

[1] Amazon Fine Food Reviews Analysis

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:

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

[Q] How to determine if a review is positive or negative?

[Ans] We could use the Score/Rating. A rating of 4 or 5 could be considered a positive review. A review of 1 or 2 could be considered negative. A review of 3 is neutral and ignored. This is an approximate and proxy way of determining the polarity (positivity/negativity) of a review.

[7.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

```
In [1]: %matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os
warnings.filterwarnings("ignore")

%matplotlib inline
# sets the backend of matplotlib to the 'inline' backend:
#With this backend, the output of plotting commands is displayed inline within from
tends like the Jupyter notebook,
#directly below the code cell that produced it. The resulting plots will then also
be stored in the notebook document.

#Functions to save objects for later use and retireve it
import pickle
def savetofile(obj,filename):
    pickle.dump(obj,open(filename+".p","wb"))
def openfromfile(filename):
    temp = pickle.load(open(filename+".p","rb"))
    return temp
```

```
C:\Users\Sai charan\Anaconda3\lib\site-packages\ipykernel\parentpoller.py:116: U
serWarning: Parent poll failed. If the frontend dies,
the kernel may be left running. Please let us know
about your system (bitness, Python, etc.) at
ipython-dev@scipy.org
ipython-dev@scipy.org"""
C:\Users\Sai charan\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarnin
g: detected Windows; aliasing chunkize to chunkize_serial
warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

Exploratory Data Analysis

[7.1.2] Data Cleaning: Deduplication

It is observed (as shown in the table below) that the reviews data had many duplicate entries. Hence it was necessary to remove duplicates in order to get unbiased results for the analysis of the data. Following is an example:

```
In [2]: #Using sqlite3 to retrieve data from sqlite file

con = sqlite3.connect("final.sqlite") #Loading Cleaned/ Preprocessed text that we did
in Text Preprocessing

#Using pandas functions to query from sql table
final = pd.read_sql_query("""
SELECT * FROM Reviews order by Time
""", con)

#Reviews is the name of the table given
#Taking only the data where score != 3 as score 3 will be neutral and it won't help
us much
final.head()
```

Out[2]:

	index	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	138706	150524	0006641040	ACITT7DI6IDDL	shari zychinski	0	
1	138683	150501	0006641040	AJ46FKXOVC7NR	Nicholas A Mesiano	2	
2	417839	451856	B00004CXX9	AIUWLEQ1ADEG5	Elizabeth Medina	0	
3	346055	374359	B00004CI84	A344SMIA5JECGM	Vincent P. Ross	1	
4	417838	451855	B00004CXX9	AJH6LUC1UT1ON	The Phantom of the Opera	0	

```
In [3]: final['Score'] = final['Score'].replace('positive',1)
final['Score'] = final['Score'].replace('negative',0)
```

```
In [4]: final.duplicated(subset={"UserId", "ProfileName", "Time", "Text"}).value_counts()
```

```
Out[4]: False      364171
dtype: int64
```

```
In [5]: final = final.drop_duplicates(subset={"UserId", "ProfileName", "Time", "Text"}, keep="first")
```

```
In [6]: final = final[final.HelpfulnessNumerator <= final.HelpfulnessDenominator]
print("Size of data", final['Id'].size, " rows ")
```

Size of data 364171 rows

```
In [7]: import re #Regex (Regular Expr Operations)
#string = r"sdfsdfd" :- r is for raw string as Regex often uses \ backslashes(\w),
so they are often raw strings(r'\d')

#####Function to remove html tags from data
def striphtml(data):
    p = re.compile('<.*?>') #Find this kind of pattern
    # print(p.findall(data)) #List of strings which follow the regex pattern
    return p.sub('', data) #Substitute nothing at the place of strings which matched
the patterns

striphtml('<a href="foo.com" class="bar">I Want This <b>text!</b></a><>')
```

Out[7]: 'I Want This text!'

```
In [8]: def strippunc(data):
    p = re.compile(r'[? ! | \ | " | # | . | , | ( | \ | / | ~ | % | * ]')
    return p.sub('', data)
strippunc("fsd*?~, , ( sdfsdfdsfvv)#")
```

Out[8]: 'fsd sdfsdfdsfvv'

```
In [9]: from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

stop = stopwords.words('english') #All the stopwords in English language
#excluding some useful words from stop words list as we doing sentiment analysis
excluding = ['against', 'not', 'don', "don't", 'ain', 'aren', "aren't", 'couldn', "cou
ldn't", 'didn', "didn't",
            'doesn', "doesn't", 'hadn', "hadn't", 'hasn', "hasn't", 'haven', "have
n't", 'isn', "isn't",
            'mightn', "mightn't", 'mustn', "mustn't", 'needn', "needn't", 'shouldn'
, "shouldn't", 'wasn',
            'wasn't', 'weren', "weren't", 'won', "won't", 'wouldn', "wouldn't"]
stop = [words for words in stop if words not in excluding]
print(stop)
```

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "
you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'hi
m', 'his', 'himself', 'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'i
ts', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what', 'which',
'who', 'whom', 'this', 'that', "that'll", 'these', 'those', 'am', 'is', 'are', '
was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does
', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as',
'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'between', 'into', '
through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up', 'do
wn', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'onc
e', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each'
, 'few', 'more', 'most', 'other', 'some', 'such', 'no', 'nor', 'only', 'own', 's
ame', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'should', "s
hould've", 'now', 'd', 'll', 'm', 'o', 're', 've', 'y', 'ma', 'shan', "shan't"]
```

```
In [10]: from nltk.stem import SnowballStemmer
snow = SnowballStemmer('english') #initialising the snowball stemmer
print("Stem/Root words of the some of the words using SnowBall Stemmer:")
print(snow.stem('tasty'))
print(snow.stem('tasteful'))
print(snow.stem('tastiest'))
print(snow.stem('delicious'))
print(snow.stem('amazing'))
print(snow.stem('amaze'))
print(snow.stem('initialize'))
print(snow.stem('fabulous'))
print(snow.stem('Honda City'))
print(snow.stem('unpleasant'

))
```

Stem/Root words of the some of the words using SnowBall Stemmer:

tasti
tast
tastiest
delici
amaz
amaz
initi
fabul
honda c
unpleas

```

In [11]: str1=' '
final_string=[]
all_positive_words=[] # store words from +ve reviews here
all_negative_words=[] # store words from -ve reviews here.
s=''
for sent in final['Text'][2:3].values: #Running only for 2nd review
    filtered_sentence=[]
    print(sent) #Each review
    sent=striphtml(sent) # remove HTML tags
    sent=strippunc(sent) # remove Punctuation Symbols
    print(sent.split())
    for w in sent.split():
        print("=====>",w)
        if((w.isalpha()) and (len(w)>2)): #If it is a numerical value or character o
f lenght less than 2
            if(w.lower() not in stop): # If it is a stopword
                s=(snow.stem(w.lower())).encode('utf8') #Stemming the word using Sn
owBall Stemmer
                print("Selected: Stem Word->",s)
                filtered_sentence.append(s)
            else:
                print("Eliminated as it is a stopword")
                continue
        else:
            print("Eliminated as it is a numerical value or character of lenght les
s than 2")
            continue
    # print(filtered_sentence)
    str1 = b" ".join(filtered_sentence) #final string of cleaned words

    final_string.append(str1)
    print("*****")
)
print("Finally selected words from the review:\n",final_string)

```

```

Beetlejuice is a well written movie ..... everything about it is excellent! From
the acting to the special effects you will be delighted you chose to view this m
ovie.
['Beetlejuice', 'is', 'a', 'well', 'written', 'movie', 'everything', 'about', 'i
t', 'is', 'excellent', 'From', 'the', 'acting', 'to', 'the', 'special', 'effects
', 'you', 'will', 'be', 'delighted', 'you', 'chose', 'to', 'view', 'this', 'movi
e']
=====> Beetlejuice
Selected: Stem Word-> b'beetlejuic'
=====> is
Eliminated as it is a numerical value or character of lenght less than 2
=====> a
Eliminated as it is a numerical value or character of lenght less than 2
=====> well
Selected: Stem Word-> b'well'
=====> written
Selected: Stem Word-> b'written'
=====> movie
Selected: Stem Word-> b'movi'
=====> everything
Selected: Stem Word-> b'everyth'
=====> about
Eliminated as it is a stopword
=====> it
Eliminated as it is a numerical value or character of lenght less than 2
=====> is
Eliminated as it is a numerical value or character of lenght less than 2
=====> excellent
Selected: Stem Word-> b'excel'
=====> From
Eliminated as it is a stopword
=====> the
Eliminated as it is a stopword
=====> acting
Selected: Stem Word-> b'act'
=====> to
Eliminated as it is a numerical value or character of lenght less than 2
=====> the
Eliminated as it is a stopword
=====> special
Selected: Stem Word-> b'special'
=====> effects
Selected: Stem Word-> b'effect'
=====> you
Eliminated as it is a stopword
=====> will
Eliminated as it is a stopword
=====> be
Eliminated as it is a numerical value or character of lenght less than 2
=====> delighted
Selected: Stem Word-> b'delight'
=====> you
Eliminated as it is a stopword
=====> chose
Selected: Stem Word-> b'chose'
=====> to
Eliminated as it is a numerical value or character of lenght less than 2
=====> view
Selected: Stem Word-> b'view'
=====> this
Eliminated as it is a stopword
=====> movie
Selected: Stem Word-> b'movi'
*****

```

```

In [12]: %%time
# Code takes a while to run as it needs to run on around 500k sentences.
i=0
str1=' '
final_string=[]
all_positive_words=[] # store words from +ve reviews here
all_negative_words=[] # store words from -ve reviews here.
s=''
for sent in final['Text'].values:
    filtered_sentence=[]
    # print(sent) #Each review
    sent=striphtml(sent)# remove HTML tags
    sent=strippunc(sent)# remove Punctuation Symbols
    # print(sent.split())
    for w in sent.split():
        # print("=====>",w)
        if((w.isalpha()) and (len(w)>2)):#If it is a numerical value or character o
f lenght less than 2
            if(w.lower() not in stop):# If it is a stopword
                s=(snow.stem(w.lower())).encode('utf8') #Stemming the word using Sn
owBall Stemmer
                #encoding as byte-string/utf-8
            # print("Selected: Stem Word->",s)
            filtered_sentence.append(s)
            if (final['Score'].values)[i] == 'Positive':
                all_positive_words.append(s) #list of all words used to describ
e positive reviews
            if(final['Score'].values)[i] == 'Negative':
                all_negative_words.append(s) #list of all words used to describ
e negative reviews reviews
            else:
                # print("Eliminated as it is a stopword")
                continue
        else:
            # print("Eliminated as it is a numerical value or character of lenght 1
ess than 2")
            continue
    # print(filtered_sentence)
    str1 = b" ".join(filtered_sentence) #final string of cleaned words
    #encoding as byte-string/utf-8

    final_string.append(str1)
    # print("*****")
    # print("Finally selected words from the review:\n",final_string)
    i+=1

```

Wall time: 7min 26s


```
In [13]: %%time
# Code takes a while to run as it needs to run on around 500k sentences.
i=0
str1=' '
final_string_nostem=[]
s=''
for sent in final['Text'].values:
    filtered_sentence=[]
    sent=striphtml(sent)# remove HTML tags
    sent=strippunc(sent)# remove Punctuation Symbols
    for w in sent.split():
        if((w.isalpha()) and (len(w)>2)):#If it is a numerical value or character o
f lenght less than 2
            if(w.lower() not in stop):# If it is a stopword
                s=w.lower().encode('utf8') #encoding as byte-string/utf-8
            else:
                continue
        else:
            continue
    str1 = b" ".join(filtered_sentence)
    final_string_nostem.append(str1)
    i+=1
```

Wall time: 47.7 s

```
In [14]: #Adding a column of CleanedText which displays the data after pre-processing of the
review
final['CleanedText']=final_string
final['CleanedText_NoStem']=final_string_nostem
final.head(3)
```

Out[14]:

	index	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	138706	150524	0006641040	ACITT7DI6IDDL	shari zychinski	0	
1	138683	150501	0006641040	AJ46FKXOVC7NR	Nicholas A Mesiano	2	
2	417839	451856	B00004CXX9	AIUWLEQ1ADEG5	Elizabeth Medina	0	

In []:

```
In [15]: final.sort_values('Time', inplace=True)
         final.head(10)
```

Out[15]:

	index	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenomin
0	138706	150524	0006641040	ACITT7DI6IDDL	shari zychinski	0	
1	138683	150501	0006641040	AJ46FKXOVC7NR	Nicholas A Mesiano	2	
2	417839	451856	B00004CXX9	AIUWLEQ1ADEG5	Elizabeth Medina	0	
3	346055	374359	B00004CI84	A344SMIA5JECGM	Vincent P. Ross	1	
4	417838	451855	B00004CXX9	AJH6LUC1UT1ON	The Phantom of the Opera	0	
5	346116	374422	B00004CI84	A1048CYU0OV4O8	Judy L. Eans	2	
6	346041	374343	B00004CI84	A1B2IZU1JLZA6	Wes	19	
7	70688	76882	B00002N8SM	A32DW342WBJ6BX	Buttersugar	0	
8	346141	374450	B00004CI84	ACJR7EQF9S6FP	Jeremy Robertson	2	
10	417883	451903	B00004CXX9	A2DEE7F9XKP3ZR	jerome	0	

```
In [16]: final=final[:100000]
```

```
In [17]: savetofile(final, "sample_lr")
```

```
In [18]: final = openfromfile("sample_lr")
```

In []:

[7.2.2] Bag of Words (BoW)

```
In [19]: from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn import preprocessing

#Breaking into Train and test
X_train, X_test, y_train, y_test = train_test_split(final['CleanedText'].values, final['Score'].values, test_size=0.3, shuffle=False)

#Text -> Uni gram Vectors
uni_gram = CountVectorizer()
X_train = uni_gram.fit_transform(X_train)
#Normalize Data
X_train = preprocessing.normalize(X_train)

X_test = uni_gram.transform(X_test)
#Normalize Data
X_test = preprocessing.normalize(X_test)
from sklearn.preprocessing import StandardScaler
sc = StandardScaler(with_mean=False)
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
print("Test Data Size: ", X_test.shape)
print("Train Data Size: ", X_train.shape)
```

```
Test Data Size: (30000, 64925)
Train Data Size: (70000, 64925)
```

```
In [20]: #To show how Time Series Split splits the data
from sklearn.model_selection import TimeSeriesSplit
tscv = TimeSeriesSplit(n_splits=10)
for train, cv in tscv.split(X_train):
#     print("%s %s" % (train, cv))
    print(X_train[train].shape, X_train[cv].shape)
```

```
(6370, 64925) (6363, 64925)
(12733, 64925) (6363, 64925)
(19096, 64925) (6363, 64925)
(25459, 64925) (6363, 64925)
(31822, 64925) (6363, 64925)
(38185, 64925) (6363, 64925)
(44548, 64925) (6363, 64925)
(50911, 64925) (6363, 64925)
(57274, 64925) (6363, 64925)
(63637, 64925) (6363, 64925)
```

GridSearchCV

```
In [21]: %time
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

clf = LogisticRegression()
#params we need to try on classifier
param_grid = {'C':[1000,500,100,50,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001,0.0005,0.0001],
              'penalty':['l1','l2']}
tscv = TimeSeriesSplit(n_splits=10) #For time based splitting
gsv = GridSearchCV(clf,param_grid,cv=tscv,verbose=1,scoring='f1_micro')
gsv.fit(X_train,y_train)
savetofile(gsv,"Log Reg/gsv_bi")
print("Best HyperParameter: ",gsv.best_params_)
print("Best Accuracy: %.2f%%"%(gsv.best_score_*100))
```

Wall time: 0 ns

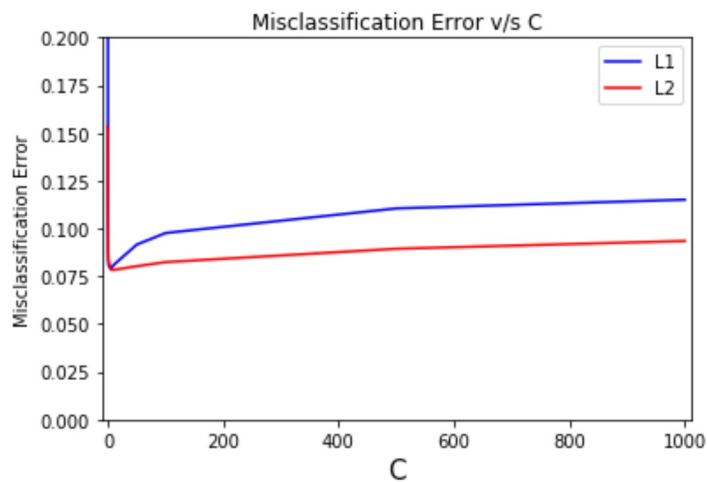
Fitting 10 folds for each of 30 candidates, totalling 300 fits

[Parallel(n_jobs=1)]: Done 300 out of 300 | elapsed: 75.3min finished

Best HyperParameter: {'C': 0.05, 'penalty': 'l1'}

Best Accuracy: 91.79%

```
In [22]: def plot_error_vs_c(gsv):  
    x1=[]  
    y1=[]  
    x2=[]  
    y2=[]  
    for a in gsv.grid_scores_:  
        if (a[0]['penalty']) == 'l1':  
            y1.append(1-a[1])  
            x1.append(a[0]['C'])  
        else:  
            y2.append(1-a[1])  
            x2.append(a[0]['C'])  
    plt.xlim(-10,1010)  
    plt.ylim(0,0.2)  
    plt.xlabel("C",fontsize=15)  
    plt.ylabel("Misclassification Error")  
    plt.title('Misclassification Error v/s C')  
    plt.plot(x1,y1,'b',label="L1")  
    plt.plot(x2,y2,'r',label="L2")  
    plt.legend()  
    plt.show()  
gsv = openfromfile("Log Reg/gsv_uni")  
plot_error_vs_c(gsv)
```



```
In [23]: def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)

    A = ((C.T) / (C.sum(axis=1))).T

    B = (C / C.sum(axis=0))

    plt.figure(figsize=(20,4))

    labels = [1,2]
    #representing A in heatmap format
    cmap=sns.light_palette("blue")
    plt.subplot(1, 3, 1)
    sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabel
s=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Confusion matrix")

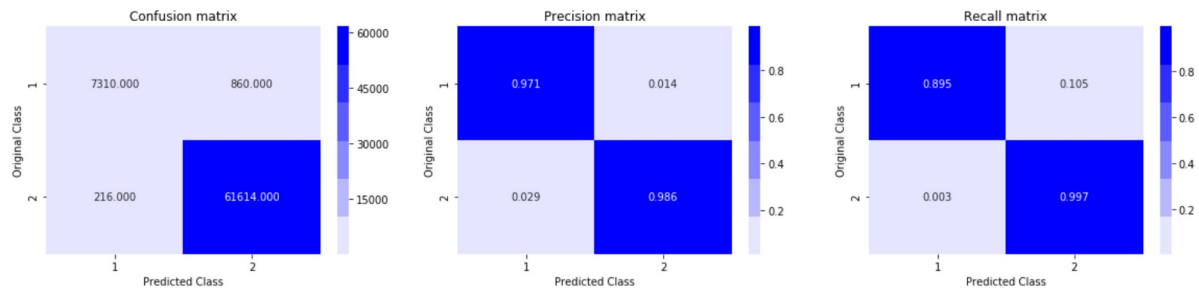
    plt.subplot(1, 3, 2)
    sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabel
s=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Precision matrix")

    plt.subplot(1, 3, 3)
    #representing B in heatmap format
    sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabel
s=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Recall matrix")

    plt.show()
```

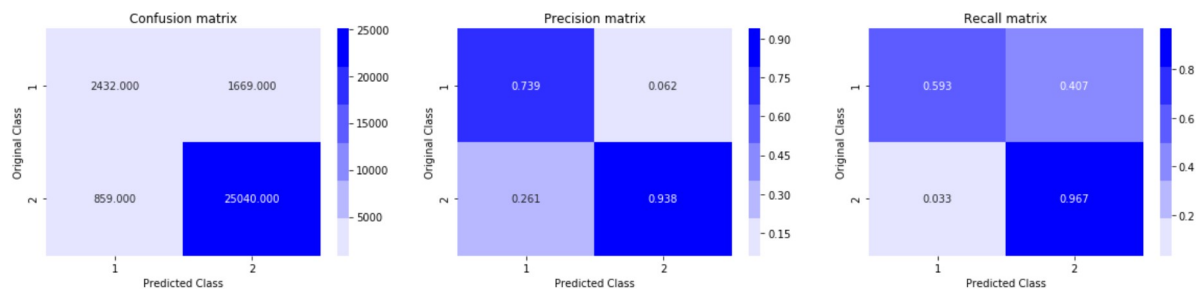
```
In [26]: from sklearn.metrics import accuracy_score
from sklearn.metrics import precision_score, recall_score, f1_score
clf = LogisticRegression(C= 0.05, penalty= 'l1')
clf.fit(X_train, y_train)
y_train_pred = clf.predict(X_train)
y_pred = clf.predict(X_test)
print("Accuracy on train set: %0.3f%%"%(accuracy_score(y_train, y_train_pred)*100))
print("Precision on train set: %0.3f"%(precision_score(y_train, y_train_pred)))
print("Recall on train set: %0.3f"%(recall_score(y_train, y_train_pred)))
print("F1-Score on train set: %0.3f"%(f1_score(y_train, y_train_pred, average='micro'
''))))
print("Non Zero weights:", np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
plot_confusion_matrix(y_train, y_train_pred)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred, average='micro'))))
print("Non Zero weights:", np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
print('Confusion matrix for the model is:')
plot_confusion_matrix(y_test, y_pred)
```

Accuracy on train set: 98.463%
 Precision on train set: 0.986
 Recall on train set: 0.997
 F1-Score on train set: 0.985
 Non Zero weights: 13254
 Confusion Matrix of test set:
 [[TN FP]
 [FN TP]]



Accuracy on test set: 91.573%
 Precision on test set: 0.938
 Recall on test set: 0.967
 F1-Score on test set: 0.916
 Non Zero weights: 13254
 Confusion Matrix of test set:
 [[TN FP]
 [FN TP]]

Confusion matrix for the model is:



Showing how sparsity increases as we increase lambda or decrease C when L1 Regularizer is used

```
In [27]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 1000, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
```

Accuracy on test set: 84.600%
 F1-Score on test set: 0.846
 Non Zero weights: 19430


```
In [28]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 100, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 84.740%
F1-Score on test set: 0.847
Non Zero weights: 17013
```

```
In [29]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 10, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 85.710%
F1-Score on test set: 0.857
Non Zero weights: 17390
```

```
In [30]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 1, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 88.003%
F1-Score on test set: 0.880
Non Zero weights: 18082
```

```
In [31]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 0.1, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 90.887%
F1-Score on test set: 0.909
Non Zero weights: 15583
```

```
In [32]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 0.01, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 91.660%
F1-Score on test set: 0.917
Non Zero weights: 5953
```

We can see how drastically the sparsity decreases from 19430 non-zero weights(@ C=1000) to only 5953 non-zero weights(@ C=0.01) when we use L1 Regularization

Using Randomized Search CV to find best parameters

```
In [33]: %time
from sklearn.model_selection import RandomizedSearchCV
from sklearn.linear_model import LogisticRegression

clf = LogisticRegression()
#params we need to try on classifier
param_grid = { 'C':[1000,500,100,50,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001,0.0005,0.001],
               'penalty':['l1','l2']}
tscv = TimeSeriesSplit(n_splits=10) #For time based splitting
gsv = RandomizedSearchCV(clf,param_grid,cv=tscv,verbose=1,scoring='f1_micro')
gsv.fit(X_train,y_train)
savetofile(gsv,"Log Reg/gsv_uni_r")
print("Best HyperParameter: ",gsv.best_params_)
print("Best Accuracy: %.2f%%"%(gsv.best_score_*100))

Wall time: 0 ns
Fitting 10 folds for each of 10 candidates, totalling 100 fits

[Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 6.5min finished

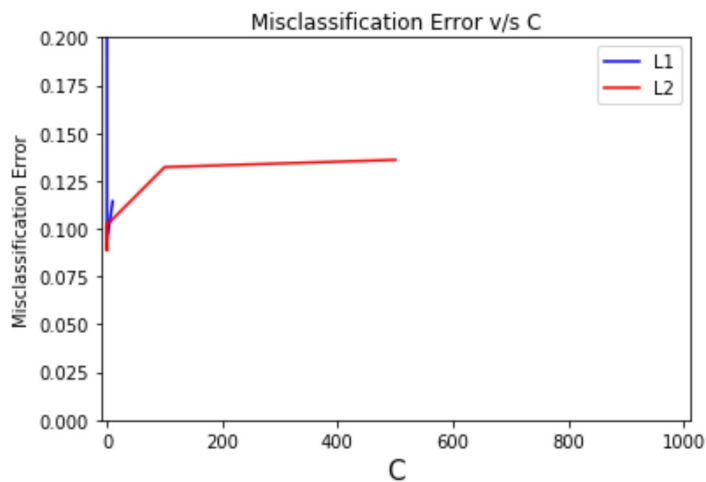
Best HyperParameter: {'penalty': 'l2', 'C': 0.001}
Best Accuracy: 91.14%
```

```

In [34]: def plot_error_vs_c_r(gsv):
    x1=[]
    y1=[]
    x2=[]
    y2=[]
    for a in gsv.grid_scores_:
        if (a[0]['penalty']) == 'l1':
            y1.append(1-a[1])
            x1.append(a[0]['C'])
        else:
            y2.append(1-a[1])
            x2.append(a[0]['C'])

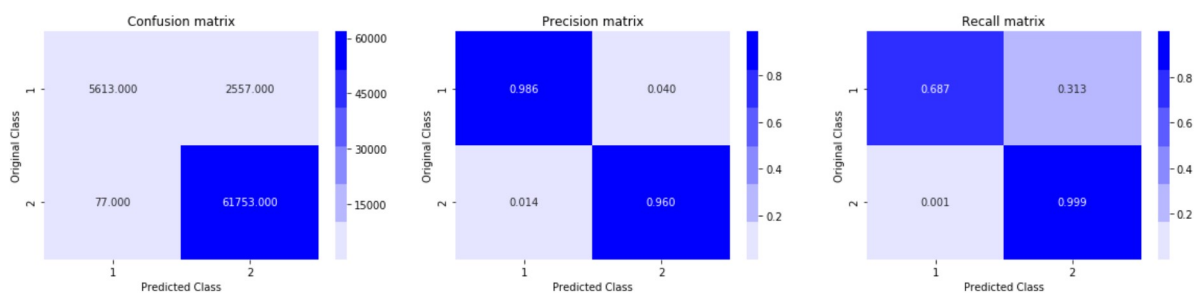
    ind1 = np.argsort(x1)
    x1=np.array(x1)
    y1=np.array(y1)
    ind2 = np.argsort(x2)
    x2=np.array(x2)
    y2=np.array(y2)
    plt.xlim(-10,1010)
    plt.ylim(0,0.2)
    plt.xlabel("C",fontsize=15)
    plt.ylabel("Misclassification Error")
    plt.title('Misclassification Error v/s C')
    plt.plot(x1[ind1],y1[ind1],'b',label="L1")
    plt.plot(x2[ind2],y2[ind2],'r',label="L2")
    plt.legend()
    plt.show()
    gsv = openfromfile("Log Reg/gsv_uni_r")
    plot_error_vs_c_r(gsv)

```



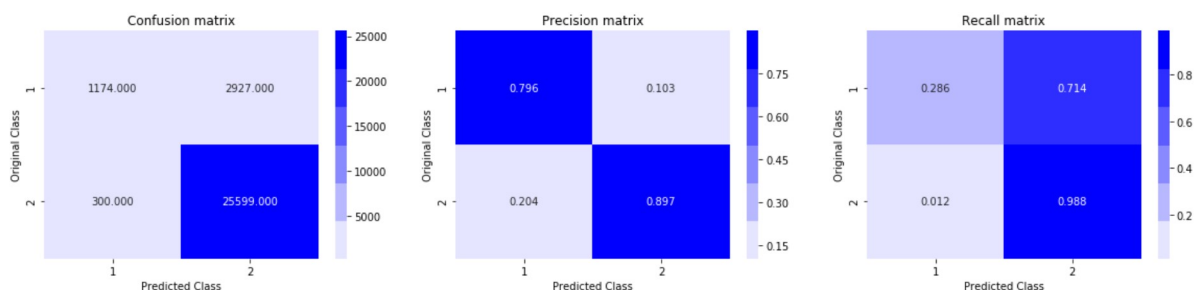
```
In [35]: clf = LogisticRegression(C= 0.0001, penalty= 'l2')
clf.fit(X_train,y_train)
y_train_pred = clf.predict(X_train)
y_pred = clf.predict(X_test)
print("Accuracy on train set: %0.3f%%"%(accuracy_score(y_train, y_train_pred)*100))
print("Precision on train set: %0.3f"%(precision_score(y_train, y_train_pred)))
print("Recall on train set: %0.3f"%(recall_score(y_train, y_train_pred)))
print("F1-Score on train set: %0.3f"%(f1_score(y_train, y_train_pred,average='micro'
'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
plot_confusion_matrix(y_train, y_train_pred)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
print('Confusion matrix for the model is:')
plot_confusion_matrix(y_test, y_pred)
```

Accuracy on train set: 96.237%
Precision on train set: 0.960
Recall on train set: 0.999
F1-Score on train set: 0.962
Non Zero weights: 64925
Confusion Matrix of test set:
[[TN FP]
[FN TP]]



Accuracy on test set: 89.243%
Precision on test set: 0.897
Recall on test set: 0.988
F1-Score on test set: 0.892
Non Zero weights: 64925
Confusion Matrix of test set:
[[TN FP]
[FN TP]]

Confusion matrix for the model is:



[7.2.4] Bi-Grams

Motivation

Now that we have our list of words describing positive and negative reviews lets analyse them.

We begin analysis by getting the frequency distribution of the words as shown below

```
In [36]: from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn import preprocessing

#Breaking into Train and test
X_train, X_test, y_train, y_test = train_test_split(final['CleanedText'].values, final['Score'].values, test_size=0.3, shuffle=False)

#taking one words and two consecutive words together
bi_gram = CountVectorizer(ngram_range=(1,2))
X_train = bi_gram.fit_transform(X_train)
#Normalize Data
X_train = preprocessing.normalize(X_train)

X_test = bi_gram.transform(X_test)
#Normalize Data
X_test = preprocessing.normalize(X_test)
from sklearn.preprocessing import StandardScaler
sc = StandardScaler(with_mean=False)
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
print("Train Data Size: ", X_train.shape)
print("Test Data Size: ", X_test.shape)

Train Data Size:  (70000, 1003102)
Test Data Size:  (30000, 1003102)
```

```
In [37]: %time
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

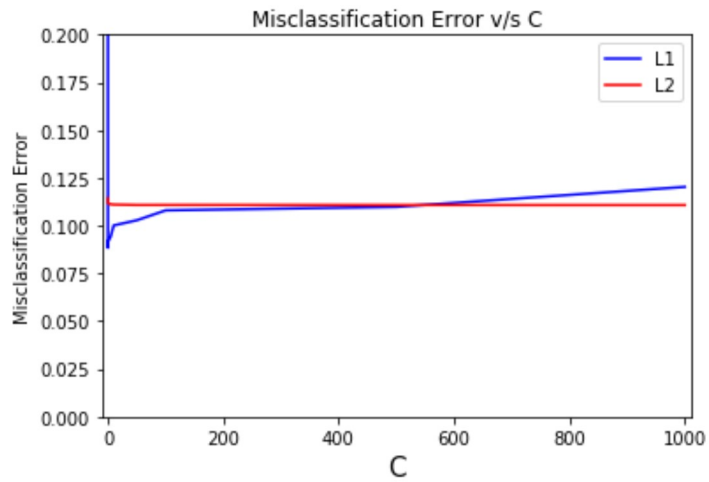
clf = LogisticRegression()
#params we need to try on classifier
param_grid = {'C': [1000, 500, 100, 50, 10, 5, 1, 0.5, 0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001],
              'penalty': ['l1', 'l2']}
tscv = TimeSeriesSplit(n_splits=10) #For time based splitting
gsv = GridSearchCV(clf, param_grid, cv=tscv, verbose=1, scoring='f1_micro')
gsv.fit(X_train, y_train)
savetofile(gsv, "Log Reg/gsv_bi")
print("Best HyperParameter: ", gsv.best_params_)
print("Best Accuracy: %.2f%%" % (gsv.best_score_*100))

Wall time: 0 ns
Fitting 10 folds for each of 30 candidates, totalling 300 fits

[Parallel(n_jobs=1)]: Done 300 out of 300 | elapsed: 42.9min finished

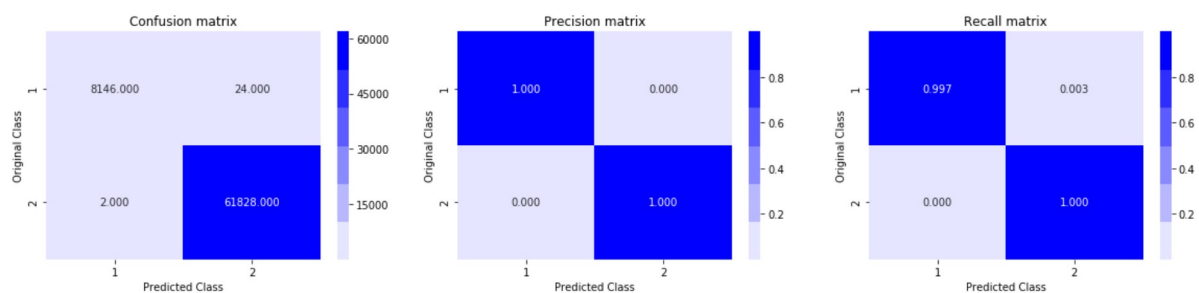
Best HyperParameter:  {'C': 0.01, 'penalty': 'l1'}
Best Accuracy: 91.16%
```

```
In [38]: gsv = openfromfile("Log Reg/gsv_bi")  
plot_error_vs_c(gsv)
```



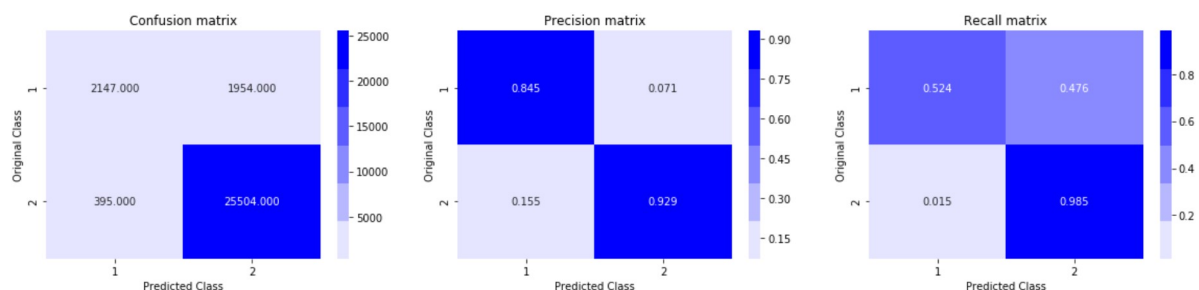
```
In [39]: clf = LogisticRegression(C= 0.01, penalty= 'l1')
clf.fit(X_train,y_train)
y_train_pred = clf.predict(X_train)
y_pred = clf.predict(X_test)
print("Accuracy on train set: %0.3f%%"%(accuracy_score(y_train, y_train_pred)*100))
print("Precision on train set: %0.3f%%"(precision_score(y_train, y_train_pred)))
print("Recall on train set: %0.3f%%"(recall_score(y_train, y_train_pred)))
print("F1-Score on train set: %0.3f%%"(f1_score(y_train, y_train_pred,average='micro'
'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
plot_confusion_matrix(y_train, y_train_pred)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("Precision on test set: %0.3f%%"(precision_score(y_test, y_pred)))
print("Recall on test set: %0.3f%%"(recall_score(y_test, y_pred)))
print("F1-Score on test set: %0.3f%%"(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
print('Confusion matrix for the model is:')
plot_confusion_matrix(y_test, y_pred)
```

Accuracy on train set: 99.963%
Precision on train set: 1.000
Recall on train set: 1.000
F1-Score on train set: 1.000
Non Zero weights: 27198
Confusion Matrix of test set:
[[TN FP]
[FN TP]]



Accuracy on test set: 92.170%
Precision on test set: 0.929
Recall on test set: 0.985
F1-Score on test set: 0.922
Non Zero weights: 27198
Confusion Matrix of test set:
[[TN FP]
[FN TP]]

Confusion matrix for the model is:



Showing how sparsity increases as we increase lambda or decrease C when L1 Regularizer is used

```
In [40]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 1000, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 84.673%
F1-Score on test set: 0.847
Non Zero weights: 467782
```

```
In [41]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 100, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 87.683%
F1-Score on test set: 0.877
Non Zero weights: 211933
```

```
In [42]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 10, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 90.337%
F1-Score on test set: 0.903
Non Zero weights: 65503
```

```
In [43]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 1, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 91.657%
F1-Score on test set: 0.917
Non Zero weights: 43566
```



```
In [44]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 0.1, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 91.650%
F1-Score on test set: 0.916
Non Zero weights: 37455
```

```
In [45]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 0.01, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 92.177%
F1-Score on test set: 0.922
Non Zero weights: 25153
```

We can see how drastically the sparsity decreases from 467782 non-zero weights(@ C=1000) to only 25153 non-zero weights(@ C=0.01) when we use L1 Regularization

Using Randomized Search CV to find best parameters

```
In [46]: %time
from sklearn.model_selection import RandomizedSearchCV
from sklearn.linear_model import LogisticRegression

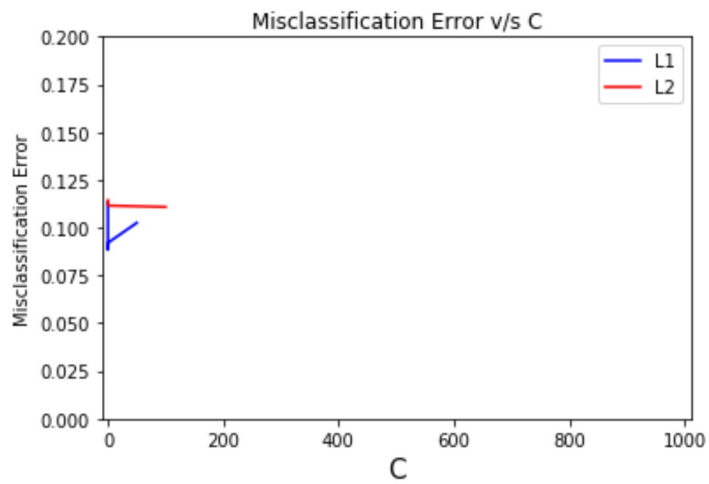
clf = LogisticRegression()
#params we need to try on classifier
param_grid = {'C':[1000,500,100,50,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001,0.0005,0.0001],
              'penalty':['l1','l2']}
tscv = TimeSeriesSplit(n_splits=10) #For time based splitting
gsv = RandomizedSearchCV(clf,param_grid,cv=tscv,verbose=1,scoring='f1_micro')
gsv.fit(X_train,y_train)
savetofile(gsv,"Log Reg/gsv_bi_r")
print("Best HyperParameter: ",gsv.best_params_)
print("Best Accuracy: %.2f%%"%(gsv.best_score_*100))

Wall time: 0 ns
Fitting 10 folds for each of 10 candidates, totalling 100 fits

[Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 18.6min finished

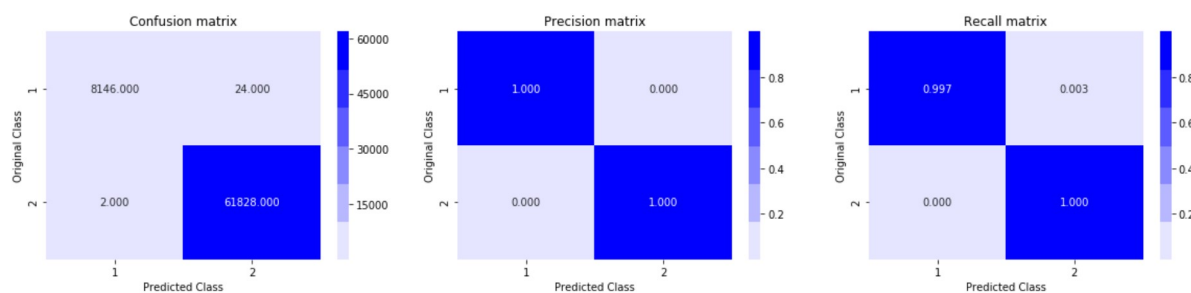
Best HyperParameter: {'penalty': 'l1', 'C': 0.01}
Best Accuracy: 91.16%
```

```
In [47]: gsv = openfromfile("Log Reg/gsv_bi_r")  
plot_error_vs_c_r(gsv)
```



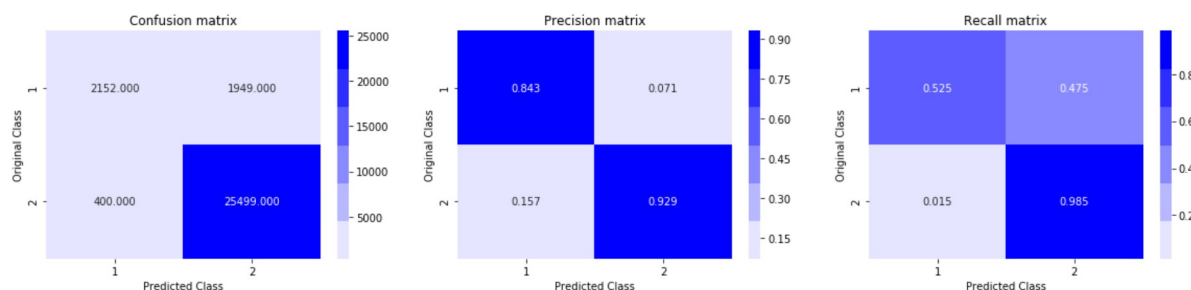
```
In [48]: clf = LogisticRegression(C= 0.01, penalty= 'l1')
clf.fit(X_train,y_train)
y_train_pred = clf.predict(X_train)
y_pred = clf.predict(X_test)
print("Accuracy on train set: %0.3f%%"%(accuracy_score(y_train, y_train_pred)*100))
print("Precision on train set: %0.3f"%(precision_score(y_train, y_train_pred)))
print("Recall on train set: %0.3f"%(recall_score(y_train, y_train_pred)))
print("F1-Score on train set: %0.3f"%(f1_score(y_train, y_train_pred,average='micro'
'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
plot_confusion_matrix(y_train, y_train_pred)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
print('Confusion matrix for the model is:')
plot_confusion_matrix(y_test, y_pred)
```

Accuracy on train set: 99.963%
Precision on train set: 1.000
Recall on train set: 1.000
F1-Score on train set: 1.000
Non Zero weights: 25528
Confusion Matrix of test set:
[[TN FP]
[FN TP]]



Accuracy on test set: 92.170%
Precision on test set: 0.929
Recall on test set: 0.985
F1-Score on test set: 0.922
Non Zero weights: 25528
Confusion Matrix of test set:
[[TN FP]
[FN TP]]

Confusion matrix for the model is:



Perturbation Test

```
In [84]: import scipy as sp
epsilon = sp.stats.distributions.norm.rvs(loc=0,scale=0.0001)
# Vector before the addition of epsilon
W_before_epsilon = clf.coef_

# Number of non zero elements in X_train_vec_standardized sparse matrix
no_of_non_zero = X_train.count_nonzero()

# Importing library to create a sparse matrix of epsilon
from scipy.sparse import csr_matrix

# Creating new sparse matrix with epsilon at same position of non-zero elements of
X_train_vec_standardized
indices_X_train = X_train.indices
indptr_X_train = X_train.indptr

# Creating a list of same element with repetition
data = [epsilon] * no_of_non_zero
Shape = X_train.shape

# Creating sparse matrix
sparse_epsilon = csr_matrix((data,indices_X_train,indptr_X_train),shape=Shape,dtype
=float)

# Add sparse_epsilon and X_train_vec_standardized to get a new sparse matrix with e
psilon added to each
# non-zero element of X_train_vec_standardized
epsilon_train = X_train + sparse_epsilon

print(X_train.shape)
print(epsilon_train.shape)

(70000, 1003102)
(70000, 1003102)
```

```
In [86]: epsilon_lr = LogisticRegression(penalty='l1', C=0.01, n_jobs=-1)
epsilon_lr.fit(epsilon_train,y_train)

# Vector after the addition of epsilon
W_after_epsilon = epsilon_lr.coef_

# Change in vectors after adding epsilon
change_vector = W_after_epsilon - W_before_epsilon

# Sort this change_vector array after making all the elements positive in ascending
order to visualize the change
sorted_change_vector = np.sort(np.absolute(change_vector))[:,::-1]

sorted_change_vector[0,0:20]
```

```
Out[86]: array([0.0177024 , 0.01069277, 0.01045934, 0.01036312, 0.00995046,
                0.00994137, 0.00987692, 0.00977656, 0.0096977 , 0.00931577,
                0.00919783, 0.00919306, 0.00918648, 0.0091204 , 0.00908224,
                0.00895077, 0.008722 , 0.00870868, 0.00864846, 0.0085583 ])
```

```
In [88]: absolute_weights = np.absolute(W_before_epsilon)
sorted_absolute_index = np.argsort(absolutely_weights)[:,-1]
top_index = sorted_absolute_index[0,0:20]

all_features = bi_gram.get_feature_names()
weight_values = clf.coef_

# Top 20 features are
print("Top 20 features with their weight values :")

for j in top_index:
    print("%12s\t--> \t%f"%(all_features[j],weight_values[0,j]))
```

```
Top 20 features with their weight values :
      great -->      0.683128
        best -->      0.472134
        love -->      0.463073
        good -->      0.369411
disappoint -->     -0.360893
      delici -->      0.353219
      perfect -->      0.307629
        excel -->      0.272745
         not -->     -0.231705
      favorit -->      0.223190
        nice -->      0.211375
       worst -->     -0.207587
      wonder -->      0.197292
not worth -->     -0.187525
        find -->      0.184926
      terribl -->     -0.176707
not disappoint -->      0.170614
high recommend -->      0.166500
        tasti -->      0.165060
         aw -->     -0.151817
```

[7.2.5] TF-IDF

```
In [90]: %%time
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn import preprocessing

#Breaking into Train and test
X_train, X_test, y_train, y_test = train_test_split(final['CleanedText'].values, final['Score'].values, test_size=0.3, shuffle=False)

tfidf = TfidfVectorizer()
X_train = tfidf.fit_transform(X_train)
#Normalize Data
X_train = preprocessing.normalize(X_train)

X_test = tfidf.transform(X_test)
#Normalize Data
X_test = preprocessing.normalize(X_test)
sc = StandardScaler(with_mean=False)
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

print("Train Data Size: ", X_train.shape)
print("Test Data Size: ", X_test.shape)

Train Data Size:  (70000, 64925)
Test Data Size:  (30000, 64925)
Wall time: 9.48 s
```

```
In [91]: %%time
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

clf = LogisticRegression()

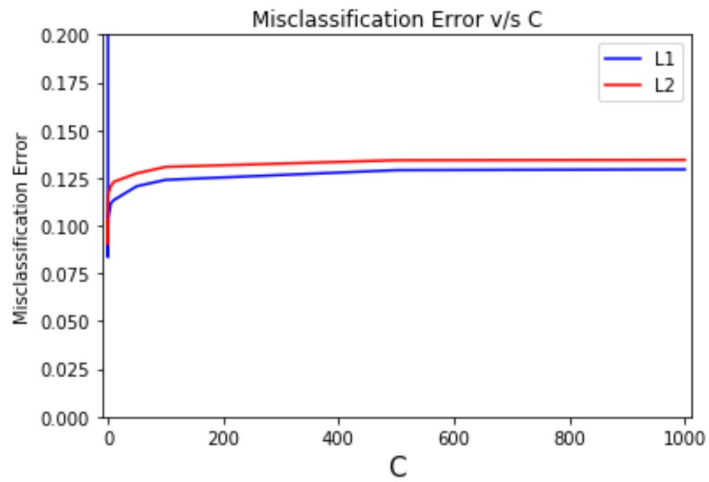
#params we need to try on classifier
param_grid = {'C': [1000, 500, 100, 50, 10, 5, 1, 0.5, 0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001],
              'penalty': ['l1', 'l2']}
tscv = TimeSeriesSplit(n_splits=10) #For time based splitting
gsv = GridSearchCV(clf, param_grid, cv=tscv, verbose=1, scoring='f1_micro')
gsv.fit(X_train, y_train)
savetofile(gsv, "Log Reg/gsv_tfidf")
print("Best HyperParameter: ", gsv.best_params_)
print("Best Accuracy: %.2f%%" % (gsv.best_score_*100))

Wall time: 0 ns
Fitting 10 folds for each of 30 candidates, totalling 300 fits

[Parallel(n_jobs=1)]: Done 300 out of 300 | elapsed: 101.7min finished

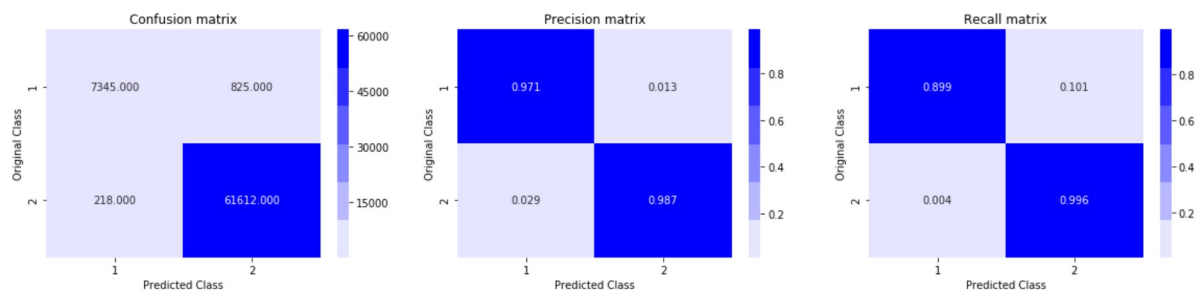
Best HyperParameter:  {'C': 0.05, 'penalty': 'l1'}
Best Accuracy: 91.65%
```

```
In [92]: gsv = openfromfile("Log Reg/gsv_tfidf")  
         plot_error_vs_c(gsv)
```



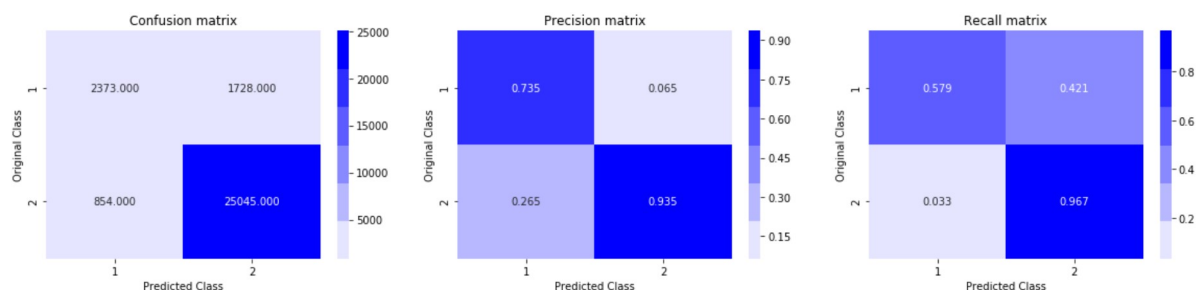
```
In [93]: clf = LogisticRegression(C= 0.05, penalty= 'l1')
clf.fit(X_train,y_train)
y_train_pred = clf.predict(X_train)
y_pred = clf.predict(X_test)
print("Accuracy on train set: %0.3f%%"%(accuracy_score(y_train, y_train_pred)*100))
print("Precision on train set: %0.3f"%(precision_score(y_train, y_train_pred)))
print("Recall on train set: %0.3f"%(recall_score(y_train, y_train_pred)))
print("F1-Score on train set: %0.3f"%(f1_score(y_train, y_train_pred,average='micro'
'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
plot_confusion_matrix(y_train, y_train_pred)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
print('Confusion matrix for the model is:')
plot_confusion_matrix(y_test, y_pred)
```

Accuracy on train set: 98.510%
Precision on train set: 0.987
Recall on train set: 0.996
F1-Score on train set: 0.985
Non Zero weights: 14008
Confusion Matrix of test set:
[[TN FP]
[FN TP]]



Accuracy on test set: 91.393%
Precision on test set: 0.935
Recall on test set: 0.967
F1-Score on test set: 0.914
Non Zero weights: 14008
Confusion Matrix of test set:
[[TN FP]
[FN TP]]

Confusion matrix for the model is:



Showing how sparsity increases as we increase lambda or decrease C when L1 Regularizer is use

```
In [103]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 1000, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 84.490%
F1-Score on test set: 0.845
Non Zero weights: 24250
```

```
In [105]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 100, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 84.697%
F1-Score on test set: 0.847
Non Zero weights: 19858
```

```
In [106]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 10, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 85.553%
F1-Score on test set: 0.856
Non Zero weights: 17594
```

```
In [107]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 1, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 87.723%
F1-Score on test set: 0.877
Non Zero weights: 18521
```

```
In [108]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 0.1, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 90.703%
F1-Score on test set: 0.907
Non Zero weights: 16211
```

```
In [109]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(C= 0.01, penalty= 'l1')
clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))

Accuracy on test set: 91.597%
F1-Score on test set: 0.916
Non Zero weights: 6124
```

We can see how drastically the sparsity increases from 24250 non-zero weights(@ C=1000) to only 3 non-zero weights(@ C=0.01) when we use L1 Regularization

Regularization Using Randomized Search CV to find best parameters

```
In [100]: %time
from sklearn.model_selection import RandomizedSearchCV
from sklearn.linear_model import LogisticRegression

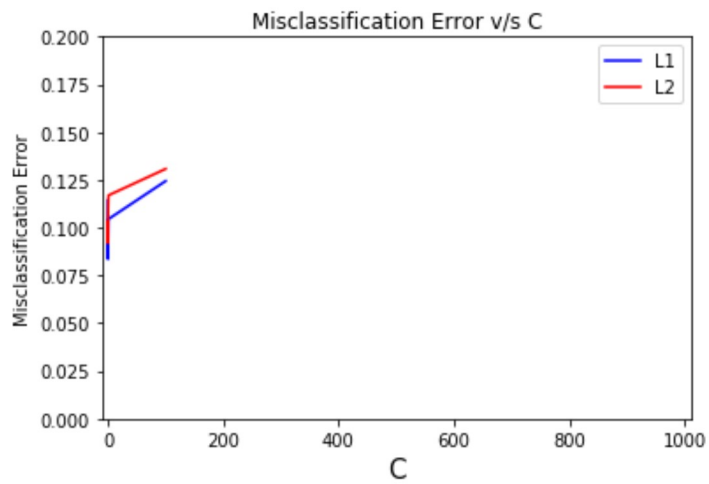
clf = LogisticRegression()
#params we need to try on classifier
param_grid = {'C':[1000,500,100,50,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001,0.0005,0.0001],
              , 'penalty':['l1','l2']}
tscv = TimeSeriesSplit(n_splits=10) #For time based splitting
gsv = RandomizedSearchCV(clf,param_grid,cv=tscv,verbose=1,scoring='f1_micro')
gsv.fit(X_train,y_train)
savetofile(gsv,"Log Reg/gsv_tfidf_r")
print("Best HyperParameter: ",gsv.best_params_)
print("Best Accuracy: %.2f%%"%(gsv.best_score_*100))

Wall time: 0 ns
Fitting 10 folds for each of 10 candidates, totalling 100 fits

[Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 7.7min finished

Best HyperParameter: {'penalty': 'l1', 'C': 0.05}
Best Accuracy: 91.65%
```

```
In [101]: gsv = openfromfile("Log Reg/gsv_tfidf_r")  
plot_error_vs_c_r(gsv)
```



```
In [102]: clf = LogisticRegression(C= 0.05, penalty= 'l1')
clf.fit(X_train,y_train)
y_train_pred = clf.predict(X_train)
y_pred = clf.predict(X_test)
print("Accuracy on train set: %0.3f%%"%(accuracy_score(y_train, y_train_pred)*100)
)
print("Precision on train set: %0.3f"%(precision_score(y_train, y_train_pred)))
print("Recall on train set: %0.3f"%(recall_score(y_train, y_train_pred)))
print("F1-Score on train set: %0.3f"%(f1_score(y_train, y_train_pred,average='micr
o'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
plot_confusion_matrix(y_train, y_train_pred)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
print('Confusion matrix for the model is:')
plot_confusion_matrix(y_test, y_pred)
```

Accuracy on train set: 98.513%

Precision on train set: 0.987

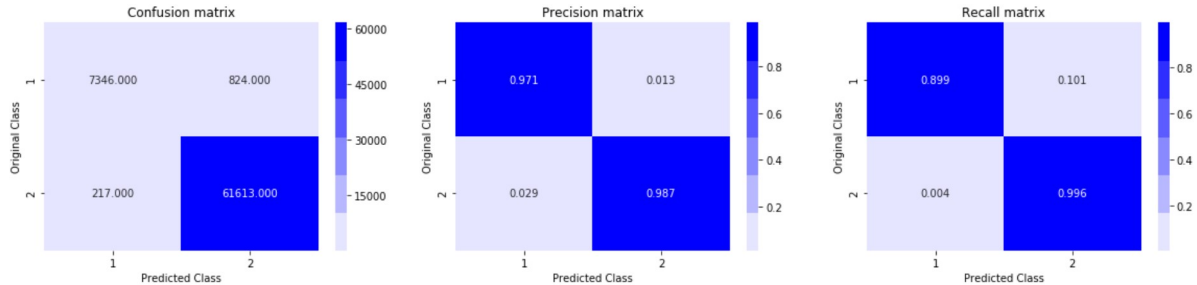
Recall on train set: 0.996

F1-Score on train set: 0.985

Non Zero weights: 13903

Confusion Matrix of test set:

```
[ [TN  FP]
 [FN TP] ]
```



Accuracy on test set: 91.390%

Precision on test set: 0.936

Recall on test set: 0.967

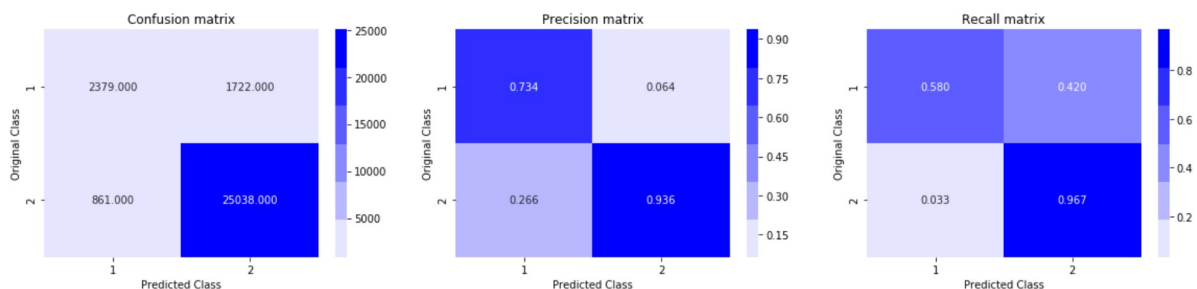
F1-Score on test set: 0.914

Non Zero weights: 13903

Confusion Matrix of test set:

```
[ [TN  FP]
 [FN TP] ]
```

Confusion matrix for the model is:



Perturbation Test

```
In [110]: import scipy as sp
epsilon = sp.stats.distributions.norm.rvs(loc=0,scale=0.0001)
# Vector before the addition of epsilon
W_before_epsilon = clf.coef_

# Number of non zero elements in X_train_vec_standardized sparse matrix
no_of_non_zero = X_train.count_nonzero()

# Importing library to create a sparse matrix of epsilon
from scipy.sparse import csr_matrix

# Creating new sparse matrix with epsilon at same position of non-zero elements of
X_train_vec_standardized
indices_X_train = X_train.indices
indptr_X_train = X_train.indptr

# Creating a list of same element with repetition
data = [epsilon] * no_of_non_zero
Shape = X_train.shape

# Creating sparse matrix
sparse_epsilon = csr_matrix((data,indices_X_train,indptr_X_train),shape=Shape,dtype=float)

# Add sparse_epsilon and X_train_vec_standardized to get a new sparse matrix with
epsilon added to each
# non-zero element of X_train_vec_standardized
epsilon_train = X_train + sparse_epsilon

print(X_train.shape)
print(epsilon_train.shape)

(70000, 64925)
(70000, 64925)
```

```
In [111]: epsilon_lr = LogisticRegression(penalty='l1', C=0.05, n_jobs=-1)
epsilon_lr.fit(epsilon_train,y_train)

# Vector after the addition of epsilon
W_after_epsilon = epsilon_lr.coef_

# Change in vectors after adding epsilon
change_vector = W_after_epsilon - W_before_epsilon

# Sort this change_vector array after making all the elements positive in ascending
order to visualize the change
sorted_change_vector = np.sort(np.absolute(change_vector))[:,::-1]

sorted_change_vector[0,0:20]
```

```
Out[111]: array([0.20545409, 0.18718968, 0.16820735, 0.16599216, 0.14544775,
0.14533518, 0.13881392, 0.13196279, 0.12740095, 0.12252001,
0.12242056, 0.1218007 , 0.11695892, 0.11260441, 0.11123762,
0.11024841, 0.10760662, 0.10624811, 0.10495961, 0.10492818])
```

```
In [112]: absolute_weights = np.absolute(W_before_epsilon)
sorted_absolute_index = np.argsort(absolutely_weights)[:,-1]
top_index = sorted_absolute_index[0,0:20]

all_features = tfidf.get_feature_names()
weight_values = clf.coef_

# Top 20 features are
print("Top 20 features with their weight values :")

for j in top_index:
    print("%12s\t--> \t%f"%(all_features[j],weight_values[0,j]))
```

Top 20 features with their weight values :

great	-->	0.684374
best	-->	0.519821
love	-->	0.461286
delici	-->	0.393307
perfect	-->	0.367119
not	-->	-0.361515
good	-->	0.331659
excel	-->	0.313476
nice	-->	0.253738
favorit	-->	0.243472
disappoint	-->	-0.234664
find	-->	0.219518
wonder	-->	0.217783
worst	-->	-0.208236
tasti	-->	0.177924
amaz	-->	0.176945
keep	-->	0.169720
addict	-->	0.156525
terribl	-->	-0.152146
yummi	-->	0.145016

In []:

In []:

[7.2.6] Word2Vec

```
In [113]: # Using Google News Word2Vectors

# in this project we are using a pretrained model by google
# its 3.3G file, once you load this into your memory
# it occupies ~9Gb, so please do this step only if you have >12G of ram
# we will provide a pickle file wich contains a dict ,
# and it contains all our courpus words as keys and model[word] as values
# To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
# from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edit
# it's 1.9GB in size.

# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.Wl7SRFAzZPY
# you can comment this whole cell
# or change these variable according to your need
is_your_ram_gt_16g=False
want_to_read_sub_set_of_google_w2v = True
want_to_read_whole_google_w2v = True
if not is_your_ram_gt_16g:
    if want_to_read_sub_set_of_google_w2v and os.path.isfile('google_w2v_for_amazon.pkl'):
        with open('google_w2v_for_amazon.pkl', 'rb') as f:
            # model is dict object, you can directly access any word vector using
            model[word]
        model = pickle.load(f)
    else:
        if want_to_read_whole_google_w2v and os.path.isfile('GoogleNews-vectors-negative300.bin'):
            model = KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)

# print("the vector representation of word 'computer'",model.wv['computer'])
# print("the similarity between the words 'woman' and 'man'",model.wv.similarity('woman', 'man'))
# print("the most similar words to the word 'woman'",model.wv.most_similar('woman'))
# this will raise an error
# model.wv.most_similar('tasti') # "tasti" is the stemmed word for tasty, tastful
```

```
In [114]: final_string = []
for sent in final['CleanedText'].values:
    sent = str(sent)
    sentence=[]
    # print(sent)
    for word in sent.split():
        # print(word)
        sentence.append(word)
    # print(sentence)
    final_string.append(sentence)
```

```
In [115]: %%time
# Train your own Word2Vec model using your own text corpus
import gensim

w2v_model=gensim.models.Word2Vec(final_string,min_count=5,size=50, workers=-1)
#min-count: Ignoring the words which occurs less than 5 times
#size:Creating vectors of size 50 for each word
#workers: Use these many worker threads to train the model (faster training with multicore machines)
```

Wall time: 1.77 s

```
In [116]: w2v_model.save('w2vmodel')
```

```
In [117]: w2v_model = gensim.models.Word2Vec.load('w2vmodel')
```

```
In [167]: w2v_words = list(w2v_model.wv.vocab)
```

```
In [118]: w2v_vocub = w2v_model.wv.vocab  
len(w2v_vocub)
```

```
Out[118]: 16909
```

```
In [119]: w2v_model.wv.most_similar('like')
```

```
Out[119]: [("spoil'", 0.5566685795783997),  
            ('southwest', 0.5187448263168335),  
            ('napolitain', 0.48771023750305176),  
            ("b'make", 0.4703911244869232),  
            ('queijo', 0.46708792448043823),  
            ('vice', 0.45246225595474243),  
            ('headach', 0.4475083351135254),  
            ('monounsatur', 0.4394482970237732),  
            ("need'", 0.4379512071609497),  
            ('arrang', 0.43731579184532166)]
```

avg w2c


```
In [120]: %%time
avg_vec = [] #List to store all the avg w2vec's
for sent in final_string[0:1]:
    cnt = 0 #to count no of words in each reviews
    sent_vec = np.zeros(50) #Initializing with zeroes
    print("sent:",sent)
    for word in sent:
        try:
            wvec = w2v_model.wv[word] #Vector of each using w2v model
            print("wvec:",wvec)
            sent_vec += wvec #Adding the vectors
            cnt += 1
        except:
            pass #When the word is not in the dictionary then do nothing
    print("sent_vec:",sent_vec)
    a_vec =sent_vec / cnt #Taking average of vectors sum of the particular review
    print("avg_vec:",a_vec)
    avg_vec.append(a_vec) #Storing the avg w2vec's for each review
    print("*****")
```

```
sent: ["b'witti'", 'littl', 'book', 'make', 'son', 'laugh', 'loud', 'recit', 'car',
', 'drive', 'along', 'alway', 'sing', 'refrain', 'hes', 'learn', 'whale', 'india',
', 'droop', 'love', 'new', 'word', 'book', 'introduc', 'silli', 'classic', 'book',
', 'will', 'bet', 'son', 'still', 'abl', 'recit', 'memori', "colleg'"]
wvec: [ 6.8155774e-03  8.0338391e-03  2.0524070e-03  5.4769060e-03
 8.6527430e-03  4.9760267e-03 -9.5563391e-03  2.7678728e-03
-4.3694251e-03  5.1570237e-03 -3.7309120e-03  1.4504878e-04
-7.8726364e-03 -7.7749914e-03  1.2366079e-03  1.6883597e-03
 3.8338881e-03 -3.7840137e-03 -4.3124352e-03  8.7198811e-03
 3.4406213e-03 -1.0414000e-04 -3.8977561e-03 -7.1372930e-04
-2.3633661e-03  8.9031989e-03 -5.4177912e-03  1.2080419e-03
-2.1027524e-03  6.2195063e-03  3.2489670e-03  3.6720117e-03
-4.7993441e-03  3.9366260e-03  8.9964475e-03 -2.9056258e-03
 8.1923390e-03 -3.3489075e-03  8.5266735e-03  5.6435908e-03
 2.7828274e-04  3.6064326e-03 -8.1068110e-03  9.1146417e-03
-9.6330056e-03 -2.1887694e-03 -2.3092022e-03 -2.4836559e-03
 7.3081910e-06  8.3241779e-03]
wvec: [ 0.00106912 -0.00334294 -0.00855349 -0.00268923 -0.00341163 -0.00180589
 0.00103952  0.00966915  0.00548612  0.00813285  0.00959688  0.00016837
-0.00987749 -0.00423928 -0.00400014 -0.00056169 -0.00464407  0.00629331
-0.00904419  0.00877909  0.00906067  0.00853112  0.00914099  0.00862897
-0.00653833  0.0055696  -0.00413966  0.00904269 -0.00580844 -0.00533257
-0.00233003 -0.00637354  0.00353411 -0.00433964  0.00361061 -0.00967616
 0.00882937 -0.00882016 -0.00194709  0.00229614 -0.00921923 -0.00432034
-0.00709985 -0.00218498 -0.00618744  0.00032735  0.00719821 -0.00983658
-0.0002708  0.0054936 ]
wvec: [-8.4407534e-03  8.6668050e-03 -8.3649727e-03 -5.5213468e-03
 9.2176832e-03  3.3547555e-03  4.4790753e-03  4.7553503e-03
-3.0217765e-04 -3.3006941e-03 -4.9324860e-03 -3.7805827e-03
 5.9023602e-03  1.2323612e-03 -3.7577774e-05  3.5210419e-03
-3.2677612e-04  1.0562962e-03 -9.1800205e-03 -2.1955410e-03
-8.9938603e-03  4.4158022e-03  2.8118223e-03  8.8645956e-03
-4.7084698e-03  7.6758317e-03  9.0797124e-03 -1.6311810e-03
 2.0808347e-03  9.0627279e-03 -2.1466543e-03  8.4444089e-03
-5.9388960e-03 -3.2090980e-03  5.6574889e-03  7.7351104e-03
-4.0460681e-03  9.5816776e-03 -5.3995569e-05 -4.5100609e-03
-2.4462973e-03  5.1194495e-03  9.1895396e-03 -8.9310240e-03
-3.2861903e-03 -6.8493411e-03  7.6813521e-03 -3.0973847e-03
-9.8433094e-03  1.6246469e-03]
wvec: [-0.00502434 -0.00369201  0.00876713 -0.00838183  0.00593766 -0.00534093
 0.0076713  0.00792619  0.00312846  0.00766176 -0.00101374 -0.0096011
 0.00628034  0.008817  -0.00168275  0.00265238 -0.00700324  0.00585293
-0.0061377 -0.00974201 -0.00052154 -0.00218882 -0.00654677  0.0089383
-0.0082328  0.00083269 -0.00730804 -0.00531235  0.00856922  0.0079186
-0.00037423  0.0095028  -0.00757354  0.00414478 -0.00948966  0.00891798
-0.00406121 -0.0075944  -0.00201797 -0.00949768 -0.00737957  0.00470853
-0.00824873 -0.00038798  0.00310391 -0.00205357 -0.00566535 -0.00440724
-0.00900015 -0.00603302]
wvec: [-0.00675152 -0.00098457 -0.00579327  0.00804078 -0.00775842 -0.00950479
-0.000615  -0.00110464 -0.00862497  0.00452825  0.00574983  0.0041848
-0.00834732  0.00243418  0.00482501  0.00774253  0.00464348 -0.00896337
-0.00113293  0.00157739 -0.00952516  0.00421187 -0.00535609  0.0001831
-0.00787035 -0.0093928  0.00627939 -0.00846677 -0.00840965  0.00010529
-0.00409689  0.00577683  0.00054423  0.00805292  0.00644578  0.00673872
-0.00968168 -0.00208418 -0.00990162 -0.00025117  0.0023344  -0.00643416
 0.00019504 -0.00060823 -0.00542158  0.00858975 -0.00427644  0.00193023
-0.00499162  0.00548009]
wvec: [-0.00598099  0.00272044  0.00820923  0.00553821 -0.00011654 -0.00328538
-0.00654433 -0.00235384  0.00179457 -0.00796135 -0.00672696  0.0050995
-0.0076665 -0.0010567  0.00247536  0.00792624 -0.00149003 -0.00733114
-0.00474387 -0.00061135 -0.00502998 -0.00028786  0.00525735  0.00621393
 0.00800868  0.00895356  0.0045481  -0.00440279 -0.00327003 -0.0094584
 0.00704338 -0.00849353 -0.00638037 -0.00085133  0.00575292  0.00507055
-0.00887217 -0.00130038 -0.00493244  0.00672765  0.0084738  0.00109222]
```

```
In [121]: %%time
np.seterr(divide='ignore', invalid='ignore')
avg_vec = [] #List to store all the avg w2vec's
for sent in final_string:
    cnt = 0 #to count no of words in each reviews
    sent_vec = np.zeros(50) #Initializing with zeroes
    for word in sent:
        try:
            wvec = w2v_model.wv[word] #Vector of each using w2v model
            sent_vec += wvec #Adding the vectors
            cnt += 1
        except:
            pass #When the word is not in the dictionary then do nothing
    sent_vec /= cnt #Taking average of vectors sum of the particular review
    avg_vec.append(sent_vec) #Storing the avg w2vec's for each review
    #print("*****")
    # Average Word2Vec
```

Wall time: 14.1 s

```
In [122]: savetofile(avg_vec,"avg_w2v_vec")
```

```
In [123]: avg_vec = openfromfile("avg_w2v_vec")
```

```
In [124]: avg_vec = np.array(avg_vec)
avg_vec.shape
```

```
Out[124]: (100000, 50)
```

```
In [125]: np.isnan(avg_vec).any()
```

```
Out[125]: True
```

```
In [126]: mask = ~np.any(np.isnan(avg_vec), axis=1)
# print(mask)
avg_vec_new = avg_vec[mask]
final_sample_new = final['Score'][mask]
print(avg_vec_new.shape)
print(final_sample_new.shape)
```

```
(99996, 50)
```

```
(99996,)
```

```
In [174]: from sklearn.model_selection import train_test_split
from sklearn import preprocessing
#Normalizing the data
avg_vec_norm = preprocessing.normalize(avg_vec_new)

#Not shuffling the data as we want it on time basis
X_train, X_test, y_train, y_test = train_test_split(avg_vec_norm,final_sample_new,
values,test_size=0.3,shuffle=False)
sc = StandardScaler(with_mean=False)
X_train = sc.fit_transform(X_train)
X_test= sc.transform(X_test)

print("Train Data Size: ",X_train.shape)
print("Test Data Size: ",X_test.shape)
```

```
Train Data Size: (69997, 50)
```

```
Test Data Size: (29999, 50)
```

```
In [131]: %time
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

clf = LogisticRegression()
#params we need to try on classifier
param_grid = {'C':[100,50,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001,0.0005,0.0001],
              'penalty':['l1','l2']}
tscv = TimeSeriesSplit(n_splits=10) #For time based splitting
gsv = GridSearchCV(clf,param_grid,cv=tscv,verbose=1,scoring='f1_micro')
gsv.fit(X_train,y_train)
savetofile(gsv,"Log Reg/gsv_w2v")
print("Best HyperParameter: ",gsv.best_params_)
print("Best Accuracy: %.2f%%"%(gsv.best_score_*100))
```

Wall time: 0 ns

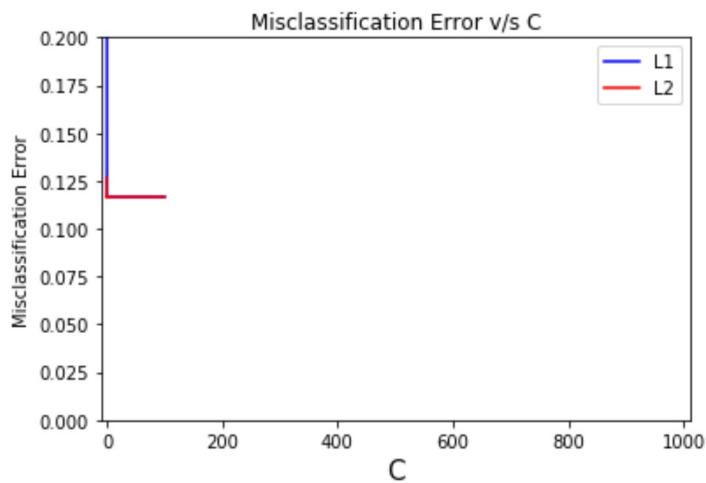
Fitting 10 folds for each of 26 candidates, totalling 260 fits

[Parallel(n_jobs=1)]: Done 260 out of 260 | elapsed: 1.2min finished

Best HyperParameter: {'C': 0.1, 'penalty': 'l2'}

Best Accuracy: 88.34%

```
In [132]: gsv = openfromfile("Log Reg/gsv_w2v")
plot_error_vs_c(gsv)
```



```
In [133]: clf = LogisticRegression(C= 0.1, penalty= 'l2')
clf.fit(X_train,y_train)
y_train_pred = clf.predict(X_train)
y_pred = clf.predict(X_test)
print("Accuracy on train set: %0.3f%%"%(accuracy_score(y_train, y_train_pred)*100)
)
print("Precision on train set: %0.3f"%(precision_score(y_train, y_train_pred)))
print("Recall on train set: %0.3f"%(recall_score(y_train, y_train_pred)))
print("F1-Score on train set: %0.3f"%(f1_score(y_train, y_train_pred,average='micro'
o))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
plot_confusion_matrix(y_train, y_train_pred)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
print('Confusion matrix for the model is:')
plot_confusion_matrix(y_test, y_pred)
```

Accuracy on train set: 88.341%

Precision on train set: 0.884

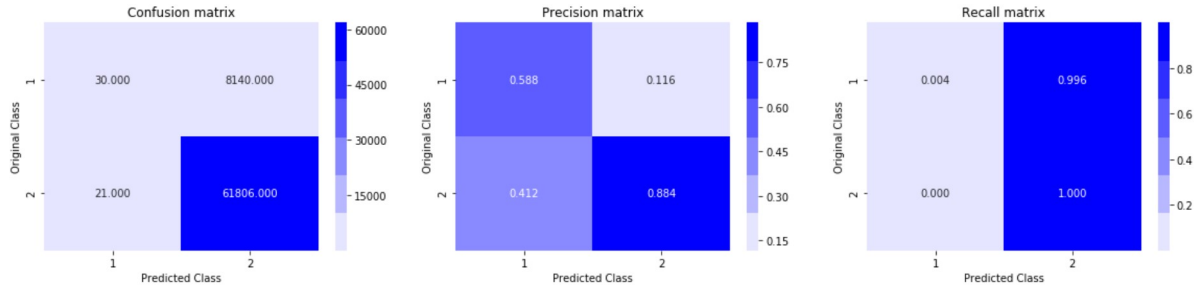
Recall on train set: 1.000

F1-Score on train set: 0.883

Non Zero weights: 50

Confusion Matrix of test set:

```
[ [TN  FP]
 [FN TP] ]
```



Accuracy on test set: 86.360%

Precision on test set: 0.864

Recall on test set: 1.000

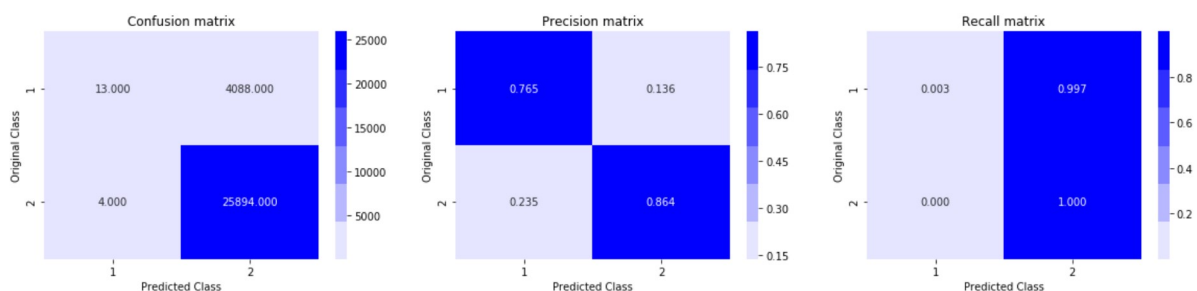
F1-Score on test set: 0.864

Non Zero weights: 50

Confusion Matrix of test set:

```
[ [TN  FP]
 [FN TP] ]
```

Confusion matrix for the model is:



Using Randomized Search CV to find best parameter

```
In [134]: %time
from sklearn.model_selection import RandomizedSearchCV
from sklearn.linear_model import LogisticRegression

clf = LogisticRegression()
#params we need to try on classifier
param_grid = {'C': [100,50,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001,0.0005,0.0001]
              , 'penalty':['l1','l2']}
tscv = TimeSeriesSplit(n_splits=10) #For time based splitting
gsv = RandomizedSearchCV(clf,param_grid,cv=tscv,verbose=1,scoring='f1_micro')
gsv.fit(X_train,y_train)
savetofile(gsv,"Log Reg/gsv_w2v_r")
print("Best HyperParameter: ",gsv.best_params_)
print("Best Accuracy: %.2f%%"%(gsv.best_score_*100))
```

Wall time: 0 ns

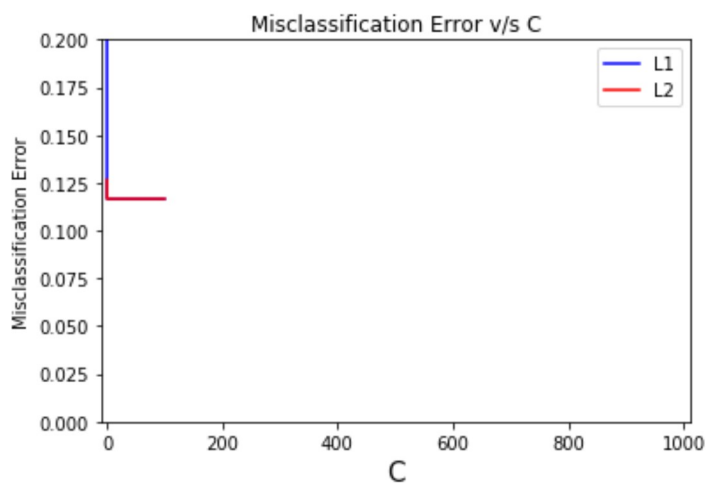
Fitting 10 folds for each of 10 candidates, totalling 100 fits

[Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 24.8s finished

Best HyperParameter: {'penalty': 'l2', 'C': 100}

Best Accuracy: 88.34%

```
In [135]: gsv = openfromfile("Log Reg/gsv_w2v_r")
plot_error_vs_c_r(gsv)
```



```
In [136]: clf = LogisticRegression(C= 100, penalty= 'l2')
clf.fit(X_train,y_train)
y_train_pred = clf.predict(X_train)
y_pred = clf.predict(X_test)
print("Accuracy on train set: %0.3f%%"%(accuracy_score(y_train, y_train_pred)*100)
)
print("Precision on train set: %0.3f"%(precision_score(y_train, y_train_pred)))
print("Recall on train set: %0.3f"%(recall_score(y_train, y_train_pred)))
print("F1-Score on train set: %0.3f"%(f1_score(y_train, y_train_pred,average='micro'
o))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
plot_confusion_matrix(y_train, y_train_pred)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred,average='micro'))))
print("Non Zero weights:",np.count_nonzero(clf.coef_))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
print('Confusion matrix for the model is:')
plot_confusion_matrix(y_test, y_pred)
```

Accuracy on train set: 88.341%

Precision on train set: 0.884

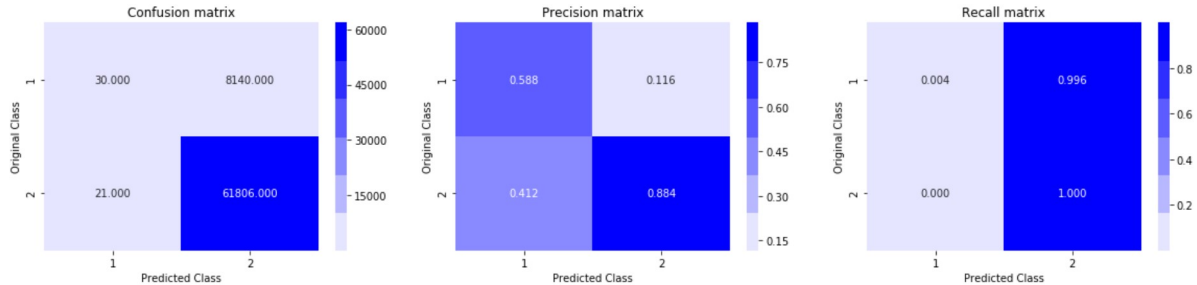
Recall on train set: 1.000

F1-Score on train set: 0.883

Non Zero weights: 50

Confusion Matrix of test set:

```
[ [TN  FP]
 [FN TP] ]
```



Accuracy on test set: 86.360%

Precision on test set: 0.864

Recall on test set: 1.000

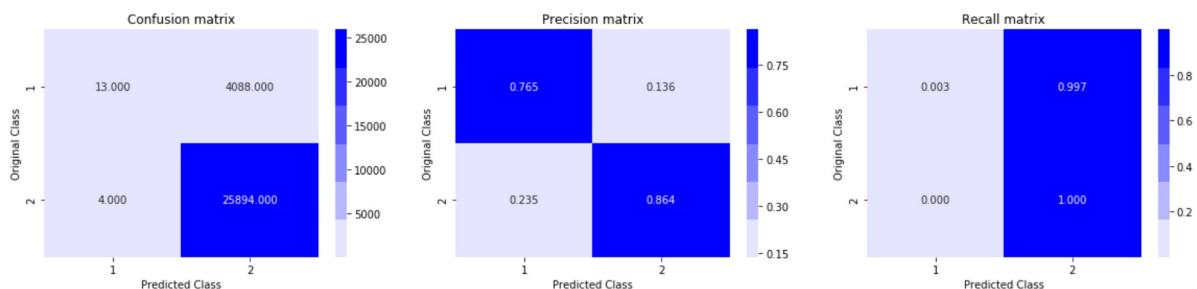
F1-Score on test set: 0.864

Non Zero weights: 50

Confusion Matrix of test set:

```
[ [TN  FP]
 [FN TP] ]
```

Confusion matrix for the model is:



Conclusion

In [181]: `from prettytable import PrettyTable`

```
# Names of models
featurization = ['Bag of Words', 'Bag of Words', 'bigram', 'bigram', 'TFIDF ', 'TFIDF ',
                 'avg w2v', 'avg w2vw']
model=['gridsearch ', 'randomsearch', 'gridsearch ', 'randomsearch', 'gridsearch ', 'ra
ndomsearch', 'gridsearch ', 'randomsearch']

# Training accuracies
Flscore= [0.916,0.892,0.922,0.922,0.967,0.914,0.864,0.864]
accuracy = [91.57,89.24,92.17,92.17,91.39,91.39,86.36,86.36]
alpha=[0.05,0.001,0.01,0.01,0.05,0.05,0.1,100]
precision=[0.938,0.897,0.929,0.929,0.935,0.936,0.864,0.864]
recall=[0.967,0.988,0.985,0.985,0.967,0.967,1.00,1.00]
numbering = [1,2,3,4,5,6,7,8]
regularization=['11', '12', '11', '11', '11', '11', '12', '12']

# Initializing prettytable
ptable = PrettyTable()

# Adding columns
ptable.add_column("S.NO.", numbering)
ptable.add_column("MODEL", featurization)
ptable.add_column("Method", model)
ptable.add_column("C", alpha)
ptable.add_column("regularization", regularization)
ptable.add_column("accuracy", accuracy)
ptable.add_column("score", Flscore)
ptable.add_column("precision", precision)
ptable.add_column("recall", recall)

# Printing the Table
print(ptable)
```

S.NO.	MODEL	Method	C	regularization	accuracy	score	precision	recall
1	Bag of Words	gridsearch	0.05	11	91.57	0.916	0.938	0.967
2	Bag of Words	randomsearch	0.001	12	89.24	0.892	0.897	0.988
3	bigram	gridsearch	0.01	11	92.17	0.922	0.929	0.985
4	bigram	randomsearch	0.01	11	92.17	0.922	0.929	0.985
5	TFIDF	gridsearch	0.05	11	91.39	0.914	0.935	0.967
6	TFIDF	randomsearch	0.05	11	91.39	0.914	0.936	0.967
7	avg w2v	gridsearch	0.1	12	86.36	0.864	0.864	1.0
8	avg w2vw	randomsearch	100	12	86.36	0.864	0.864	1.0

In []: