

NB 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 NB 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
```

```
C:\Users\deepak\Anaconda3\lib\site-packages\gensim\utils.py:1209: UserWarning: detected Windows; aliasing chunkize to chunk
ize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

```
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.naive_bayes import BernoulliNB
2 from sklearn.naive_bayes import MultinomialNB
3 from sklearn.model_selection import TimeSeriesSplit
4 from sklearn.model_selection import GridSearchCV
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
```

Using Bernauli NB

Running Gridsearch CV

```

1 alpha_range = [0.000001,0.00001,0.0001,0.001,0.01,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0,10,20,30,40,50]
2 T= TimeSeriesSplit(n_splits=5)
3
4 param_grid = dict(alpha=alpha_range)
5 print(param_grid)
6
7 # instantiate and fit the grid
8 grid = GridSearchCV(BernoulliNB(), param_grid, cv=T, scoring='accuracy', return_train_score=False,n_jobs=5)

{'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, 10, 20, 30, 40, 50]}

```

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

Fitting Grid 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_)
```

```
0.8992424810620265
{'alpha': 0.1}
```

```
1 # Naive bayes
2 Bnb_optimal=BernoulliNB(alpha=0.1, binarize=0.0, fit_prior=True)
3
4 # fitting the model
5 Bnb_optimal.fit(Bow_x_train, y_train)
6
7 # predict the response
8 pred_bow = Bnb_optimal.predict(Bow_x_test)
9
10 # evaluate accuracy
11 acc = accuracy_score(y_test, pred_bow) * 100
12 print('\nThe accuracy of the Bernaulli NB classifier for a = %d is %f%%' % (0.1, acc))
```

```
The accuracy of the Bernaulli NB classifier for a = 0 is 90.045000%
```

```
1 # Train & Test Error
2 print("The overall accuracy score for the Train Data is : ", metrics.accuracy_score(y_train,Bnb_optimal
3 print("The overall accuracy score for the Test Data is : ", metrics.accuracy_score(y_test,pred_bow))
```

```
The overall accuracy score for the Train Data is : 0.9273625
The overall accuracy score for the Test Data is : 0.90045
```

```
1 #Feature importance
2 neg_class_prob_sorted = Bnb_optimal.feature_log_prob_[1, :].argsort()
3 pos_class_prob_sorted = Bnb_optimal.feature_log_prob_[0, :].argsort()
4
5 neg_bow_bnb=neg_class_prob_sorted[::-1]
6 pos_bow_bnb=pos_class_prob_sorted[::-1]
```

```
1 #Top 20 negative words
2 n=(np.take(count_vect.get_feature_names(), neg_bow_bnb[:20]))
3 sn = ""
4 for i in n:
5     sn += str(i)+", "
6 print(sn)
```

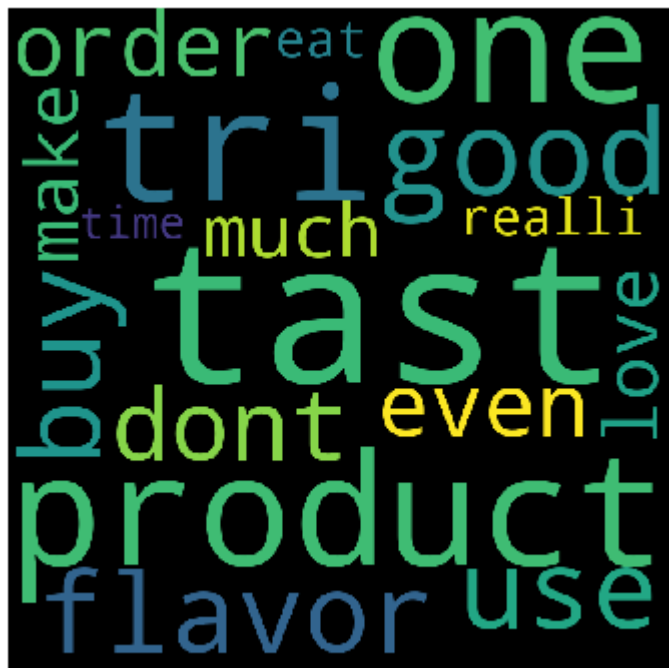
tast,like,product,one,would,tri,good,flavor,buy,get,use,dont,order,even,much,make,love,realli,eat,time,

```
1 #Top 20 positive words
2 p=(np.take(count_vect.get_feature_names(), pos_bow_bnb[:20]))
3 sp = ""
4 for i in p:
5     sp += str(i)+", "
6 print(sp)
```

tast,like,great,good,love,flavor,one,use,product,tri,make,get,best,time,buy,find,amazon,eat,realli,also,

```
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 *****



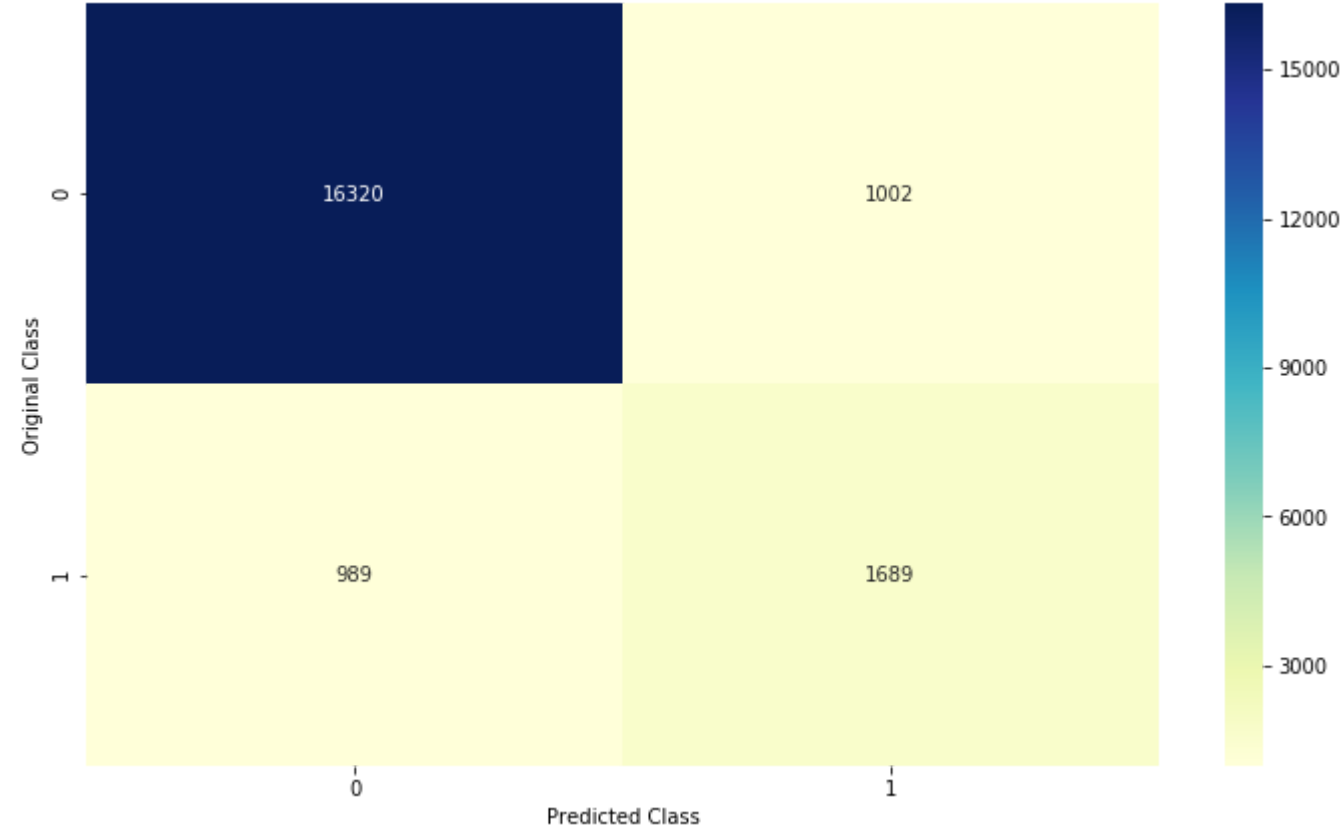
***** Top 20 Positive words *****



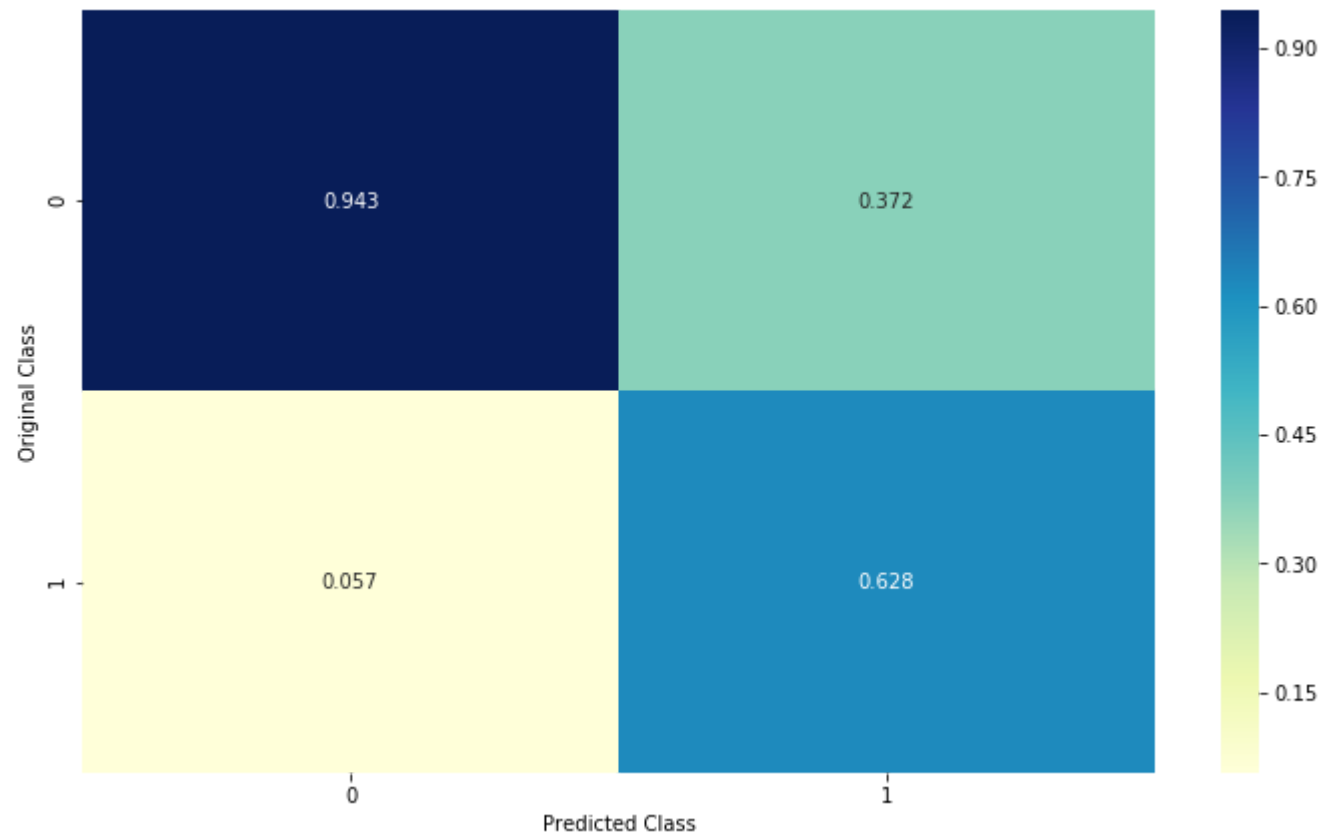
```
1 #confusion matrix, Precision score & Recall score
2 C = confusion_matrix(y_test, pred_bow)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

```
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()
```

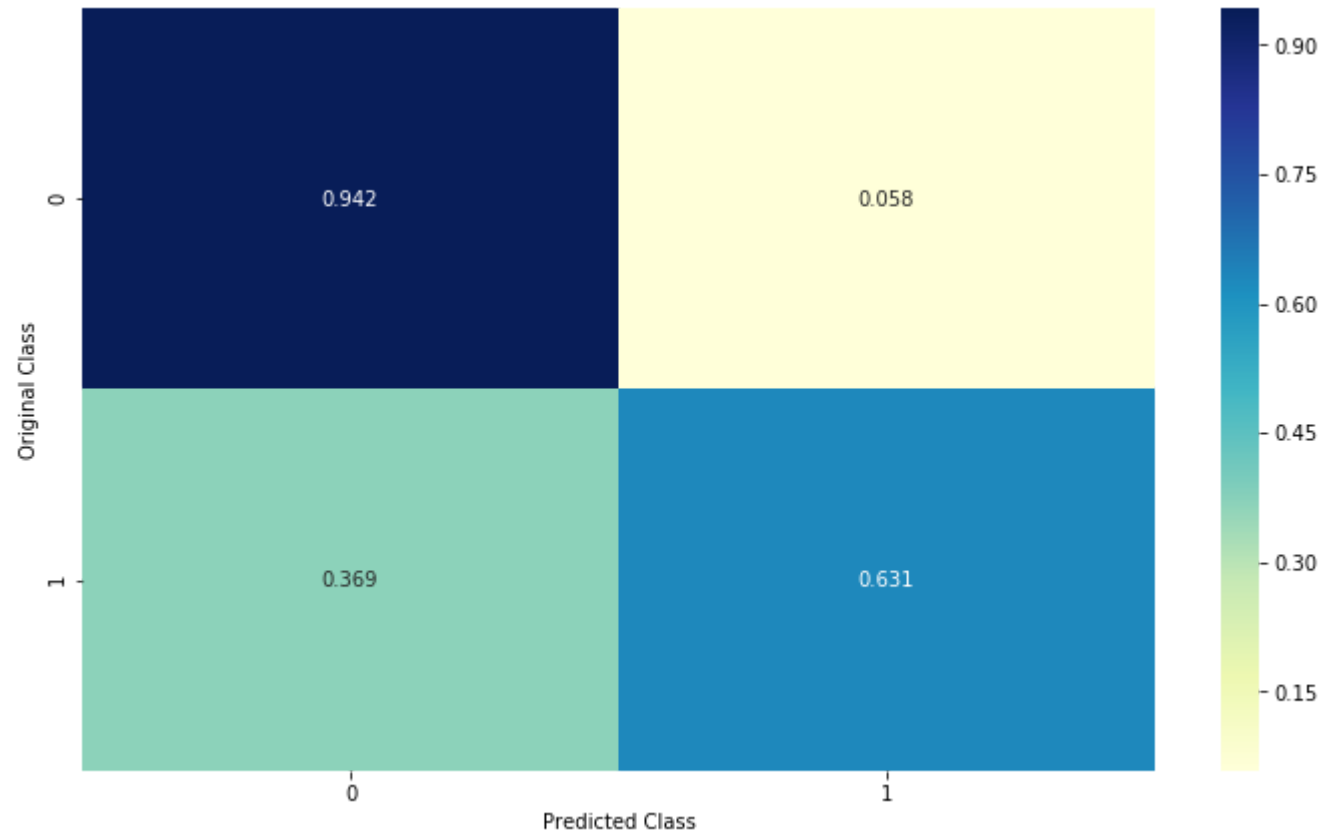
----- Confusion matrix -----



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



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



Tf-Idf

```
1 #Initiating Vectorizer
2 count_vect = CountVectorizer(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

Fitting Grid Search CV 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_)
```

0.8943373584339609
{'alpha': 0.0001}

```
1 # Naive bayes
2 Bnb_optimal=BernoulliNB(alpha=0.0001, binarize=0.0, fit_prior=True)
3
4 # fitting the model
5 Bnb_optimal.fit(Tfidf_x_train, y_train)
6
7 # predict the response
8 pred_tfidf = Bnb_optimal.predict(Tfidf_x_test)
9
10 # evaluate accuracy
11 acc = accuracy_score(y_test, pred_tfidf) * 100
12 print('\nThe accuracy of the Bernoulli NB classifier for a = %d is %f%%' % (0.0001, acc))
```

The accuracy of the Bernoulli NB classifier for a = 0 is 88.180000%

```
1 # Train & Test Error
2 print("The overall accuracy score for the Train Data is : ", metrics.accuracy_score(y_train,Bnb_optimal
3 print("The overall accuracy score for the Test Data is : ", metrics.accuracy_score(y_test,pred_tfidf))
```

The overall accuracy score for the Train Data is : 0.998775

The overall accuracy score for the Test Data is : 0.8818

```
1 #Feature importance
2 neg_class_prob_sorted = Bnb_optimal.feature_log_prob_[1, :].argsort()
3 pos_class_prob_sorted = Bnb_optimal.feature_log_prob_[0, :].argsort()
4
5 neg_tfidf_bnb=neg_class_prob_sorted[::-1]
6 pos_tfidf_bnb=pos_class_prob_sorted[::-1]
```



```
1 #Top 20 negative words
2 n=(np.take(count_vect.get_feature_names(), neg_tfidf_bnb[:20]))
3 sn = ""
4 for i in n:
5     sn += str(i)+", "
6 print(sn)
```

tast,like,product,one,would,tri,good,flavor,buy,get,use,dont,order,even,much,make,love,realli,eat,time,

```
1 #Top 20 positive words
2 p=(np.take(count_vect.get_feature_names(), pos_tfidf_bnb[:20]))
3 sp = ""
4 for i in p:
5     sp += str(i)+", "
6 print(sp)
```

tast,like,great,good,love,flavor,one,use,product,tri,make,get,best,time,buy,find,amazon,eat,realli,also,

```
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 *****



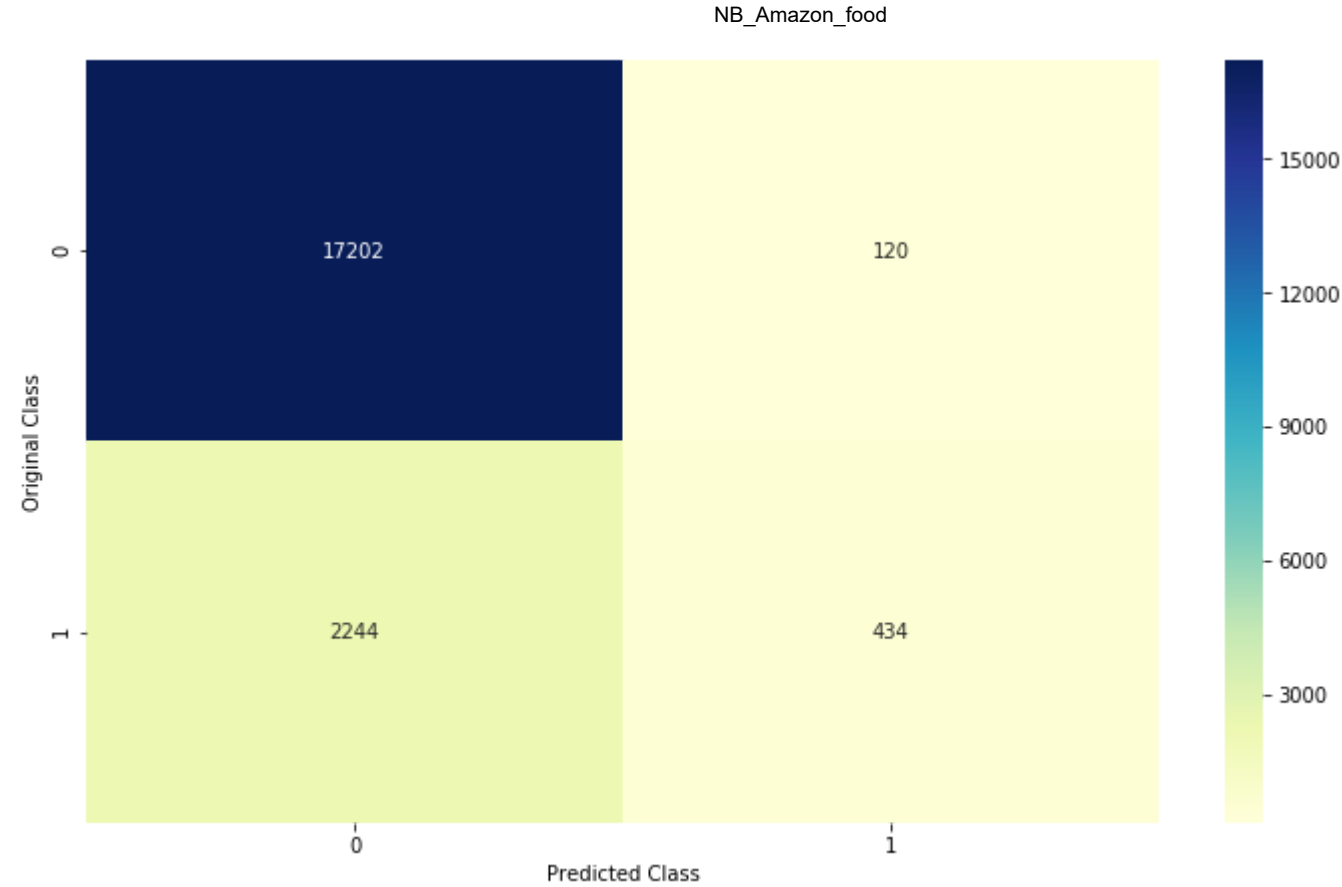
***** Top 20 Positive words *****



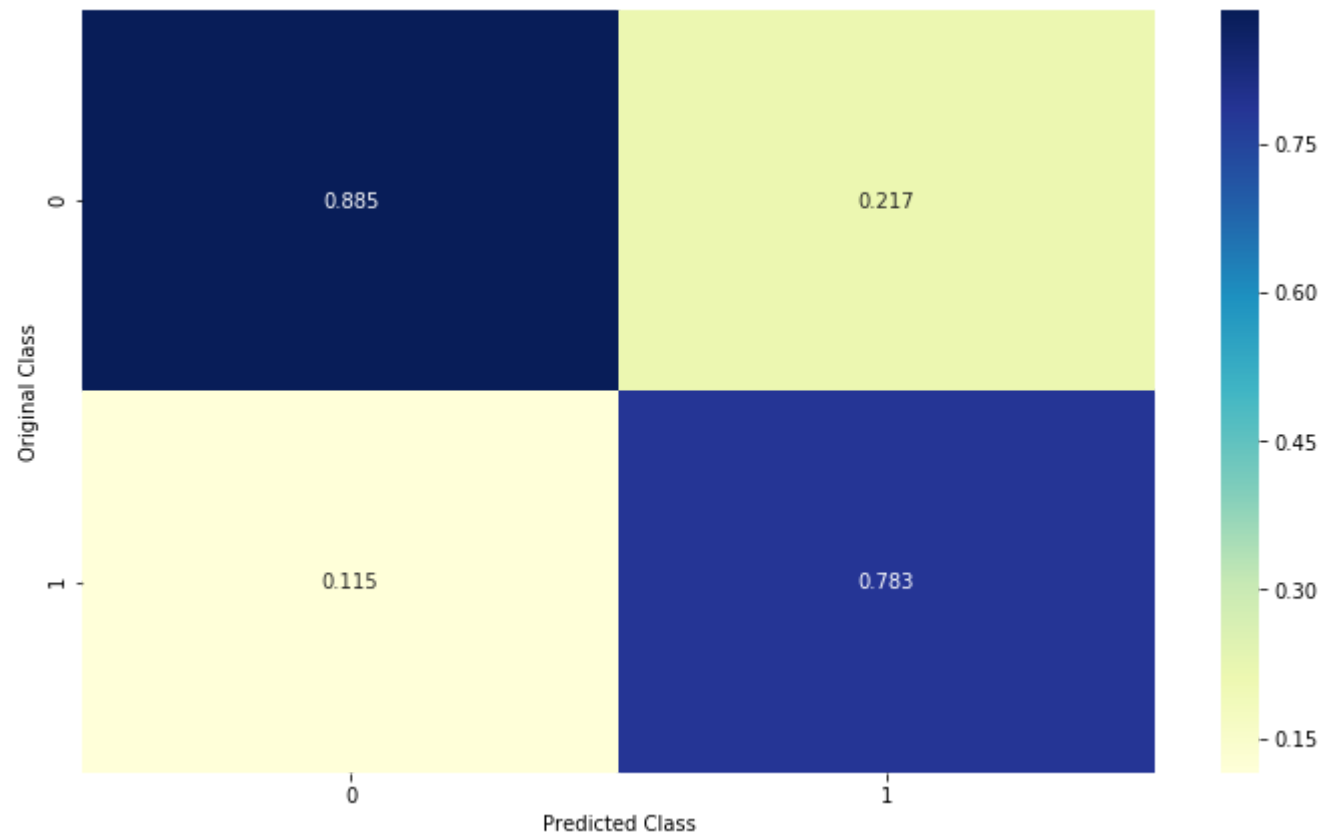
```
1 #confusion matrix, Precision score & Recall score
2 C = confusion_matrix(y_test, pred_tfidf)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

```
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()
```

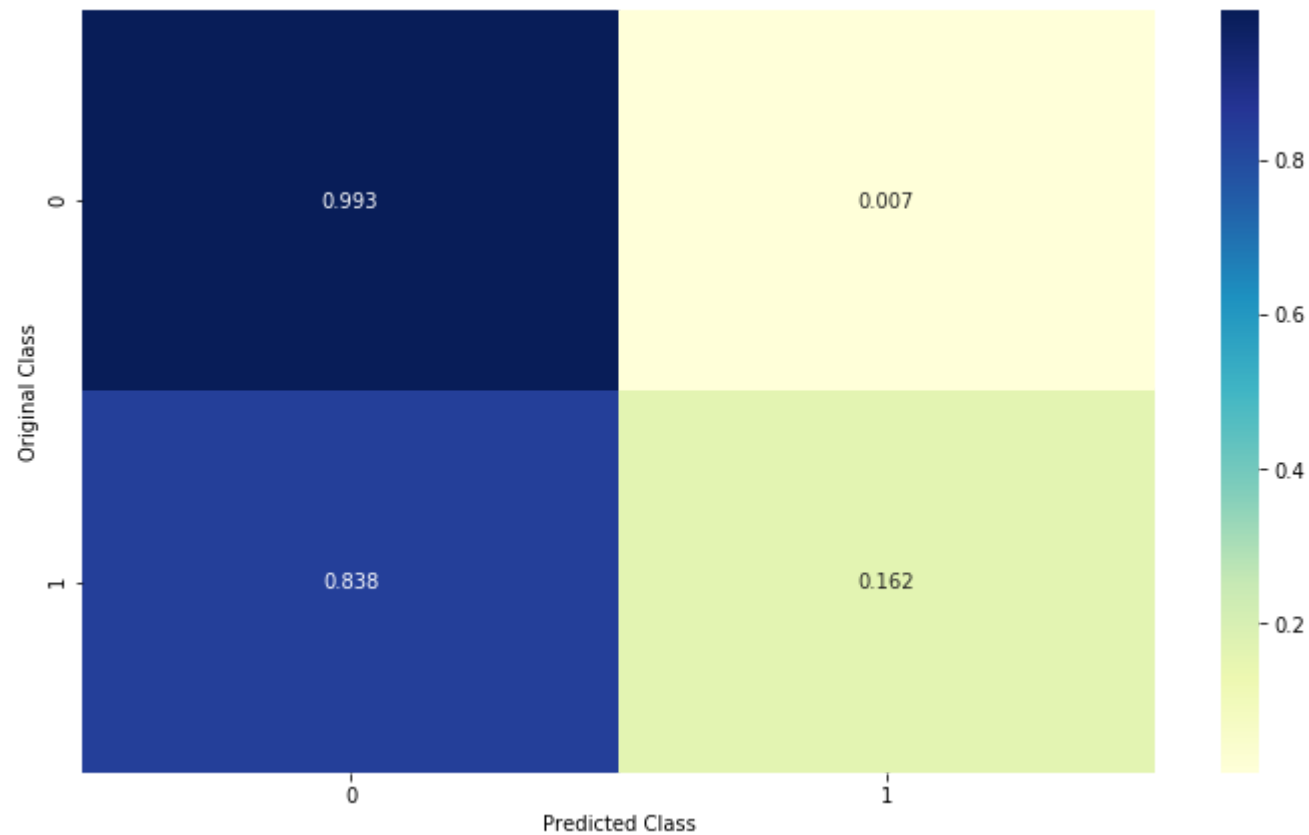
----- Confusion matrix -----



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



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



Using MultinomialNB

Running Gridsearch CV


```
1 alpha_range = [0.000001,0.00001,0.0001,0.001,0.01,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0,10,20,30,40,50]
2 T= TimeSeriesSplit(n_splits=5)
3
4 param_grid = dict(alpha=alpha_range)
5 print(param_grid)
6
7 # instantiate and fit the grid
8 grid = GridSearchCV(MultinomialNB(), param_grid, cv=T, scoring='accuracy', return_train_score=False,n_jobs=5)

{'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, 10, 20, 30, 40, 50]}
```

Fitting Grid 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_)

0.9072976824420611
{'alpha': 0.5}
```

```
1 # Naive bayes
2 Mnb_optimal=MultinomialNB(alpha=0.5,fit_prior=True)
3
4 # fitting the model
5 Mnb_optimal.fit(Bow_x_train, y_train)
6
7 # predict the response
8 pred_bow = Mnb_optimal.predict(Bow_x_test)
9
10 # evaluate accuracy
11 acc = accuracy_score(y_test, pred_bow) * 100
12 print('\nThe accuracy of the Multinomial NB classifier for a = %d is %f%%' % (0.5, acc))
```

The accuracy of the Multinomial NB classifier for a = 0 is 90.920000%

```
1 # Train & Test Error
2 print("The overall accuracy score for the Train Data is : ", metrics.accuracy_score(y_train,Mnb_optimal
3 print("The overall accuracy score for the Test Data is : ", metrics.accuracy_score(y_test,pred_bow))
```

The overall accuracy score for the Train Data is : 0.9290625

The overall accuracy score for the Test Data is : 0.9092

```
1 #Feature importance
2 neg_class_prob_sorted = Mnb_optimal.feature_log_prob_[1, :].argsort()
3 pos_class_prob_sorted = Mnb_optimal.feature_log_prob_[0, :].argsort()
4
5 neg_bow_Mnb=neg_class_prob_sorted[::-1]
6 pos_bow_Mnb=pos_class_prob_sorted[::-1]
```

```
1 #Top 20 negative words
2 n=(np.take(count_vect.get_feature_names(), neg_bow_Mnb[:20]))
3 sn = ""
4 for i in n:
5     sn += str(i)+", "
6 print(sn)
```

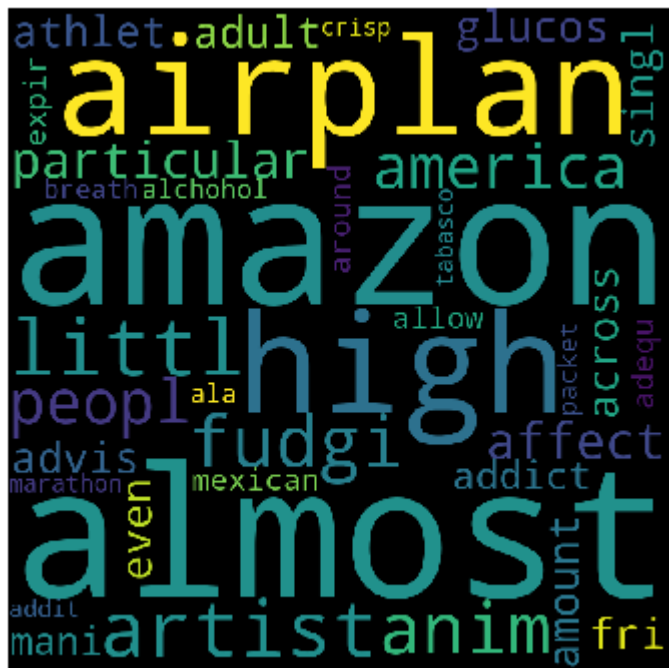
amazon high,airplan littl,also artist,almost fudgi,anim peopl,america particular,affect glucos,adult athlet,across fri,advis singl,amount even,addict mani,almost mexican,adequ expir,allow around,alcohol breath,ala packet,also marathon,addit crisp,amazon tabasco,

```
1 #Top 20 positive words
2 p=(np.take(count_vect.get_feature_names(), pos_bow_Mnb[:20]))
3 sp = ""
4 for i in p:
5     sp += str(i)+", "
6 print(sp)
```

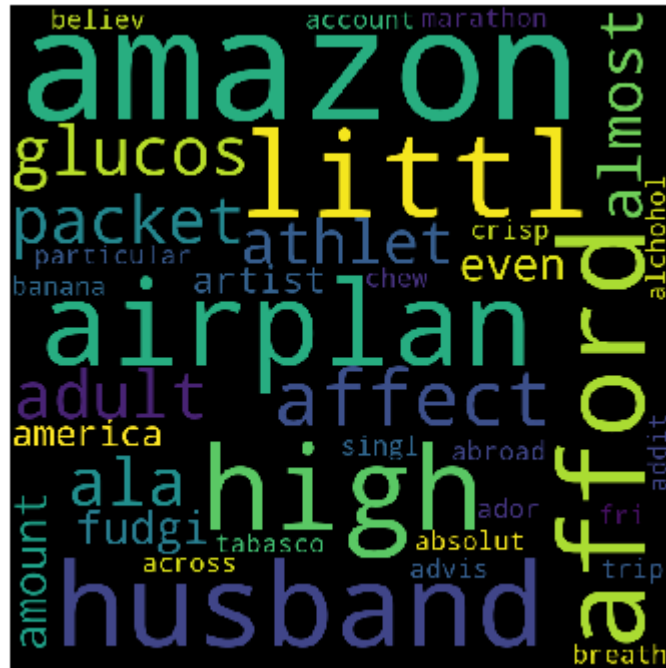
amazon high,airplan littl,afford husband,affect glucos,ala packet,adult athlet,almost fudgi,amount even,also artist,america particular,alcohol breath,advis singl,account believ,amazon tabasco,across fri,ador chew,absolut banana,addit crisp,also marathon,abroad trip,

```
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 *****



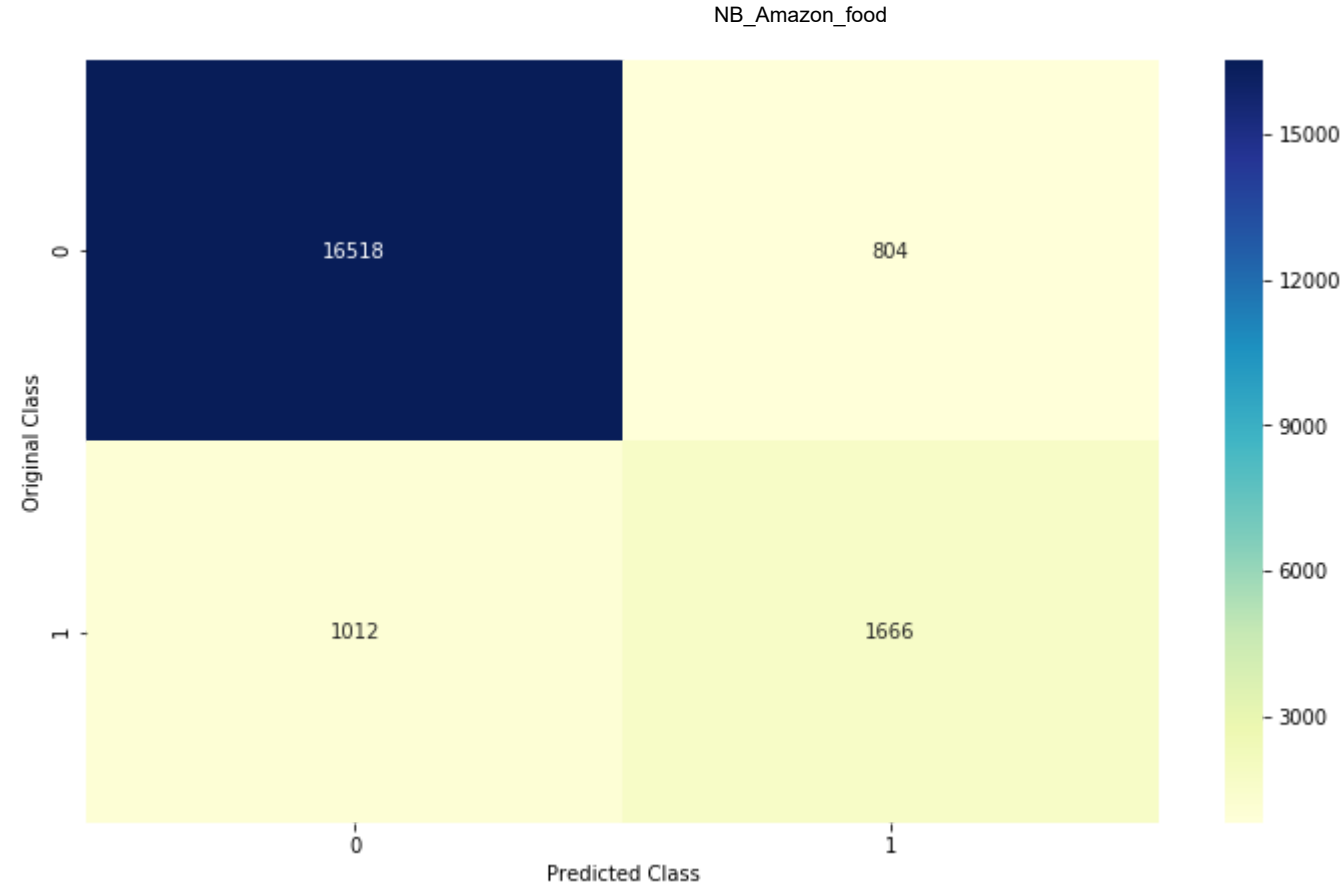
***** Top 20 Positive words *****



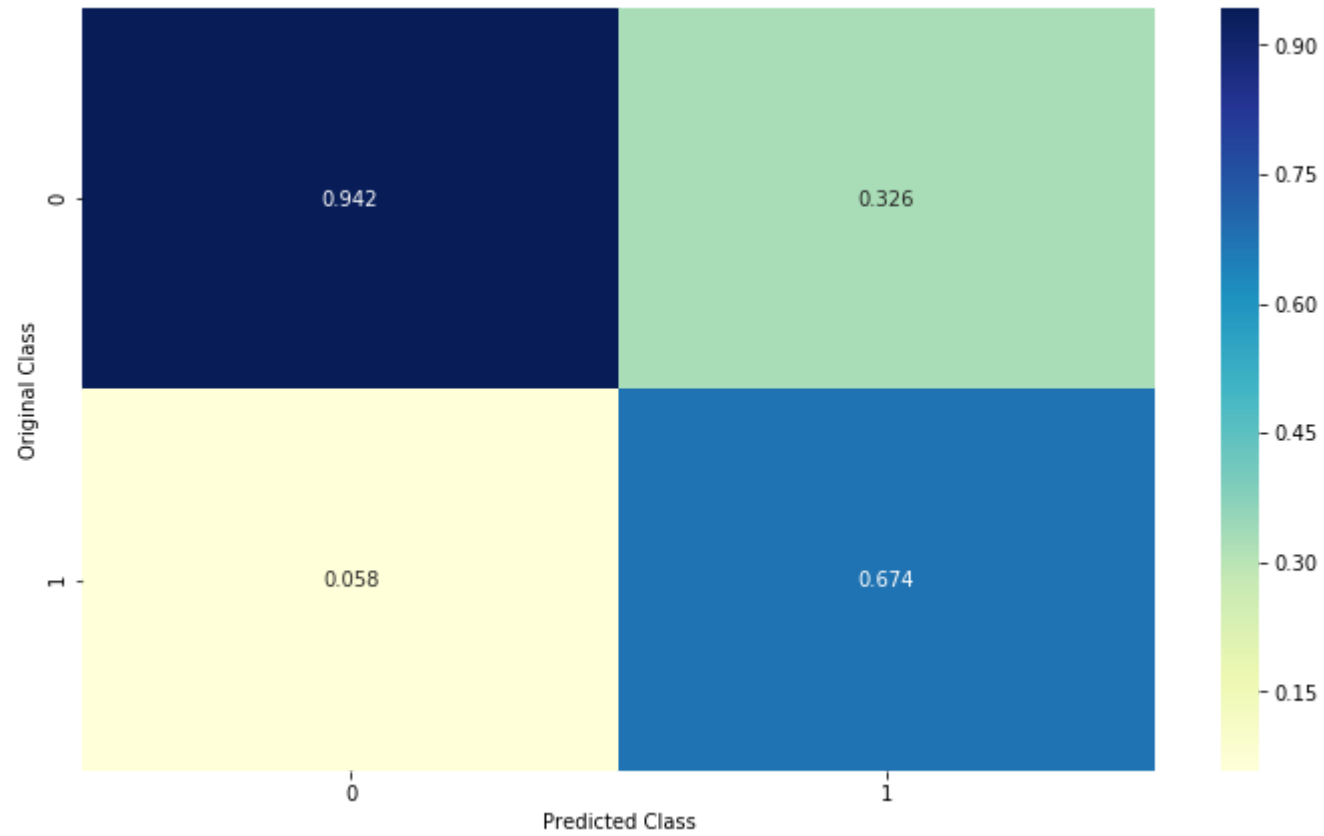
```
1 #confusion matrix, Precision score & Recall score
2 C = confusion_matrix(y_test, pred_bow)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

```
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()
```

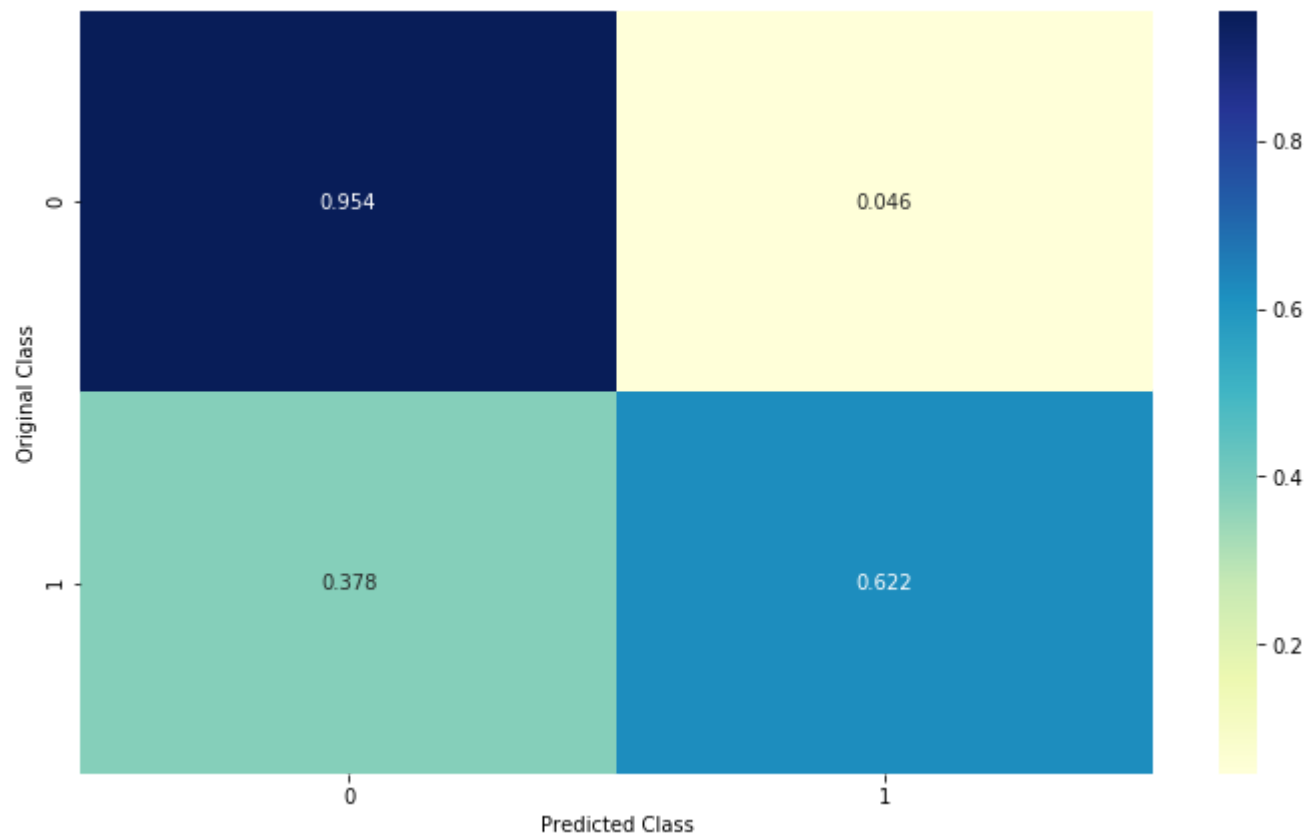
----- Confusion matrix -----



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



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



Fitting gridsearch cv on Tfidf

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

```
0.8972174304357609
{'alpha': 0.5}
```

```
1  # Naive bayes
2  mnb_optimal=BernoulliNB(alpha=0.5, binarize=0.0, fit_prior=True)
3
4  # fitting the model
5  mnb_optimal.fit(Tfidf_x_train, y_train)
6
7  # predict the response
8  pred_tfidf = mnb_optimal.predict(Tfidf_x_test)
9
10 # evaluate accuracy
11 acc = accuracy_score(y_test, pred_tfidf) * 100
12 print('\nThe accuracy of the Multinomial NB classifier for a = %d is %f%%' % (0.0001, acc))
```

The accuracy of the Multinomial NB classifier for a = 0 is 87.305000%

```
1  # Train & Test Error
2  print("The overall accuracy score for the Train Data is : ", metrics.accuracy_score(y_train,mnb_optimal
3  print("The overall accuracy score for the Test Data is : ", metrics.accuracy_score(y_test,pred_tfidf))
```

The overall accuracy score for the Train Data is : 0.9253375

The overall accuracy score for the Test Data is : 0.87305

```
1  #Feature importance
2  neg_class_prob_sorted = Mnb_optimal.feature_log_prob_[1, :].argsort()
3  pos_class_prob_sorted = Mnb_optimal.feature_log_prob_[0, :].argsort()
4
5  neg_bow_Mnb=neg_class_prob_sorted[::-1]
6  pos_bow_Mnb=pos_class_prob_sorted[::-1]
```

```
1 #Top 20 negative words
2 n=(np.take(count_vect.get_feature_names(), neg_bow_Mnb[:20]))
3 sn = ""
4 for i in n:
5     sn += str(i)+", "
6 print(sn)
```

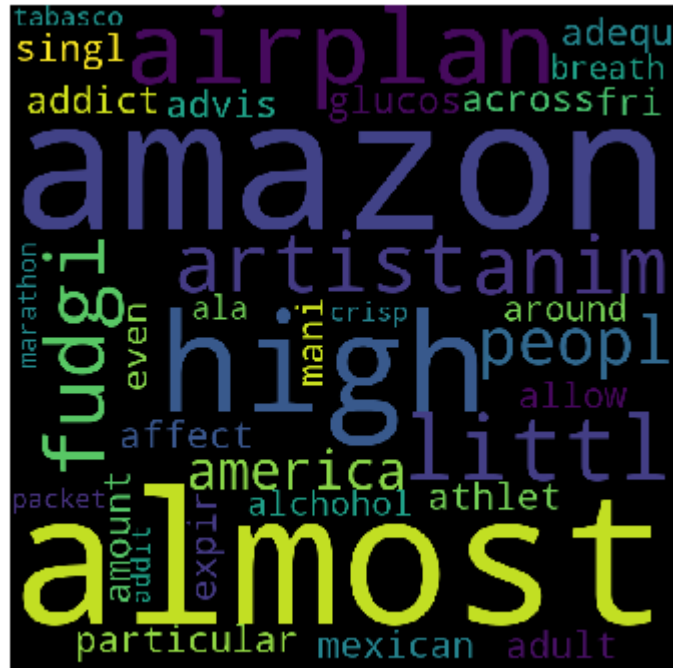
amazon high,airplan littl,also artist,almost fudgi,anim peopl,america particular,affect glucos,adult athlet,across fri,advis singl,amount even,addict mani,almost mexican,adequ expir,allow around,alcohol breath,ala packet,also marathon,addit crisp,amazon tabasco,

```
1 #Top 20 positive words
2 p=(np.take(count_vect.get_feature_names(), pos_bow_Mnb[:20]))
3 sp = ""
4 for i in p:
5     sp += str(i)+", "
6 print(sp)
```

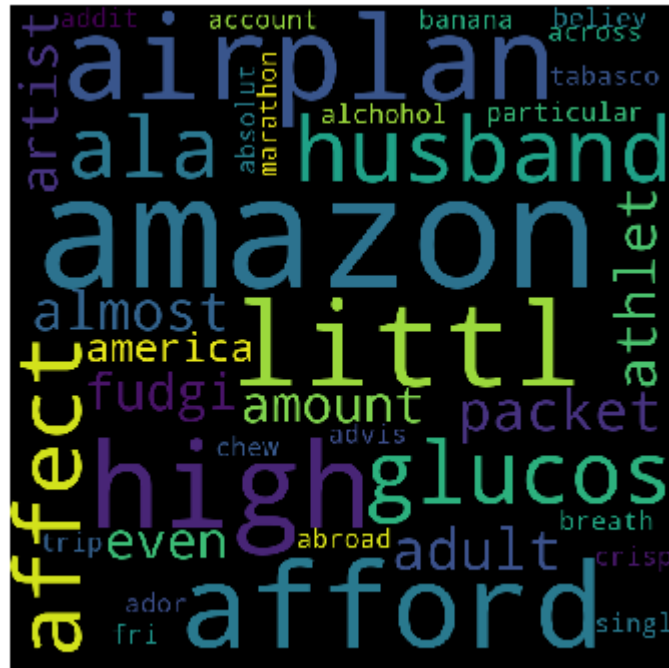
amazon high,airplan littl,afford husband,affect glucos,ala packet,adult athlet,almost fudgi,amount even,also artist,america particular,alcohol breath,advis singl,account believ,amazon tabasco,across fri,ador chew,absolut banana,addit crisp,also marathon,abroad trip,

```
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 *****



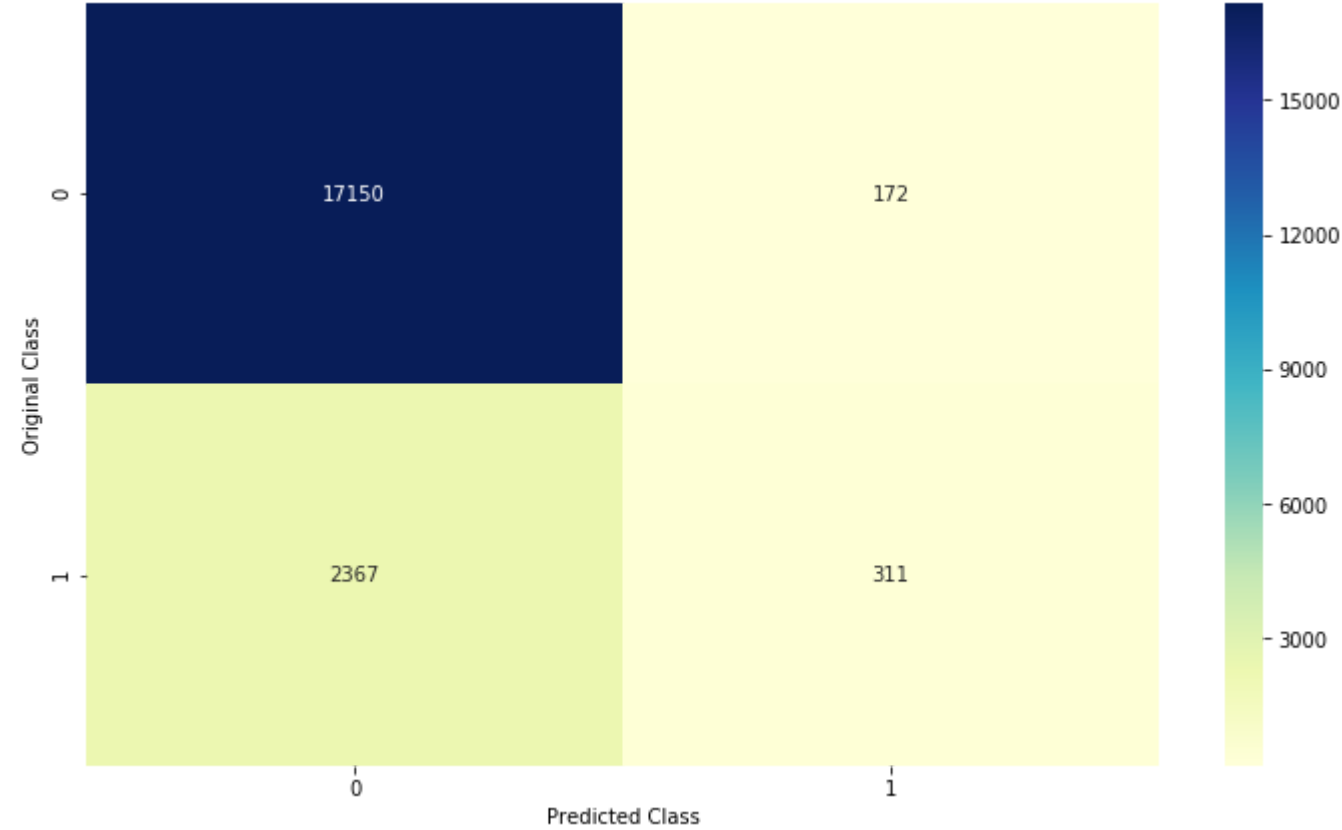
```
***** Top 20 Positive words *****
```



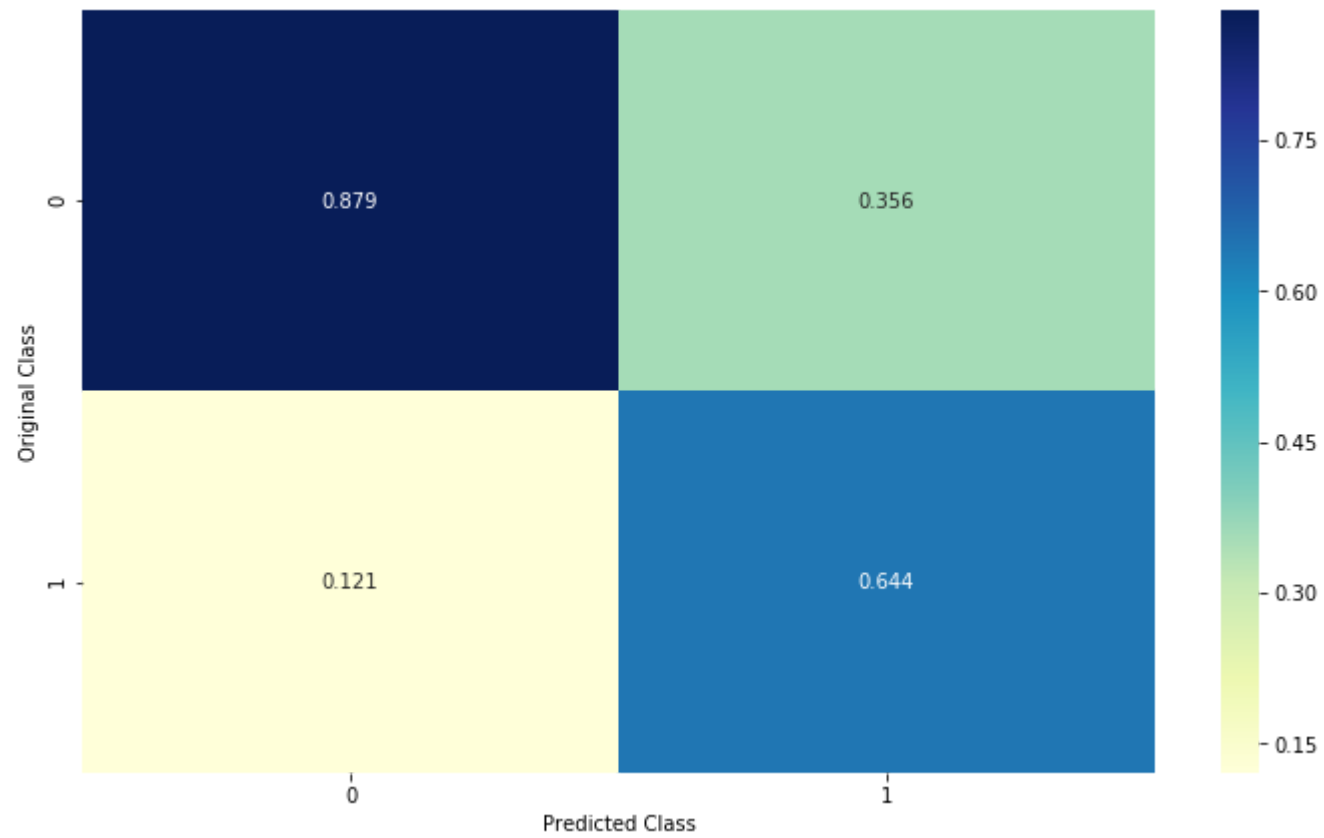
```
1 #confusion matrix, Precision score & Recall score
2 C = confusion_matrix(y_test, pred_tfidf)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

```
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()
```

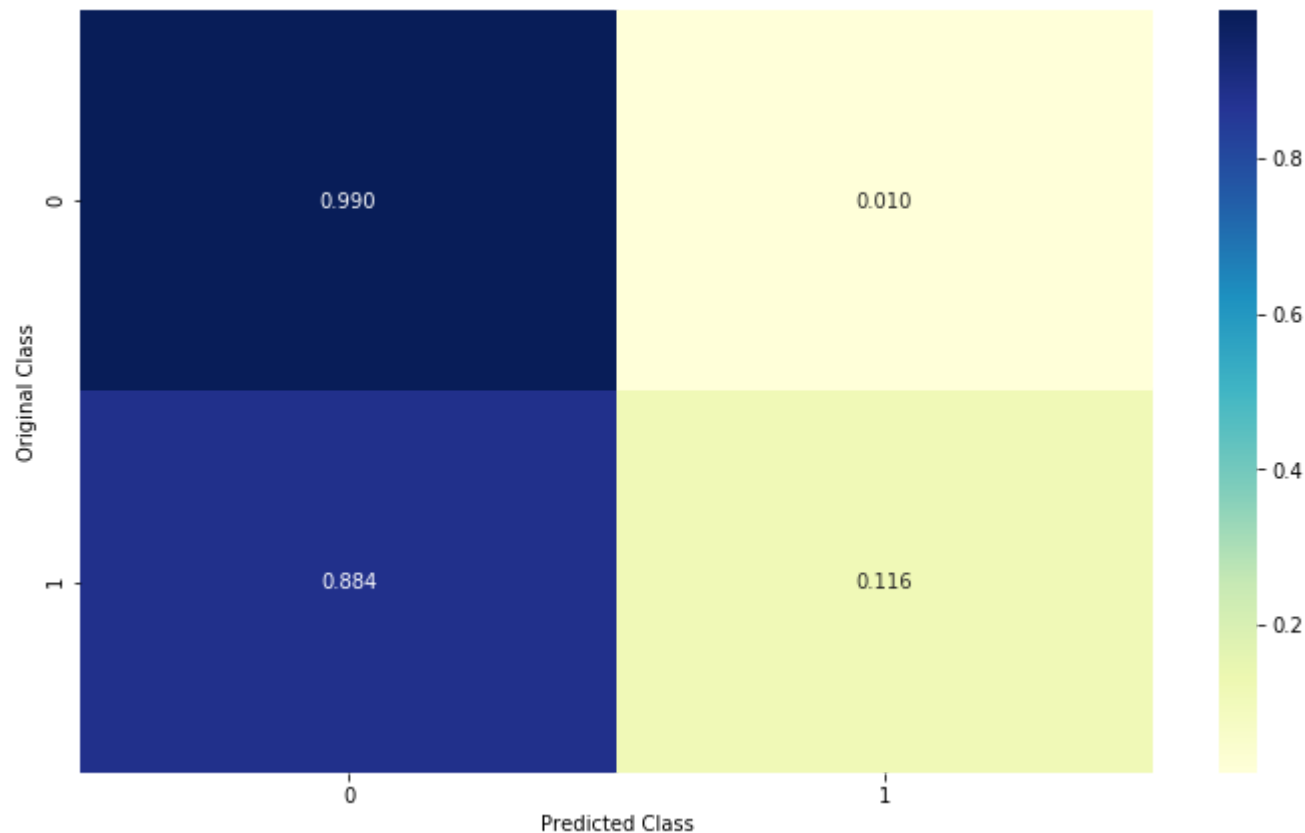
----- Confusion matrix -----



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



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



Bernaulli NB v/s Multinomial NB

```
1 !pip install prettytable
```

Collecting prettytable

Downloading <https://files.pythonhosted.org/packages/ef/30/4b0746848746ed5941f052479e7c23d2b56d174b82f4fd34a25e389831f5/prettytable-0.7.2.tar.bz2> (https://files.pythonhosted.org/packages/ef/30/4b0746848746ed5941f052479e7c23d2b56d174b82f4fd34a25e389831f5/prettytable-0.7.2.tar.bz2)

Building wheels for collected packages: prettytable

Running setup.py bdist_wheel for prettytable: started

Running setup.py bdist_wheel for prettytable: finished with status 'done'

Stored in directory: C:\Users\deepak\AppData\Local\pip\Cache\wheels\80\34\1c\3967380d9676d162cb59513bd9dc862d0584e045a162095606

Successfully built prettytable

Installing collected packages: prettytable

Successfully installed prettytable-0.7.2

```
1 from prettytable import PrettyTable
```

```
1 x=PrettyTable()
2 x.field_names = ["Model", "BernuallNB_BOW", "BernuallNB_tfidf", "MultinomialNB_BOW", "MultinomialNB_Tf:
3 x.add_row(["Hyperparameter", 0.1, 0.0001, 0.5, 0.5])
4 x.add_row(["Train Error", 0.927, 0.998, 0.929, 0.925])
5 x.add_row(["Test Error", 0.900, 0.881, 0.909, 0.873])
6
7 print(x)
```

Model	BernuallNB_BOW	BernuallNB_tfidf	MultinomialNB_BOW	MultinomialNB_Tfidf
Hyperparameter	0.1	0.0001	0.5	0.5
Train Error	0.927	0.998	0.929	0.925
Test Error	0.9	0.881	0.909	0.873