

# Logistic Regression 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 Logistic regression using L1 regularization 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
33 import os
```

```
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.linear_model import LogisticRegression
2 from sklearn.model_selection import GridSearchCV
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
```

## Applying L1 regularization

### Gridsearch CV using L2

```

1 gamma_range = [0.000001,0.00001,0.0001,0.001,0.01,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0,5,10,15,20,30,
2 T= TimeSeriesSplit(n_splits=5)
3 weight=[None,'balanced']
4
5 param_grid = dict(C=gamma_range,class_weight=weight)
6 print(param_grid)
7
8 # instantiate and fit the grid
9 grid = GridSearchCV(LogisticRegression(penalty='l1'), param_grid, cv=T, scoring='f1', return_train_score=

```

{'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, 5, 10, 15, 20, 30, 40, 50], '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 Scaler=StandardScaler(with_mean=False)
4 Bow_x_train = Scaler.fit_transform(Bow_x_train)
5 Bow_x_test = Scaler.transform(Bow_x_test)
6
7 print(Bow_x_train.shape)
8 print(Bow_x_test.shape)
```

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

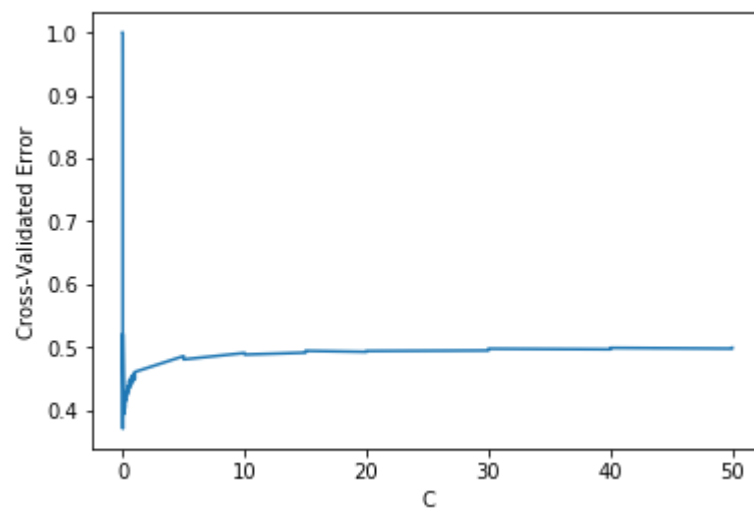
### Fitting Grid search cv on BOW

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

```
0.6296566961732983
{'C': 0.01, 'class_weight': 'balanced'}
```

```
1 #Plotting C 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
2  #{'C': 0.01, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.01,class_weight='balanced')
4
5  # fitting the model
6  LR_optimal.fit(Bow_x_train, y_train)
7
8  # predict the response
9  pred_bow = LR_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,LR_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 : 0.7461583187972424

The overall f1\_score for the Test Data is : 0.6732643118148599

### 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
2 #{'C': 0.01, 'class_weight': 'balanced'}
3 LR_optimal_noise=LogisticRegression(penalty='l1',C=0.01,class_weight='balanced')
4
5 # fitting the model
6 LR_optimal_noise.fit(Noise_Bow_x_train, y_train)
7
8 # predict the response
9 pred_bow = LR_optimal_noise.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,LR_optimal.predict(Noise_Bow_x_train)))
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 : 0.7461893092993311

The overall f1\_score for the Test Data is : 0.6737707413609377

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

(33433,)

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

(1, 33433)

```
1 #Weights after adding noise  
2 LR_optimal_noise.coef_.shape  
  
(1, 33433)
```

```

1 merge_arr = np.concatenate([LR_optimal.coef_, LR_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	0.000000	0.000000	NaN
aaaaaaaaagghh	0.000000	0.000000	NaN
aaaaah	0.000000	0.000000	NaN
aaaaahhhhhhhhhhhhhhhh	0.000000	0.000000	NaN
aaaah	0.000000	0.000000	NaN
aaah	0.000000	0.000000	NaN
aachen	0.000000	0.000000	NaN
aad	0.000000	0.000000	NaN
aadp	0.000000	0.000000	NaN
aafco	0.000000	0.000000	NaN
aagh	-0.000872	-0.000104	-88.105923
aah	0.000000	0.000000	NaN
aahh	0.000000	0.000000	NaN
aand	0.000000	0.000000	NaN
aardvark	-0.006501	-0.001462	-77.508067
aarrgh	0.009030	0.009029	-0.002622
ab	0.000000	0.000000	NaN
aback	-0.000381	-0.000386	1.460521
abandon	0.000000	0.000000	NaN
abaolut	0.000000	0.000000	NaN
abattoir	0.000000	0.000000	NaN
abba	0.000000	0.000000	NaN
abbey	0.000000	0.000000	NaN
abbi	0.000000	0.000000	NaN

	0	1	2
abbott	0.000000	0.000000	NaN
abbrevi	0.000000	0.000000	NaN
abc	0.000000	0.000000	NaN
abcstor	0.000000	0.000000	NaN
abd	0.000000	0.000000	NaN
abdomen	0.000000	0.000000	NaN
...	...	...	...
zot	-0.011494	-0.011496	0.014257
zotz	0.000000	0.000000	NaN
zour	0.000000	0.000000	NaN
zout	0.000000	0.000000	NaN
zowi	0.000000	0.000000	NaN
zreport	0.000000	0.000000	NaN
zsweet	0.000000	0.000000	NaN
zuc	0.000000	0.000000	NaN
zucchini	0.000000	0.000000	NaN
zuccini	0.000000	0.000000	NaN
zuccnini	0.000000	0.000000	NaN
zuchinni	0.000000	0.000000	NaN
zuke	0.000000	0.000000	NaN
zulu	0.000000	0.000000	NaN
zum	0.000000	0.000000	NaN
zummi	0.000000	0.000000	NaN
zune	0.000000	0.000000	NaN
zupreem	0.000000	0.000000	NaN
zurich	0.000000	0.000000	NaN
zwar	0.000000	0.000000	NaN
zwieback	0.000646	0.000641	-0.704890
zwiebeck	0.000000	0.000000	NaN

	0	1	2
zydeco	0.000000	0.000000	NaN
zzzzzs	0.000000	0.000000	NaN
zzzzzz	0.000000	0.000000	NaN
zzzzzzzz	0.000759	0.000756	-0.283037
zzzzzzzzzz	0.000000	0.000000	NaN
zzzzzzzzzzzz	0.000000	0.000000	NaN
zzzzzzzzzzzzzz	0.000000	0.000000	NaN
çay	0.000000	0.000000	NaN

33433 rows x 3 columns

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

(253, 3)

253 features out of 33433 shows percentage change > 30 post pertubation test i.e 0.75%

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

```
1 feature_names = np.array(vocabulary.get_feature_names())
2 sorted_coef_index = LR_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,best,love,perfect,delici,excel,good,favorit,amaz,nice,wonder,addict,tasti,find,awesom,thank,keep,happi,smooth,year,

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

disappoint,worst,terribl,thought,tast,aw,unfortun,horribl,bad,bland,would,return,stale,money,didnt,hope,product,stick,thre  
w,weak,

```
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 \*\*\*\*\*



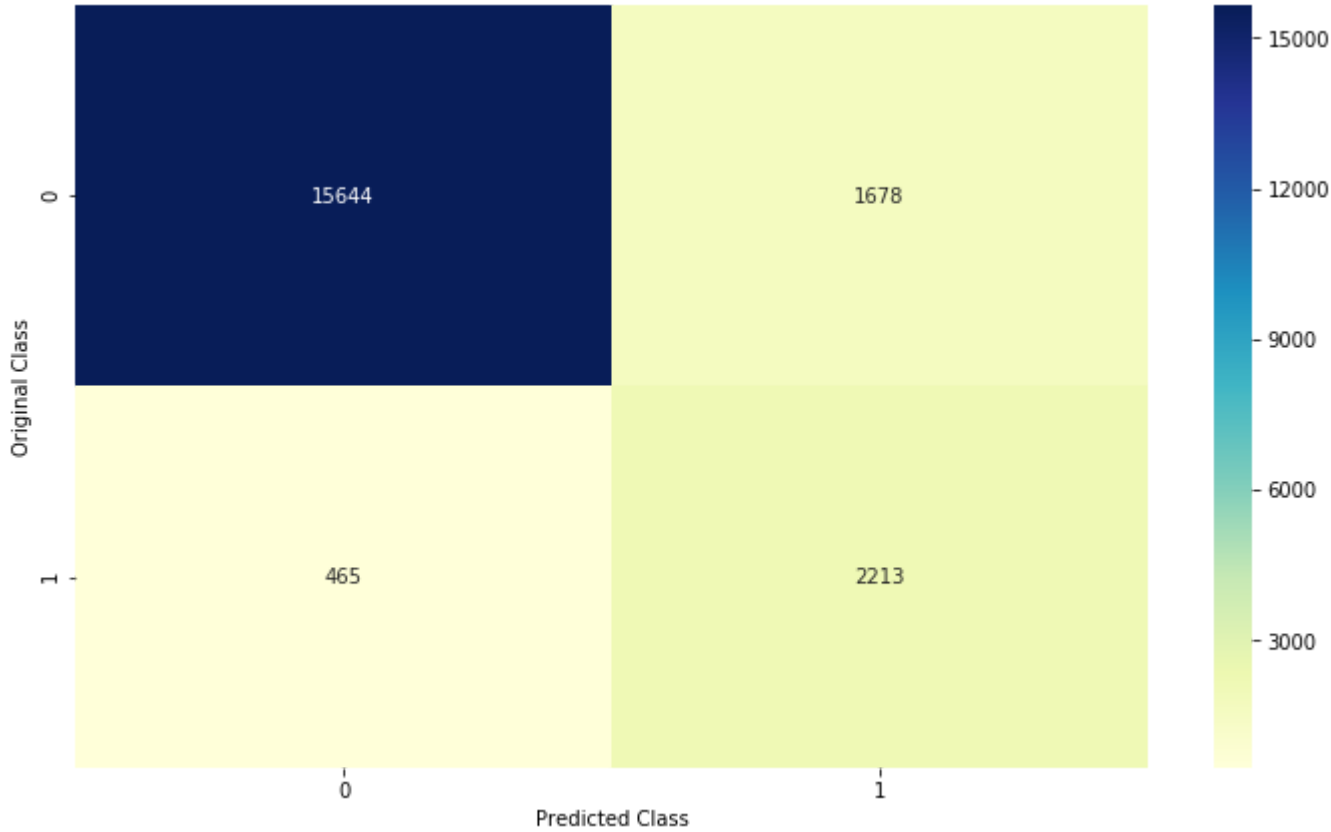




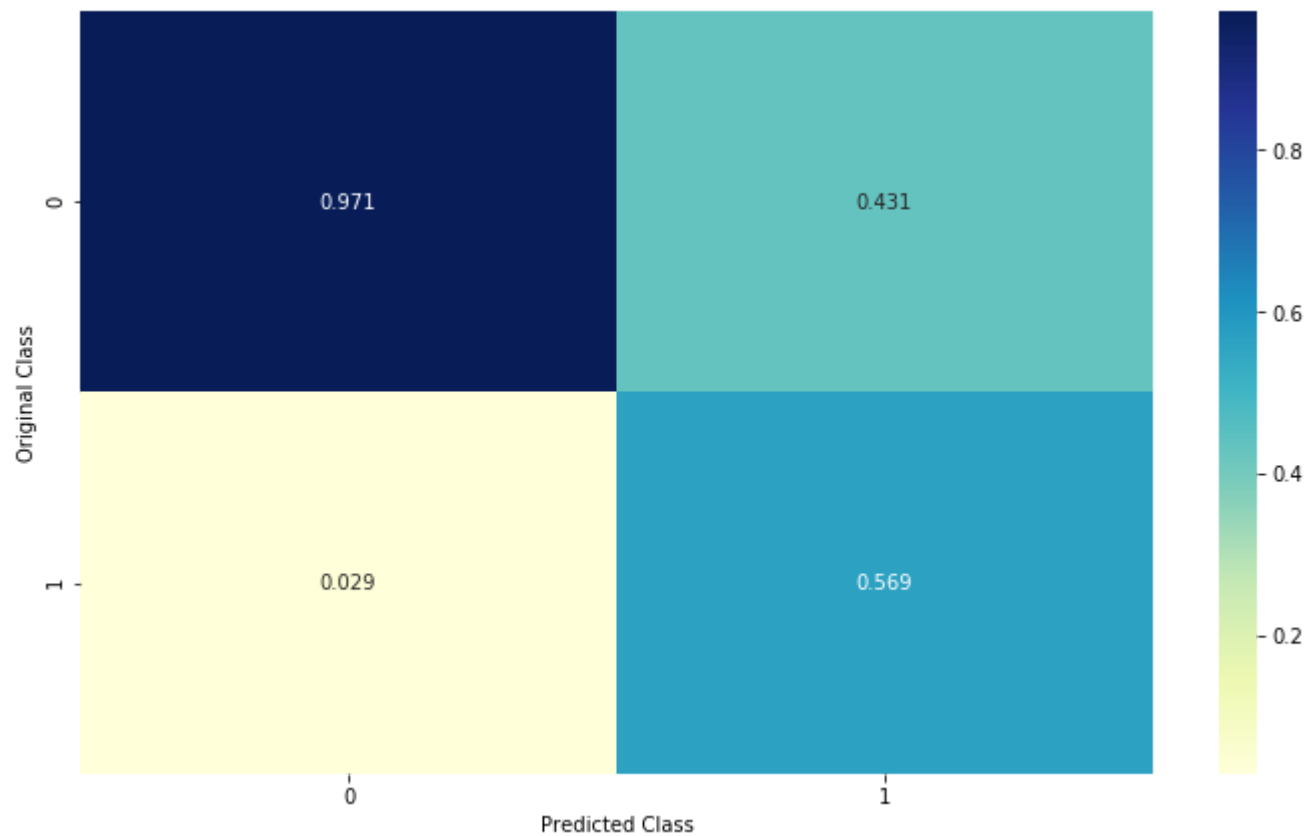
```
1 #Confusion matrix
2 C = confusion_matrix(y_test, pred_bow)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

```
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()
```

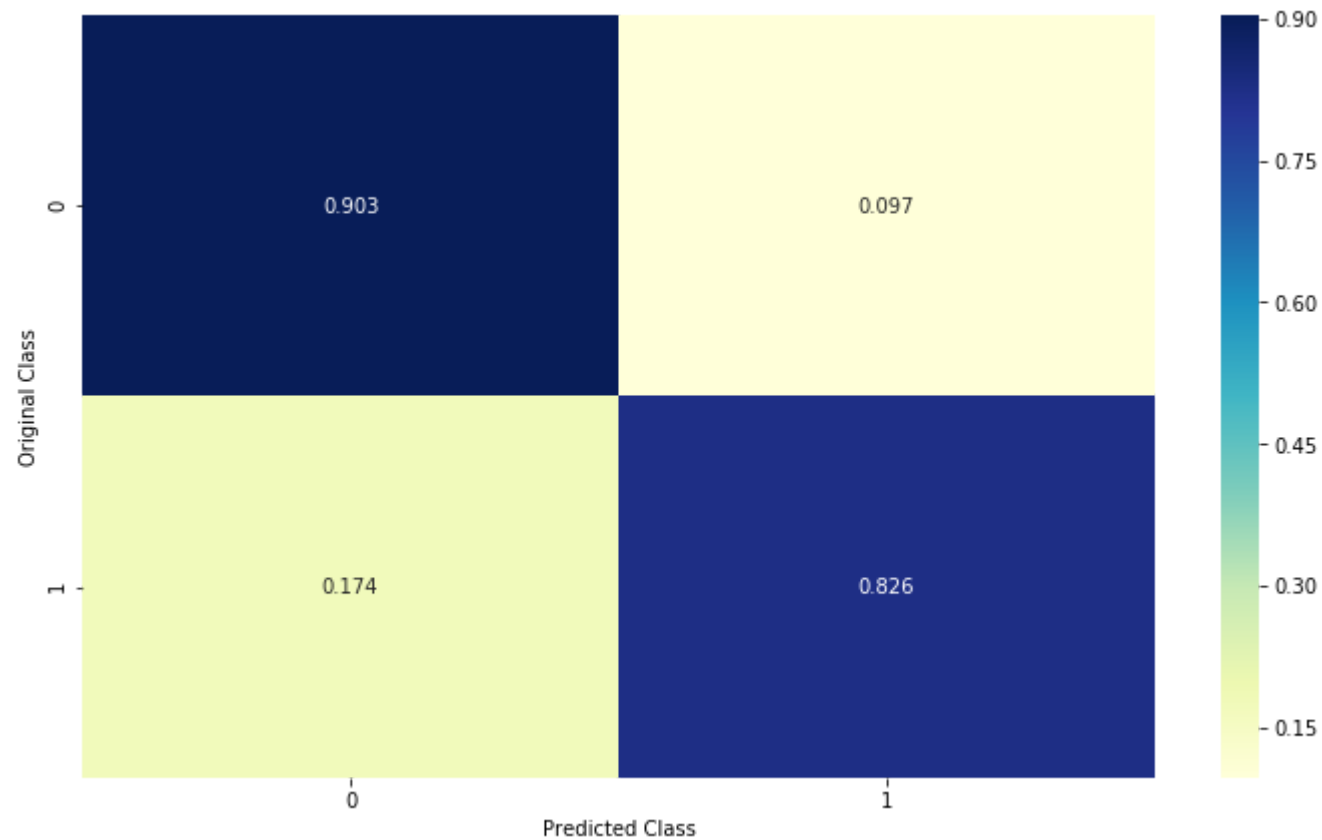
----- Confusion matrix -----



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



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



What if we decrease the value of C? How does it effect Sparsity and Error?

$C=1/\lambda$

```
1 #Checking current non-zero elements
2 a=np.array(LR_optimal.coef_)
3 print(np.count_nonzero(a))
```

5293

```
1  #Decreasing the value of C
2  #{'C': 0.01, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.001,class_weight='balanced')
4
5  # fitting the model
6  LR_optimal.fit(Bow_x_train, y_train)
7
8  # predict the response
9  pred_bow = LR_optimal.predict(Bow_x_test)
10
11 # evaluate f1_score
12 f1_score = f1_score(y_test, pred_bow)
```

```
1  #Checking non-zero elements
2  a=np.array(LR_optimal.coef_)
3  print(np.count_nonzero(a))
```

245

```
1  #Decreasing the value of C
2  #{'C': 0.01, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.001,class_weight='balanced')
4
5  # fitting the model
6  LR_optimal.fit(Bow_x_train, y_train)
7
8  # predict the response
9  pred_bow = LR_optimal.predict(Bow_x_test)
10
11 # evaluate f1_score
12 f1_score = f1_score(y_test, pred_bow)
```

```
1 #Checking non-zero elements
2 a=np.array(LR_optimal.coef_)
3 print(np.count_nonzero(a))
```

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_)
```

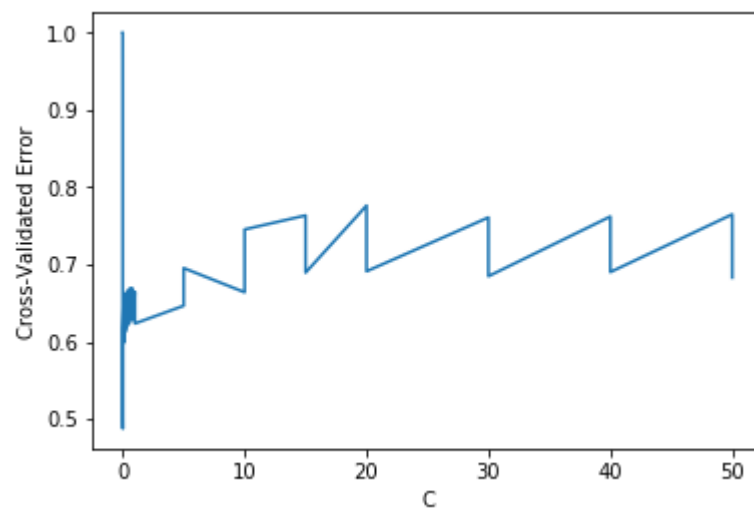
```
0.5125440355342478
```

```
{'C': 0.01, 'class_weight': 'balanced'}
```



```
1 #Plotting C 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
2 #{'C': 0.01, 'class_weight': 'balanced'}
3 LR_optimal=LogisticRegression(penalty='l1',C=0.01,class_weight='balanced')
4
5 # fitting the model
6 LR_optimal.fit(Tfidf_x_train, y_train)
7
8 # predict the response
9 pred_tfidf = LR_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,LR_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 : 0.954717356687898

The overall f1\_score for the Test Data is : 0.679509632224168

### 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 from sklearn.metrics import f1_score
2 #{'C': 0.01, 'class_weight': 'balanced'}
3 LR_optimal_noise=LogisticRegression(penalty='l1',C=0.01,class_weight='balanced')
4
5 # fitting the model
6 LR_optimal_noise.fit(Noise_Tfidf_x_train, y_train)
7
8 # predict the response
9 pred_tfidf = LR_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,LR_optimal.predict(Noise_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 : 0.9545273631840796

The overall f1\_score for the Test Data is : 0.6779925795743019

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

(1013943,)

```
1 LR_optimal.coef_.shape
```

(1, 1013943)

```
1 LR_optimal_noise.coef_.shape
```

(1, 1013943)

```

1 merge_arr = np.concatenate([LR_optimal.coef_, LR_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	0.0	0.000000	NaN
aaa condit	0.0	0.000000	NaN
aaa perfect	0.0	0.000000	NaN
aaaaaaaaagghh	0.0	0.000000	NaN
aaaaah	0.0	0.000000	NaN
aaaaah awak	0.0	0.000000	NaN
aaaaah satisfi	0.0	0.000000	NaN
aaaaahhhhhhhhhhhhhhhhh	0.0	0.000000	NaN
aaaaahhhhhhhhhhhhhhhhh angel	0.0	0.000000	NaN
aaaah	0.0	0.000000	NaN
aaaah snob	0.0	0.000000	NaN
aaah	0.0	0.000000	NaN
aaah inhal	0.0	0.000000	NaN
aaah miss	0.0	0.000000	NaN
aaah sip	0.0	0.000000	NaN
aachen	0.0	0.000000	NaN
aachen munich	0.0	0.000000	NaN
aad	0.0	0.000000	NaN
aad sausag	0.0	0.000000	NaN
aadp	0.0	0.000000	NaN
aafco	0.0	0.000000	NaN
aafco also	0.0	0.000000	NaN
aafco certifi	0.0	0.000168	inf
aafco countri	0.0	0.000000	NaN

	0	1	2
aafco definit	0.0	0.000000	NaN
aafco dog	0.0	0.000000	NaN
aafco guidelin	0.0	0.000000	NaN
aafco requir	0.0	0.000000	NaN
aagh	0.0	0.000000	NaN
aagh yelp	0.0	0.000000	NaN
...	...	...	...
zum heal	0.0	0.000000	NaN
zummi	0.0	0.000000	NaN
zummi love	0.0	0.000000	NaN
zummi tast	0.0	0.000000	NaN
zummi tri	0.0	0.000000	NaN
zune	0.0	0.000000	NaN
zune video	0.0	0.000000	NaN
zupreem	0.0	0.000000	NaN
zupreem ferret	0.0	0.000000	NaN
zurich	0.0	0.000000	NaN
zurich schnatzlet	0.0	0.000000	NaN
zwar	0.0	0.000000	NaN
zwar billig	0.0	0.000000	NaN
zwieback	0.0	0.000000	NaN
zwieback toast	0.0	0.000000	NaN
zwiebeck	0.0	0.000000	NaN
zwiebeck toast	0.0	0.000000	NaN
zydeco	0.0	0.000000	NaN
zydeco saturday	0.0	0.000000	NaN
zzzzzs	0.0	0.000000	NaN
zzzzzs larg	0.0	0.000000	NaN
zzzzzz	0.0	0.000000	NaN

	0	1	2
zzzzzz say	0.0	0.000000	NaN
zzzzzzzz	0.0	0.000000	NaN
zzzzzzzz high	0.0	0.000000	NaN
zzzzzzzzzz	0.0	0.000000	NaN
zzzzzzzzzzzz	0.0	0.000000	NaN
zzzzzzzzzzzz final	0.0	0.000000	NaN
zzzzzzzzzzzzzz	0.0	0.000000	NaN
çay	0.0	0.000000	NaN

1013943 rows x 3 columns

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

(12588, 3)

12588 features out of 1013943 shows percentage change > 30 post pertubation test i.e 1.24%

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

```
1 feature_names = np.array(vocabulary.get_feature_names())
2 sorted_coef_index = LR_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,best,love,delici,good,perfect,excel,high recommend,favorit,wonder,nice,find,tasti,amaz,keep,always,enjoy,addict,thank,easi,

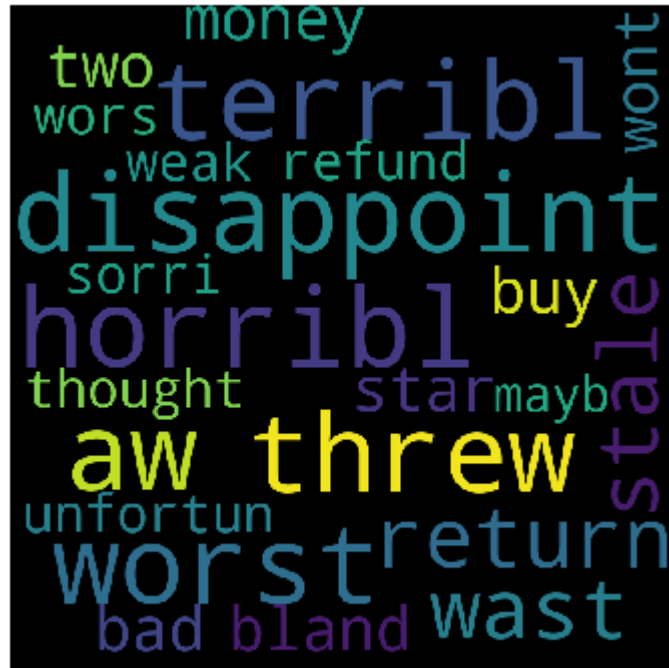
```
1 n=feature_names[sorted_coef_index[: -21: -1]]
2
3 sn = ""
4 for i in n:
5     sn += str(i) + ", "
6 print(sn)
```

disappoint, worst, terrible, aw, threw, horrible, return, stale, waste money, two star, won't buy, bad, bland, unfortunate, thought, refund, sorry, worse, weak, maybe,

```
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 \*\*\*\*\*





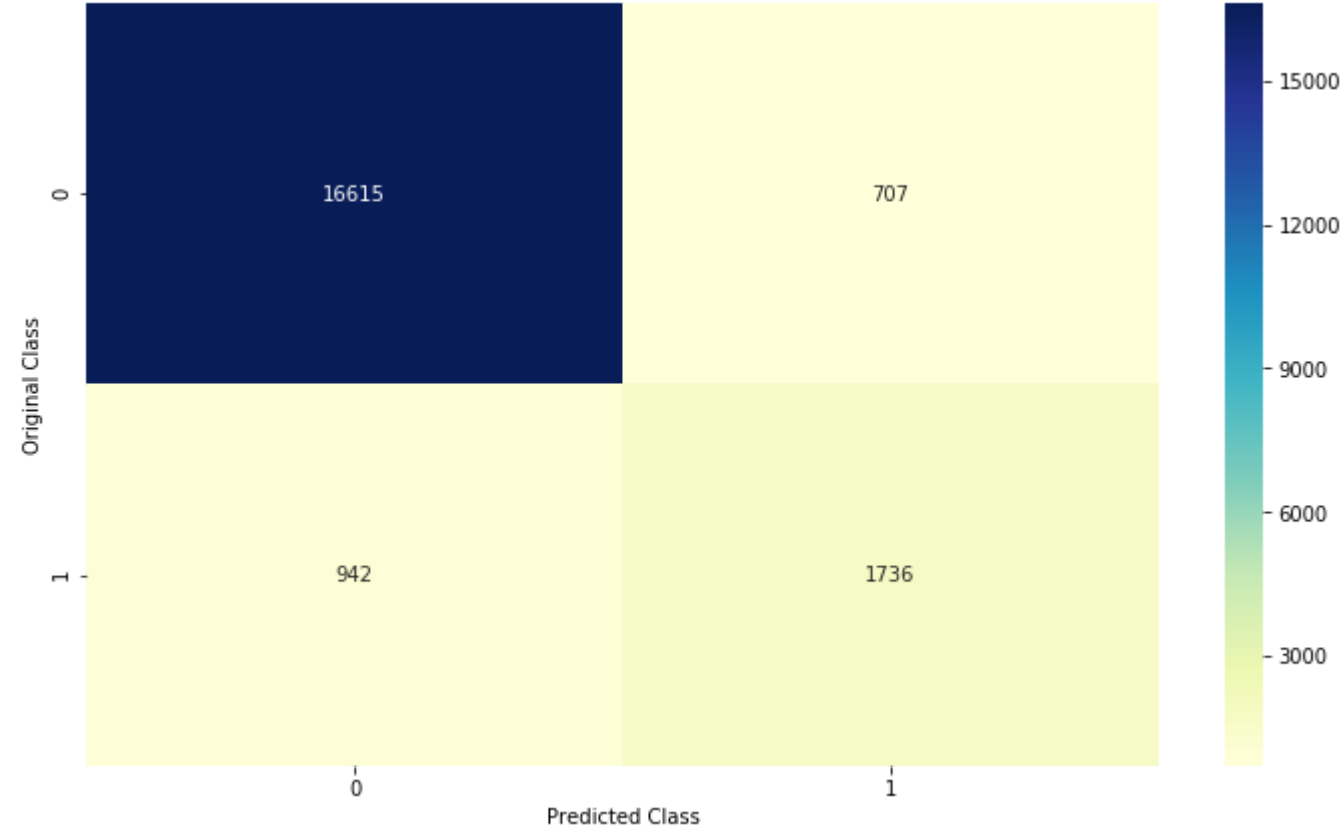
\*\*\*\*\* Top 20 Positive words \*\*\*\*\*



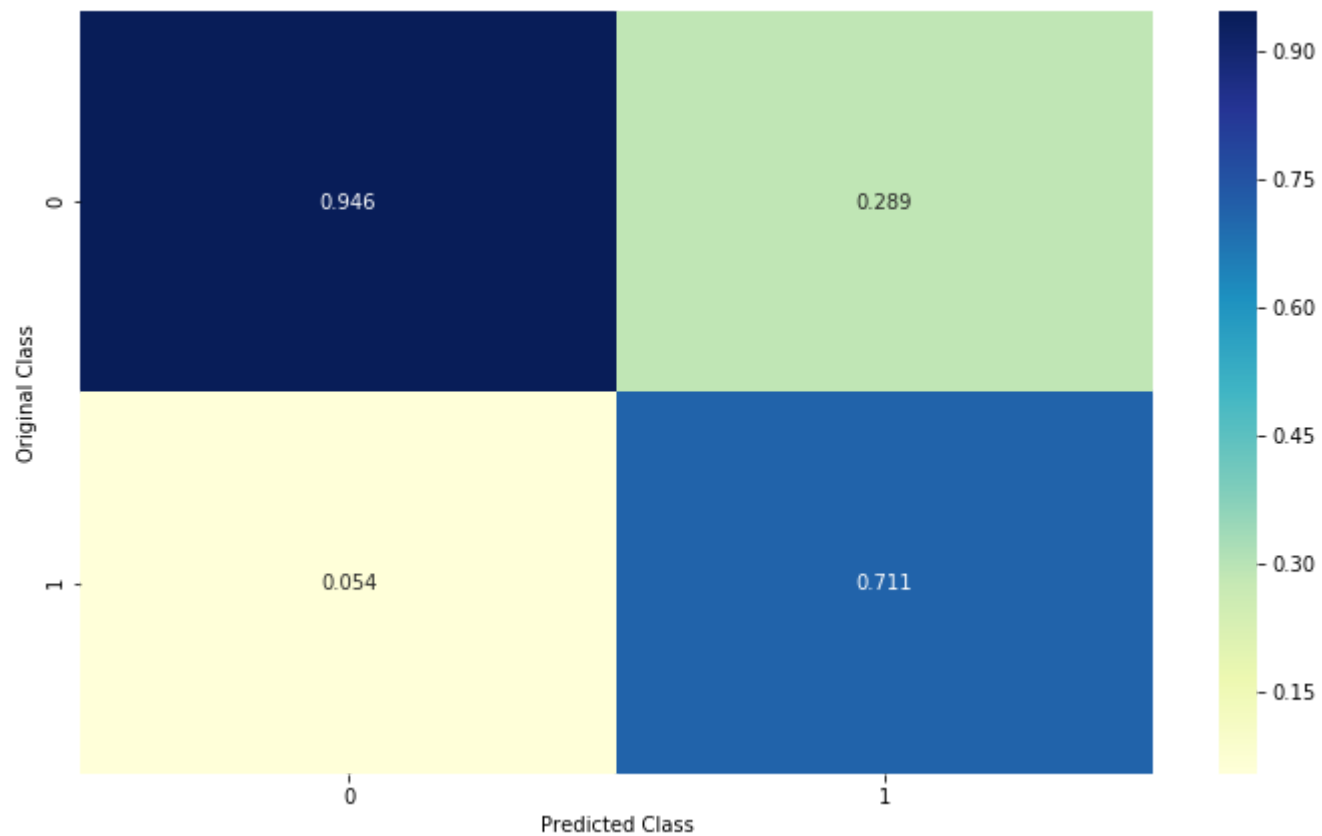
```
1 #Confusion matrix
2 C = confusion_matrix(y_test, pred_tfidf)
3 A = (((C.T)/(C.sum(axis=1))).T)
4 B = (C/C.sum(axis=0))
5 labels = [0,1]
```

```
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()
```

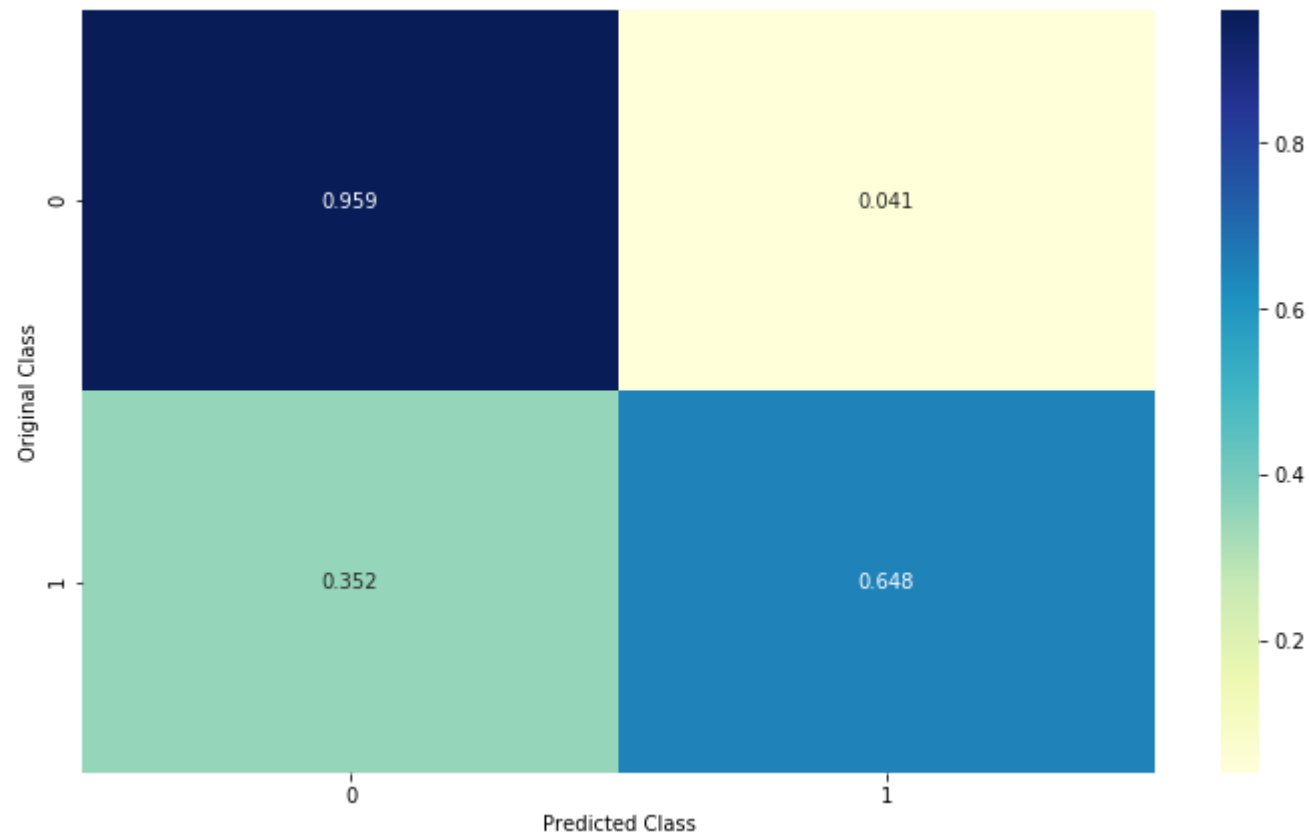
----- Confusion matrix -----



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



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



What if we decrease the value of C? How does it effect Sparsity and Error?

$C=1/\lambda$

```
1 #Checking current non-zero elements
2 a=np.array(LR_optimal.coef_)
3 print(np.count_nonzero(a))
```

31591

```
1  #Decreasing value of c
2  #{'C': 0.01, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.001,class_weight='balanced')
4
5  # fitting the model
6  LR_optimal.fit(Tfidf_x_train, y_train)
7
8  # predict the response
9  pred_tfidf = LR_optimal.predict(Tfidf_x_test)
10
11 # evaluate accuracy
12 f1_score = f1_score(y_test, pred_tfidf)
```

```
1  #Checking current non-zero elements
2  a=np.array(LR_optimal.coef_)
3  print(np.count_nonzero(a))
```

577

```
1  #Decreasing value of c
2  #{'C': 0.01, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.0001,class_weight='balanced')
4
5  # fitting the model
6  LR_optimal.fit(Tfidf_x_train, y_train)
7
8  # predict the response
9  pred_tfidf = LR_optimal.predict(Tfidf_x_test)
10
11 # evaluate accuracy
12 f1_score = f1_score(y_test, pred_tfidf)
```

```

1  #Checking current non-zero elements
2  a=np.array(LR_optimal.coef_)
3  print(np.count_nonzero(a))

```

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', '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(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
4 Standard=StandardScaler()
5 sent_vectors_train_avgw2v = Standard.fit_transform(sent_vectors_train_avgw2v)
6 sent_vectors_test_avgw2v = Standard.transform(sent_vectors_test_avgw2v)
7
8 print(sent_vectors_train_avgw2v.shape)
9 print(sent_vectors_test_avgw2v.shape)
```

(80000, 50)

(20000, 50)

### Fitting grid search on Avg-W2V

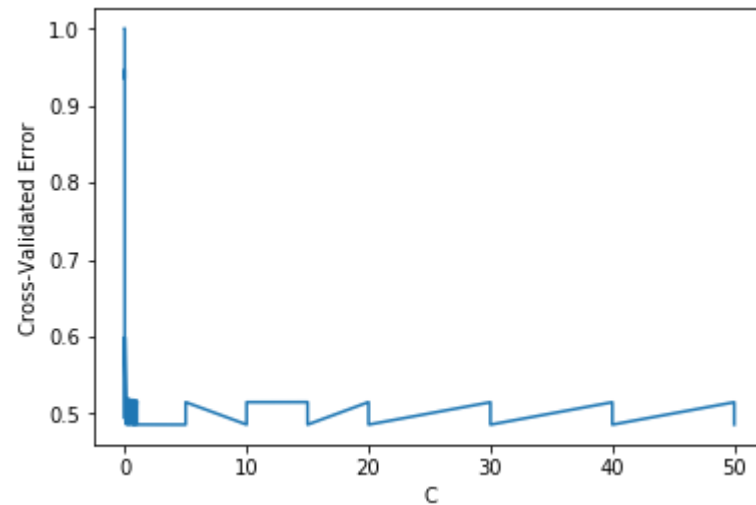
```
1 grid.fit(sent_vectors_train_avgw2v, y_train)
2
3 # examine the best model
4 print(grid.best_score_)
5 print(grid.best_params_)
```

0.5141808414956328

{'C': 5, 'class\_weight': 'balanced'}

```
1 #Plotting C 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
2 #{'C': 5, 'class_weight': 'balanced'}
3 LR_optimal=LogisticRegression(penalty='l1',C=5,class_weight='balanced')
4
5 # fitting the model
6 LR_optimal.fit(sent_vectors_train_avgw2v, y_train)
7
8 # predict the response
9 pred_avg_w2v = LR_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,LR_optimal.predict(sent_
16 print("The overall f1_score for the Test Data is : ", metrics.f1_score(y_test,pred_avg_w2v))
```

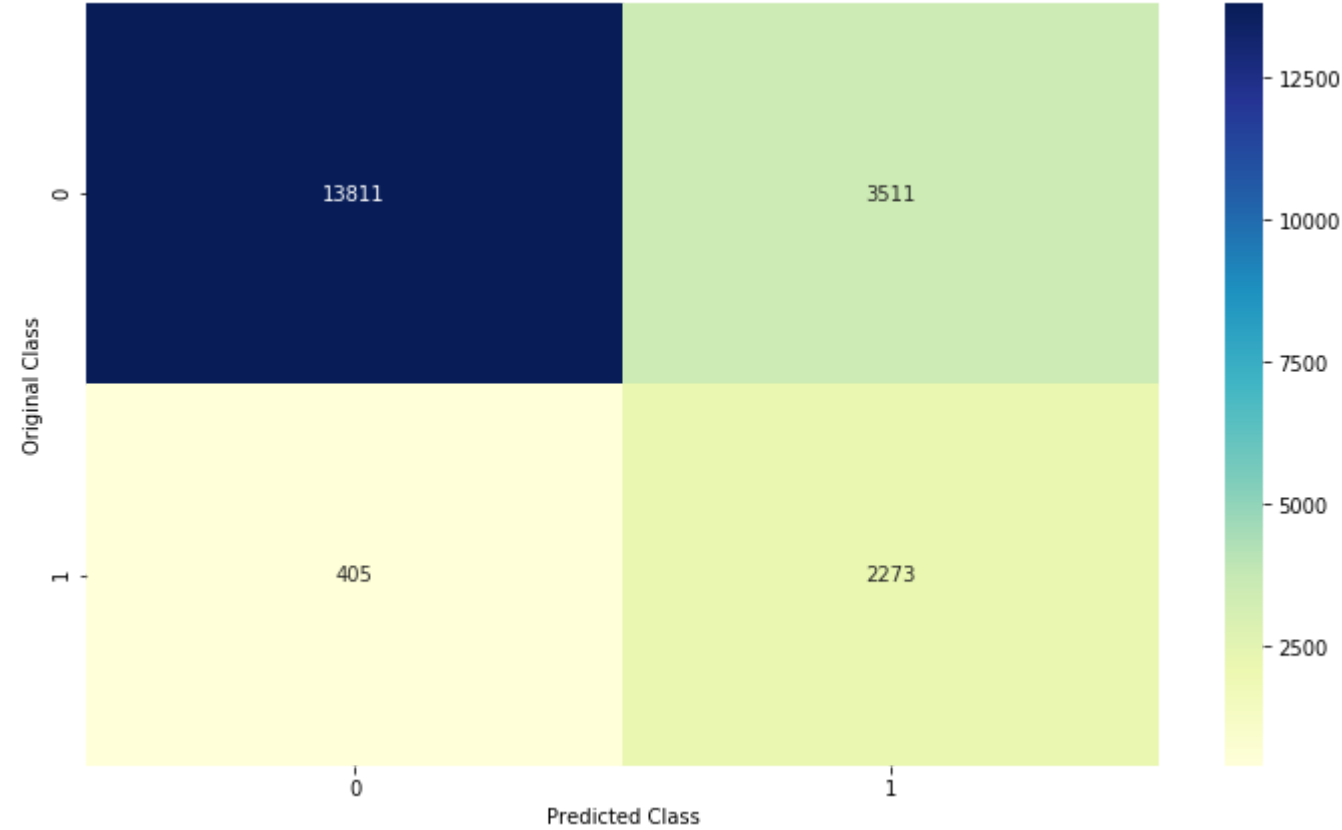
The overall f1\_score for the Train Data is : 0.5156973751930005

The overall f1\_score for the Test Data is : 0.5372252422595132

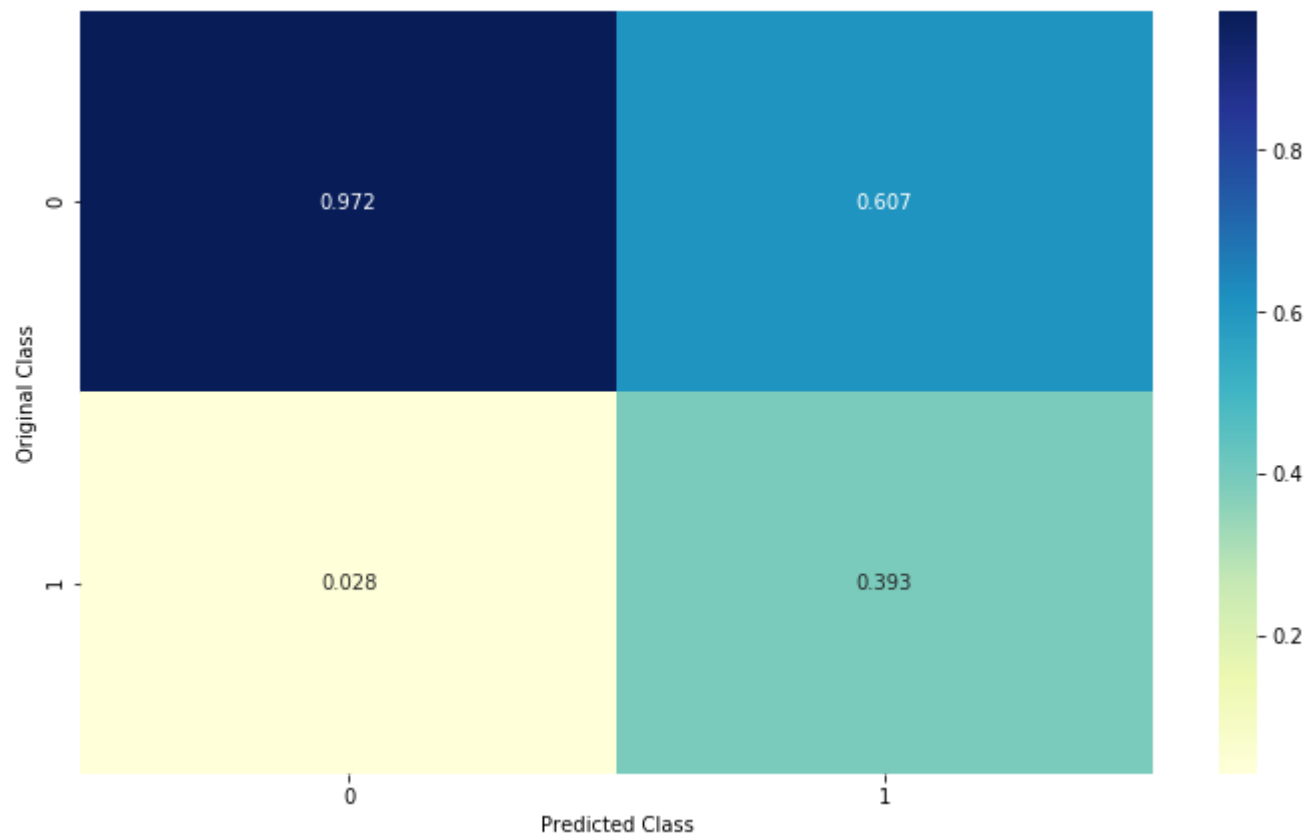
```
1 #Confusion matrix
2 C = confusion_matrix(y_test, pred_avg_w2v)
3 A =(((C.T)/(C.sum(axis=1))).T)
4 B =(C/C.sum(axis=0))
5 labels = [0,1]
```

```
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 -----

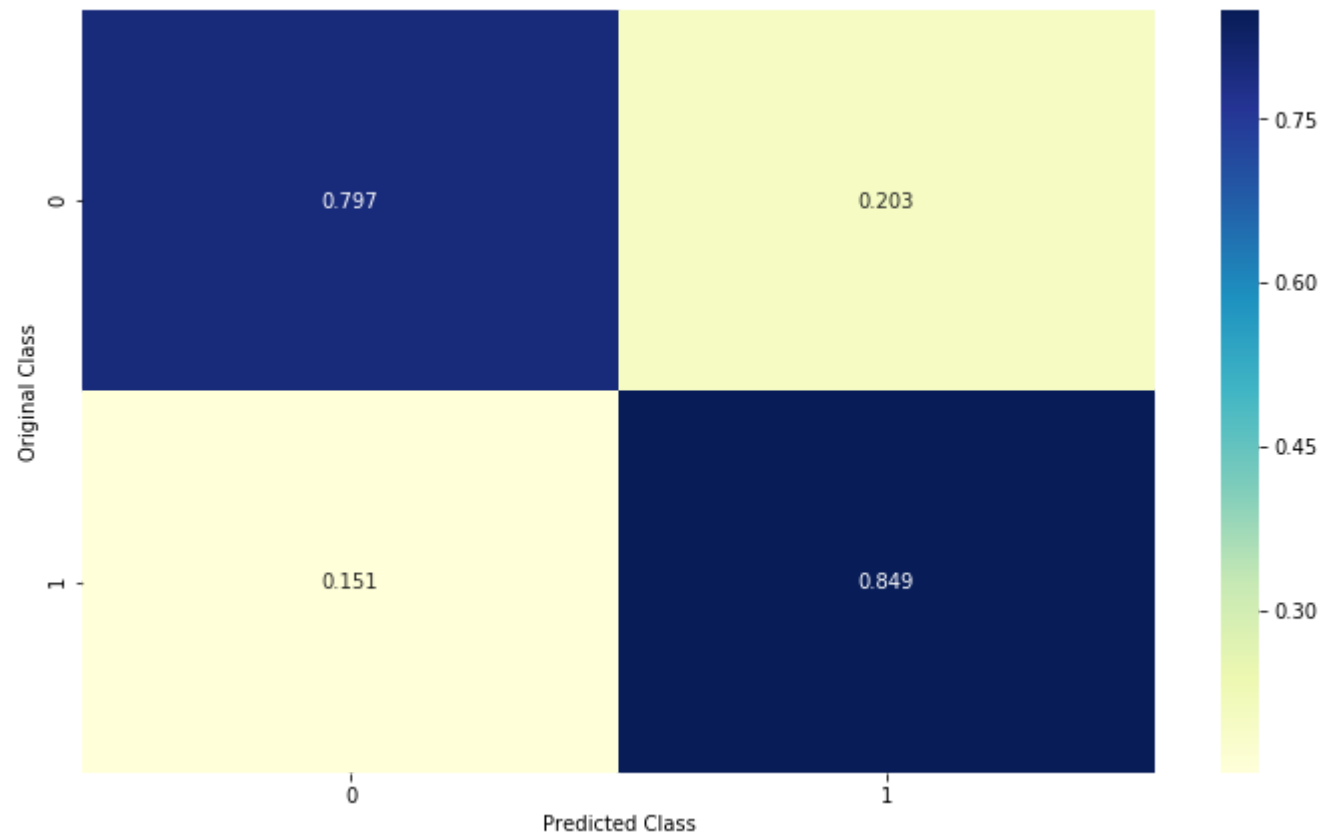


----- Precision matrix (Column Sum=1) -----



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





What if we decrease the value of C? How does it effect Sparsity and Error?

$C=1/\lambda$

```
1 #Checking current non-zero elements
2 a=np.array(LR_optimal.coef_)
3 print(np.count_nonzero(a))
```

50

```
1  #Decrese the value of C
2  #{'C': 5, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.1,class_weight='balanced')
4
5  # fitting the model
6  LR_optimal.fit(sent_vectors_train_avg2v, y_train)
7
8  # predict the response
9  pred_avg_w2v = LR_optimal.predict(sent_vectors_test_avg2v)
10
11 # evaluate f1_score
12 f1_score = f1_score(y_test, pred_avg_w2v)
```

```
1  a=np.array(LR_optimal.coef_)
2  print(np.count_nonzero(a))
```

50

```
1  #Decrese the value of C
2  #{'C': 5, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.001,class_weight='balanced')
4
5  # fitting the model
6  LR_optimal.fit(sent_vectors_train_avg2v, y_train)
7
8  # predict the response
9  pred_avg_w2v = LR_optimal.predict(sent_vectors_test_avg2v)
10
11 # evaluate f1_score
12 f1_score = f1_score(y_test, pred_avg_w2v)
```

```
1 a=np.array(LR_optimal.coef_)
2 print(np.count_nonzero(a))
```

33

```
1 #Decrese the value of C
2 #{'C': 5, 'class_weight': 'balanced'}
3 LR_optimal=LogisticRegression(penalty='l1',C=0.0001,class_weight='balanced')
4
5 # fitting the model
6 LR_optimal.fit(sent_vectors_train_avg2v, y_train)
7
8 # predict the response
9 pred_avg_w2v = LR_optimal.predict(sent_vectors_test_avg2v)
10
11 # evaluate f1_score
12 f1_score = f1_score(y_test, pred_avg_w2v)
```

```
1 a=np.array(LR_optimal.coef_)
2 print(np.count_nonzero(a))
```

3

## Tf-idf W2V

```
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_)))
```





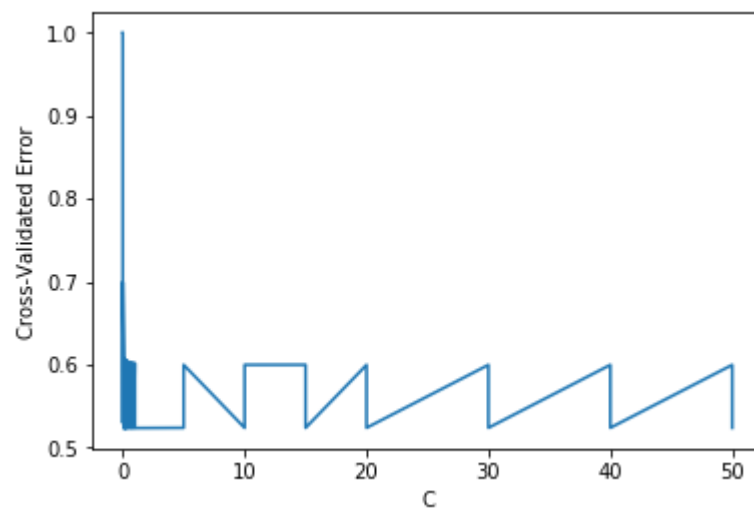
```
1 grid.fit(tfidf_w2v_sent_vectors_train, y_train)
2
3 # examine the best model
4 print(grid.best_score_)
5 print(grid.best_params_)
```

0.47690283662194616

{'C': 0.1, 'class\_weight': 'balanced'}

```
1 #Plotting C 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  #{'C': 0.1, 'class_weight': 'balanced'}
2  LR_optimal=LogisticRegression(penalty='l1',C=0.1,class_weight=None)
3
4  # fitting the model
5  LR_optimal.fit(tfidf_w2v_sent_vectors_train, y_train)
6
7  # predict the response
8  pred_tfidf_w2v_sent_vectors_test = LR_optimal.predict(tfidf_w2v_sent_vectors_test)
9
10 # evaluate f1_score
11 f1_score = f1_score(y_test, pred_tfidf_w2v_sent_vectors_test)
12
13 # Train & Test Error
14 print("The overall f1_score for the Train Data is : ", metrics.f1_score(y_train,LR_optimal.predict(tfidf_w2v_sent_vectors_train)))
15 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 : 0.41298833079654995

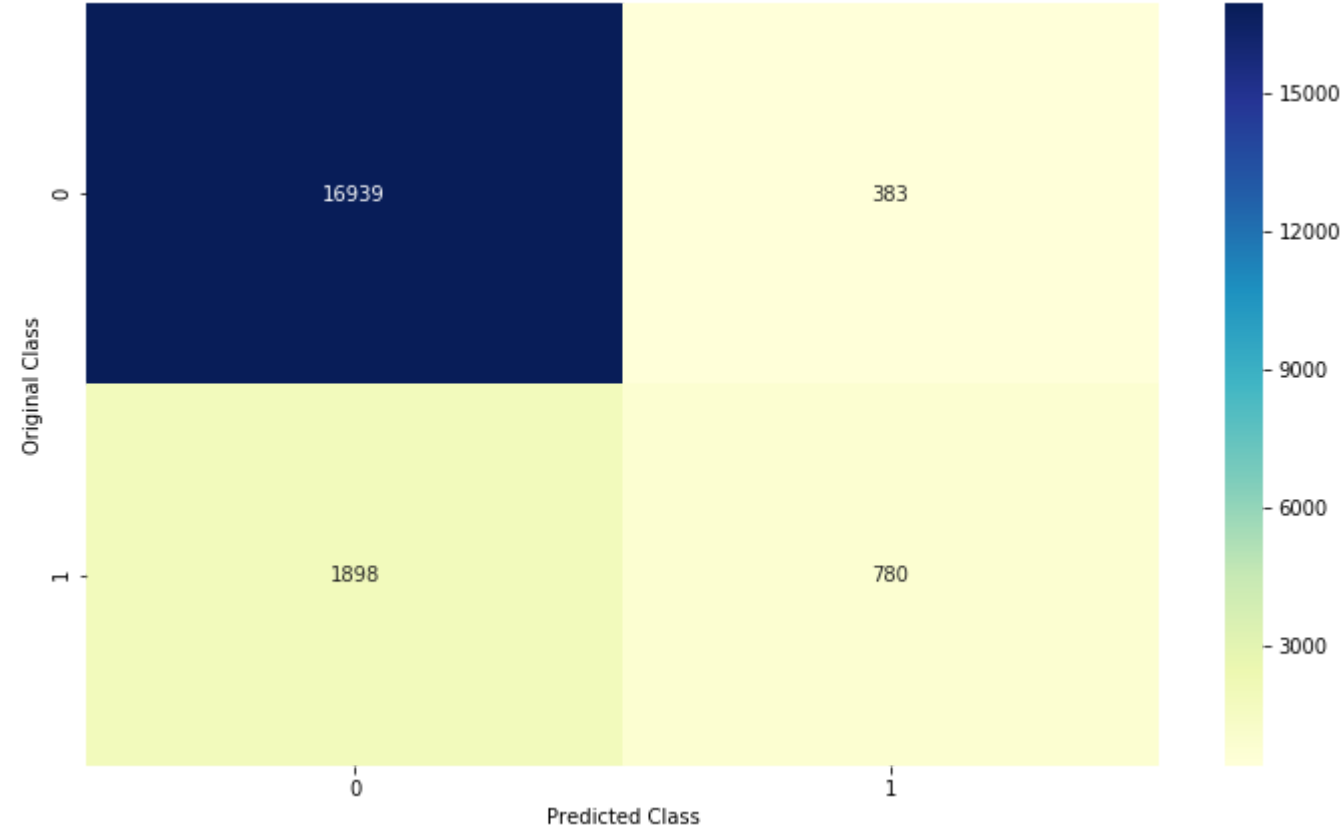
The overall f1\_score for the Test Data is : 0.40614423327258525

```
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]
```

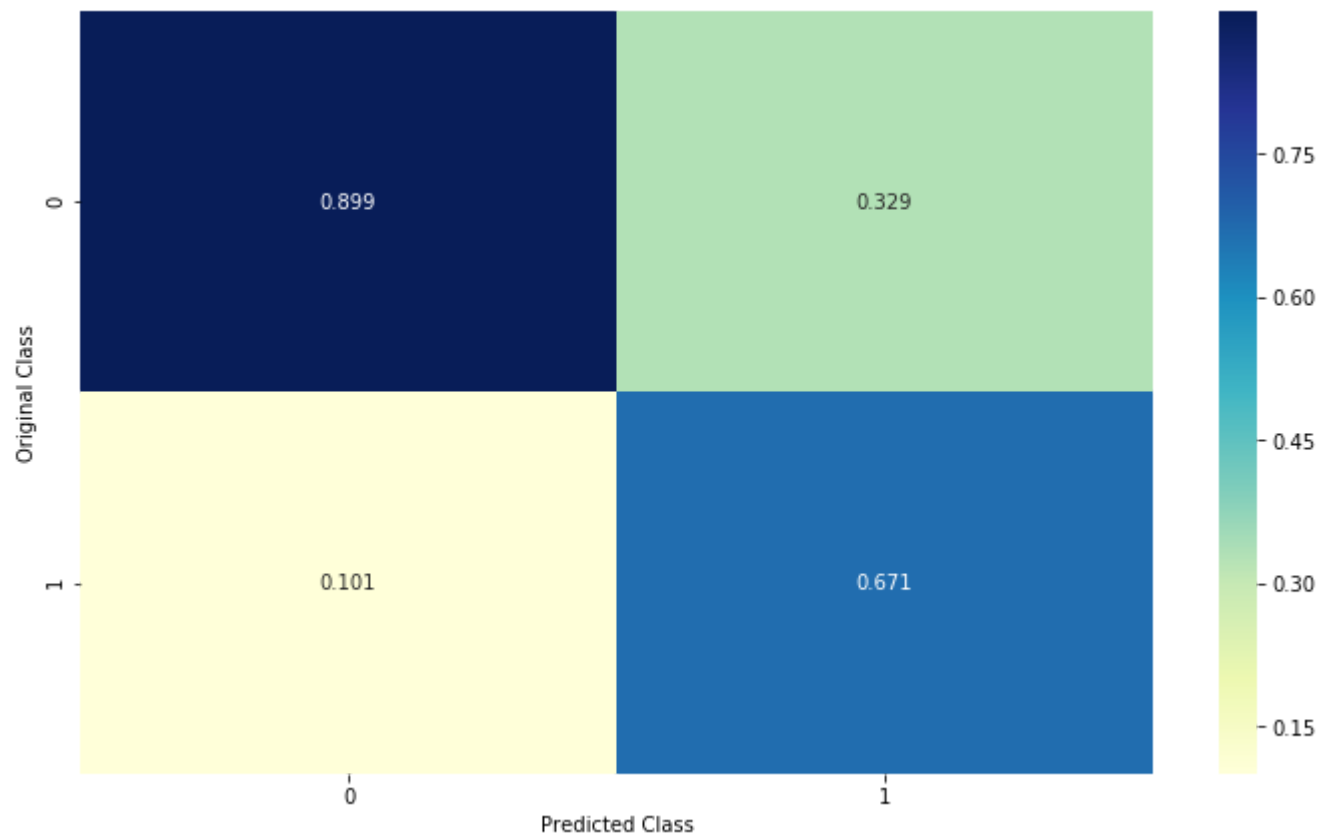


```
1  print("-"*20, "Confusion matrix", "-"*20)
2  plt.figure(figsize=(12,7))
3  sns.heatmap(C, annot=True, cmap="YlGnBu", fmt="g", xticklabels=labels, yticklabels=labels)
4  plt.xlabel('Predicted Class')
5  plt.ylabel('Original Class')
6  plt.show()
7
8  print("-"*20, "Precision matrix (Column Sum=1)",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

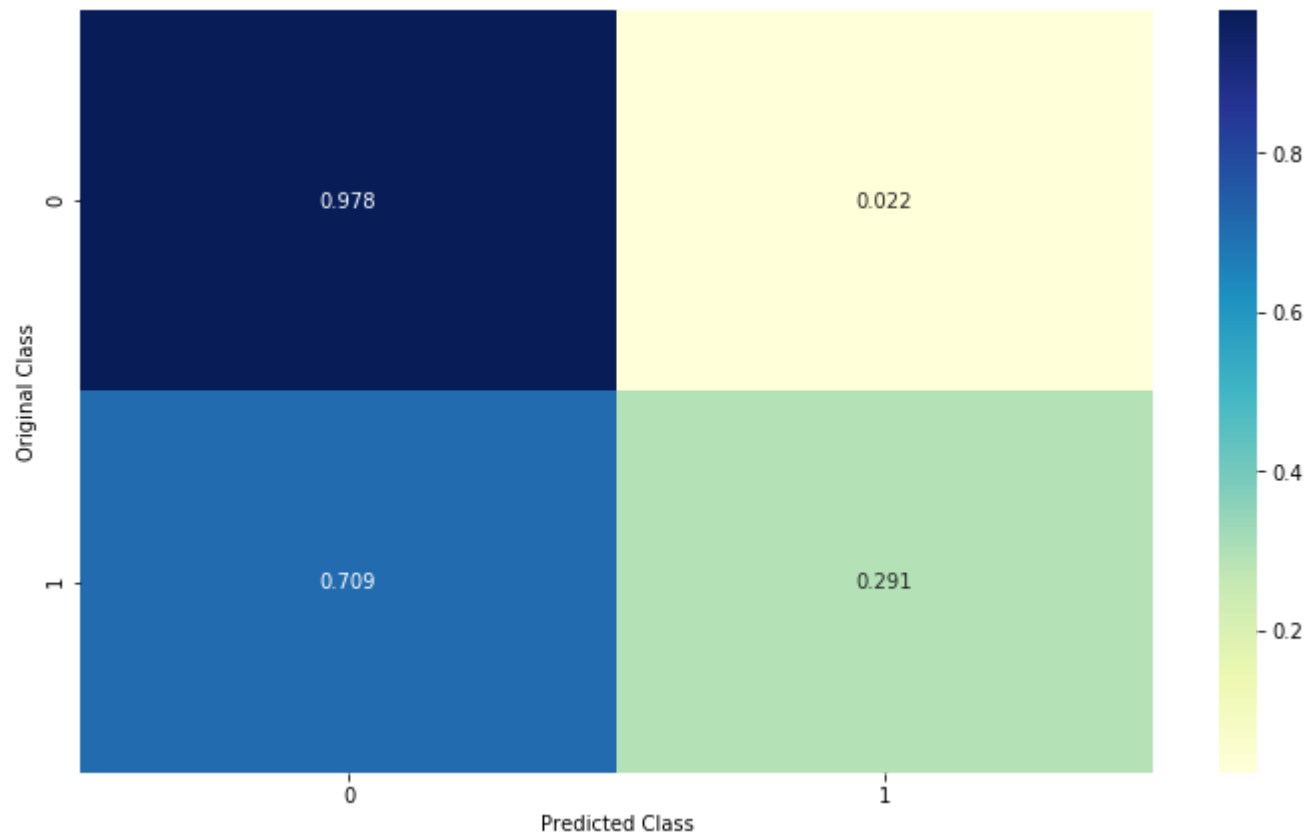
----- Confusion matrix -----



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



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



What if we decrease the value of C? How does it effect Sparsity and Error?

$C=1/\lambda$

```
1 #Checking current non-zero elements
2 a=np.array(LR_optimal.coef_)
3 print(np.count_nonzero(a))
```

50

```
1  #Decreasing value of C
2  #{'C': 0.1, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.001,class_weight=None)
4
5  # fitting the model
6  LR_optimal.fit(tfidf_w2v_sent_vectors_train, y_train)
7
8  # predict the response
9  pred_tfidf_w2v_sent_vectors_test = LR_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
```

```
1  #Checking current non-zero elements
2  a=np.array(LR_optimal.coef_)
3  print(np.count_nonzero(a))
```

19

```
1  #Decreasing value of C
2  #{'C': 0.1, 'class_weight': 'balanced'}
3  LR_optimal=LogisticRegression(penalty='l1',C=0.0001,class_weight=None)
4
5  # fitting the model
6  LR_optimal.fit(tfidf_w2v_sent_vectors_train, y_train)
7
8  # predict the response
9  pred_tfidf_w2v_sent_vectors_test = LR_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
```

```

1  #Checking current non-zero elements
2  a=np.array(LR_optimal.coef_)
3  print(np.count_nonzero(a))

```

0

Reporting f1\_score for above featurization with L1 regularizer

```

1  from prettytable import PrettyTable

```

```

1  x=PrettyTable()
2  x.field_names = ["Model", "Bow", "Tfidf", "Avg-W2V", "Tfidf-W2V"]
3  x.add_row(["C", 0.01, 0.01, 5, 0.1])
4  x.add_row(["Train f1_score", 0.74, 0.95, 0.51, 0.41])
5  x.add_row(["Test f1_score", 0.67, 0.67, 0.53, 0.40])
6
7  print(x)

```

```

+-----+-----+-----+-----+-----+
|   Model   | Bow | Tfidf | Avg-W2V | Tfidf-W2V |
+-----+-----+-----+-----+-----+
|      C      | 0.01 | 0.01 | 5       | 0.1       |
| Train f1_score | 0.74 | 0.95 | 0.51    | 0.41     |
| Test f1_score  | 0.67 | 0.67 | 0.53    | 0.4      |
+-----+-----+-----+-----+-----+

```