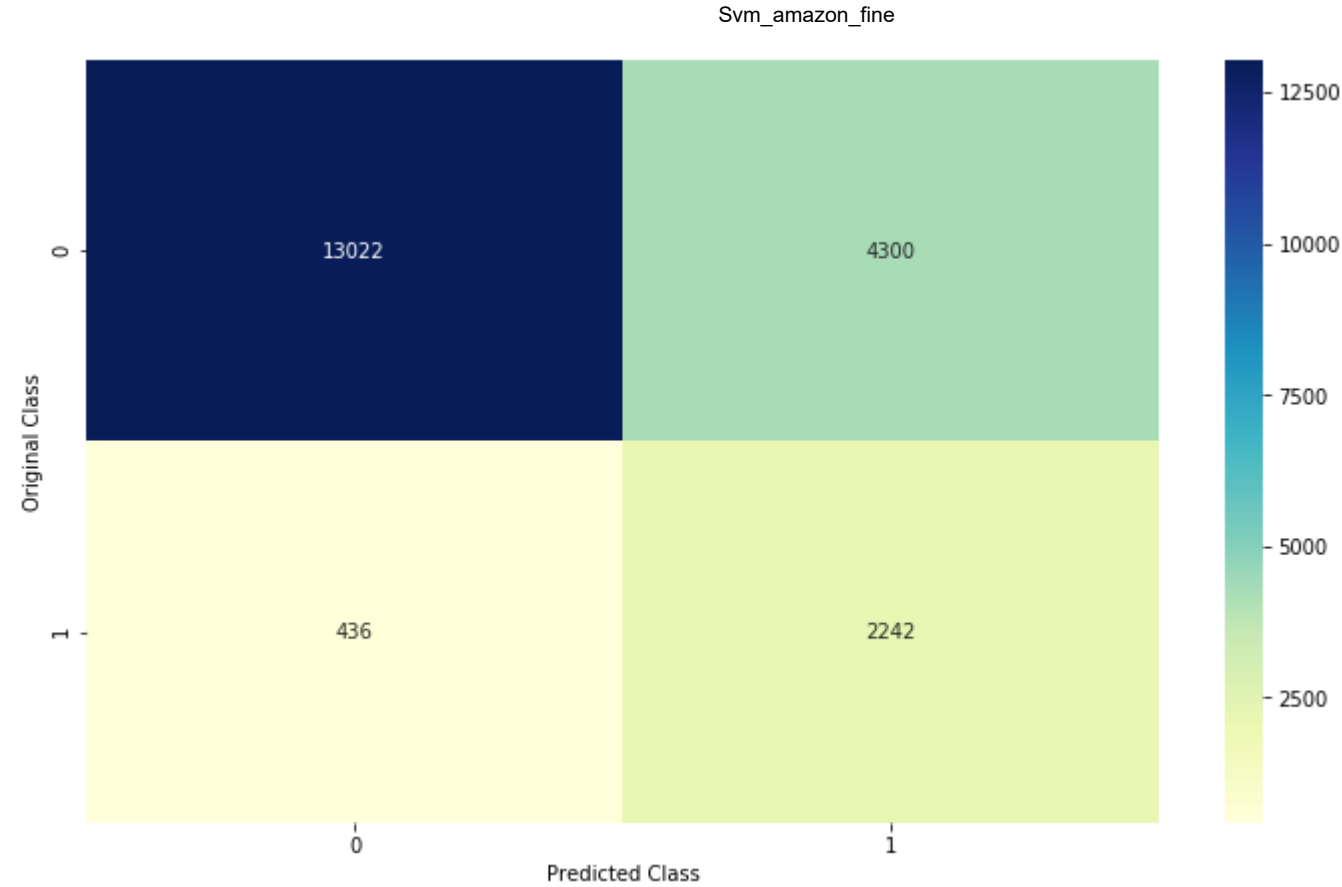
The overall f1_score for the Train Data is : 0.458559030811346

The overall f1_score for the Test Data is : 0.4863340563991323

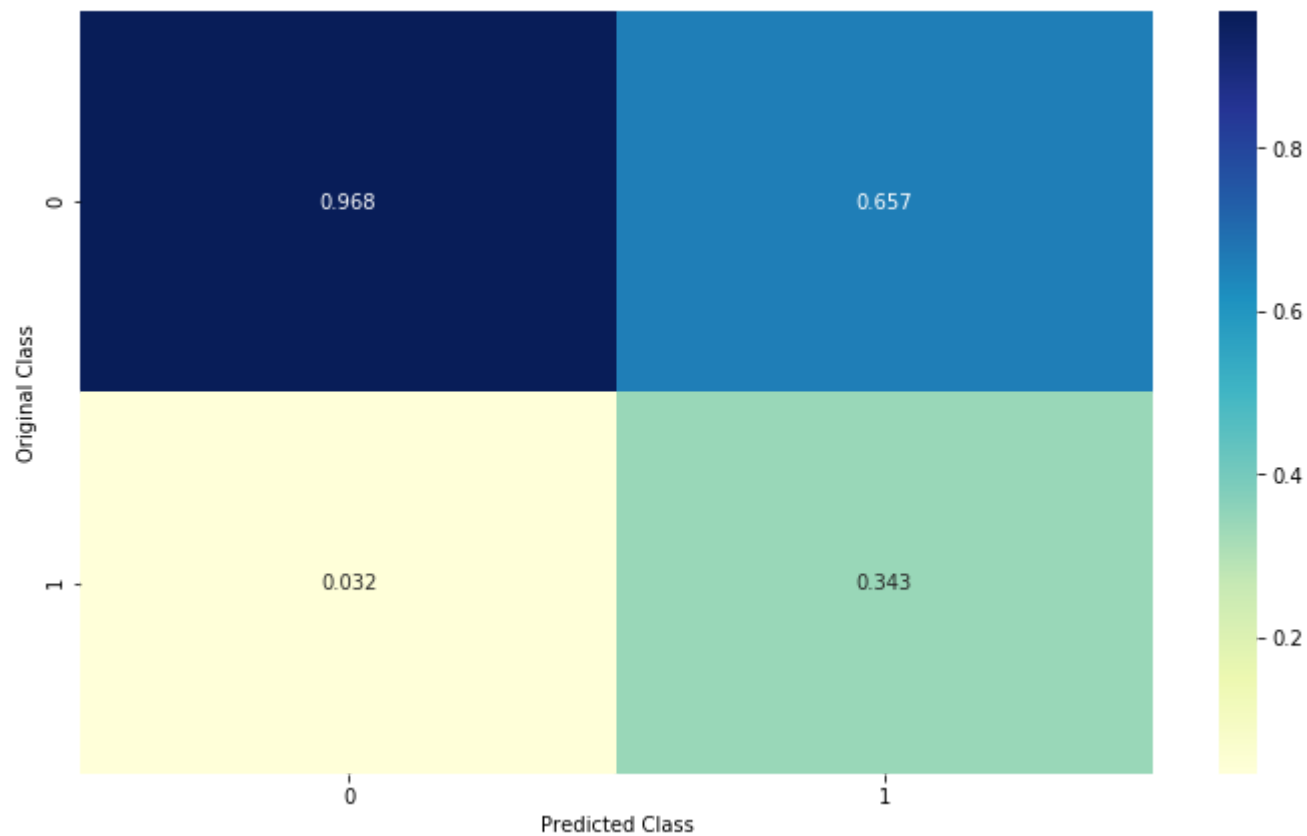
```
1 #Confusion matrix
2 C = confusion_matrix(y_test, pred_tfidf_w2v_sent_vectors_test)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

```
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 (Column 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()
```

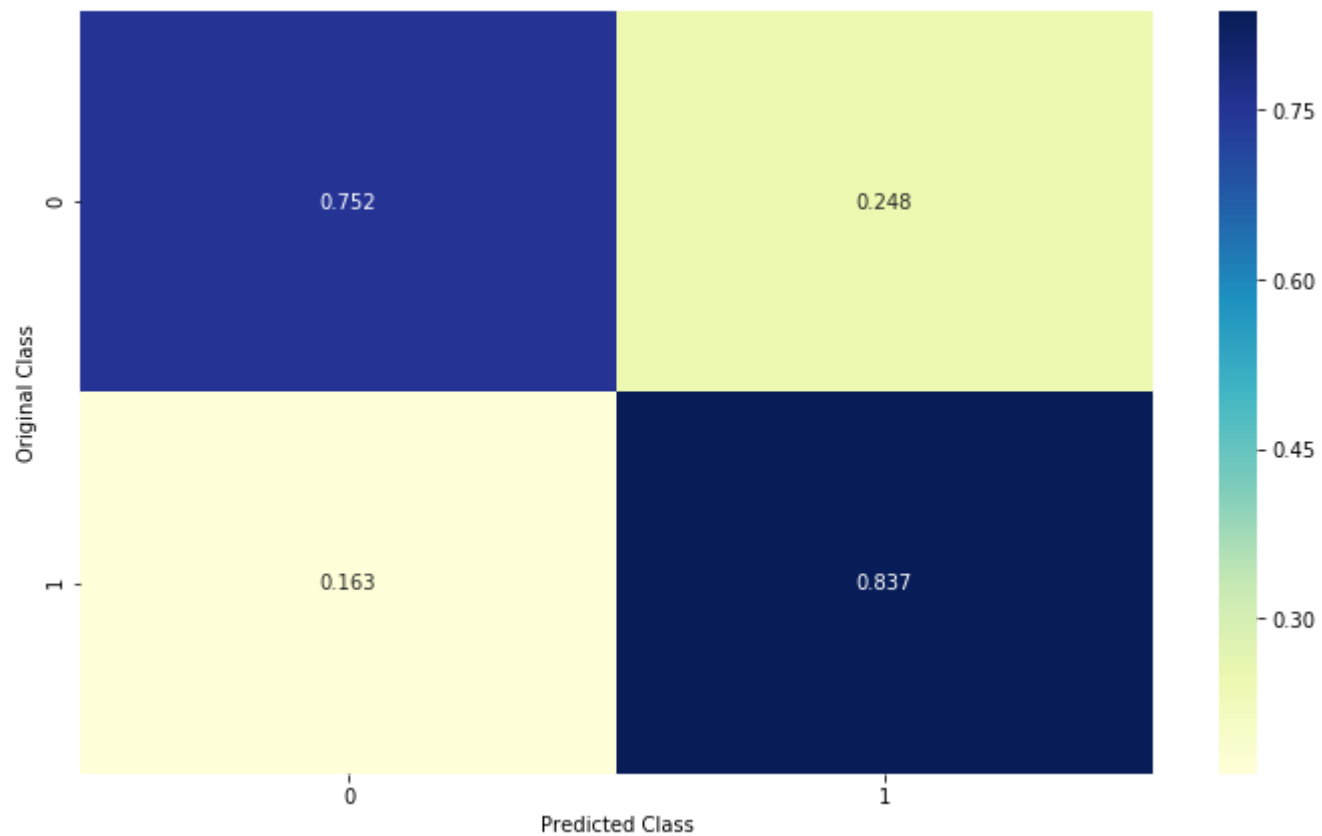
----- Confusion matrix -----



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



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



Reporting f1_score after running SGD classifier on above featurization

```
1 from prettytable import PrettyTable
```

```

1 x=PrettyTable()
2 x.field_names = ["Model", "Bow", "Tfidf", "Avg-W2V", "Tfidf-W2V"]
3 x.add_row(["Alpha", 3, 5, 0.01, 0.01])
4 x.add_row(["Train f1_score", 0.748, 0.998, 0.478, 0.458])
5 x.add_row(["Test f1_score", 0.693, 0.582, 0.502, 0.486])
6
7 print(x)

```

```

+-----+-----+-----+-----+-----+
|   Model   |  Bow  | Tfidf | Avg-W2V | Tfidf-W2V |
+-----+-----+-----+-----+-----+
|   Alpha   |    3   |    5   |    0.01  |    0.01   |
| Train f1_score | 0.748 | 0.998 |    0.478 |    0.458   |
| Test f1_score  | 0.693 | 0.582 |    0.502 |    0.486   |
+-----+-----+-----+-----+-----+

```

We can see that BOW has the best f1_score therefore we will run SVC classifier on it.

Randomised search cv

```

1 C = [0.000001, 0.00001, 0.0001, 0.001, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15,
2 T= TimeSeriesSplit(n_splits=3)
3 weight=[None, 'balanced']
4
5 param_distributions = dict(C=C, class_weight=weight)
6 print(param_distributions)
7
8 # instantiate and fit the grid
9 grid = RandomizedSearchCV(SVC(), param_distributions, cv=T, scoring='f1', return_train_score=False, n_jo

```

```

{'C': [1e-06, 1e-05, 0.0001, 0.001, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10,
15, 20, 30, 40, 50], 'class_weight': [None, 'balanced']}

```

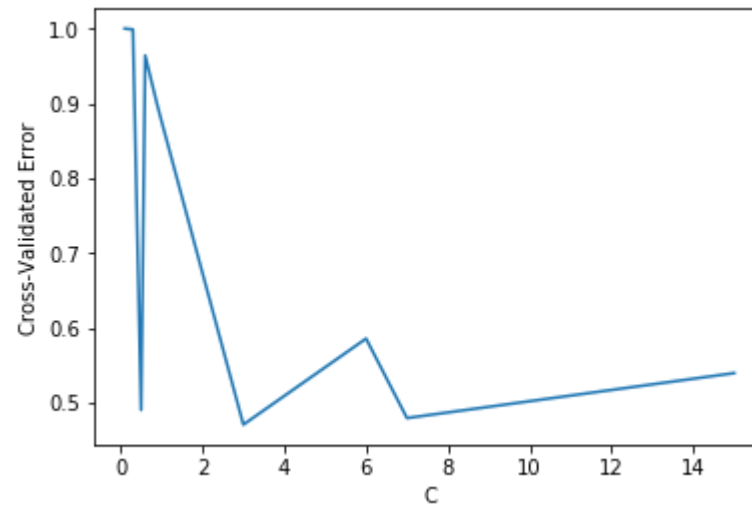
```
1 grid.fit(Bow_x_train, y_train)
2
3 # examine the best model
4 print(grid.best_score_)
5 print(grid.best_params_)
```

0.5289631785421901

{'class_weight': 'balanced', 'C': 3}

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

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




```
1  #{'class_weight': 'balanced', 'C': 3}
2  svm_optimal=SVC(C=3,class_weight='balanced')
3
4  # fitting the model
5  svm_optimal.fit(Bow_x_train, y_train)
6
7  # predict the response
8  pred_bow = svm_optimal.predict(Bow_x_test)
9
10 # evaluate f1_score
11 f1_score = f1_score(y_test, pred_bow)
12
13 # Train & Test Error
14 print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,svm_optimal.predict(Bow_
15 print("The overall f1_score for the Test Data is : ", metrics.f1_score(y_test,pred_bow))
```

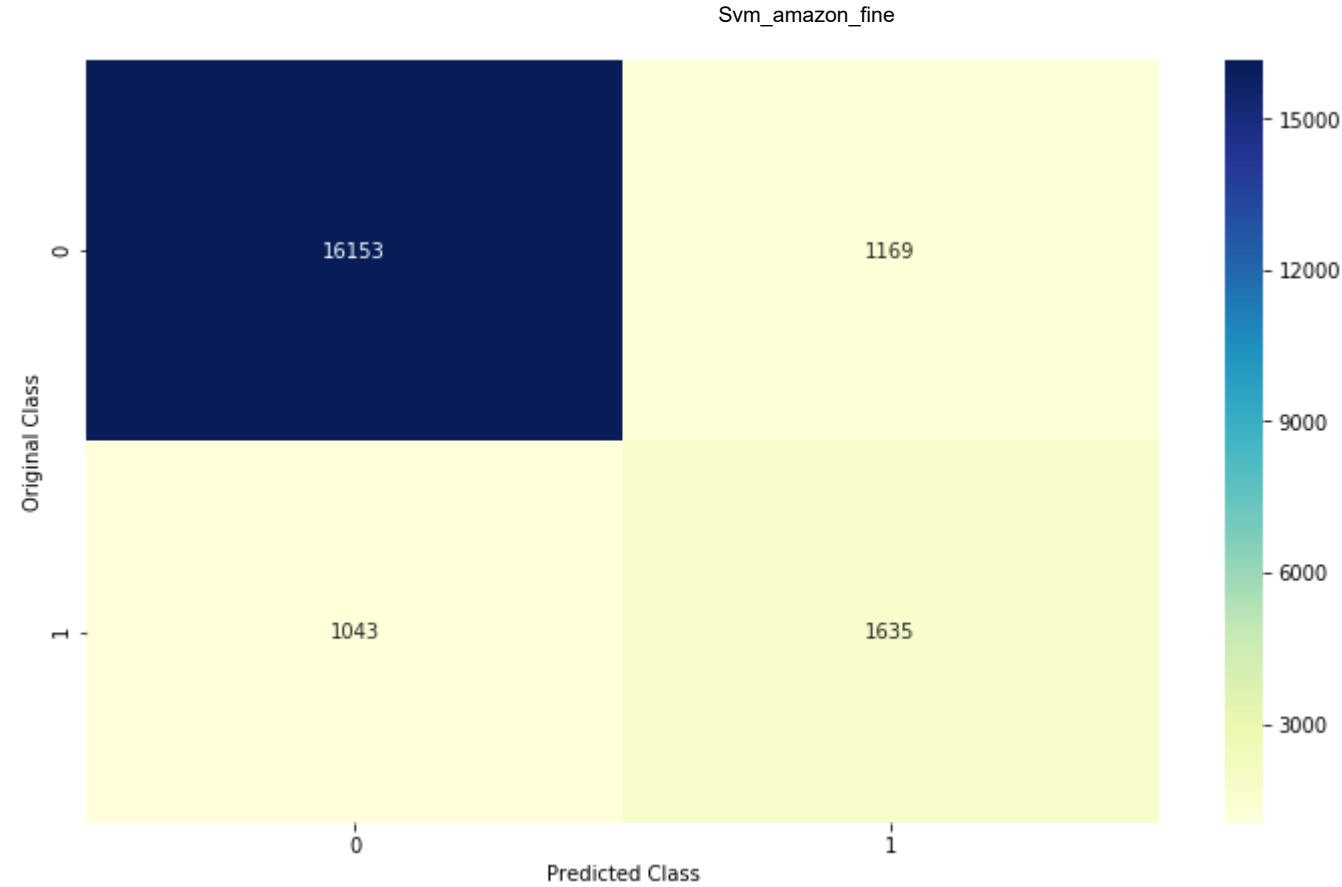
The overall f1_score for the Train Data is : 0.8612301957129542

The overall f1_score for the Test Data is : 0.5964976286026997

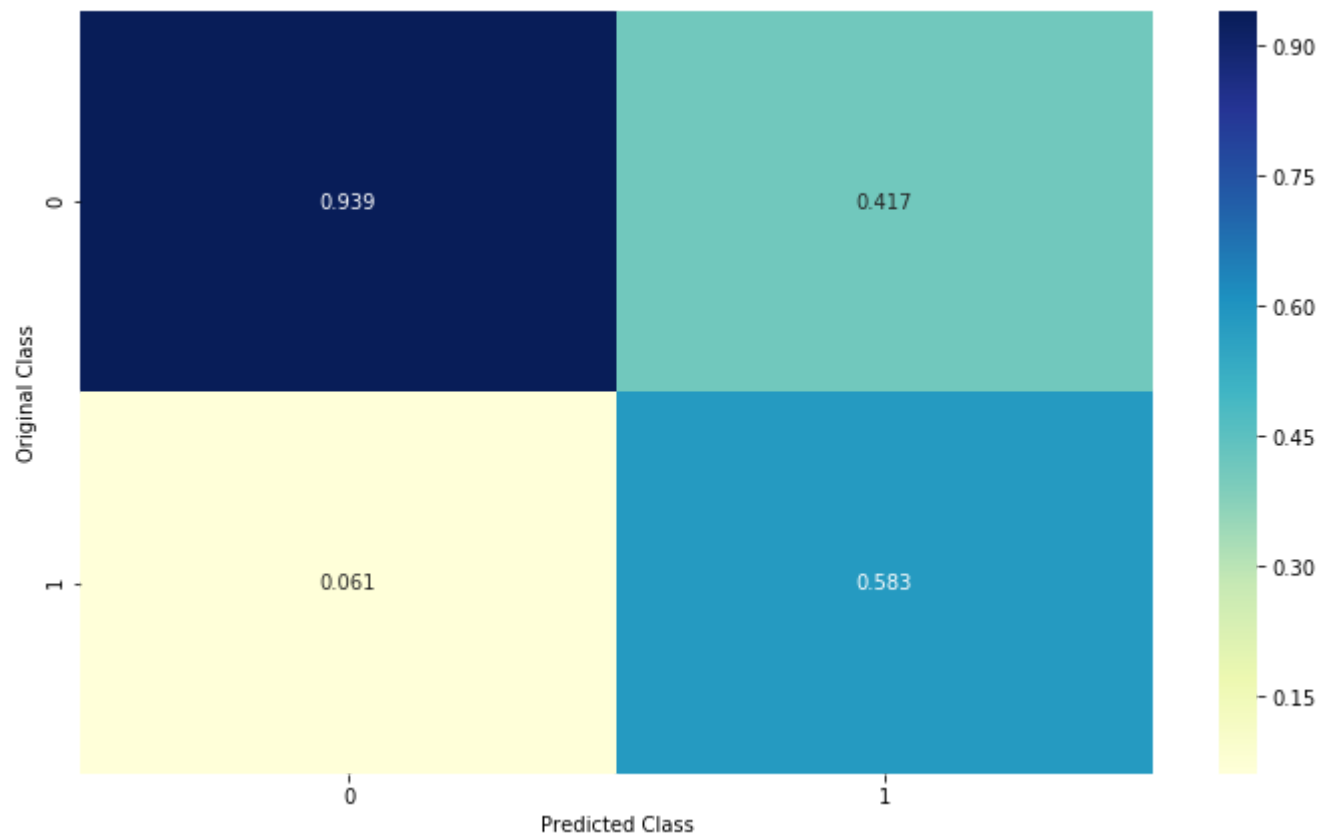
```
1  #Confusion matrix
2  C = confusion_matrix(y_test, pred_bow)
3  A =(((C.T)/(C.sum(axis=1))).T)
4  B =(C/C.sum(axis=0))
5  labels = [0,1]
```

```
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 (Column 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 (Column Sum=1) -----



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

