

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>)

EDA: <https://nycdatascience.com/blog/student-works/amazon-fine-foods-visualization/> (<https://nycdatascience.com/blog/student-works/amazon-fine-foods-visualization/>)

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:

We want to train a model such that it is able to classify the incoming data point which is a review text into positive review and negative review. For this task we will take in consideration of the review text and will work on it using different vectorizers like Bag of words, tfidf, and word to vector to generate features that can be feeded to our model for making predictions

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

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

[1]. Reading Data

[1.1] 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 [0]: 1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5
6 import sqlite3
7 import pandas as pd
8 import numpy as np
9 import nltk
10 import string
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 from sklearn.feature_extraction.text import TfidfTransformer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.metrics import confusion_matrix
18 from sklearn import metrics
19 from sklearn.metrics import roc_curve, auc
20 from nltk.stem.porter import PorterStemmer
21
22 import re
23 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
24 import string
25 from nltk.corpus import stopwords
26 from nltk.stem import PorterStemmer
27 from nltk.stem.wordnet import WordNetLemmatizer
28
29 from gensim.models import Word2Vec
30 from gensim.models import KeyedVectors
31 import pickle
32
33 from tqdm import tqdm
34 import os
```

```
In [2]: 1 from google.colab import drive
        2 drive.mount('/content/gdrive')
```

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awg%3Aoauth%3A2.0%3Aoob&scope=email%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdocs.test%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fpeopleapi.readonly&response_type=code (https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awg%3Aoauth%3A2.0%3Aoob&scope=email%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdocs.test%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fpeopleapi.readonly&response_type=code)

Enter your authorization code:

.....

Mounted at /content/gdrive

```

In [3]: 1 # using SQLite Table to read data.
2 con = sqlite3.connect('gdrive/My Drive/database.sqlite')
3 filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3""", con)
4 # Give reviews with Score>3 a positive rating(1), and reviews with a score<3 a negative rating(0).
5 def partition(x):
6     if x < 3:
7         return 0
8     return 1
9 #changing reviews with score less than 3 to be positive and vice-versa
10 actualScore = filtered_data['Score']
11 positiveNegative = actualScore.map(partition)
12 filtered_data['Score'] = positiveNegative
13 print("Number of data points in our data", filtered_data.shape)
14 filtered_data.head(3)

```

Number of data points in our data (525814, 10)

Out[3]:

		Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian		1	1	1	1303862400	Good Quality Dog Food	Good Quality Dog Food
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa		0	0	0	1346976000	Not as Advertised	Not as Advertised
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres	"Natalia Corres"	1	1	1	1219017600	"Delight" says it all	"Delight" says it all

```
In [0]: 1 display = pd.read_sql_query("""
2 SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
3 FROM Reviews
4 GROUP BY UserId
5 HAVING COUNT(*)>1
6 """, con)
```

```
In [5]: 1 print(display.shape)
2 display.head()
```

(80668, 7)

```
Out[5]:
```

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
0	#oc-R115TNMSPFT9I7	B007Y59HVM	Breyton	1331510400	2	Overall its just OK when considering the price...	2
1	#oc-R11D9D7SHXIJB9	B005HG9ET0	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms, u...	3
2	#oc-R11DNU2NBKQ23Z	B007Y59HVM	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately not ...	2
3	#oc-R11O5J5ZVQE25C	B005HG9ET0	Penguin Chick	1346889600	5	This will be the bottle that you grab from the...	3
4	#oc-R12KPBODL2B5ZD	B007OSBE1U	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2

```
In [6]: 1 display[display['UserId']=='AZY10LLTJ71NX']
```

```
Out[6]:
```

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
80638	AZY10LLTJ71NX	B006P7E5ZI	undertheshrine "undertheshrine"	1334707200	5	I was recommended to try green tea extract to ...	5

```
In [7]: 1 display['COUNT(*)'].sum()
```

Out[7]: 393063

[2] Exploratory Data Analysis

[2.1] 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 [8]:

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5 ORDER BY ProductID
6 """, con)
7 display.head()
```

Out[8]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS

As it can be seen above that same user has multiple reviews with same values for HelpfulnessNumerator, HelpfulnessDenominator, Score,

Time, Summary and Text and on doing analysis it was found that

ProductId=B000HDOPZG was Loacker Quadratini Vanilla Wafer Cookies, 8.82-Ounce Packages (Pack of 8)

ProductId=B000HDL1RQ was Loacker Quadratini Lemon Wafer Cookies, 8.82-Ounce Packages (Pack of 8) and so on

It was inferred after analysis that reviews with same parameters other than ProductId belonged to the same product just having different flavour or quantity. Hence in order to reduce redundancy it was decided to eliminate the rows having same parameters.

The method used for the same was that we first sort the data according to ProductId and then just keep the first similar product review and delete the others. for eg. in the above just the review for ProductId=B000HDL1RQ remains. This method ensures that there is only one representative for each product and deduplication without sorting would lead to possibility of different representatives still existing for the same product.

```
In [0]: 1 #Sorting data according to ProductId in ascending order
        2 sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False, kind='quicksort',
```

```
In [10]: 1 #Deduplication of entries
        2 final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first', inplace=False)
        3 final.shape
```

Out[10]: (364173, 10)

```
In [11]: 1 #Checking to see how much % of data still remains
        2 (final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

Out[11]: 69.25890143662969

Observation:- It was also seen that in two rows given below the value of HelpfulnessNumerator is greater than HelpfulnessDenominator which is not practically possible hence these two rows too are removed from calculations


```
In [12]: 1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND Id=44737 OR Id=64422
5 ORDER BY ProductID
6 """, con)
7
8 display.head()
```

```
Out[12]:
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5	1224892800	Bought This for My Son at College
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2	4	1212883200	Pure cocoa taste with crunchy almonds inside

```
In [0]: 1 final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [14]: 1 #Before starting the next phase of preprocessing Lets see the number of entries Left
2 print(final.shape)
3
4 #How many positive and negative reviews are present in our dataset?
5 final['Score'].value_counts()
```

```
(364171, 10)
```

```
Out[14]: 1 307061
0 57110
Name: Score, dtype: int64
```

Taking 100k datapoints

```
In [0]: 1 final = final.sample(100000)#sampling 100k datapoints
```

```
In [16]: 1 final = final.sort_values('Time',ascending = True)
        2 final
```

```
Out[16]:
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
346094	374400	B00004CI84	A2DEE7F9XKP3ZR	jerome	0	3	1	95999040
1145	1244	B00002Z754	A3B8RCEI0FXFI6	B G Chase	10	10	1	96223680
138001	149770	B00004S1C5	A1KXONFPU2XQ5K	Stephanie Manley	8	8	1	96577920
346115	374421	B00004CI84	A1FJOY14X3MUHE	Justin Howard	2	2	1	96629760

[3] Preprocessing

[3.1]. Preprocessing Review Text

Now that we have finished deduplication our data requires some preprocessing before we go on further with analysis and making the prediction model.

Hence in the Preprocessing phase we do the following in the order below:-

1. Begin by removing the html tags
2. Remove any punctuations or limited set of special characters like , or . or # etc.
3. Check if the word is made up of english letters and is not alpha-numeric

4. Check to see if the length of the word is greater than 2 (as it was researched that there is no adjective in 2-letters)
5. Convert the word to lowercase
6. Remove Stopwords
7. Finally Snowball Stemming the word (it was observed to be better than Porter Stemming)

After which we collect the words used to describe positive and negative reviews

```
In [17]: 1 # printing some random reviews
2 sent_0 = final['Text'].values[0]
3 print(sent_0)
4 print("="*50)
5
6 sent_1000 = final['Text'].values[1000]
7 print(sent_1000)
8 print("="*50)
9
10 sent_1500 = final['Text'].values[1500]
11 print(sent_1500)
12 print("="*50)
13
14 sent_4900 = final['Text'].values[4900]
15 print(sent_4900)
16 print("="*50)
```

I'm getting crazy.I'm looking for Beatlejuice french version video.<p>Is it really impossible today not to find the French VHS version of this film ?<p>Could U please tell me something about it ? Tks

=====

This mixer is easy to use and makes using natural peanut butter so much pleasurable. It mixed the peanut butter thoroughly and clean up was easy. The peanut butter was much better mixed than when we attempted to mix it ourselves with utensils.

=====

These delicious mini-crackers are small enough to grab a satisfying handful. Yep, it's got all the organic requisites, organically-fed cows whose lips never taste hormones or antibiotics (well, maybe these are injected, I don't know), and the crackers contain non-hydrogenated, trans-fat free organic palm oil.

The taste is mild but you can taste the cheddar, and the texture is more flakey than hard or brittle. If you grab 30 crackers (one ounce, there are 5 ounces, or 150 crackers per box), you'll get 310 mg of sodium, or 13% of the Recommended Daily Value (RDV) for a 2,000 calorie diet. Now that's by any means low-sodium, but you can easily eat half of that and feel you've had a snack. Per 30 crackers, there are 4.5 grams (g) of total fat, of which 1.5 are saturated, and none are trans fat. Total carbohydrates are 19g, including less than 1g of sugar, but also less than 1g of dietary fiber. Hey, no one said this is health food, but compare this to what you're snacking on now (and that cheese gives it 3g of protein).

Full ingredients and nutritional information may be found at latejuly.com/products, but other good stuff includes organic wheat flour, organic whey, sea salt, and organic evaporated cane juice. Nothing artificial, no corn syrup, preservatives, hydrogenated oils, and they're both Kosher and Lacto Vegetarian (gotta make sure there's no meat in your crackers...). If that's not enough, "Late July" is independently owned and family operated.

Mainly, though, it's the light taste, the organic ingredients, and that enticing small size you can grab by the fistful. Full information at the above-mentioned website.

=====

My puppy goes crazy for this stuff! Its kinda gross to us humans but the dogs love the stuff!

=====

```
In [18]: 1 # remove urls from text python: https://stackoverflow.com/a/40823105/4084039
2 sent_0 = re.sub(r"http\S+", "", sent_0)
3 sent_1000 = re.sub(r"http\S+", "", sent_1000)
4 sent_150 = re.sub(r"http\S+", "", sent_1500)
5 sent_4900 = re.sub(r"http\S+", "", sent_4900)
6
7 print(sent_0)
```

I'm getting crazy.I'm looking for Beatlejuice french version video.<p>Is it really impossible today not to find the French VHS version of this film ?<p>Could U please tell me something about it ? Tks

```
In [19]: 1 # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all-tags-from-an-element
2 from bs4 import BeautifulSoup
3
4 soup = BeautifulSoup(sent_0, 'lxml')
5 text = soup.get_text()
6 print(text)
7 print("="*50)
8
9 soup = BeautifulSoup(sent_1000, 'lxml')
10 text = soup.get_text()
11 print(text)
12 print("="*50)
13
14 soup = BeautifulSoup(sent_1500, 'lxml')
15 text = soup.get_text()
16 print(text)
17 print("="*50)
18
19 soup = BeautifulSoup(sent_4900, 'lxml')
20 text = soup.get_text()
21 print(text)
```

I'm getting crazy.I'm looking for Beatlejuice french version video.Is it really impossible today not to find the French VHS version of this film ?Could U please tell me something about it ? Tks

=====

This mixer is easy to use and makes using natural peanut butter so much pleasurable. It mixed the peanut butter thoroughly and clean up was easy. The peanut butter was much better mixed than when we attempted to mix it ourselves with utensils.

=====

These delicious mini-crackers are small enough to grab a satisfying handful. Yep, it's got all the organic requisites, organically-fed cows whose lips never taste hormones or antibiotics (well, maybe these are injected, I don't know), and the crackers contain non-hydrogenated, trans-fat free organic palm oil. The taste is mild but you can taste the cheddar, and the texture is more flakey than hard or brittle. If you grab 30 crackers (one ounce, there are 5 ounces, or 150 crackers per box), you'll get 310 mg of sodium, or 13% of the Recommended Daily Value (RDV) for a 2,000 calorie diet. Now that's by any means low-sodium, but you can easily eat half of that and feel you've had a snack. Per 30 crackers, there are 4.5 grams (g) of total fat, of which 1.5 are saturated, and none are trans fat. Total carbohydrates are 19g, including less than 1g of sugar, but also less than 1g of dietary fiber. Hey, no one said this is health food, but compare this to what you're snacking on now (and that cheese gives it 3g of protein). Full ingredients and nutritional information may be found at latejuly.com/products, but other good stuff includes organic wheat flour, organic whey, sea salt, and organic evaporated cane juice. Nothing artificial, no corn syrup, preservatives, hydrogenated oils, and they're both Kosher and Lacto Vegetarian (gotta make sure there's no meat in your crackers...). If that's not enough, "Late July" is independently owned

ed and family operated. Mainly, though, it's the light taste, the organic ingredients, and that enticing small size you can grab by the fistful. Full information at the above-mentioned website.

=====

My puppy goes crazy for this stuff! Its kinda gross to us humans but the dogs love the stuff!

```
In [0]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\ 're", " are", phrase)
12    phrase = re.sub(r"\ 's", " is", phrase)
13    phrase = re.sub(r"\ 'd", " would", phrase)
14    phrase = re.sub(r"\ 'll", " will", phrase)
15    phrase = re.sub(r"\ 't", " not", phrase)
16    phrase = re.sub(r"\ 've", " have", phrase)
17    phrase = re.sub(r"\ 'm", " am", phrase)
18    return phrase
```

```
In [21]: 1 sent_1500 = decontracted(sent_1500)
          2 print(sent_1500)
          3 print("="*50)
```

These delicious mini-crackers are small enough to grab a satisfying handful. Yep, it is got all the organic requisites, organically-fed cows whose lips never taste hormones or antibiotics (well, maybe these are injected, I do not know), and the crackers contain non-hydrogenated, trans-fat free organic palm oil.

The taste is mild but you can taste the cheddar, and the texture is more flakey than hard or brittle. If you grab 30 crackers (one ounce, there are 5 ounces, or 150 crackers per box), you will get 310 mg of sodium, or 13% of the Recommended Daily Value (RDV) for a 2,000 calorie diet. Now that is by any means low-sodium, but you can easily eat half of that and feel you have had a snack. Per 30 crackers, there are 4.5 grams (g) of total fat, of which 1.5 are saturated, and none are trans fat. Total carbohydrates are 19g, including less than 1g of sugar, but also less than 1g of dietary fiber. Hey, no one said this is health food, but compare this to what you are snacking on now (and that cheese gives it 3g of protein).

Full ingredients and nutritional information may be found at [latejuly.com/products](https://www.latejuly.com/products), but other good stuff includes organic wheat flour, organic whey, sea salt, and organic evaporated cane juice. Nothing artificial, no corn syrup, preservatives, hydrogenated oils, and they are both Kosher and Lacto Vegetarian (gotta make sure there is no meat in your crackers...). If that is not enough, "Late July" is independently owned and family operated.

Mainly, though, it is the light taste, the organic ingredients, and that enticing small size you can grab by the fistful. Full information at the above-mentioned website.

=====

```
In [22]: 1 #remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
          2 sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
          3 print(sent_0)
```

I'm getting crazy. I'm looking for Beatlejuice french version video.<p>Is it really impossible today not to find the French VHS version of this film ?<p>Could U please tell me something about it ? Tks


```
In [23]: 1 #remove spacial character: https://stackoverflow.com/a/5843547/4084039
2 sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
3 print(sent_1500)
```

These delicious mini crackers are small enough to grab a satisfying handful Yep it is got all the organic requi sites organically fed cows whose lips never taste hormones or antibiotics well maybe these are injected I do no t know and the crackers contain non hydrogenated trans fat free organic palm oil br br The taste is mild but yo u can taste the cheddar and the texture is more flakey than hard or brittle If you grab 30 crackers one ounce t here are 5ounces or 150 crackers per box you will get 310 mg of sodium or 13 of the Recommended Daily Value RDV for a 2 000 calorie diet Now that is by any means low sodium but you can easily eat half of that and feel you h ave had a snack Per 30 crackers there are 4 5 grams g ofttotal fat of which 1 5 are saturated and none are trans fat Total carbohydrates are 19g including less than 1g of sugar but also less than 1g of dietary fiber Hey noon e said this is health food but compare this to what you are snacking on now and that cheese gives it 3g of prot ein br br Full ingredients and nutritional information may be found at latejuly com products but other good stu ff includes organic wheat flour organic whey sea salt and organic evaporated cane juice Nothing artificial no c orn syrup preservatives hydrogenated oils and they are both Kosher and Lacto Vegetarian gotta make sure there i s no meat in your crackers If that is not enough Late July is independently owned and family operated br br Mai nly though it is the light taste the organic ingredients and that enticing small size you can grab by the fistf ul Full information at the above mentioned website

```
In [0]: 1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'
3 # <br /><br /> ==> after the above steps, we are getting "br br"
4 # we are including them into stop words list
5 # instead of <br /> if we have <br/> these tags would have revmoved in the 1st step
6
7 stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "
8     "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
9     'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their'
10    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'tho
11    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', '
12    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', '
13    'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before',
14    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again'
15    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'f
16    'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
17    's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', '
18    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't",
19    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mus
20    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'were
21    'won', "won't", 'wouldn', "wouldn't"])
```

```
In [25]: 1 # Combining all the above students
2 from tqdm import tqdm
3 preprocessed_reviews = []
4 # tqdm is for printing the status bar
5 for sentence in tqdm(final['Text'].values):
6     sentence = re.sub(r"http\S+", "", sentence)
7     sentence = BeautifulSoup(sentence, 'lxml').get_text()
8     sentence = decontracted(sentence)
9     sentence = re.sub("\S*\d\S*", "", sentence).strip()
10    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
11    # https://gist.github.com/sebleier/554280
12    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
13    preprocessed_reviews.append(sentence.strip())
```

100%|██████████| 100000/100000 [00:44<00:00, 2241.36it/s]

```
In [26]: 1 preprocessed_reviews[1500]
```

Out[26]: 'delicious mini crackers small enough grab satisfying handful yep got organic requisites organically fed cows w hose lips never taste hormones antibiotics well maybe injected not know crackers contain non hydrogenated trans fat free organic palm oil taste mild taste cheddar texture flakey hard brittle grab crackers one ounce crackers per box get mg sodium recommended daily value rdv calorie diet means low sodium easily eat half feel snack per crackers grams g of total fat saturated none trans fat total carbohydrates including less sugar also less dietary fiber hey noone said health food compare snacking cheese gives protein full ingredients nutritional information may found latejuly com products good stuff includes organic wheat flour organic whey sea salt organic evaporated cane juice nothing artificial no corn syrup preservatives hydrogenated oils kosher lacto vegetarian gotta make sure no meat crackers not enough late july independently owned family operated mainly though light taste organic ingredients enticing small size grab fistful full information mentioned website'

[3.2] Preprocessing Review Summary

```
In [27]: 1 from sklearn.model_selection import train_test_split
2
3 X_train,X_test,Y_train,Y_test = train_test_split(preprocessed_reviews,final['Score'],test_size=0.3,random_st
4 print('After Splitting size of:\n')
5 print('Train Dataset is:',np.array(X_train).shape)
6 print('Test Dataset is:',np.array(X_test).shape)
```

After Splitting size of:

Train Dataset is: (70000,)

Test Dataset is: (30000,)

[4] Featurization

READING THE FILES

BAG OF WORDS

```
In [28]: 1 bow_vect = CountVectorizer()#initiating the vectorizer
2 bow_vect.fit(X_train)#fittin the data into vectorizer makes it learn from the corpus
3
4 train_bow = bow_vect.transform(X_train)
5 test_bow = bow_vect.transform(X_test)
6
7 print('TRAINING SET FOR BAG OF WORDS',train_bow.shape)
8 print('TEST SET FOR BAG OF WORDS',test_bow.shape)
```

TRAINING SET FOR BAG OF WORDS (70000, 50640)

TEST SET FOR BAG OF WORDS (30000, 50640)

```
In [0]: 1 import pickle
2 pickle.dump(train_bow,open('train_bow.p','wb'))
3 pickle.dump(test_bow,open('test_bow.p','wb'))
```

TFIDF

```
In [0]: 1 tfidf_vect = TfidfVectorizer(ngram_range = (1,2),min_df = 10)
2 tfidf_vect.fit(X_train)
3
4 train_tfidf = tfidf_vect.transform(X_train)
5 test_tfidf = tfidf_vect.transform(X_test)
6
7
```

```
In [0]: 1 pickle.dump(train_tfidf,open('train_tfidf.p','wb'))
2 pickle.dump(test_tfidf,open('test_tfidf.p','wb'))
```

Word 2 vector

```
In [0]: 1 s_train = []
2 for sent in X_train:
3     s_train.append(sent.split())
4 #preparing the training data for word to vector vectorization
5
6 s_test = []
7 for sent in X_test:
8     s_test.append(sent.split())
9 #preparing the test data for word to vector fatorization
```

```
In [33]: 1 w2v_model=Word2Vec(s_train,min_count=5,size=50, workers=4)# min_count = 5 considers only words that occurred
2
3 w2v_words = list(w2v_model.wv.vocab)
4 print("number of words that occurred minimum 5 times ",len(w2v_words))
5
```

number of words that occurred minimum 5 times 16090

Average Word 2 Vector

In [34]:

```
1  #Average word2vec
2  #computing average word to vector for training data
3
4  train_set = [] # the avg-w2v for each sentence/review is stored in this list
5  for sent in tqdm(s_train):
6      sent_vec = np.zeros(50)
7      cnt_words = 0; # num of words with a valid vector in the sentence/review
8      for word in sent: #
9          if word in w2v_words:
10             vec = w2v_model.wv[word]
11             sent_vec += vec
12             cnt_words += 1
13     if cnt_words != 0:
14         sent_vec /= cnt_words
15     train_set.append(sent_vec)
16
17 print(len(train_set))#number of datapoints in training set
```

100%|██████████| 70000/70000 [02:20<00:00, 499.37it/s]

70000

```
In [35]: 1 #computing average word to vector for test data
2
3 test_set = [] # the avg-w2v for each sentence/review is stored in this list
4 for sent in s_test:
5     sent_vec = np.zeros(50)
6     cnt_words = 0; # num of words with a valid vector in the sentence/review
7     for word in sent: #
8         if word in w2v_words:
9             vec = w2v_model.wv[word]
10            sent_vec += vec
11            cnt_words += 1
12        if cnt_words != 0:
13            sent_vec /= cnt_words
14        test_set.append(sent_vec)
15
16 print(len(test_set))#number of datapoints in test set
```

30000

```
In [0]: 1 pickle.dump(train_set,open('train_avgw2v.p','wb'))
2        pickle.dump(test_set,open('test_avgw2v.p','wb'))
```

TFIDF WORD 2 VECTOR

```
In [0]: 1 #tfidf word 2 vec
2
3 vect = TfidfVectorizer()#initializing the tfidf vectorizer
4
5 tf_idf = vect.fit_transform(X_train)#fitting the training data
6 dictionary = dict(zip(vect.get_feature_names(), list(vect.idf_)))#zipping both of the feature names and vect
```

```
In [38]: 1 import itertools
          2 dict(itertools.islice(dictionary.items(),20))
          3 #printing first 20 elements of the dictionary
```

```
Out[38]: {'aa': 9.853679713649695,
          'aaa': 10.076823264963906,
          'aaaa': 11.057652517975633,
          'aaaaa': 10.769970445523851,
          'aaaaaaaaaaaa': 11.463117626083797,
          'aaaaaaaaaaaaaa': 11.463117626083797,
          'aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa': 11.463117626083797,
          'aaaaaaahhhhhh': 11.463117626083797,
          'aaaaaah': 11.057652517975633,
          'aaaaaahhh': 11.463117626083797,
          'aaaaaahhhh': 11.463117626083797,
          'aaaaaawwwwwwwww': 11.463117626083797,
          'aaaah': 11.463117626083797,
          'aaaahhhhhhhhhhhh': 11.463117626083797,
          'aaaannnnddd': 11.463117626083797,
          'aaah': 11.057652517975633,
          'aadults': 11.463117626083797,
          'aafco': 10.210354657588429,
          'aafes': 10.769970445523851,
          'aah': 11.463117626083797}
```

```
In [39]: 1 tfidf_feat = vect.get_feature_names() # tfidf words/col-names
          2 print(tfidf_feat[:20])
```

```
['aa', 'aaa', 'aaaa', 'aaaaa', 'aaaaaaaaaaaa', 'aaaaaaaaaaaaaa', 'aaaaaaaaaaaaaaaaaaaaaaaaaaaaa', 'aaaaaaahhhhhh', 'aaaaaah', 'aaaaaahhh', 'aaaaaahhhh', 'aaaaaawwwwwwwww', 'aaaah', 'aaaahhhhhhhhhhhh', 'aaaannnnddd', 'aaah', 'aadults', 'aafco', 'aafes', 'aah']
```

```
In [40]: 1 train_set_tfidfw2v = []; # the tfidf-w2v for each sentence/review in training set is stored in this list
2 row=0;
3 for sent in tqdm(s_train): # for each review/sentence
4     sent_vec = np.zeros(50) # as word vectors are of zero length
5     weight_sum =0; # num of words with a valid vector in the sentence/review
6     for word in sent: # for each word in a review/sentence
7         if word in w2v_words and word in tfidf_feat:
8             vec = w2v_model.wv[word]
9             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
10            sent_vec += (vec * tf_idf)
11            weight_sum += tf_idf
12        if weight_sum != 0:
13            sent_vec /= weight_sum
14        train_set_tfidfw2v.append(sent_vec)
15        row += 1
16    print(len(train_set_tfidfw2v))
```

100%|██████████| 70000/70000 [30:52<00:00, 37.79it/s]

70000


```
In [41]: 1 test_set_tfidfw2v = []; # the tfidf-w2v for each sentence/review in test set is stored in this list
2 row=0;
3 for sent in tqdm(s_test): # for each review/sentence
4     sent_vec = np.zeros(50) # as word vectors are of zero length
5     weight_sum =0; # num of words with a valid vector in the sentence/review
6     for word in sent: # for each word in a review/sentence
7         if word in w2v_words and word in tfidf_feat:
8             vec = w2v_model.wv[word]
9             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
10            sent_vec += (vec * tf_idf)
11            weight_sum += tf_idf
12        if weight_sum != 0:
13            sent_vec /= weight_sum
14        test_set_tfidfw2v.append(sent_vec)
15        row += 1
16
17 print(len(test_set_tfidfw2v))
18
```

100%|██████████| 30000/30000 [13:20<00:00, 37.47it/s]

30000

```
In [0]: 1 import pickle
2 pickle.dump(train_set_tfidfw2v,open('train_tfidfw2v.p','wb'))
3 pickle.dump(test_set_tfidfw2v,open('test_tfidfw2v.p','wb'))
```

1. Apply SVM on these feature sets

- **SET 1:**Review text, preprocessed one converted into vectors using (BOW)
- **SET 2:**Review text, preprocessed one converted into vectors using (TFIDF)
- **SET 3:**Review text, preprocessed one converted into vectors using (AVG W2v)
- **SET 4:**Review text, preprocessed one converted into vectors using (TFIDF W2v)

2. Procedure

- You need to work with 2 versions of SVM
 - Linear kernel

- RBF kernel
- When you are working with linear kernel, use SGDClassifier' with hinge loss because it is computationally less expensive.
- When you are working with 'SGDClassifier' with hinge loss and trying to find the AUC score, you would have to use [CalibratedClassifierCV \(https://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html)
- Similarly, like kdtree or knn, when you are working with RBF kernel it's better to reduce the number of dimensions. You can put min_df = 10, max_features = 500 and consider a sample size of 40k points.

3. Hyper parameter tuning (find best alpha in range $[10^{-4}$ to 10^4], and the best penalty among 'l1', 'l2')

- Find the best hyper parameter which will give the maximum [AUC \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value
- Find the best hyper parameter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning



4. Feature importance

- When you are working on the linear kernel with BOW or TFIDF please print the top 10 best features for each of the positive and negative classes.

5. Feature engineering

- To increase the performance of your model, you can also experiment with feature engineering like :
 - Taking length of reviews as another feature.
 - Considering some features from review summary as well.

6. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
 
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.
 
- Along with plotting ROC curve, you need to print the [confusion matrix \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).



(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

7. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library [link \(http://zetcode.com/python/prettytable/\)](http://zetcode.com/python/prettytable/)



Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakag, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on you train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link. \(https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf\)](https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf)

Applying SVM

[5.1] Linear SVM

[5.1.1] Applying Linear SVM on BOW, SET 1

```
In [0]: 1 from sklearn.model_selection import GridSearchCV
2 from sklearn.linear_model import SGDClassifier
3 from sklearn.calibration import CalibratedClassifierCV
4 from sklearn.calibration import calibration_curve
5 from sklearn.model_selection import TimeSeriesSplit
6 from sklearn.metrics import roc_auc_score
7 from sklearn.metrics import confusion_matrix
8 from sklearn.metrics import f1_score
9 from sklearn.metrics import recall_score
10 from sklearn.metrics import precision_score
11 from sklearn.metrics import roc_curve
12
13
14
15
```

```
In [0]: 1 #opening adn reading all the saved files\
2
3 #BAG OF VECTORS
4 train_bow = pickle.load(open('train_bow.p', 'rb'))
5 test_bow = pickle.load(open('test_bow.p', 'rb'))
6
7
8
9
10
```

In [0]:

```

1  #function for tuning the hyperparameters
2  alpha = [10**i for i in range(-4,5,1)]
3  tscv = TimeSeriesSplit(n_splits = 5)# for times series cross validation
4  params = {'alpha':alpha}
5
6  def svm(regularization,train_set,train_y):
7      clf = SGDClassifier(penalty = regularization,loss = 'hinge',random_state = 42)
8      #we will be checking for both l1 and l2 regularizations
9      gsm = GridSearchCV(clf,param_grid = params,verbose = 1,cv = tscv,scoring = 'roc_auc',return_train_score
10     #cv = tscv does cross validation according to time series split
11     gsm.fit(train_set,train_y)
12     return gsm
13  #*****
14
15  def error(regularization,train_set,train_y):
16      print('Corresponding hyperparameter will give best auc on CV data')
17      gsm = svm(regularization,train_set,train_y)
18      best_alpha = gsm.best_params_
19      t_auc = gsm.cv_results_['mean_train_score']
20      cv_auc = gsm.cv_results_['mean_test_score']
21      print('best hyperparameter is',best_alpha)
22      sns.set_style('darkgrid')
23      plt.figure(figsize=(15,6))
24      plt.plot(alpha,t_auc,'g',label = 'training AUC')#t_auc refers to the auc on training data
25      plt.plot(alpha,cv_auc,'r',label='validation AUC')# c_auc refers to the auc on cross validation data
26      #plotting the graph between AUC and hyperparameter for tuning
27
28      plt.xscale('log')#taking log scale for x axis for better analysing the results
29      plt.xlabel('hyperparameter Alpha',fontsize=18)
30      plt.ylabel('Area under curve',fontsize=18)
31      #plt.xticks([])
32      #plt.yticks([])
33      plt.legend(loc = 'best')
34      plt.title('AUC vs hyperparameter ',fontsize=18)
35      return best_alpha
36
37
38  #*****
39
40  def best_classifier(best_alpha,regularization,train_set,train_y,test_set,test_y):
41      clf_optimal = SGDClassifier(alpha = best_alpha,penalty = regularization,loss = 'hinge',random_state = 42)
42      clf_optimal.fit(train_set,train_y)

```

```

43 clb = CalibratedClassifierCV(clf_optimal,cv = 5,method = 'sigmoid')
44 clb.fit(train_set,train_y)
45 train_proba = clb.predict_proba(train_set)[: ,1]
46 test_proba = clb.predict_proba(test_set)[: ,1]
47
48 pred_tr = clb.predict(train_set)
49 pred_test = clb.predict(test_set)
50
51 train_auc = roc_auc_score(train_y,train_proba)
52 test_auc = roc_auc_score(test_y,test_proba)
53
54 print('AUC on training data is',train_auc)
55 print('AUC on test data is',test_auc)
56 return train_auc,test_auc,train_proba,test_proba,pred_tr,pred_test
57
58 *****
59
60 #computing function for reliability curve
61 #referred to :https://machinelearningmastery.com/calibrated-classification-model-in-scikit-learn/
62 def reliability_curve(best_alpha,regularization,trainX, testX, train_y,test_y):
63     def uncalibrated(best_alpha,regularization,trainX, testX, train_y):
64         # fit a model
65         model = SGDClassifier(alpha = best_alpha,penalty = regularization,loss = 'hinge',random_state = 42)
66         model.fit(trainX, train_y)
67         # predict probabilities
68         return model.decision_function(testX)
69
70     #predict calibrated probabilities
71     def calibrated(best_alpha,regularization,trainX, testX, train_y):
72         # define model
73         model = SGDClassifier(alpha = best_alpha,penalty = regularization,loss = 'hinge',random_state = 42)
74         # define and fit calibration model
75         calibrated = CalibratedClassifierCV(model, method='sigmoid', cv=5)
76         calibrated.fit(trainX, train_y)
77         # predict probabilities
78         return calibrated.predict_proba(testX)[: , 1]
79     # generate 2 class dataset
80
81     # uncalibrated predictions
82     yhat_uncalibrated = uncalibrated(best_alpha,regularization,trainX, testX, train_y)
83     # calibrated predictions
84     yhat_calibrated = calibrated(best_alpha,regularization,trainX, testX, train_y)
85     # reliability diagrams

```

```
86 fop_uncalibrated, mpv_uncalibrated = calibration_curve(test_y, yhat_uncalibrated, n_bins=10, normalize=1
87 fop_calibrated, mpv_calibrated = calibration_curve(test_y, yhat_calibrated, n_bins=10)
88 # plot perfectly calibrated
89 plt.figure(figsize = (15,6))
90 plt.plot([0, 1], [0, 1], linestyle='--', color='black')
91 # plot model reliabilities
92 plt.plot(mpv_uncalibrated, fop_uncalibrated, marker='.',label = 'Uncalibrated')
93 plt.plot(mpv_calibrated, fop_calibrated, marker='.',label = 'Calibrated')
94 plt.xlabel('Expected Probabilities',fontsize =18)
95 plt.ylabel('Predicted Probabilities',fontsize=18)
96 plt.legend(loc = 'best')
97 plt.title('Calibrated vs Uncalibrated',fontsize =18)
98 plt.show()
99
100 #*****
101
102
```

```

In [0]: 1
2 def plot_confusion_matrix(test_y, pred_y):
3     print('Confusion Matrix')
4     C = confusion_matrix(test_y, pred_y)
5
6     A = (((C.T)/(C.sum(axis=1))).T)#for recall matrix
7
8     B = (C/C.sum(axis=0))#for precision matrix
9     plt.figure(figsize=(20,4))
10
11     labels = [0,1]
12     # representing A in heatmap format
13     cmap=sns.light_palette("blue")
14     plt.subplot(1, 3, 1)
15     sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
16     plt.xlabel('Predicted Class')
17     plt.ylabel('Original Class')
18     plt.title("Confusion matrix")
19
20     plt.subplot(1, 3, 2)
21     sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
22     plt.xlabel('Predicted Class')
23     plt.ylabel('Original Class')
24     plt.title("Precision matrix")
25
26     plt.subplot(1, 3, 3)
27     # representing B in heatmap format
28     sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
29     plt.xlabel('Predicted Class')
30     plt.ylabel('Original Class')
31     plt.title("Recall matrix")
32
33     plt.show()
34
35
36     #*****
37
38
39 def plot_roc(train_y,train_proba,test_y,test_proba,auc_train,auc_test):
40     print('plotting ROC on Test data')
41     fpr_tr, tpr_tr, _ = roc_curve(train_y,train_proba)
42     fpr_test, tpr_test, _ = roc_curve(test_y,test_proba)

```



```

43 #calculating the fpr,tpr and thresholds for each training and test dataset
44 sns.set_style('darkgrid')
45 plt.figure(figsize=(15,8))
46 plt.plot(np.linspace(0,1,100),np.linspace(0,1,100),"g--")#this plots the roc curve for AUC = 0.5
47 plt.plot(fpr_tr,tpr_tr,'r',linewidth=2,label="train auc="+str(auc_train))
48 plt.plot(fpr_test,tpr_test,'b',linewidth=1,label=" test auc="+str(auc_test))
49 plt.xlabel('False positive rate(1-specificity)',fontsize=18)
50 plt.ylabel('True positive rate(sensitivity)',fontsize=18)
51 plt.title('Reciever operating characteristics curve',fontsize=18)
52 plt.legend(loc='best')
53 plt.show()
54
55 *****
56
57 def imp_features(best_alpha,regularization,train_set,train_y,vect):
58     clf = SGDClassifier(alpha = best_alpha,penalty = regularization,loss = 'hinge',random_state = 42)
59     clf.fit(train_set,train_y)
60     w = clf.coef_[0]#finding the coefficients of all features
61     #print(w)
62
63     features = vect.get_feature_names()#getting name of the features after fitting and transforming by count
64     negative_indices = np.argsort(w)
65     positive_indices = np.argsort(w)[::-1]
66     pos_dict = {}
67     neg_dict = {}
68
69
70
71     print('TOP 20 important features for positive class and their coefficients in this featurization are:\n')
72     for i in (positive_indices[0:20]):
73         pos_dict[features[i]] = w[i]
74         #print("%s\t --> \t%f"%(features[i],w[i]))
75     pos_df = pd.DataFrame.from_dict(pos_dict,orient = 'index',columns=['Coefficients'])
76     print(pos_df)
77     print("*****")
78
79
80     print('TOP 20 important features for negative class and their coefficients in this featurization are:\n')
81     for i in (negative_indices[0:20]):
82         neg_dict[features[i]] = w[i]
83     neg_df = pd.DataFrame.from_dict(neg_dict,orient = 'index',columns=['Coefficients'])
84     print(neg_df)
85

```

Using L2 Regularization

In [50]:

```

1  """classifier for tuning"""
2  gsm = svm('l2',train_bow,Y_train)
3  print(gsm)
4  print('*****\n')
5
6
7  """best hyperparameter"""
8  best_alpha_l2_bow = error('l2',train_bow,Y_train)['alpha']
9  print('*****\n')
10
11
12 """best Classifier fitted with tuned Hyperparameter"""
13 train_auc_l2_bow,test_auc_l2_bow,train_proba,test_proba,train_pred,test_pred = best_classifier(best_alpha_l2
14                                             'l2',train_bow,Y_train,test_bow,Y_test)
15 print('*****')
16
17
18 """Reliability Curve"""
19 print('The Calibration Curve')
20 reliability_curve(best_alpha_l2_bow,'l2',train_bow,test_bow,Y_train,Y_test)
21 print('*****\n')
22
23 """Confusion Matrix"""
24 print('for Training data:\n')
25 plot_confusion_matrix(Y_train,train_pred)
26
27 print('for Test data')
28 plot_confusion_matrix(Y_test,test_pred)
29 print('*****\n')
30
31 """ROC CURVE"""
32 plot_roc(Y_train,train_proba,Y_test,test_proba,train_auc_l2_bow,test_auc_l2_bow)
33
34
35

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 8.5s finished

```

GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=5),
             error_score='raise-deprecating',

```

```

estimator=SGDClassifier(alpha=0.0001, average=False,
                        class_weight=None, early_stopping=False,
                        epsilon=0.1, eta0=0.0, fit_intercept=True,
                        l1_ratio=0.15, learning_rate='optimal',
                        loss='hinge', max_iter=1000,
                        n_iter_no_change=5, n_jobs=None,
                        penalty='l2', power_t=0.5, random_state=42,
                        shuffle=True, tol=0.001,
                        validation_fraction=0.1, verbose=0,
                        warm_start=False),
iid='warn', n_jobs=None,
param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000,
                      10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)

```

Corresponding hypereparameter will give best auc on CV data
Fitting 5 folds for each of 9 candidates, totalling 45 fits

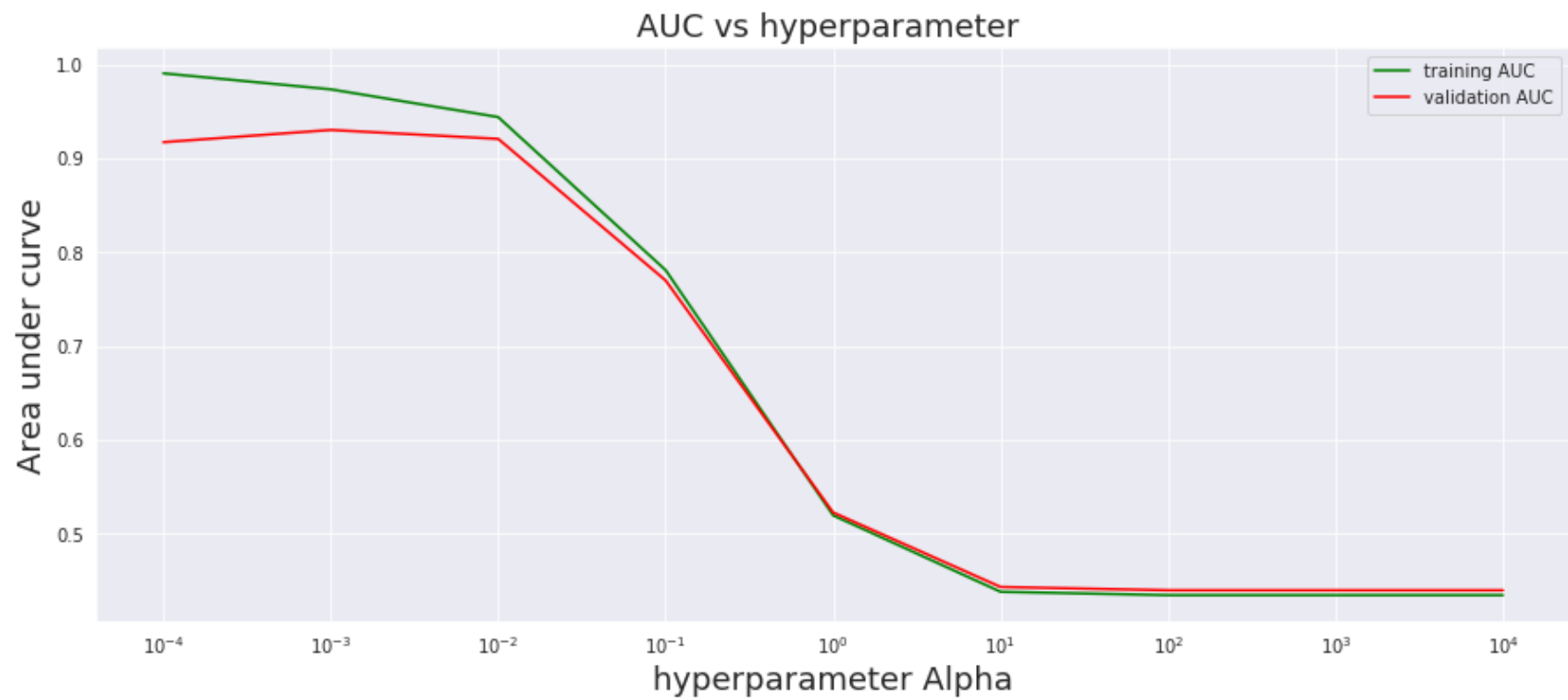
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 8.8s finished

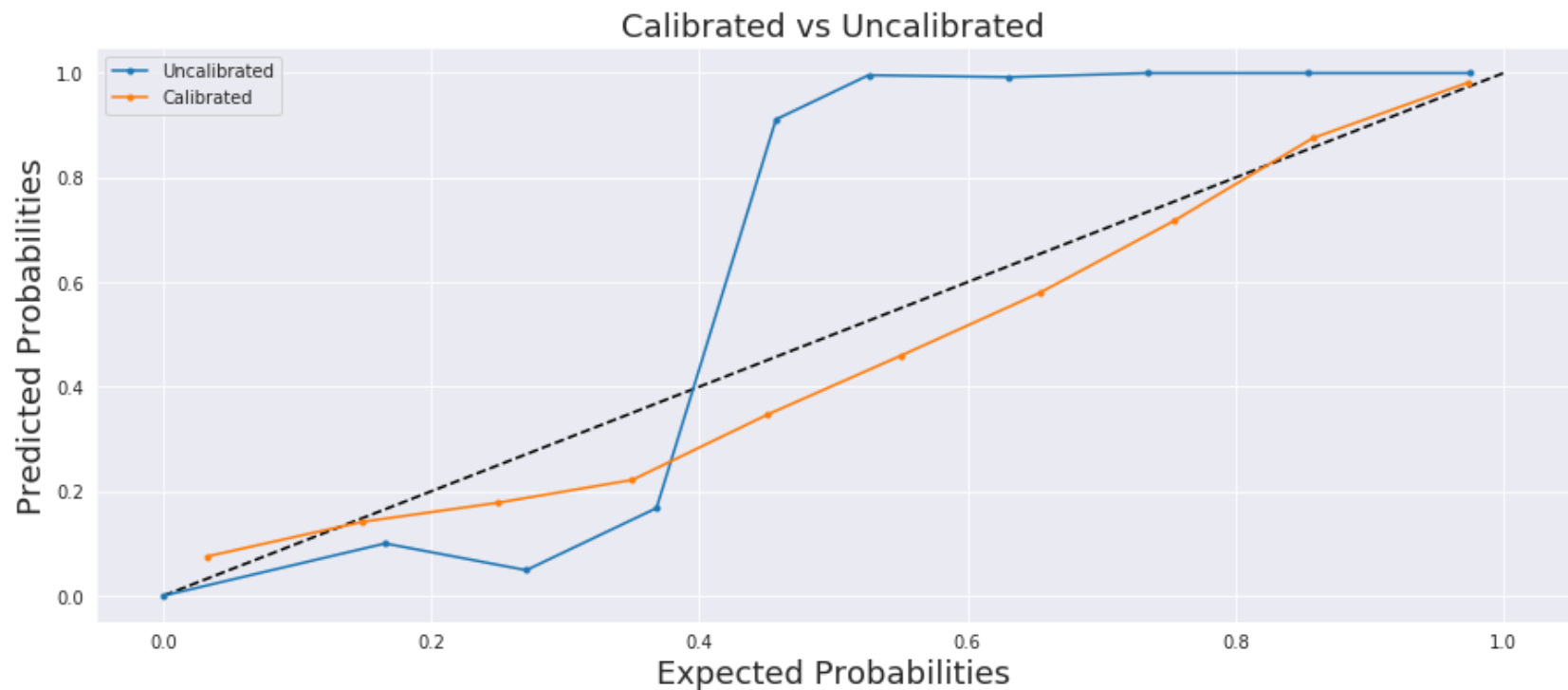
best hyperparameter is {'alpha': 0.001}

AUC on training data is 0.9635540335920942

AUC on test data is 0.9402083193654796

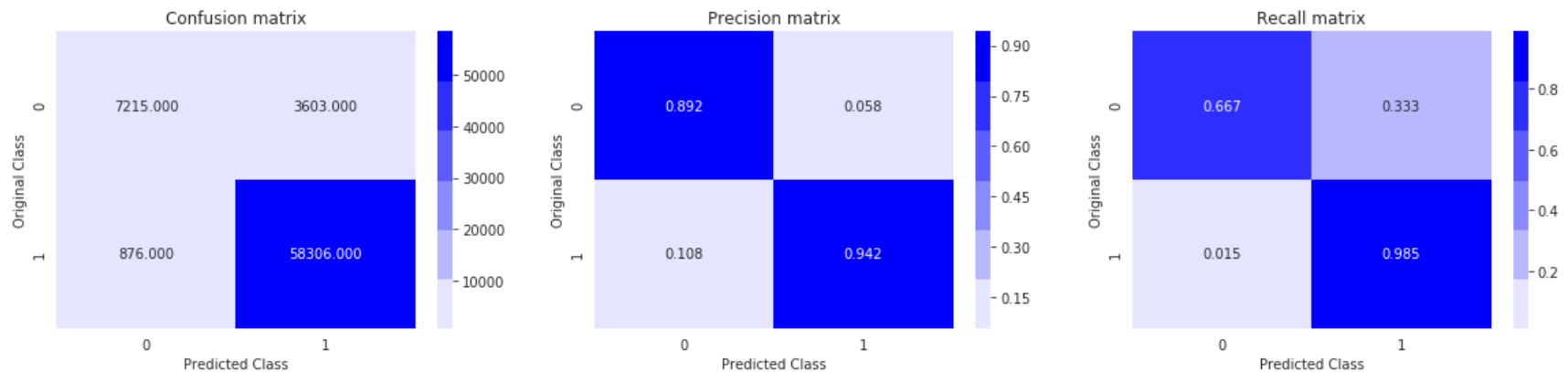
The Calibration Curve



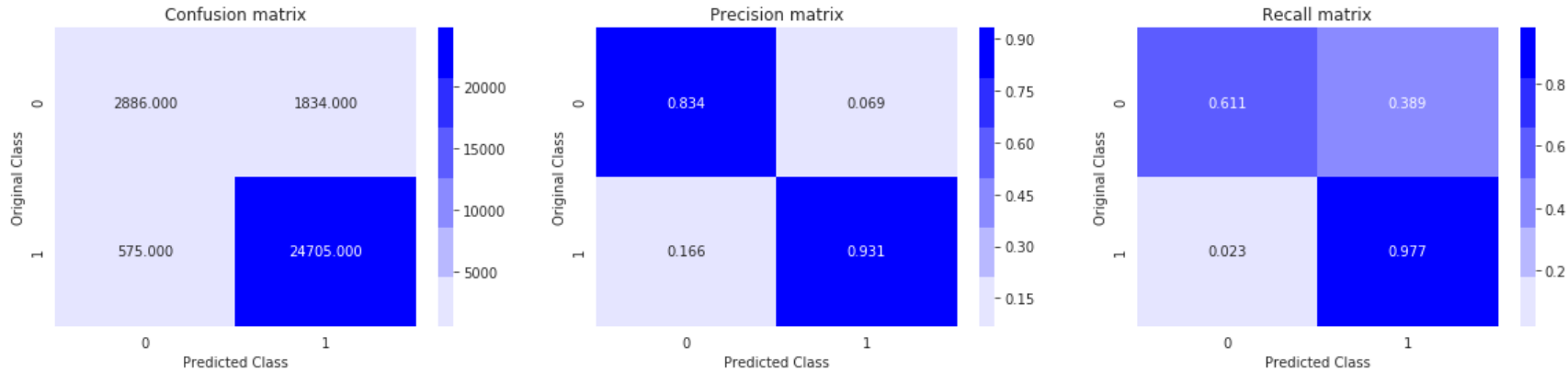


for Training data:

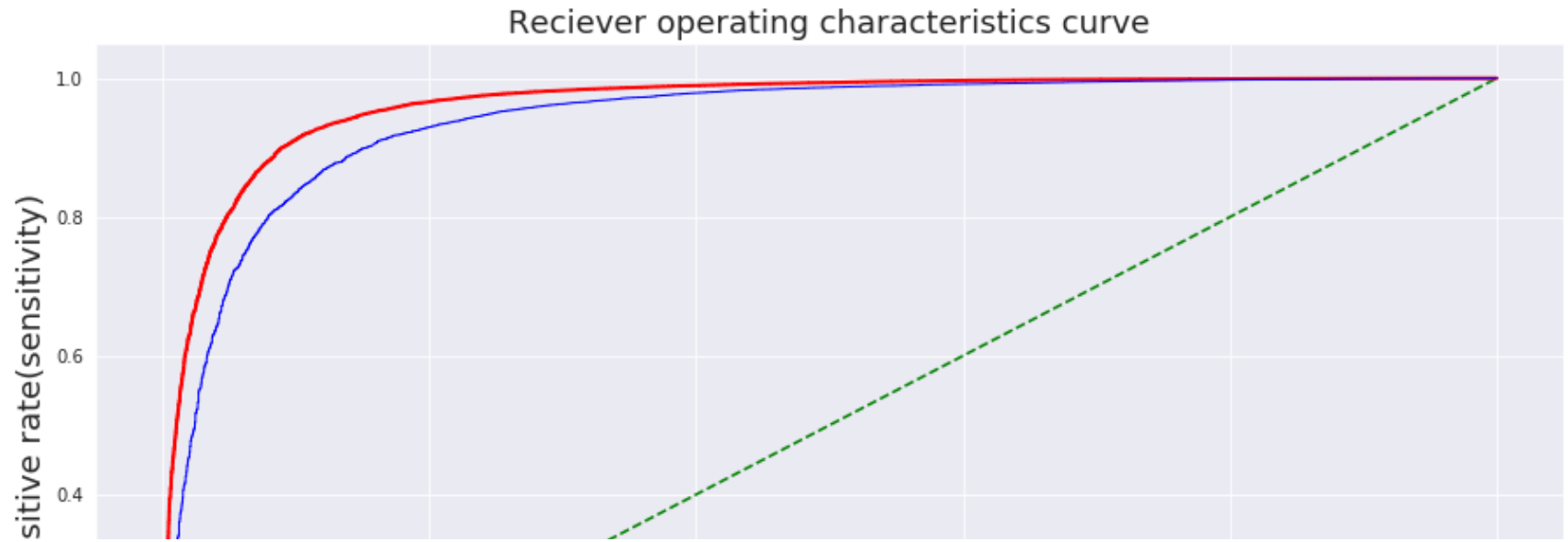
Confusion Matrix



for Test data
Confusion Matrix



plotting ROC on Test data




```
In [0]: 1 """TOP Features for both classes"""
        2 imp_features(best_alpha_l2_bow, 'l2', train_bow, Y_train, bow_vect)
```

TOP 20 important features for positive class and their coefficients in Bag of Words featurization are:

	Coefficients
delicious	0.628254
perfect	0.591130
excellent	0.565429
great	0.508315
highly	0.508315
awesome	0.505459
smooth	0.474046
amazing	0.468335
yummy	0.465479
best	0.456912
wonderful	0.431211
happy	0.422644
loves	0.422644
nice	0.408365
thank	0.399798
tasty	0.379808
definitely	0.376952
complaint	0.376952
fantastic	0.376952
pleased	0.374097

TOP 20 important features for negative class and their coefficients in Bag of Words featurization are:

	Coefficients
worst	-0.899546
awful	-0.842432
disappointing	-0.796740
disappointed	-0.699647
disappointment	-0.696791
threw	-0.696791
terrible	-0.688224
unfortunately	-0.668234
stale	-0.642533
horrible	-0.633965
tasteless	-0.585419
weak	-0.582563

waste	-0.573996
poor	-0.571140
disgusting	-0.565429
bland	-0.548294
sorry	-0.516882
rip	-0.514026
return	-0.511170
unpleasant	-0.499748

Using L1 Regularization

In [51]:

```

1  """classifier for tuning"""
2  gsm = svm('l1',train_bow,Y_train)
3  print(gsm)
4  print('*****\n')
5
6
7  """best hyperparameter"""
8  best_alpha_l1_bow = error('l1',train_bow,Y_train)['alpha']
9  print('*****\n')
10
11
12 """best Classifier fitted with tuned Hyperparameter"""
13 train_auc_l1_bow,test_auc_l1_bow,train_proba,test_proba,train_pred,test_pred = best_classifier(best_alpha_l1
14                                             'l1',train_bow,Y_train,test_bow,Y_test)
15 print('*****')
16
17
18 """Reliability Curve"""
19 print('The Calibration Curve')
20 reliability_curve(best_alpha_l1_bow,'l1',train_bow,test_bow,Y_train,Y_test)
21 print('*****\n')
22
23 """Confusion Matrix"""
24 print('for Training data:\n')
25 plot_confusion_matrix(Y_train,train_pred)
26
27 print('for Test data')
28 plot_confusion_matrix(Y_test,test_pred)
29 print('*****\n')
30
31 """ROC CURVE"""
32 plot_roc(Y_train,train_proba,Y_test,test_proba,train_auc_l1_bow,test_auc_l1_bow)
33
34
35

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 29.0s finished

```

GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=5),
             error_score='raise-deprecating',

```

```

estimator=SGDClassifier(alpha=0.0001, average=False,
                        class_weight=None, early_stopping=False,
                        epsilon=0.1, eta0=0.0, fit_intercept=True,
                        l1_ratio=0.15, learning_rate='optimal',
                        loss='hinge', max_iter=1000,
                        n_iter_no_change=5, n_jobs=None,
                        penalty='l1', power_t=0.5, random_state=42,
                        shuffle=True, tol=0.001,
                        validation_fraction=0.1, verbose=0,
                        warm_start=False),
iid='warn', n_jobs=None,
param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000,
                      10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)

```

Corresponding hypereparameter will give best auc on CV data

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

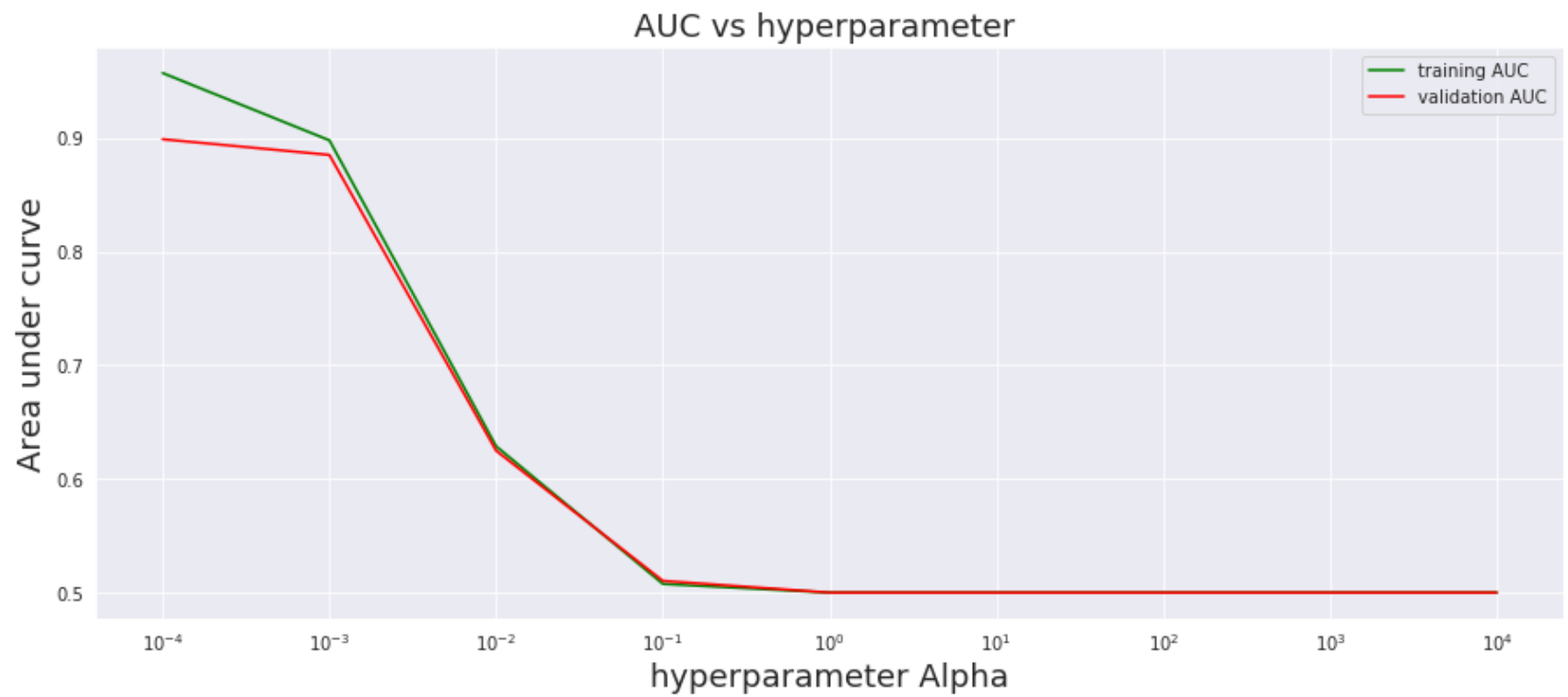
[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 28.5s finished

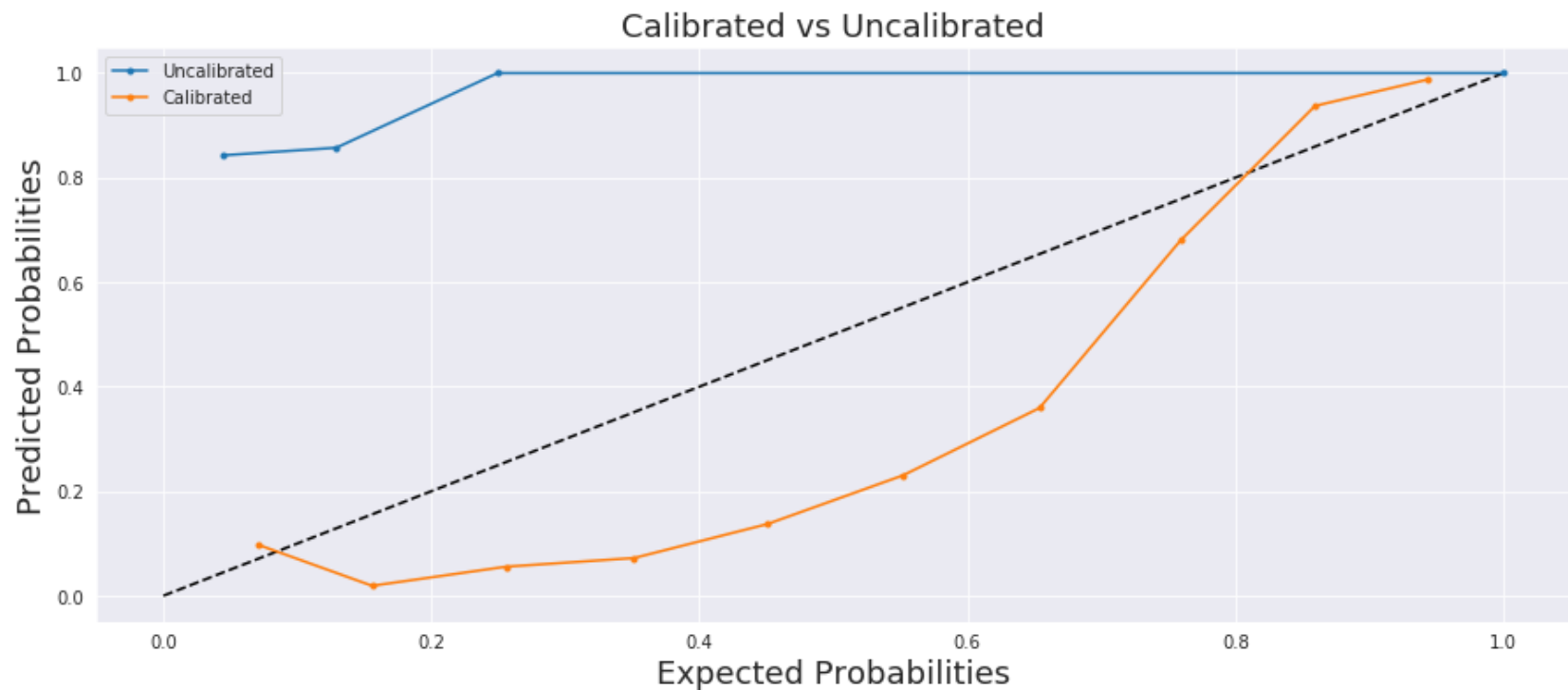
best hyperparameter is {'alpha': 0.0001}

AUC on training data is 0.956162146262999

AUC on test data is 0.9303229130350246

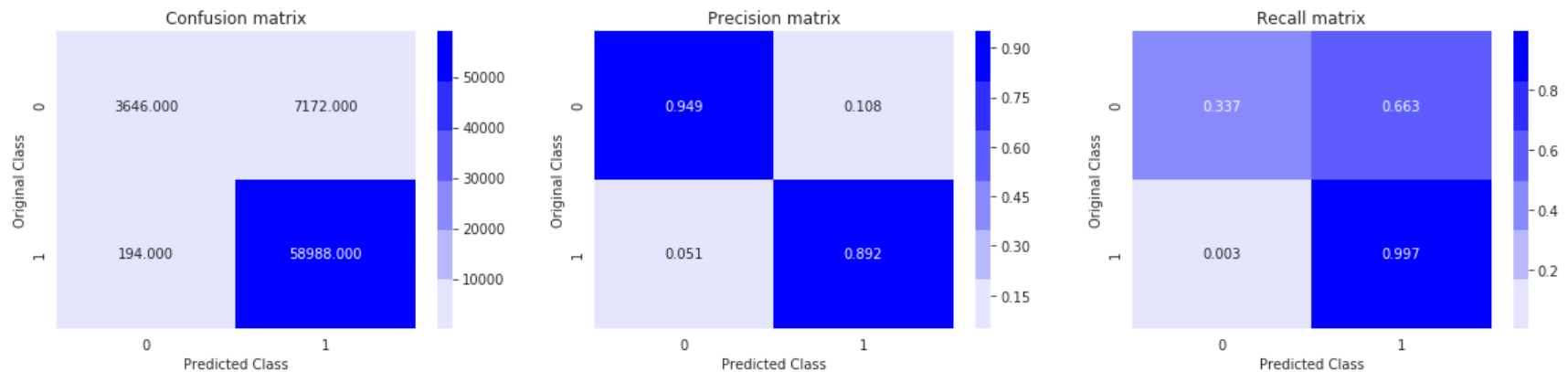
The Calibration Curve



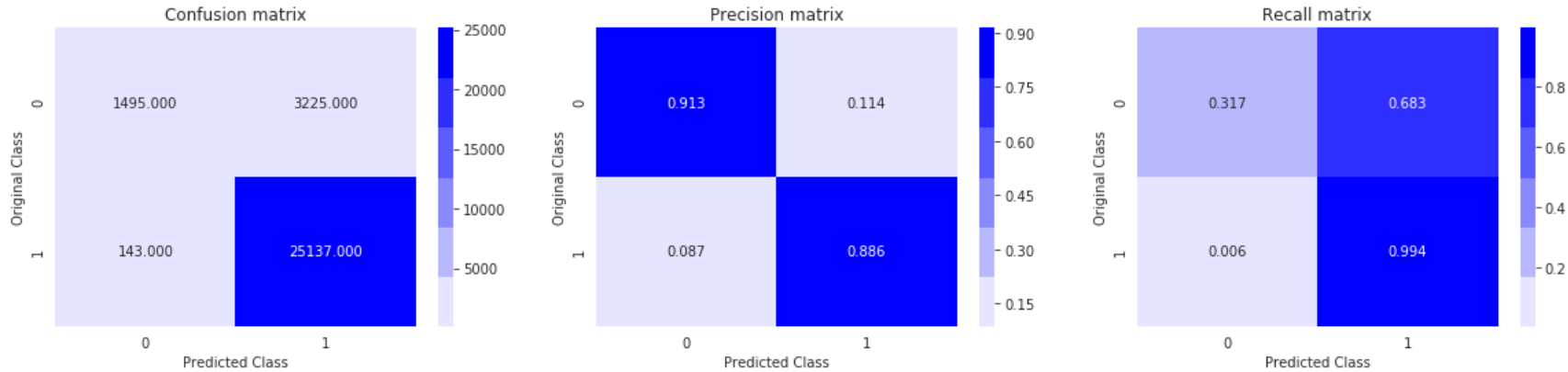


for Training data:

Confusion Matrix



for Test data
Confusion Matrix



plotting ROC on Test data




```
In [0]: 1 """TOP Features for both classes"""
        2 imp_features(best_alpha_l1_bow, '11', train_bow, Y_train, bow_vect)
```

TOP 20 important features for positive class and their coefficients in Bag of Words featurization are:

	Coefficients
knife	104.474975
corner	76.963398
edge	76.236796
gently	74.148043
point	53.998964
proceed	44.844656
tofu	42.977610
slowly	41.747194
withdraw	38.607035
izze	38.357197
slit	38.241826
gentle	35.820080
tip	33.461279
slits	29.697758
insulin	29.464263
membrane	28.668077
procedure	28.603500
sulfate	26.594956
vegetable	23.110526
constipation	22.643934

TOP 20 important features for negative class and their coefficients in Bag of Words featurization are:

	Coefficients
iams	-43.559593
disappointing	-32.887452
pureed	-32.127471
milka	-29.474204
disappointment	-28.367887
worst	-27.275682
awful	-27.102317
marshmallow	-25.113100
hopes	-24.536193
disgusting	-23.996222
threw	-23.979561
molars	-23.841928

worse	-21.691247
unfortunately	-21.194842
september	-20.181955
rip	-19.520538
terrible	-19.431291
ruined	-19.276387
sad	-18.732239
bland	-18.728568

[5.1.2] Applying Linear SVM on TFIDF, SET 2

```
In [0]: 1 #reading tfidf files
        2 train_tfidf = pickle.load(open('train_tfidf.p','rb'))
        3 test_tfidf = pickle.load(open('test_tfidf.p','rb'))
        4
```

Using L2 Regularization

In [53]:

```

1  """classifier for tuning"""
2  gsm = svm('l2',train_tfidf,Y_train)
3  print(gsm)
4  print('*****\n')
5
6
7  """best hyperparameter"""
8  best_alpha_l2_tfidf = error('l2',train_tfidf,Y_train)['alpha']
9  print('*****\n')
10
11
12  """best Classifier fitted with tuned Hyperparameter"""
13  train_auc_l2_tfidf,test_auc_l2_tfidf,train_proba,test_proba,train_pred,test_pred = best_classifier(best_alpha_l2_tfidf,
14  'l2',train_tfidf,Y_train,test_tfidf,Y_test)
15  print('*****')
16
17
18  """Reliability Curve"""
19  print('The Calibration Curve')
20  reliability_curve(best_alpha_l2_tfidf,'l2',train_tfidf,test_tfidf,Y_train,Y_test)
21  print('*****\n')
22
23  """Confusion Matrix"""
24  print('for Training data:\n')
25  plot_confusion_matrix(Y_train,train_pred)
26
27  print('for Test data')
28  plot_confusion_matrix(Y_test,test_pred)
29  print('*****\n')
30
31  """ROC CURVE"""
32  plot_roc(Y_train,train_proba,Y_test,test_proba,train_auc_l2_tfidf,test_auc_l2_tfidf)
33
34

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 10.5s finished

```

GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=5),
             error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,

```

```

class_weight=None, early_stopping=False,
epsilon=0.1, eta0=0.0, fit_intercept=True,
l1_ratio=0.15, learning_rate='optimal',
loss='hinge', max_iter=1000,
n_iter_no_change=5, n_jobs=None,
penalty='l2', power_t=0.5, random_state=42,
shuffle=True, tol=0.001,
validation_fraction=0.1, verbose=0,
warm_start=False),
iid='warn', n_jobs=None,
param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000,
                      10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)
*****

```

Corresponding hypereparameter will give best auc on CV data
Fitting 5 folds for each of 9 candidates, totalling 45 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 10.3s finished

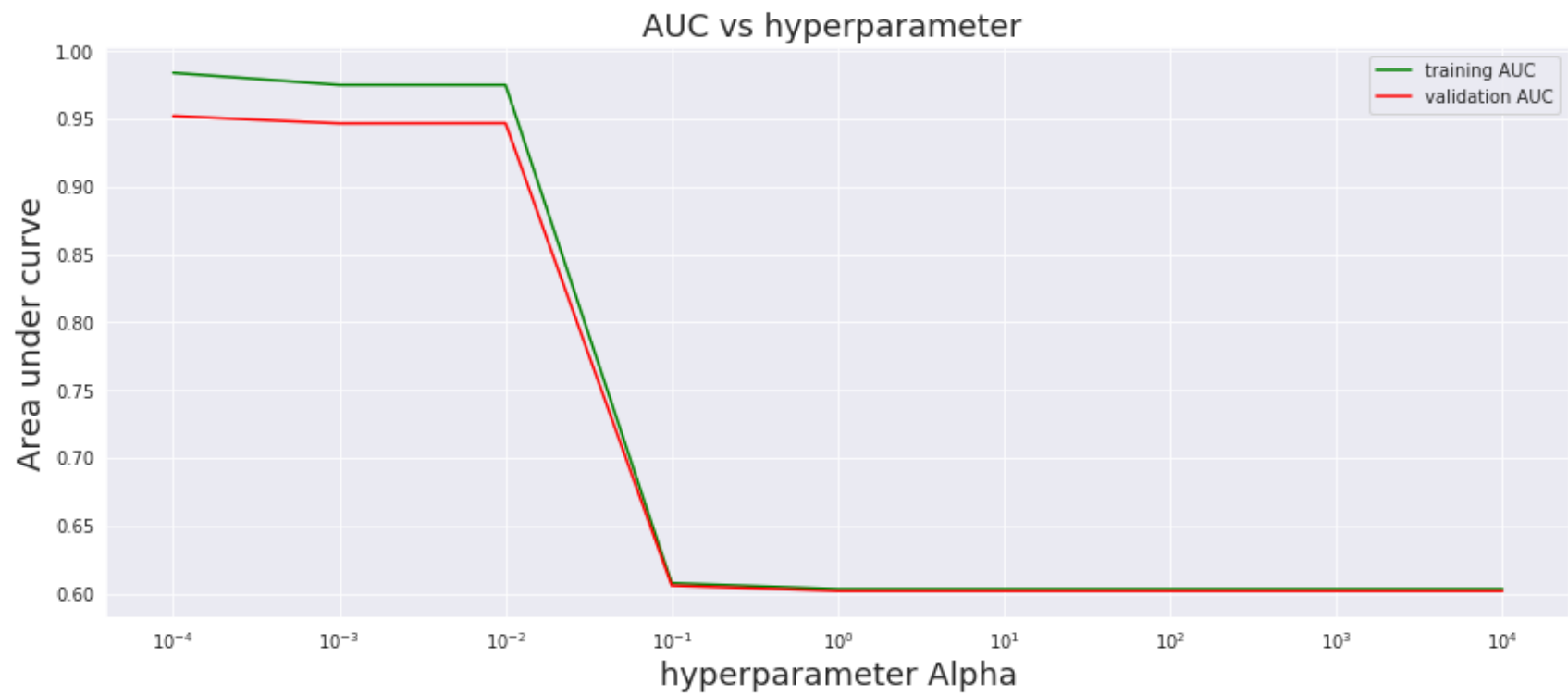
```

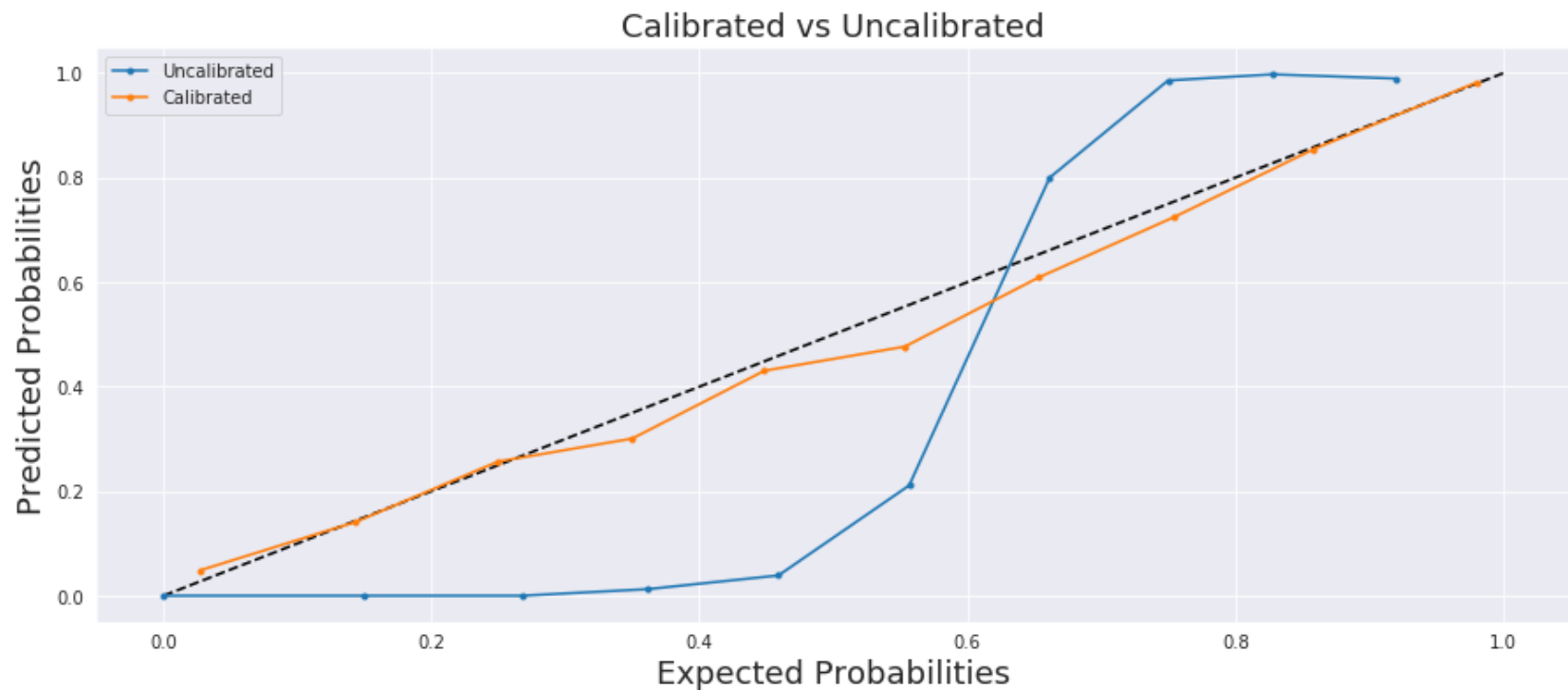
best hyperparameter is {'alpha': 0.0001}

AUC on training data is 0.9736325736030264

AUC on test data is 0.9557599378486376

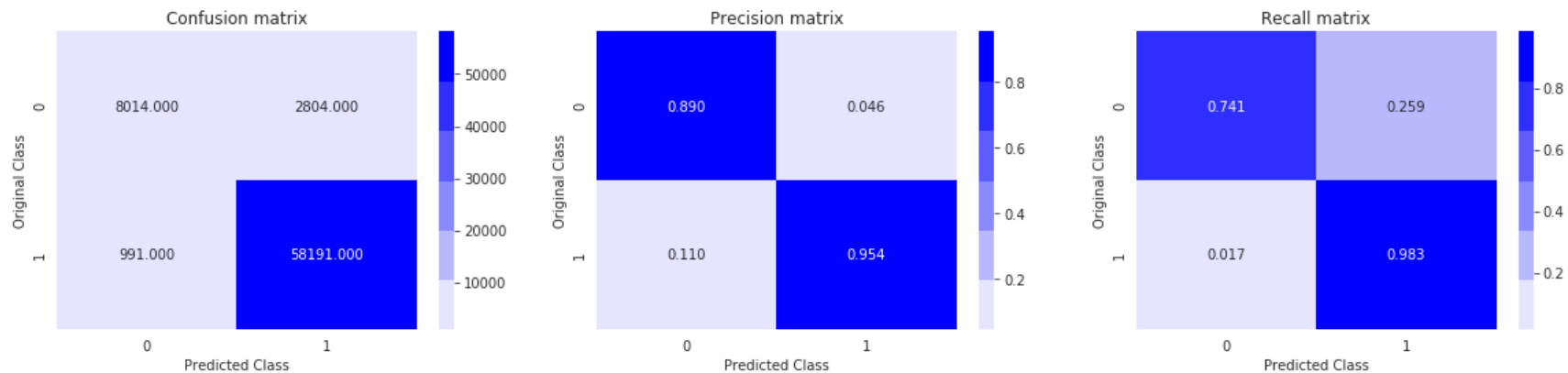
The Calibration Curve



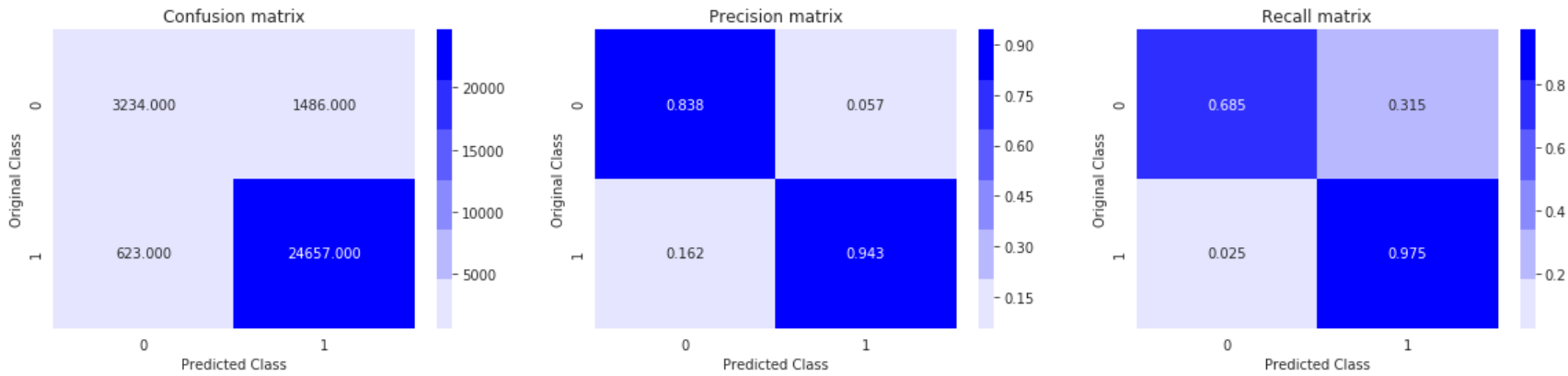


for Training data:

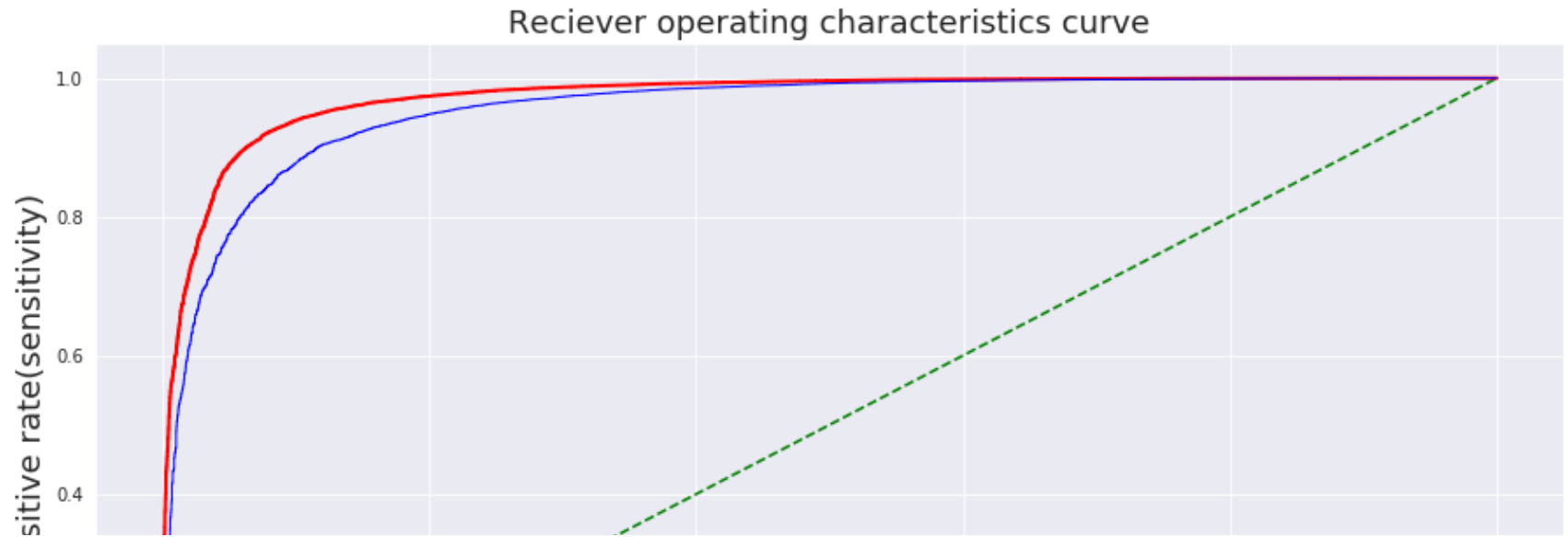
Confusion Matrix



for Test data
Confusion Matrix



plotting ROC on Test data




```
In [0]: 1 """TOP Features for both classes"""
        2 imp_features(best_alpha_l2_tfidf, 'l2', train_tfidf, Y_train, tfidf_vect)
```

TOP 20 important features for positive class and their coefficients in this featurization are:

	Coefficients
great	3.524112
best	2.738388
delicious	2.585764
not disappointed	2.428233
good	2.368122
love	2.120625
perfect	2.036168
excellent	1.951569
nice	1.904471
loves	1.776171
happy	1.586472
not bad	1.552323
wonderful	1.540727
amazing	1.493950
awesome	1.490634
tasty	1.484353
yummy	1.462984
favorite	1.450721
definitely	1.333967
smooth	1.324058

TOP 20 important features for negative class and their coefficients in this featurization are:

	Coefficients
disappointed	-3.998880
worst	-3.453786
awful	-3.165934
not worth	-3.026346
not buy	-2.999386
not recommend	-2.981703
not	-2.816763
not good	-2.785361
terrible	-2.702545
disappointing	-2.616888
horrible	-2.572291
stale	-2.528569

return	-2.522660
unfortunately	-2.491579
threw	-2.466702
bad	-2.321640
money	-2.319694
disgusting	-2.232970
waste money	-2.183431
not order	-2.100434

Using L1 Regularization

In [54]:

```

1  """classifier for tuning"""
2  gsm = svm('l1',train_tfidf,Y_train)
3  print(gsm)
4  print('*****\n')
5
6
7  """best hyperparameter"""
8  best_alpha_l1_tfidf = error('l1',train_tfidf,Y_train)['alpha']
9  print('*****\n')
10
11
12 """best Classifier fitted with tuned Hyperparameter"""
13 train_auc_l1_tfidf,test_auc_l1_tfidf,train_proba,test_proba,train_pred,test_pred = best_classifier(best_alpha_l1_tfidf,train_tfidf,Y_train,test_tfidf,Y_test)
14
15 print('*****')
16
17
18 """Reliability Curve"""
19 print('The Calibration Curve')
20 reliability_curve(best_alpha_l1_tfidf,'l1',train_tfidf,test_tfidf,Y_train,Y_test)
21 print('*****\n')
22
23 """Confusion Matrix"""
24 print('for Training data:\n')
25 plot_confusion_matrix(Y_train,train_pred)
26
27 print('for Test data')
28 plot_confusion_matrix(Y_test,test_pred)
29 print('*****\n')
30
31 """ROC CURVE"""
32 plot_roc(Y_train,train_proba,Y_test,test_proba,train_auc_l1_tfidf,test_auc_l1_tfidf)
33
34

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 16.8s finished

```

GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=5),
             error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,

```

```

class_weight=None, early_stopping=False,
epsilon=0.1, eta0=0.0, fit_intercept=True,
l1_ratio=0.15, learning_rate='optimal',
loss='hinge', max_iter=1000,
n_iter_no_change=5, n_jobs=None,
penalty='l1', power_t=0.5, random_state=42,
shuffle=True, tol=0.001,
validation_fraction=0.1, verbose=0,
warm_start=False),
iid='warn', n_jobs=None,
param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000,
                      10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)
*****

```

Corresponding hypereparameter will give best auc on CV data
 Fitting 5 folds for each of 9 candidates, totalling 45 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 17.0s finished

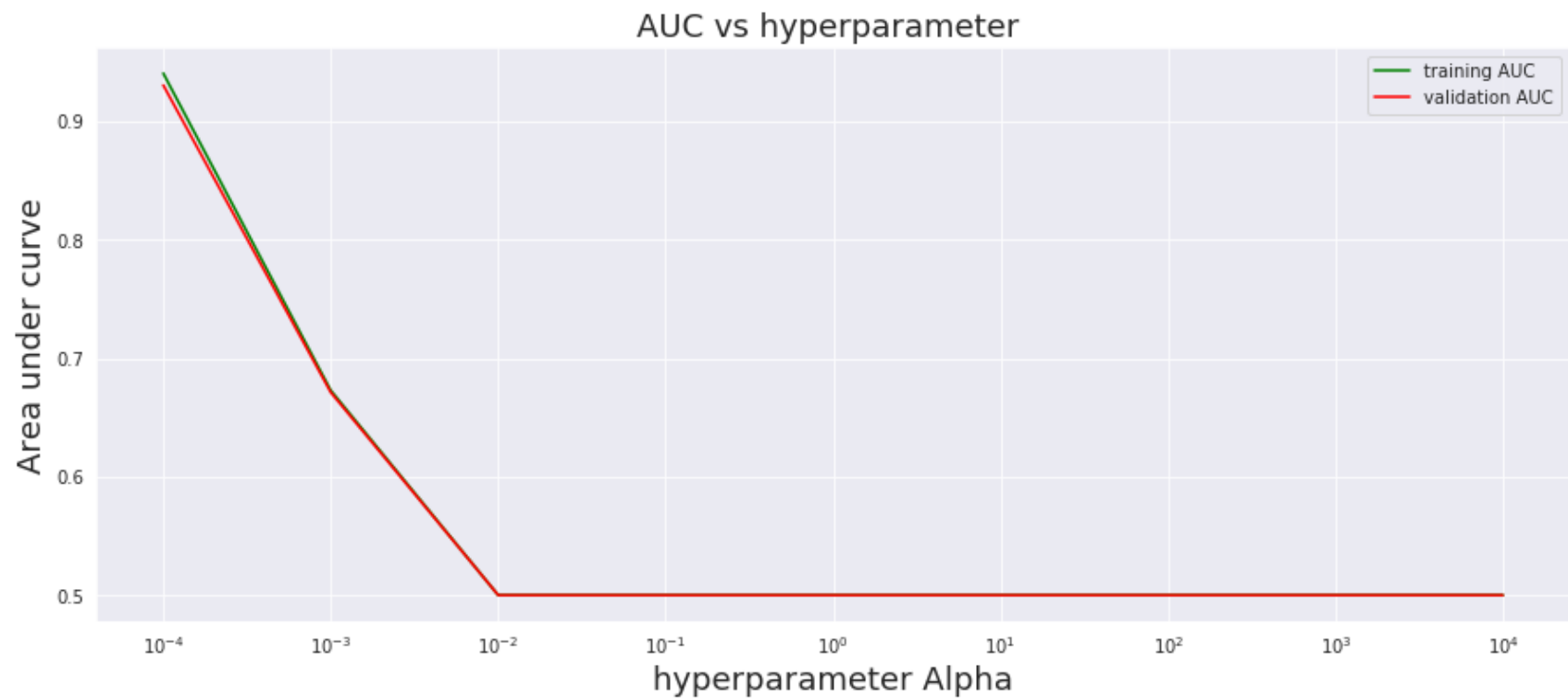
```

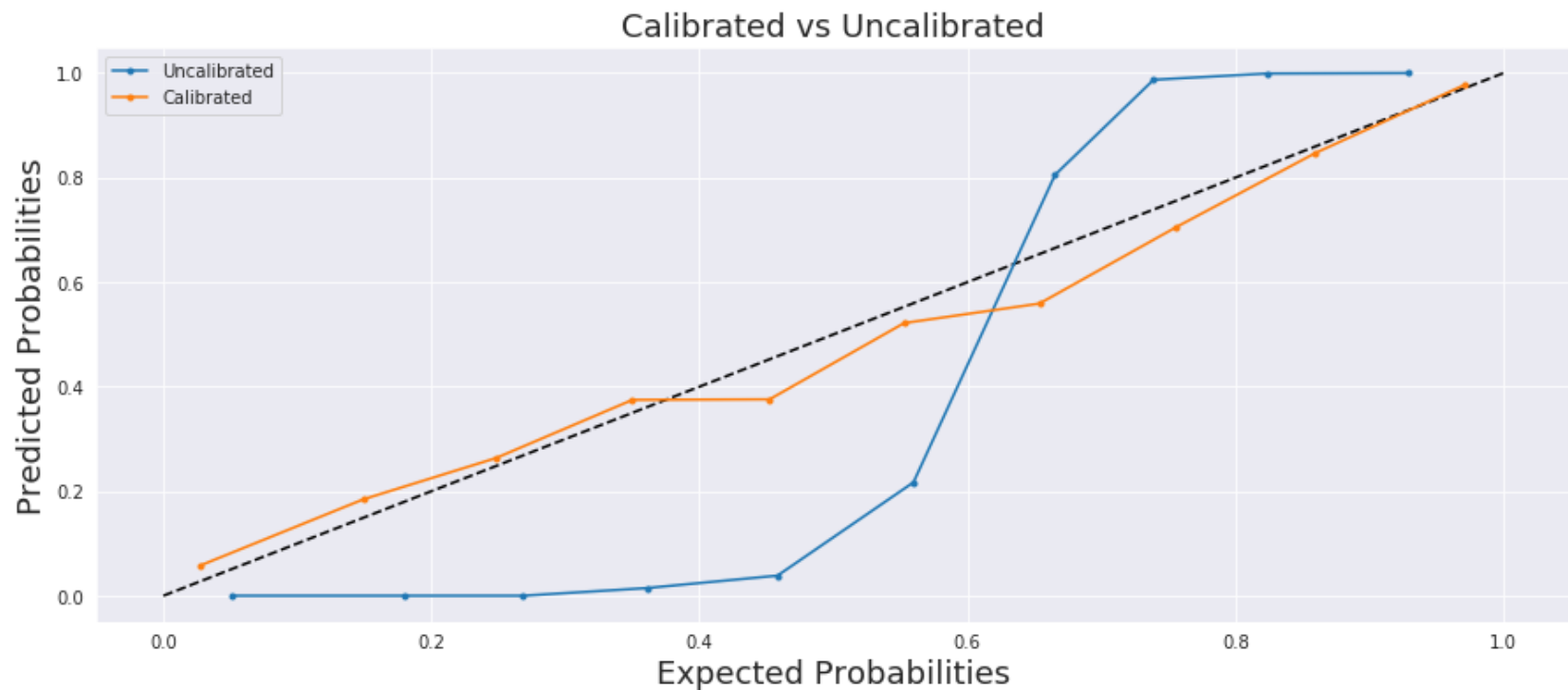
best hyperparameter is {'alpha': 0.0001}

AUC on training data is 0.9369532257610144

AUC on test data is 0.9330814202960738

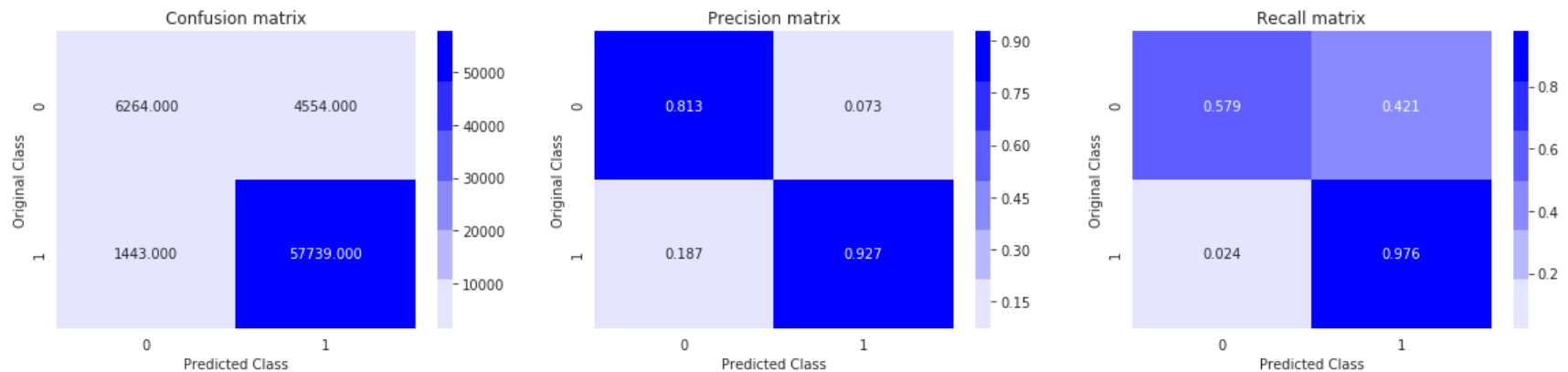
The Calibration Curve





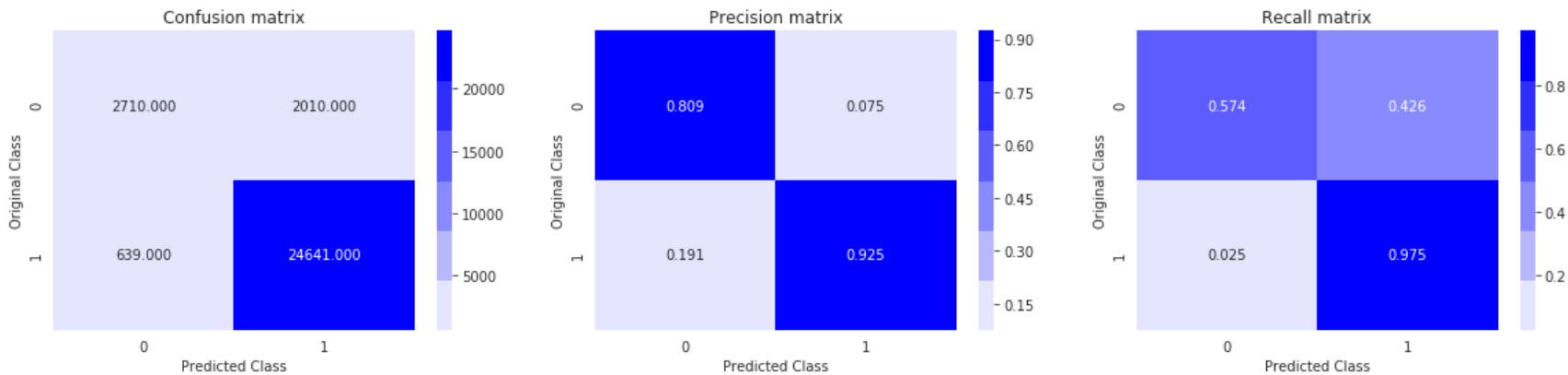
for Training data:

Confusion Matrix

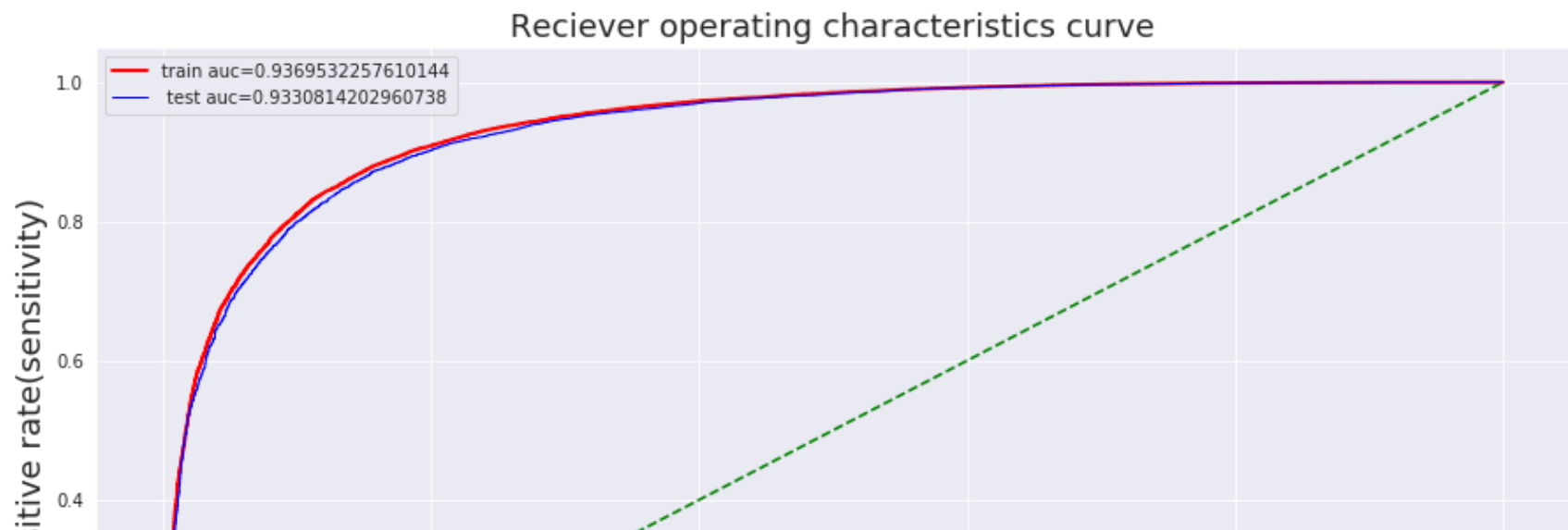


for Test data

Confusion Matrix



plotting ROC on Test data




```
In [0]: 1 """TOP Features for both classes"""
        2 imp_features(best_alpha_l1_tfidf, '12', train_tfidf, Y_train, tfidf_vect)
```

TOP 20 important features for positive class and their coefficients in this featurization are:

	Coefficients
great	3.524112
best	2.738388
delicious	2.585764
not disappointed	2.428233
good	2.368122
love	2.120625
perfect	2.036168
excellent	1.951569
nice	1.904471
loves	1.776171
happy	1.586472
not bad	1.552323
wonderful	1.540727
amazing	1.493950
awesome	1.490634
tasty	1.484353
yummy	1.462984
favorite	1.450721
definitely	1.333967
smooth	1.324058

TOP 20 important features for negative class and their coefficients in this featurization are:

	Coefficients
disappointed	-3.998880
worst	-3.453786
awful	-3.165934
not worth	-3.026346
not buy	-2.999386
not recommend	-2.981703
not	-2.816763
not good	-2.785361
terrible	-2.702545
disappointing	-2.616888
horrible	-2.572291
stale	-2.528569

return	-2.522660
unfortunately	-2.491579
threw	-2.466702
bad	-2.321640
money	-2.319694
disgusting	-2.232970
waste money	-2.183431
not order	-2.100434

[5.1.3] Applying Linear SVM on AVG W2V, SET 3

```
In [0]: 1 #Average WORD 2 VECTOR
        2 train_avgw2v = pickle.load(open('train_avgw2v.p', 'rb'))
        3 test_avgw2v = pickle.load(open('test_avgw2v.p', 'rb'))
```

Using L2 Regularization

In [57]:

```

1  """classifier for tuning"""
2  gsm = svm('l2',train_avgw2v,Y_train)
3  print(gsm)
4  print('*****\n')
5
6
7  """best hyperparameter"""
8  best_alpha_l2_avgw2v = error('l2',train_avgw2v,Y_train)['alpha']
9  print('*****\n')
10
11
12  """best Classifier fitted with tuned Hyperparameter"""
13  train_auc_l2_avgw2v,test_auc_l2_avgw2v,train_proba,test_proba,train_pred,test_pred = best_classifier(best_al
14  'l2',train_avgw2v,Y_train,test_avgw2v,Y_test)
15  print('*****')
16
17
18  """Reliability Curve"""
19  print('The Calibration Curve')
20  reliability_curve(best_alpha_l2_avgw2v,'l2',train_avgw2v,test_avgw2v,Y_train,Y_test)
21  print('*****\n')
22
23  """Confusion Matrix"""
24  print('for Training data:\n')
25  plot_confusion_matrix(Y_train,train_pred)
26
27  print('for Test data')
28  plot_confusion_matrix(Y_test,test_pred)
29  print('*****\n')
30
31  """ROC CURVE"""
32  plot_roc(Y_train,train_proba,Y_test,test_proba,train_auc_l2_avgw2v,test_auc_l2_avgw2v)
33

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 10.1s finished

```

GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=5),
             error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,
                                     class_weight=None, early_stopping=False,

```

```

epsilon=0.1, eta0=0.0, fit_intercept=True,
l1_ratio=0.15, learning_rate='optimal',
loss='hinge', max_iter=1000,
n_iter_no_change=5, n_jobs=None,
penalty='l2', power_t=0.5, random_state=42,
shuffle=True, tol=0.001,
validation_fraction=0.1, verbose=0,
warm_start=False),
iid='warn', n_jobs=None,
param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000,
                      10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)

```

Corresponding hypereparameter will give best auc on CV data

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

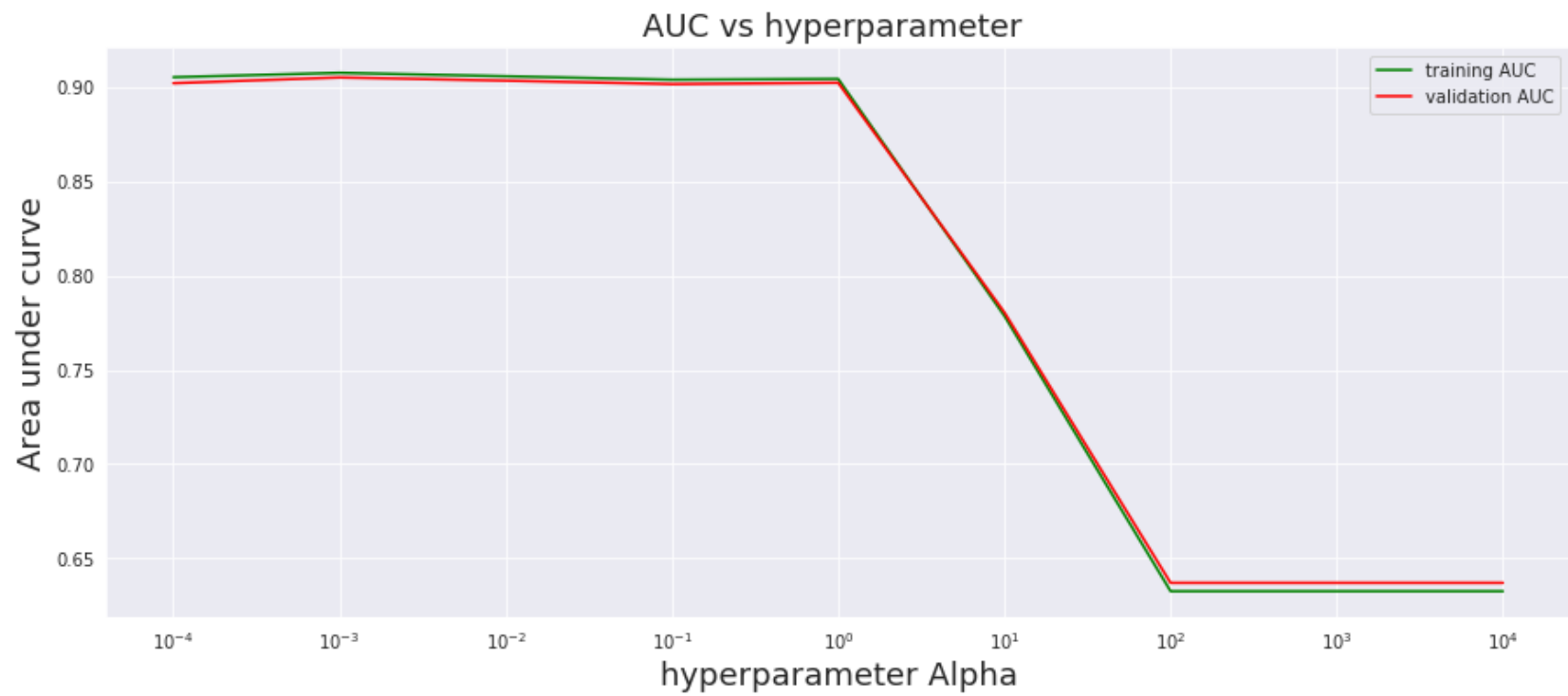
[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 10.3s finished

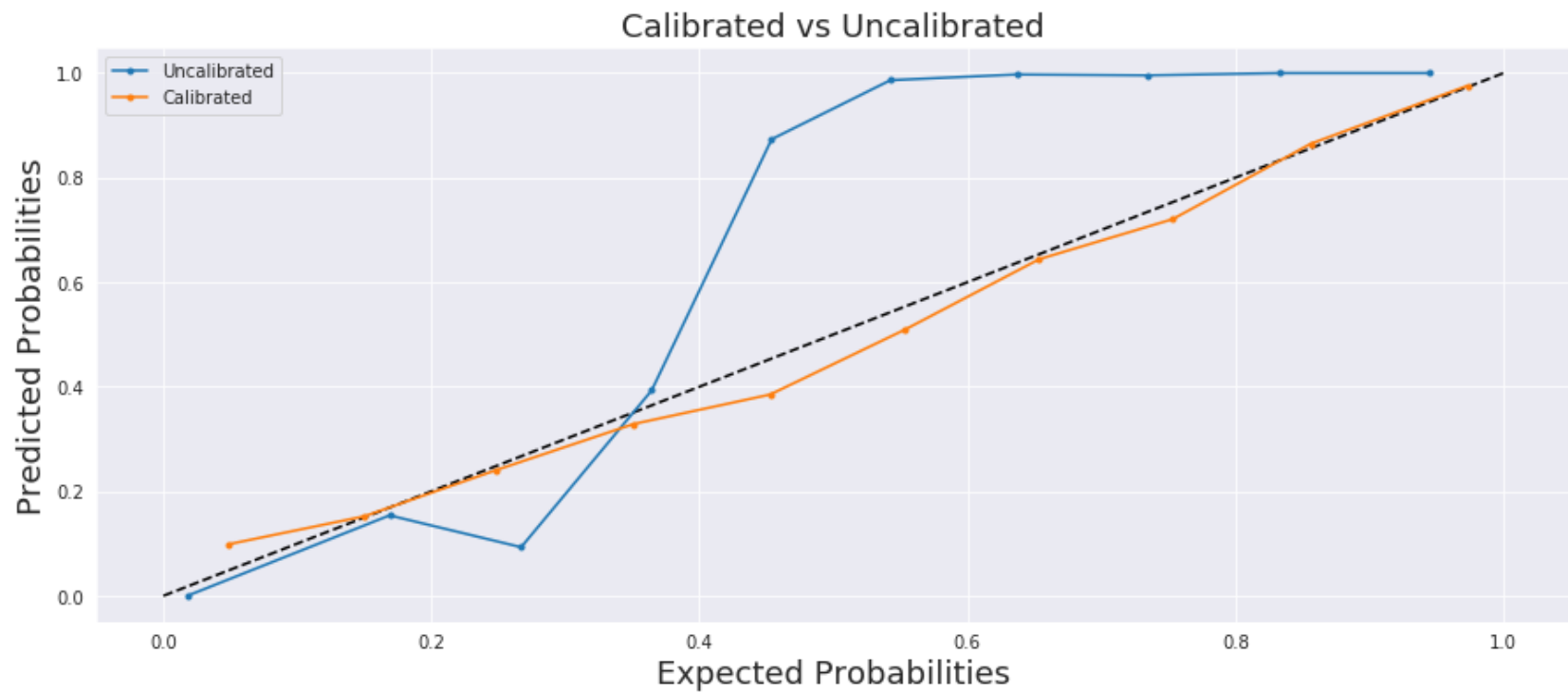
best hyperparameter is {'alpha': 0.001}

AUC on training data is 0.9074093991368202

AUC on test data is 0.9058933923111994

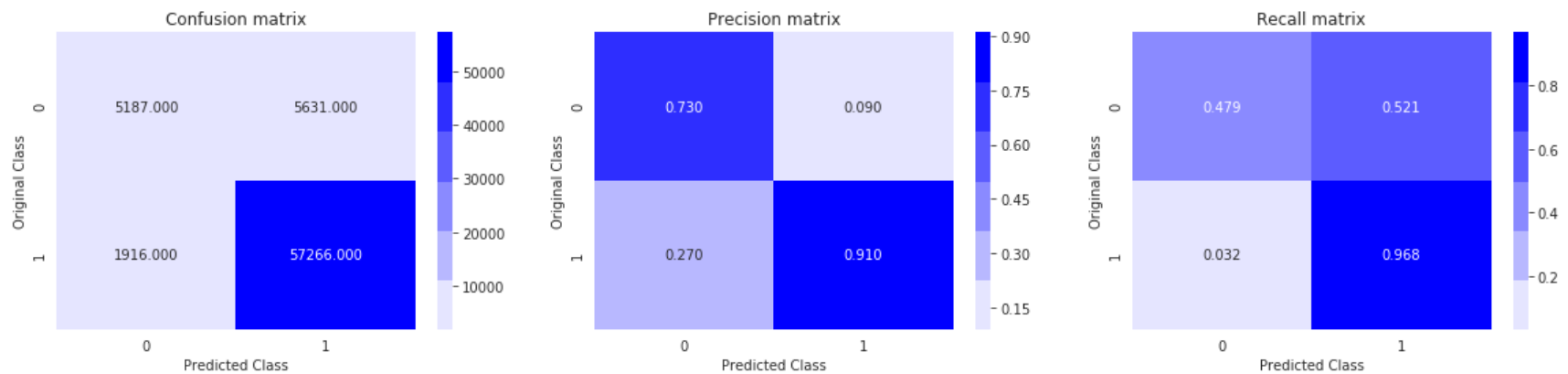
The Calibration Curve





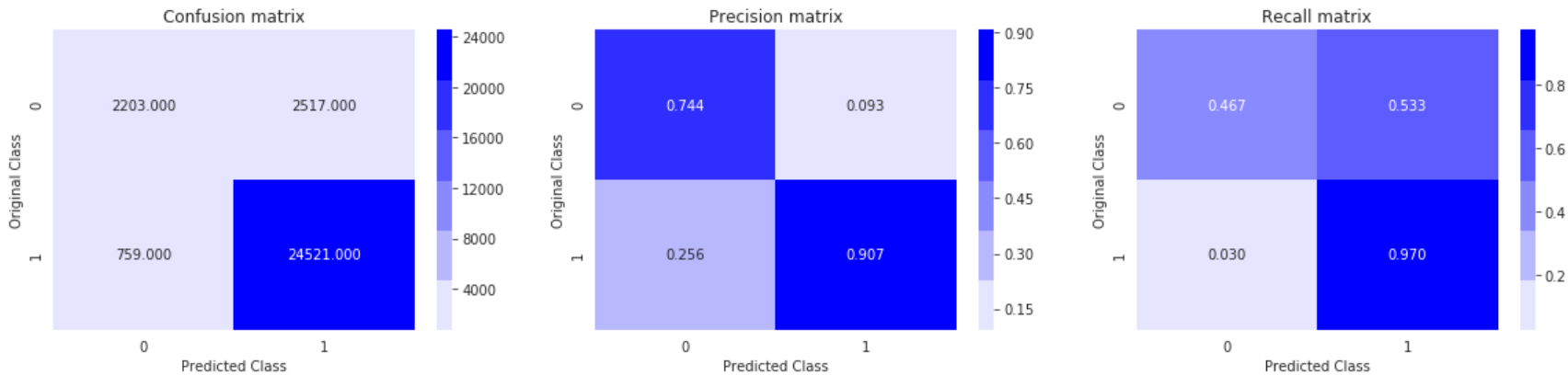
for Training data:

Confusion Matrix



for Test data

Confusion Matrix



plotting ROC on Test data



Using L1 Regularization

In [58]:

```

1  """classifier for tuning"""
2  gsm = svm('l1',train_avg2v,Y_train)
3  print(gsm)
4  print('*****\n')
5
6
7  """best hyperparameter"""
8  best_alpha_l1_avg2v = error('l1',train_avg2v,Y_train)['alpha']
9  print('*****\n')
10
11
12  """best Classifier fitted with tuned Hyperparameter"""
13  train_auc_l1_avg2v,test_auc_l1_avg2v,train_proba,test_proba,train_pred,test_pred = best_classifier(best_al
14  'l1',train_avg2v,Y_train,test_avg2v,Y_test)
15  print('*****')
16
17
18  """Reliability Curve"""
19  print('The Calibration Curve')
20  reliability_curve(best_alpha_l1_avg2v,'l1',train_avg2v,test_avg2v,Y_train,Y_test)
21  print('*****\n')
22
23  """Confusion Matrix"""
24  print('for Training data:\n')
25  plot_confusion_matrix(Y_train,train_pred)
26
27  print('for Test data')
28  plot_confusion_matrix(Y_test,test_pred)
29  print('*****\n')
30
31  """ROC CURVE"""
32  plot_roc(Y_train,train_proba,Y_test,test_proba,train_auc_l1_avg2v,test_auc_l1_avg2v)
33

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 14.0s finished

```

GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=5),
             error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,
                                     class_weight=None, early_stopping=False,

```

```

epsilon=0.1, eta0=0.0, fit_intercept=True,
l1_ratio=0.15, learning_rate='optimal',
loss='hinge', max_iter=1000,
n_iter_no_change=5, n_jobs=None,
penalty='l1', power_t=0.5, random_state=42,
shuffle=True, tol=0.001,
validation_fraction=0.1, verbose=0,
warm_start=False),
iid='warn', n_jobs=None,
param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000,
                      10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)
*****

```

Corresponding hypereparameter will give best auc on CV data
 Fitting 5 folds for each of 9 candidates, totalling 45 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 14.3s finished

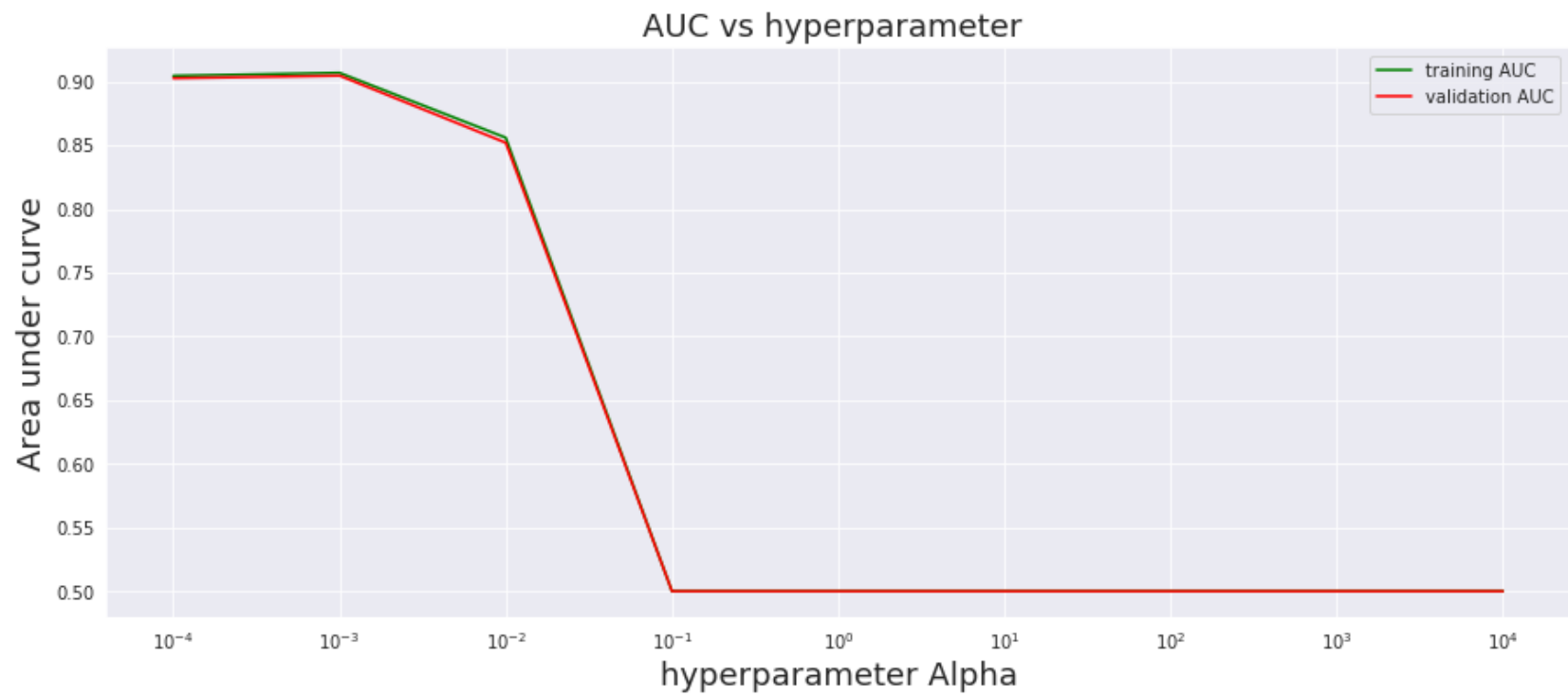
```

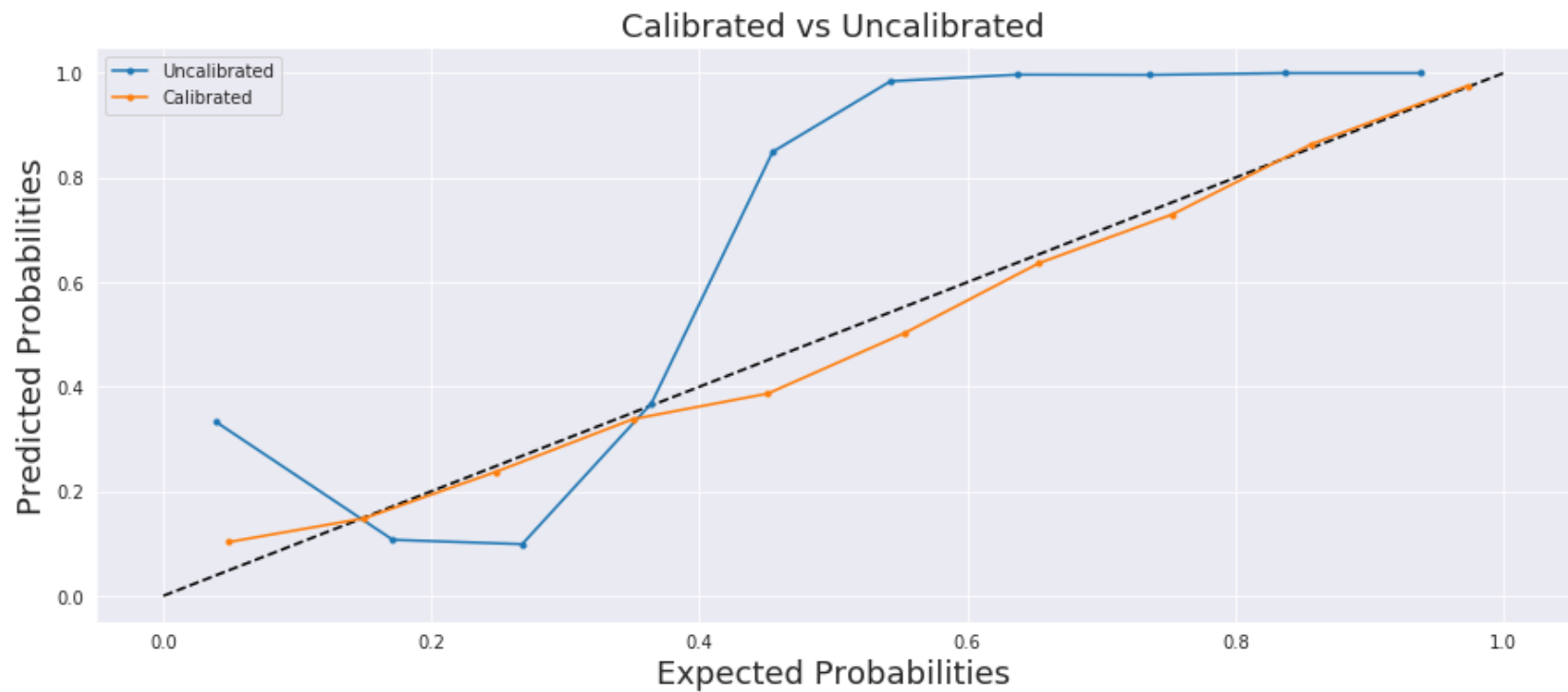
best hyperparameter is {'alpha': 0.001}

AUC on training data is 0.9069544023365721

AUC on test data is 0.9056091185502038

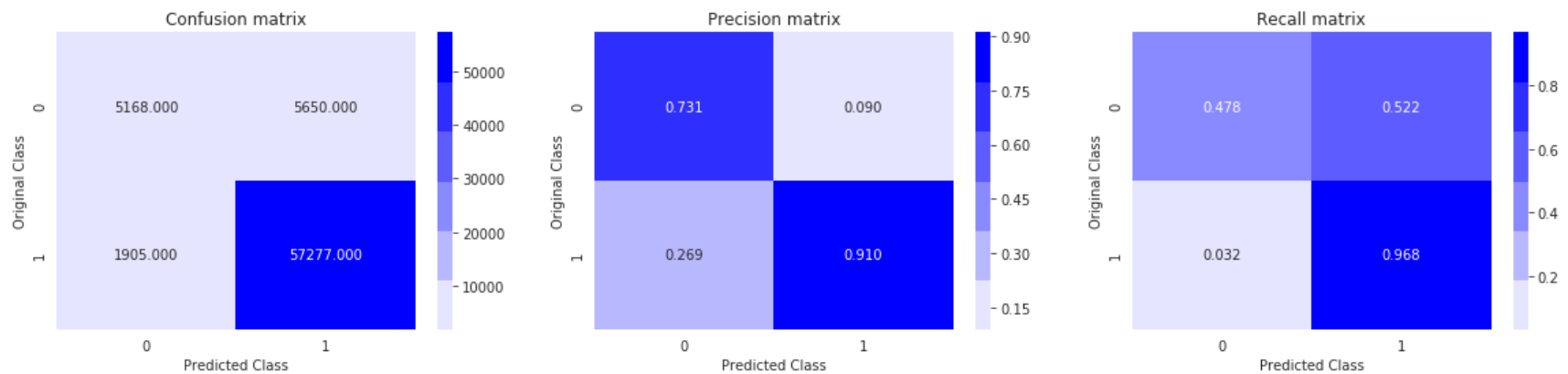
The Calibration Curve





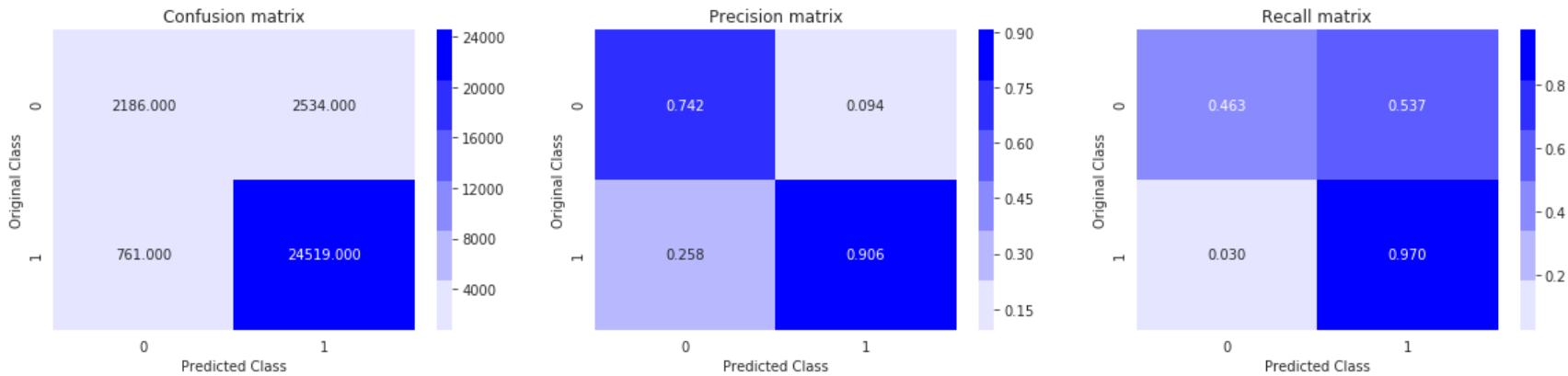
for Training data:

Confusion Matrix

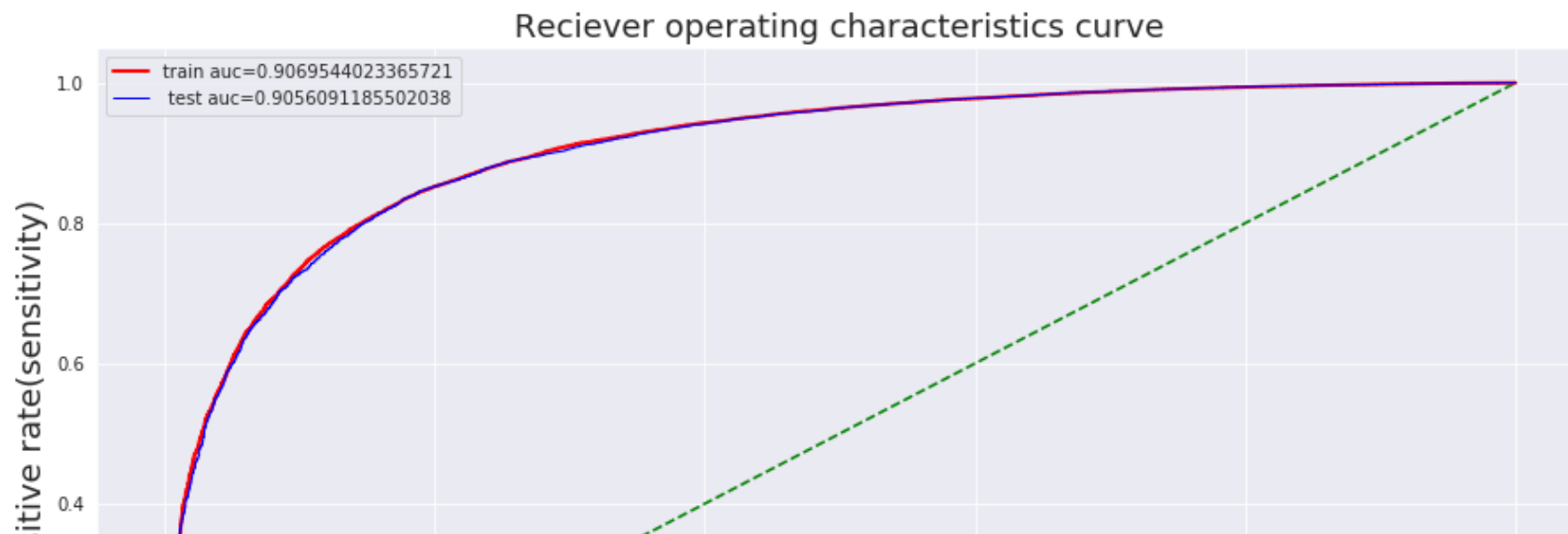


for Test data

Confusion Matrix



plotting ROC on Test data



[5.1.4] Applying Linear SVM on TFIDF W2V, SET 4

In [60]:

```
1 #TFIDF WORD 2 VECTOR
2 train_tfidf2v = pickle.load(open('train_tfidf2v.p', 'rb'))
3 test_tfidf2v = pickle.load(open('test_tfidf2v.p', 'rb'))
4
5 print(len(train_tfidf2v))
6 print(len(test_tfidf2v))
```

70000

30000

Using L2 Regularization

In [61]:

```

1  """classifier for tuning"""
2  gsm = svm('l2',train_tfidfv2v,Y_train)
3  print(gsm)
4  print('*****\n')
5
6
7  """best hyperparameter"""
8  best_alpha_l2_tfidfv2v = error('l2',train_tfidfv2v,Y_train)['alpha']
9  print('*****\n')
10
11
12  """best Classifier fitted with tuned Hyperparameter"""
13  train_auc_l2_tfidfv2v,test_auc_l2_tfidfv2v,train_proba,test_proba,train_pred,test_pred = best_classifier(bes
14  print('*****')
15
16
17  """Reliability Curve"""
18  print('The Calibration Curve')
19  reliability_curve(best_alpha_l2_tfidfv2v,'l2',train_tfidfv2v,test_tfidfv2v,Y_train,Y_test)
20  print('*****\n')
21
22  """Confusion Matrix"""
23  print('for Training data:\n')
24  plot_confusion_matrix(Y_train,train_pred)
25
26  print('for Test data')
27  plot_confusion_matrix(Y_test,test_pred)
28  print('*****\n')
29
30  """ROC CURVE"""
31  plot_roc(Y_train,train_proba,Y_test,test_proba,train_auc_l2_tfidfv2v,test_auc_l2_tfidfv2v)
32

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 10.1s finished

```

GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=5),
             error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,
                                     class_weight=None, early_stopping=False,

```

```

epsilon=0.1, eta0=0.0, fit_intercept=True,
l1_ratio=0.15, learning_rate='optimal',
loss='hinge', max_iter=1000,
n_iter_no_change=5, n_jobs=None,
penalty='l2', power_t=0.5, random_state=42,
shuffle=True, tol=0.001,
validation_fraction=0.1, verbose=0,
warm_start=False),
iid='warn', n_jobs=None,
param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000,
                      10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)

```

Corresponding hypereparameter will give best auc on CV data

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

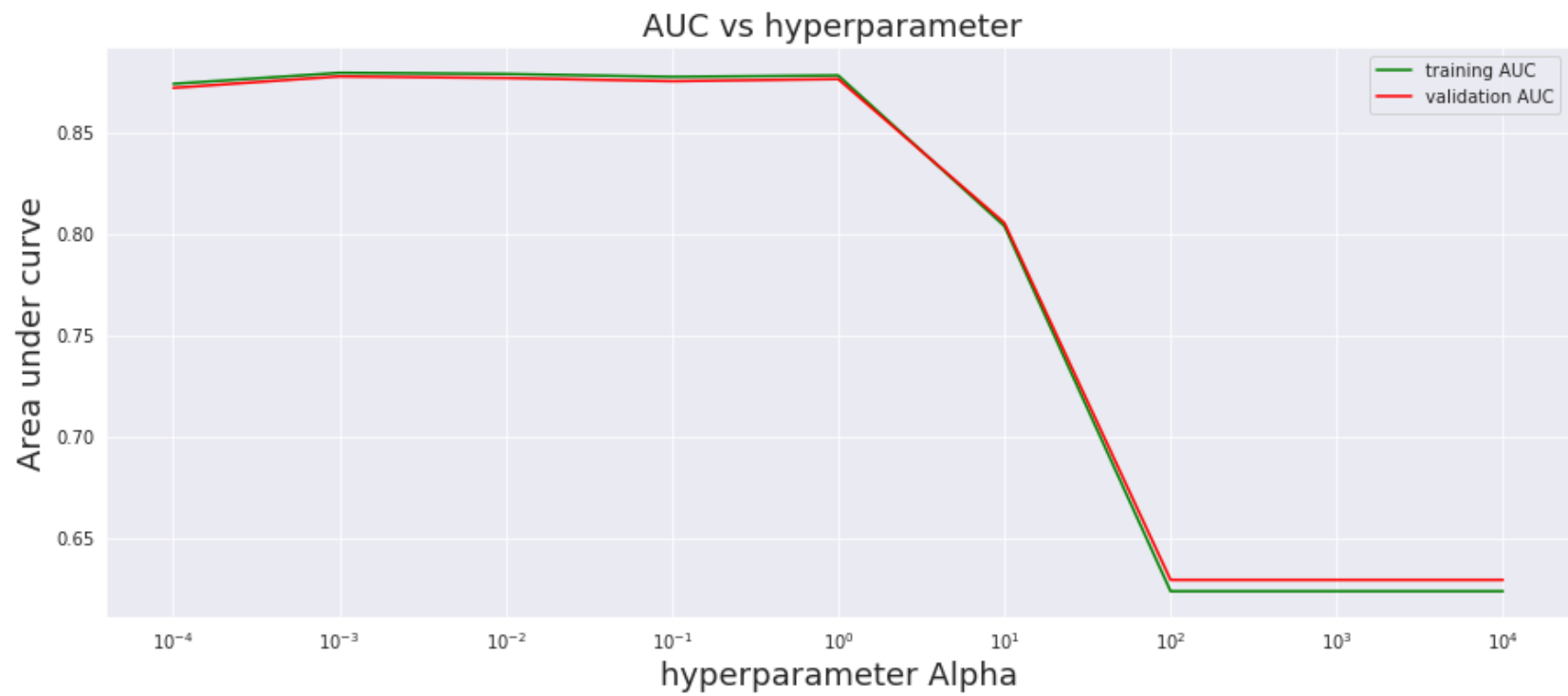
[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 11.3s finished

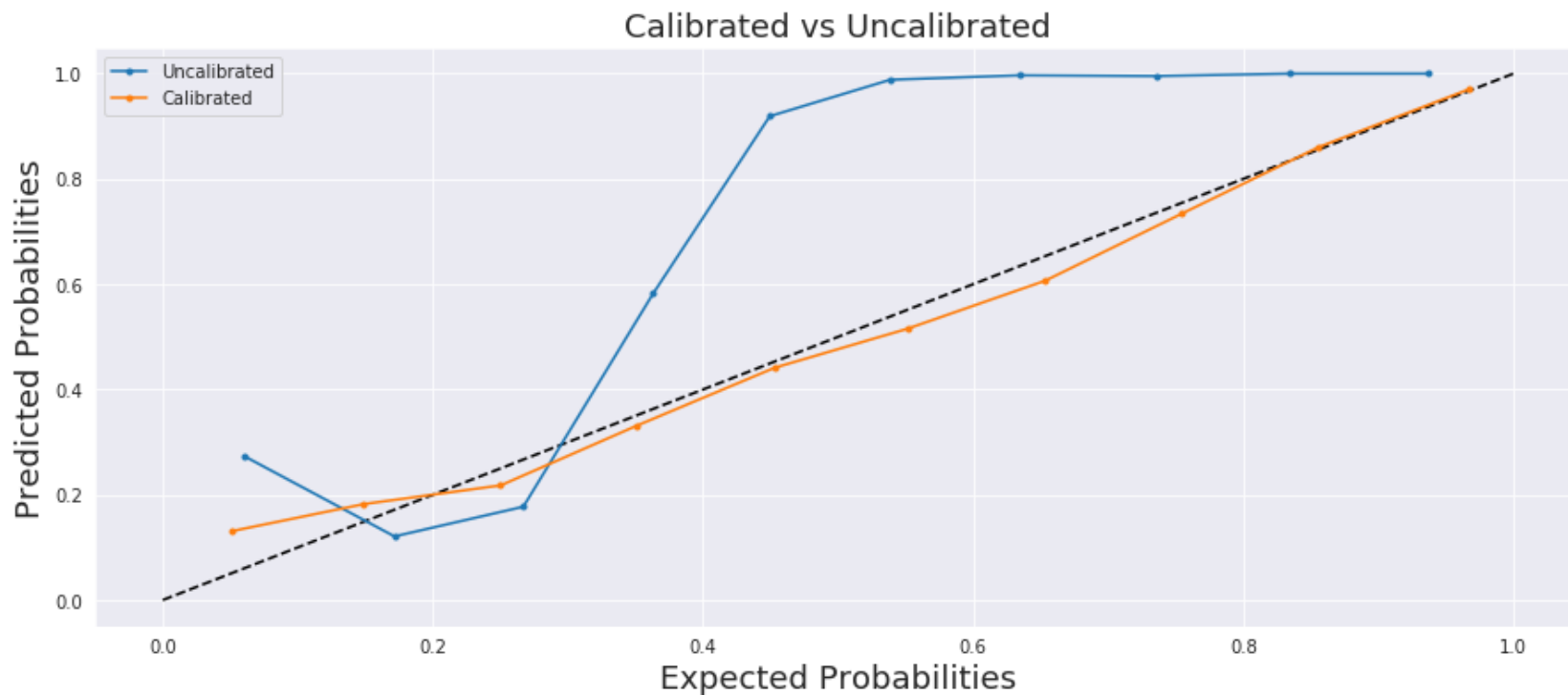
best hyperparameter is {'alpha': 0.001}

AUC on training data is 0.881121875946514

AUC on test data is 0.8774580042506973

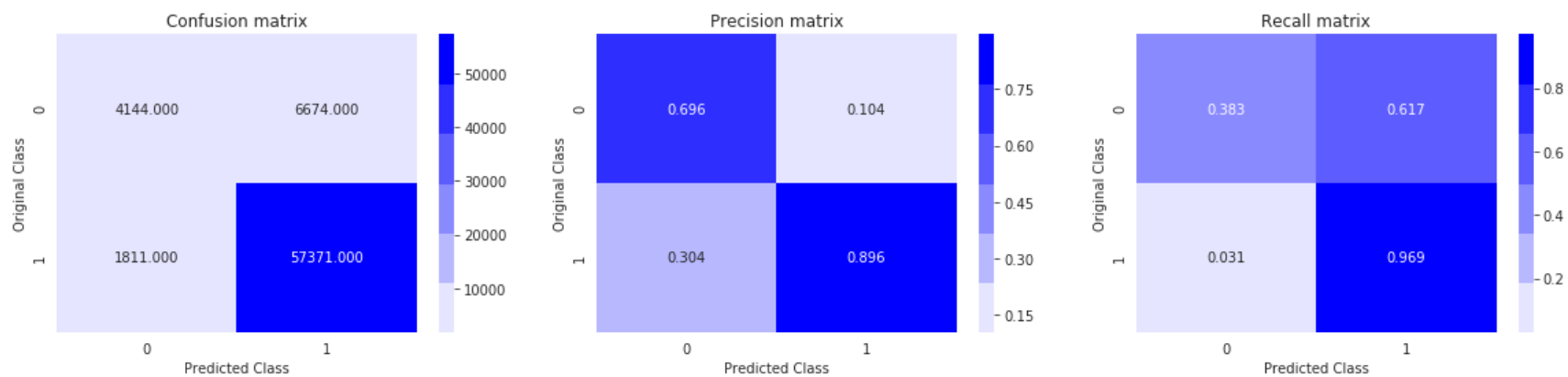
The Calibration Curve





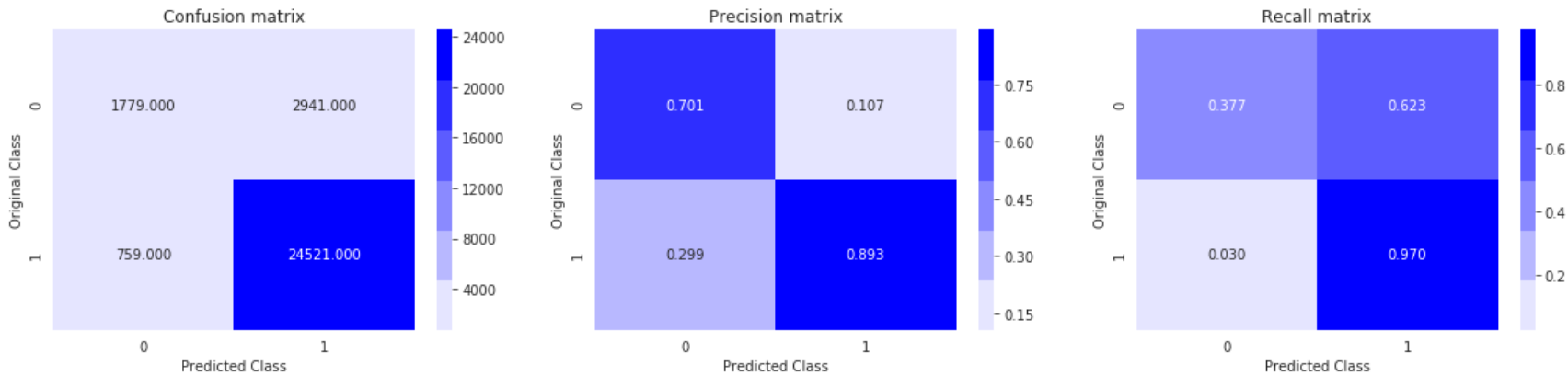
for Training data:

Confusion Matrix

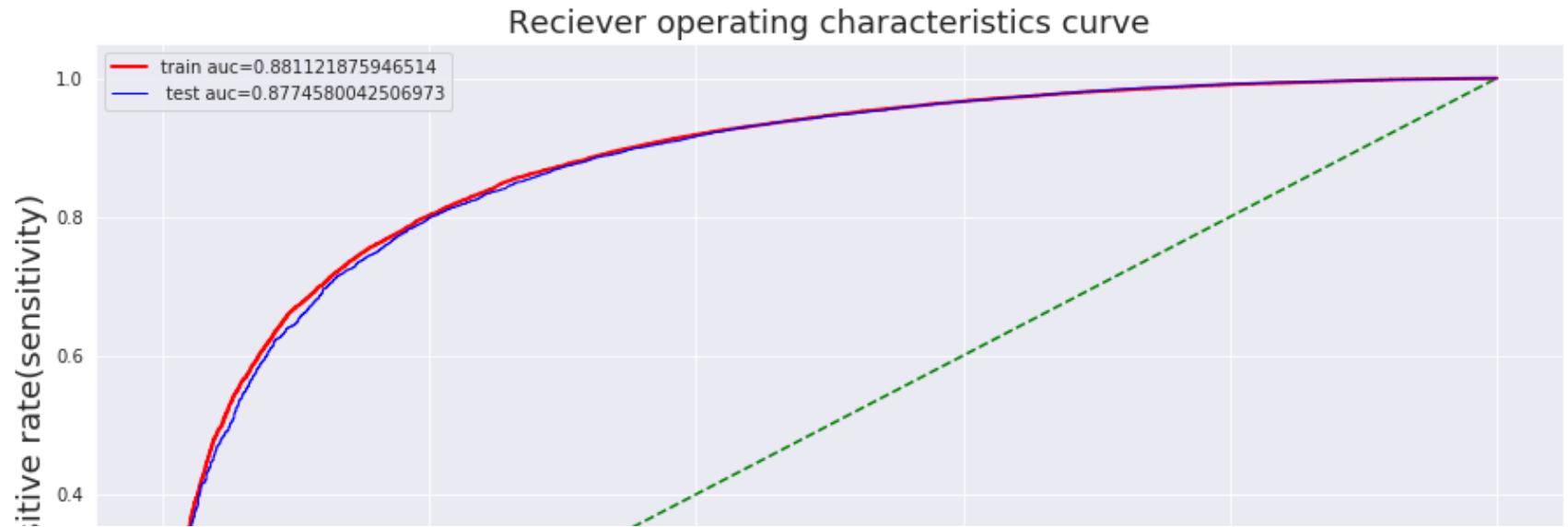


for Test data

Confusion Matrix



plotting ROC on Test data



In [0]:

1

Using L1 Regularization

In [62]:

```

1  """classifier for tuning"""
2  gsm = svm('l1',train_tfidfw2v,Y_train)
3  print(gsm)
4  print('*****\n')
5
6
7  """best hyperparameter"""
8  best_alpha_l1_tfidfw2v = error('l1',train_tfidfw2v,Y_train)['alpha']
9  print('*****\n')
10
11
12  """best Classifier fitted with tuned Hyperparameter"""
13  train_auc_l1_tfidfw2v,test_auc_l1_tfidfw2v,train_proba,test_proba,train_pred,test_pred = best_classifier(bes
14  'l1',train_tfidfw2v,Y_train,test_tfidfw2v,Y_test)
15  print('*****')
16
17
18  """Reliability Curve"""
19  print('The Calibration Curve')
20  reliability_curve(best_alpha_l1_tfidfw2v,'l1',train_tfidfw2v,test_tfidfw2v,Y_train,Y_test)
21  print('*****\n')
22
23  """Confusion Matrix"""
24  print('for Training data:\n')
25  plot_confusion_matrix(Y_train,train_pred)
26
27  print('for Test data')
28  plot_confusion_matrix(Y_test,test_pred)
29  print('*****\n')
30
31  """ROC CURVE"""
32  plot_roc(Y_train,train_proba,Y_test,test_proba,train_auc_l1_tfidfw2v,test_auc_l1_tfidfw2v)
33

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 16.6s finished

```

GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=5),
             error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,

```

```

class_weight=None, early_stopping=False,
epsilon=0.1, eta0=0.0, fit_intercept=True,
l1_ratio=0.15, learning_rate='optimal',
loss='hinge', max_iter=1000,
n_iter_no_change=5, n_jobs=None,
penalty='l1', power_t=0.5, random_state=42,
shuffle=True, tol=0.001,
validation_fraction=0.1, verbose=0,
warm_start=False),
iid='warn', n_jobs=None,
param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000,
                      10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)
*****

```

Corresponding hypereparameter will give best auc on CV data
 Fitting 5 folds for each of 9 candidates, totalling 45 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 45 out of 45 | elapsed: 15.9s finished

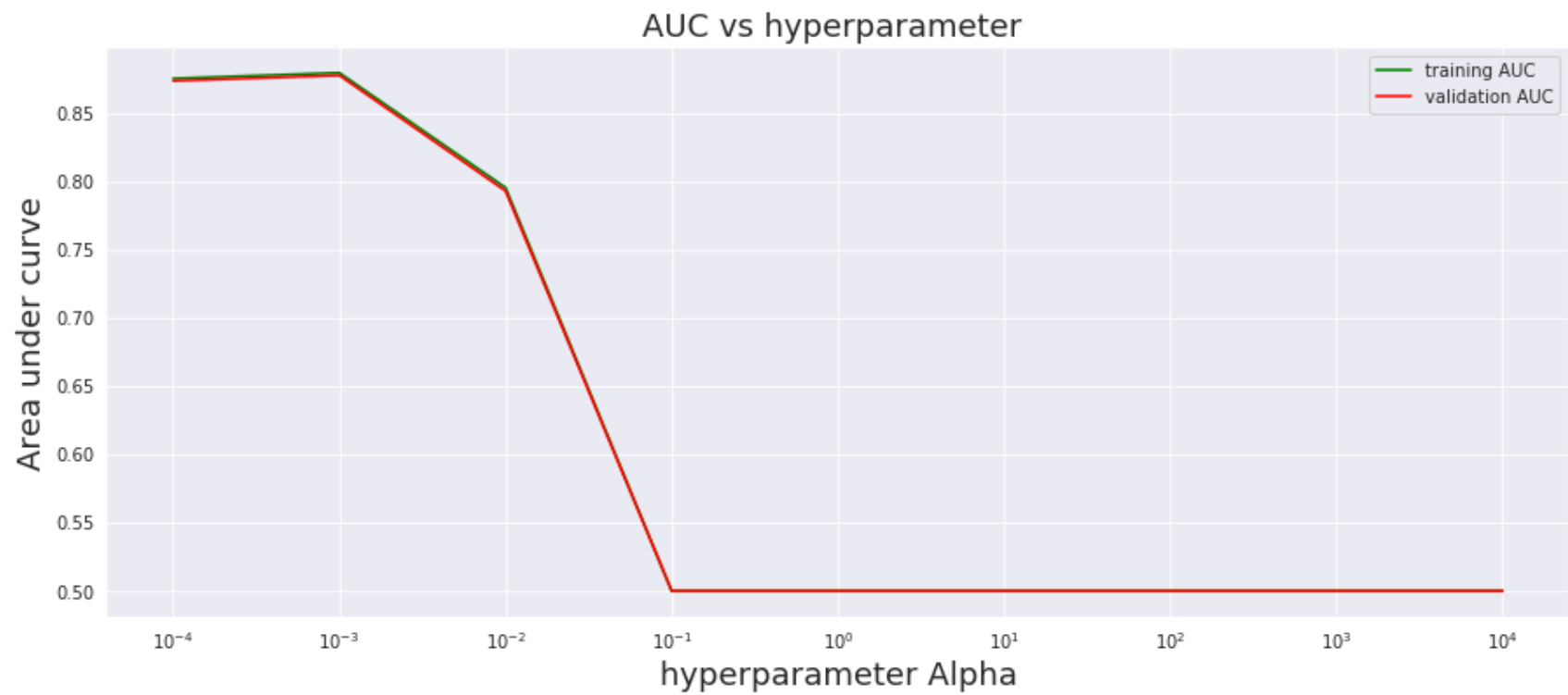
```

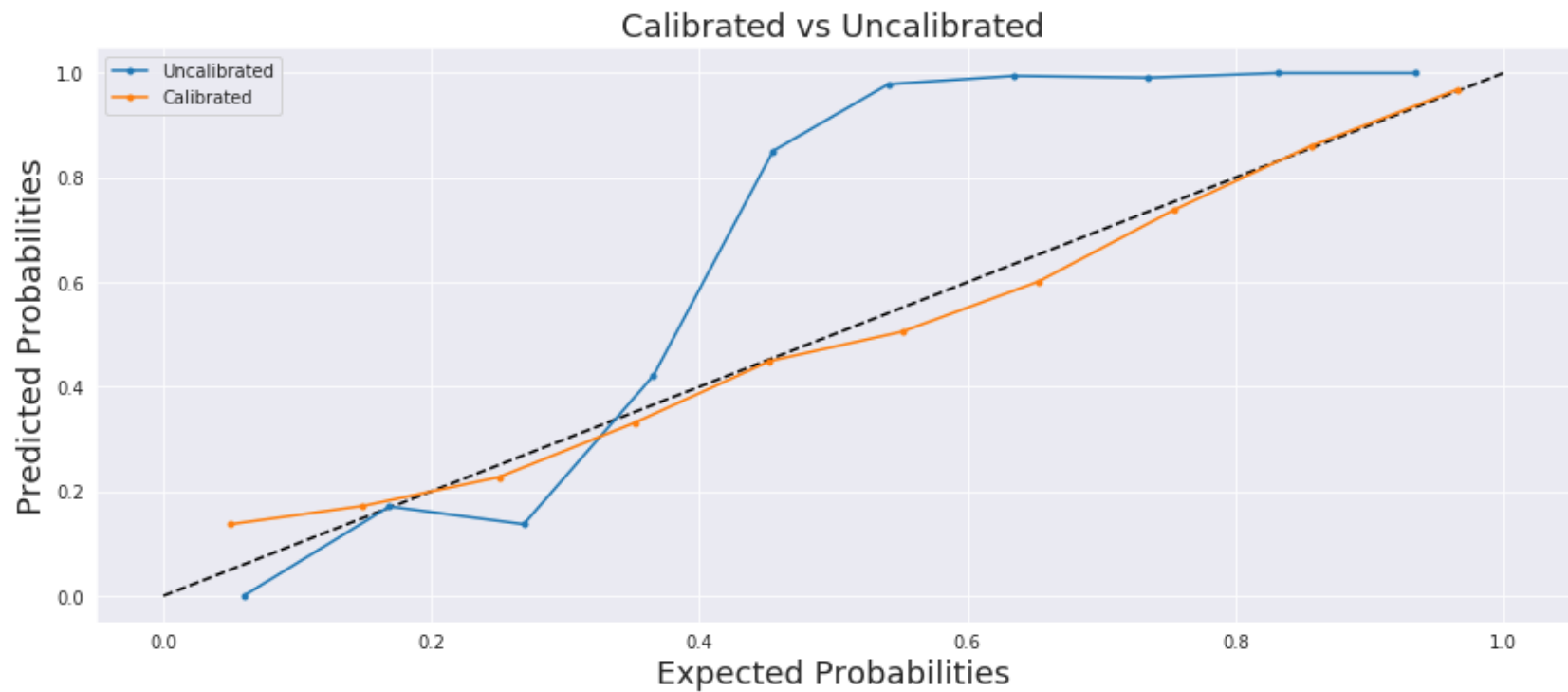
best hyperparameter is {'alpha': 0.001}

AUC on training data is 0.8799045916367207

AUC on test data is 0.8762896407691483

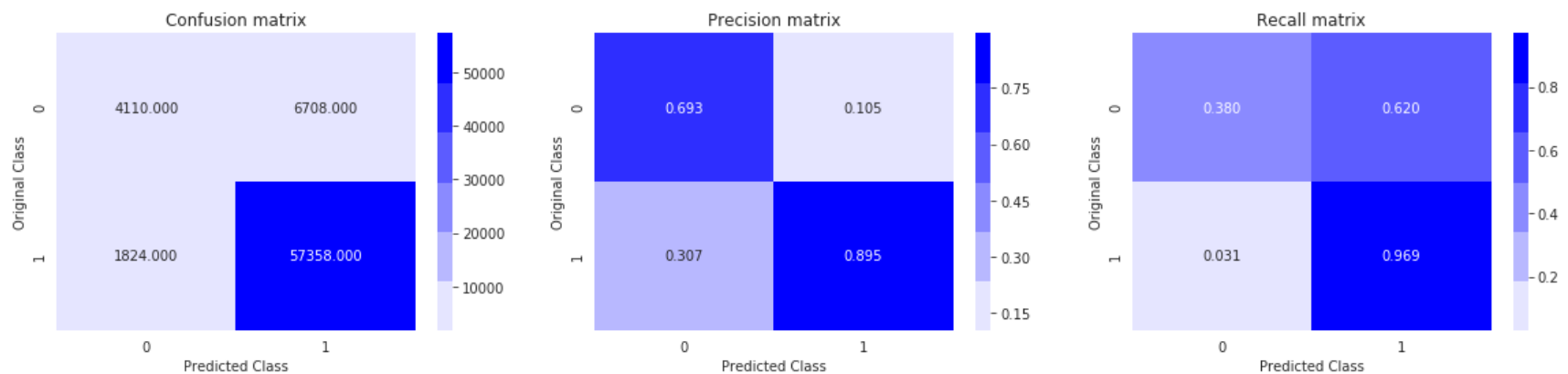
The Calibration Curve



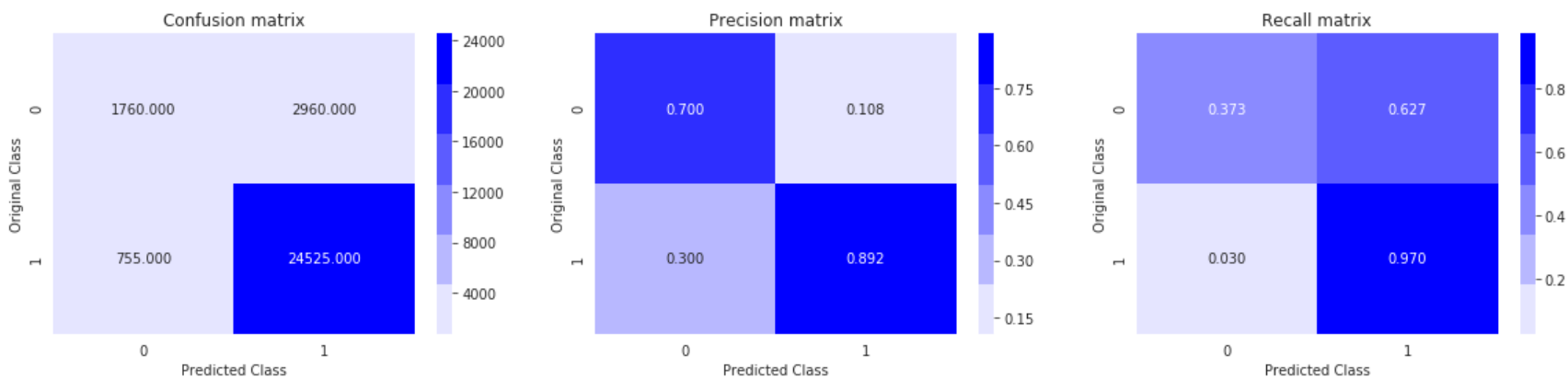


for Training data:

Confusion Matrix



for Test data
Confusion Matrix



plotting ROC on Test data



In [0]:

```
1 pickle.dump(best_alpha_l2_tfidf2v,open('best_alpha_l2_tfidf2v.p','wb'))
2 pickle.dump(test_auc_l2_tfidf2v,open('test_auc_l2_tfidf2v.p','wb'))
3 pickle.dump(best_alpha_l1_tfidf2v,open('best_alpha_l1_tfidf2v.p','wb'))
4 pickle.dump(test_auc_l1_tfidf2v,open('test_auc_tfidf2v.p','wb'))
```

[6] Conclusions

```
In [63]: 1 from prettytable import PrettyTable
2
3 #table for random forest
4 table = PrettyTable()
5 no = [1,2,3,4,5,6,7,8]
6 vectorizers = ['Bag of vectors','Bag of Vectors','TFIDF','TFIDF','Average Word 2 vector','Average Word 2 Vec
7 regularization = ['l2','l1','l2','l1','l2','l1','l2','l1']
8 alphas = [best_alpha_l2_bow,best_alpha_l1_bow,best_alpha_l2_tfidf,best_alpha_l1_tfidf,best_alpha_l2_avgw2v,b
9           best_alpha_l2_tfidfw2v,'0.001']
10 AUC = [test_auc_l2_bow,test_auc_l1_bow,test_auc_l2_tfidf,test_auc_l1_tfidf,test_auc_l2_avgw2v,test_auc_l1_av
11        test_auc_l2_tfidfw2v,test_auc_l1_tfidfw2v]#their respective auc scores
12
13 table.add_column("SNo",no)
14 table.add_column('Vectorizers',vectorizers)
15 table.add_column('Regularization',regularization)
16 table.add_column('Hyperparameter(alpha)',alphas)
17 table.add_column('AUC on test',AUC)
18 print('\t\t\t Table for SGDClassifier using Hinge Loss')
19 print(table)
```

Table for SGDClassifier using Hinge Loss

SNo	Vectorizers	Regularization	Hyperparameter(alpha)	AUC on test
1	Bag of vectors	l2	0.001	0.9402083193654796
2	Bag of Vectors	l1	0.0001	0.9303229130350246
3	TFIDF	l2	0.0001	0.9557599378486376
4	TFIDF	l1	0.0001	0.9330814202960738
5	Average Word 2 vector	l2	0.001	0.9058933923111994
6	Average Word 2 Vector	l1	0.001	0.9056091185502038
7	TFIDF Word 2 Vector	l2	0.001	0.8774580042506973
8	TFIDF Word 2 Vector	l1	0.001	0.8762896407691483

- We observed that the L2 regularization seems to work better than l1 regularization in almost all the cases.
- The L1 regularization because of creating sparsity, increase the value of coefficients of most important features.
- Calibrated Classifier tells more realistic picture of expected and predicted probabilities, hence working greatly in this case of linear SVM. We plotted the reliability curve to actually see how closely the SGDClassifier fits the sigmoid function, thus using 'sigmoid' method rather than isotonic.
- Though Word 2 Vector gives semantic understanding, then also simple TFIDF Scores gave best AUC Score

