

# Amazon Fine Food Reviews Analysis Using SVM

Data Source: <https://www.kaggle.com/snap/amazon-fine-food-reviews> (<https://www.kaggle.com/snap/amazon-fine-food-reviews>)

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

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.

## 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 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 id above 3, then the recommendation wil be set to "positive". Otherwise, it will be set to "negative".

```
In [76]: %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

from sklearn.model_selection import train_test_split
from sklearn.metrics import roc_auc_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import accuracy_score
from sklearn.cross_validation import cross_val_score
from collections import Counter
from sklearn import cross_validation
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
from sklearn.calibration import CalibratedClassifierCV
from sklearn.svm import SVC
from sklearn.linear_model import SGDClassifier
from sklearn.svm import LinearSVC

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

## [1]. Reading Data

```
In [77]: # using SQLite Table to read data.
con = sqlite3.connect('D:\\TGM\\ML\\AmazonFineFoodReviews\\database.sqlite')

# filtering only positive and negative reviews i.e.
# not taking into consideration those reviews with Score=3
# SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 500000 data points
# you can change the number to any other number based on your computing power

# filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000""", con)
# for tsne assignment you can take 5k data points

filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 100000""", con)

# Give reviews with Score>3 a positive rating(1), and reviews with a score<3 a negative rating(0).
def partition(x):
    if x < 3:
        return 0
    return 1

#changing reviews with score less than 3 to be positive and vice-versa
actualScore = filtered_data['Score']
positiveNegative = actualScore.map(partition)
filtered_data['Score'] = positiveNegative
print("Number of data points in our data", filtered_data.shape)
filtered_data.head(3)
```

Number of data points in our data (100000, 10)

Out[77]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpfulne
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1

```
In [78]: display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

```
In [79]: print(display.shape)
display.head()
```

```
(80668, 7)
```

```
Out[79]:
```

	UserId	ProductId	ProfileName	Time	Score	Text	Cou
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 [80]: display[display['UserId']=='AZY10LLTJ71NX']
```

```
Out[80]:
```

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

```
In [81]: display['COUNT(*)'].sum()
```

```
Out[81]: 393063
```

## Exploratory Data Analysis

### [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 [82]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[82]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpful
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2

As can be seen above the same user has multiple reviews of the with the 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 [83]: #Sorting data according to ProductId in ascending order
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False, kind='quicksort', na_position='last')
```

```
In [84]: #Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first', inplace=False)
final.shape
```

```
Out[84]: (87775, 10)
```

```
In [85]: #Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[85]: 87.775
```

**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 [86]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)

display.head()
```

Out[86]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpfulr
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2

```
In [87]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [88]: #Before starting the next phase of preprocessing Lets see the number of entrie
s left
print(final.shape)

#How many positive and negative reviews are present in our dataset?
print(final['Score'].value_counts())

(87773, 10)
1    73592
0    14181
Name: Score, dtype: int64
```

## [3]. Text Preprocessing.

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 [89]: # printing some random reviews
sent_0 = final['Text'].values[0]
print(sent_0)
print("="*50)

sent_1000 = final['Text'].values[1000]
print(sent_1000)
print("="*50)

sent_1500 = final['Text'].values[1500]
print(sent_1500)
print("="*50)

sent_4900 = final['Text'].values[4900]
print(sent_4900)
print("="*50)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought were eaten and I threw the rest away. I would not buy the candy again.

=====

was way to hot for my blood, took a bite and did a jig lol

=====

My dog LOVES these treats. They tend to have a very strong fish oil smell. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of these without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's definitely worth it to buy a big bag if your dog eats them a lot.

=====

```
In [90]: # https://stackoverflow.com/a/47091490/4084039
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"'re", " are", phrase)
    phrase = re.sub(r"'s", " is", phrase)
    phrase = re.sub(r"'d", " would", phrase)
    phrase = re.sub(r"'ll", " will", phrase)
    phrase = re.sub(r"'t", " not", phrase)
    phrase = re.sub(r"'ve", " have", phrase)
    phrase = re.sub(r"'m", " am", phrase)
    return phrase
```

```
In [91]: sent_4900 = decontracted(sent_4900)
print(sent_4900)
print("="*50)
```

My dog LOVES these treats. They tend to have a very strong fish oil smell. So if you are afraid of the fishy smell, do not get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of these without worrying about him over eating. Amazon is price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It is definitely worth it to buy a big bag if your dog eats them a lot.

=====

```
In [92]: #remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

```
In [93]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

was way to hot for my blood took a bite and did a jig lol

```
In [94]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have revmoved in the 1st
step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours',
'ourselves', 'you', "you're", "you've",\
    "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
, 'him', 'his', 'himself', \
    'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'it
self', 'they', 'them', 'their',\
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 't
hat', "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', 'becau
se', 'as', 'until', 'while', 'of', \
    'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
'through', 'during', 'before', 'after',\
    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on',
'off', 'over', 'under', 'again', 'further',\
    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'a
ll', 'any', 'both', 'each', 'few', 'more',\
    'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'tha
n', 'too', 'very', \
    's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shoul
d've", 'now', 'd', 'll', 'm', 'o', 're', \
    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
"didn't", 'doesn', "doesn't", 'hadn',\
    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'm
a', 'mightn', "mightn't", 'mustn',\
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shoul
dn't", 'wasn', "wasn't", 'weren', "weren't", \
    'won', "won't", 'wouldn', "wouldn't"])
```

```
In [95]: # fitting all the above students
from bs4 import BeautifulSoup
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentence in tqdm(final['Text'].values):
    sentence = re.sub(r"http\S+", "", sentence)
    sentence = BeautifulSoup(sentence, 'lxml').get_text()
    sentence = decontracted(sentence)
    sentence = re.sub("\S*\d\S*", "", sentence).strip()
    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
    preprocessed_reviews.append(sentence.strip())
```

```
100%|███████████████████████████████████████████████████████████████████████████████  
██████████| 87773/87773 [00:56<00:00, 1566.45it/s]
```

```
In [96]: preprocessed_reviews[1500]
```

```
Out[96]: 'way hot blood took bite jig lol'
```

```
In [97]: final['cleaned_text']=preprocessed_reviews
```

```
In [98]: final.shape
```

```
Out[98]: (87773, 11)
```

```
In [99]: final["Score"].value_counts()
```

```
Out[99]: 1    73592  
         0    14181  
         Name: Score, dtype: int64
```



```
In [100]: #Sorted the data based on time and took 100k data points
final["Time"] = pd.to_datetime(final["Time"], unit = "s")
final = final.sort_values(by = "Time")
final.head()
```

Out[100]:

	<b>Id</b>	<b>ProductId</b>	<b>UserId</b>	<b>ProfileName</b>	<b>HelpfulnessNumerator</b>	<b>He</b>
<b>70688</b>	76882	B00002N8SM	A32DW342WBJ6BX	Buttersugar	0	0
<b>1146</b>	1245	B00002Z754	A29Z5PI9BW2PU3	Robbie	7	7
<b>1145</b>	1244	B00002Z754	A3B8RCEI0FXFI6	B G Chase	10	10
<b>28086</b>	30629	B00008RCMI	A19E94CF5O1LY7	Andrew Arnold	0	0
<b>28087</b>	30630	B00008RCMI	A284C7M23F0APC	A. Mendoza	0	0

```
In [101]: Y = final['Score'].values
X = final['cleaned_text'].values
#X = preprocessed_reviews
#Y = np.array(final['Score'])
print(Y.shape)
print(type(Y))
print(X.shape)
print(type(X))
```

```
(87773,)
<class 'numpy.ndarray'>
(87773,)
<class 'numpy.ndarray'>
```

```
In [102]: # split the data set into train and test
X_Train, X_Test, Y_Train, Y_Test = train_test_split(X,Y,test_size=0.3, random_
state=0)

# split the train data set into cross validation train and cross validation te
st
X_tr, X_cv, Y_tr, Y_cv = train_test_split(X,Y, test_size=0.3, random_state=0)

print('='*100)
print("After splitting")
print("X_Train Shape:",X_Train.shape, "Y_Train Shape:",Y_Train.shape)
print("X_cv Shape:",X_cv.shape, "Y_cv Shape",Y_cv.shape)
print("X_Test Shape",X_Test.shape, "Y_Test Shape",Y_Test.shape)

=====
=====
After splitting
X_Train Shape: (61441,) Y_Train Shape: (61441,)
X_cv Shape: (26332,) Y_cv Shape (26332,)
X_Test Shape (26332,) Y_Test Shape (26332,)
```

## [3.2] Preprocess Summary

## [4] Featurization

### [4.1] BAG OF WORDS

```
In [103]: #Bow
count_vect = CountVectorizer(ngram_range=(1,2)) #in scikit-learn
count_vect.fit(X_Train)
print("some feature names ", count_vect.get_feature_names()[:10])
X_Train_Bow = count_vect.transform(X_Train)
X_Test_Bow = count_vect.transform(X_Test)
X_CV_Bow = count_vect.transform(X_cv)

print('='*50)

#final_counts = count_vect.transform(X_Test)

print("the type of X Train : ",type(X_Train_Bow))
print("the shape of Train BOW vectorizer ",X_Train_Bow.get_shape())
print("the shape of Test BOW vectorizer ",X_Test_Bow.get_shape())
print("the shape of CV BOW vectorizer ",X_CV_Bow.get_shape())
#print("the number of unique words ", final_counts.get_shape()[1])

some feature names  ['aa', 'aa caffene', 'aa coffee', 'aa cups', 'aa dark',
'aa extra', 'aa favorite', 'aa kona', 'aa may', 'aa not']
=====
the type of X Train :  <class 'scipy.sparse.csr.csr_matrix'>
the shape of Train BOW vectorizer  (61441, 1076376)
the shape of Test BOW vectorizer  (26332, 1076376)
the shape of CV BOW vectorizer  (26332, 1076376)
```

```
In [104]: import warnings
warnings.filterwarnings('ignore')
scalar = StandardScaler(with_mean=False)
X_Train_Bow = scalar.fit_transform(X_Train_Bow)
X_Test_Bow = scalar.transform(X_Test_Bow)
X_CV_Bow = scalar.transform(X_CV_Bow)

print("the type of X Train : ",type(X_Train_Bow))
print("the shape of Train BOW vectorizer ",X_Train_Bow.get_shape())
print("the shape of Test BOW vectorizer ",X_Test_Bow.get_shape())
print("the shape of CV BOW vectorizer ",X_CV_Bow.get_shape())

the type of X Train :  <class 'scipy.sparse.csr.csr_matrix'>
the shape of Train BOW vectorizer  (61441, 1076376)
the shape of Test BOW vectorizer  (26332, 1076376)
the shape of CV BOW vectorizer  (26332, 1076376)
```

## [4.1] Linear Kernel

```

In [109]: import math
def Optimal_Lamda_L1(X_Train,Y_Train,X_CV,Y_CV):
    train_AUC_L1 = []
    CV_AUC_L1 = []
    tuned_parameters=[10**-4, 10**-3, 10**-2, 10**-1, 1,10**1, 10**2, 10**3, 1
0**4]
    for j in tqdm(tuned_parameters):
        clf = SGDClassifier(alpha=j, class_weight='balanced', penalty='l1', lo
ss='hinge')
        calibrated_clf = CalibratedClassifierCV(clf, cv=5)
        calibrated_clf.fit(X_Train, Y_Train)
        y_train_pred = calibrated_clf.predict_proba(X_Train)[:,-1]
        y_cv_pred = calibrated_clf.predict_proba(X_CV)[:,-1]
        train_AUC_L1.append(roc_auc_score(Y_Train,y_train_pred))
        CV_AUC_L1.append(roc_auc_score(Y_CV, y_cv_pred))

    #Error plots with penalty L1
    plt.plot(np.log(tuned_parameters), train_AUC_L1, label='Train AUC-L1')
    plt.plot(np.log(tuned_parameters), CV_AUC_L1, label='CV AUC-L1')
    plt.legend()
    plt.xlabel("Hyperparameter (Lambda)")
    plt.ylabel("AUC")
    plt.title("AUC PLOT")
    plt.show()

    #Cv auc scores with penalty L1
    print("CV AUS Scores with Penalty=? Cv auc scores with penalty L1")
    print(CV_AUC_L1)
    print("Maximun AUC value :",max(CV_AUC_L1))
    print("Index",CV_AUC_L1.index(max(CV_AUC_L1)))

```

```

In [110]: import math
def Optimal_Lamda_L2(X_Train,Y_Train,X_CV,Y_CV):
    train_AUC_L2 = []
    CV_AUC_L2 = []
    cv_scores = []
    tuned_parameters=[10**-4, 10**-3, 10**-2, 10**-1, 1,10**1, 10**2, 10**3, 1
0**4]
    for j in tqdm(tuned_parameters):
        clf = SGDClassifier(alpha=j, class_weight='balanced', penalty='l2', lo
ss='hinge')
        calibrated_clf = CalibratedClassifierCV(clf, cv=5)
        calibrated_clf.fit(X_Train, Y_Train)
        y_train_pred = calibrated_clf.predict_proba(X_Train)[:,-1]
        y_cv_pred = calibrated_clf.predict_proba(X_CV)[:,-1]
        train_AUC_L2.append(roc_auc_score(Y_Train,y_train_pred))
        CV_AUC_L2.append(roc_auc_score(Y_CV, y_cv_pred))

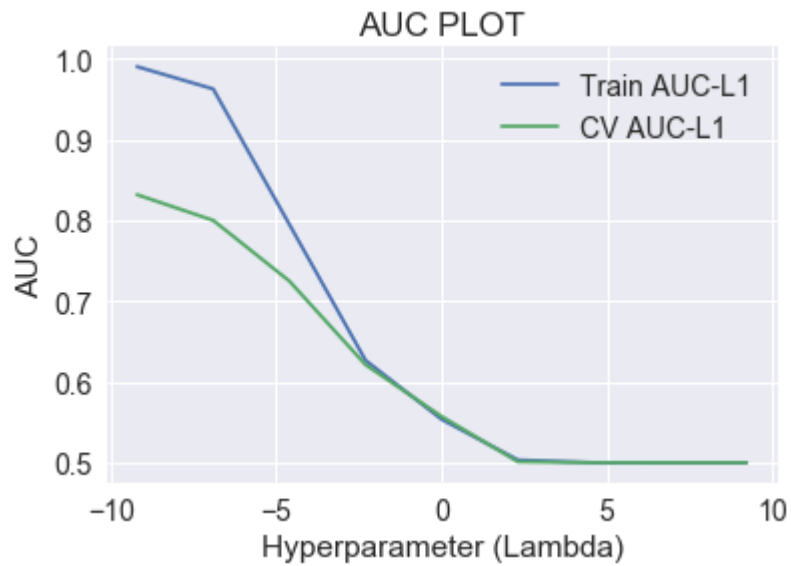
    #Error plots with penalty L2
    plt.plot(np.log(tuned_parameters), train_AUC_L2, label='Train AUC-L2')
    plt.plot(np.log(tuned_parameters), CV_AUC_L2, label='CV AUC-L2')
    plt.legend()
    plt.xlabel("Hyperparameter (Lambda)")
    plt.ylabel("AUC")
    plt.title("AUC PLOT")
    plt.show()

    #Cv auc scores with penalty L2
    print("CV AUS Scores with Penalty=? Cv auc scores with penalty L2")
    print(CV_AUC_L2)
    print("Maximun AUC value :",max(CV_AUC_L2))
    print("Index",CV_AUC_L2.index(max(CV_AUC_L2)))

```

```
In [107]: Optimal_Lamda_L1(X_Train_Bow, Y_Train, X_CV_Bow,Y_cv)
```

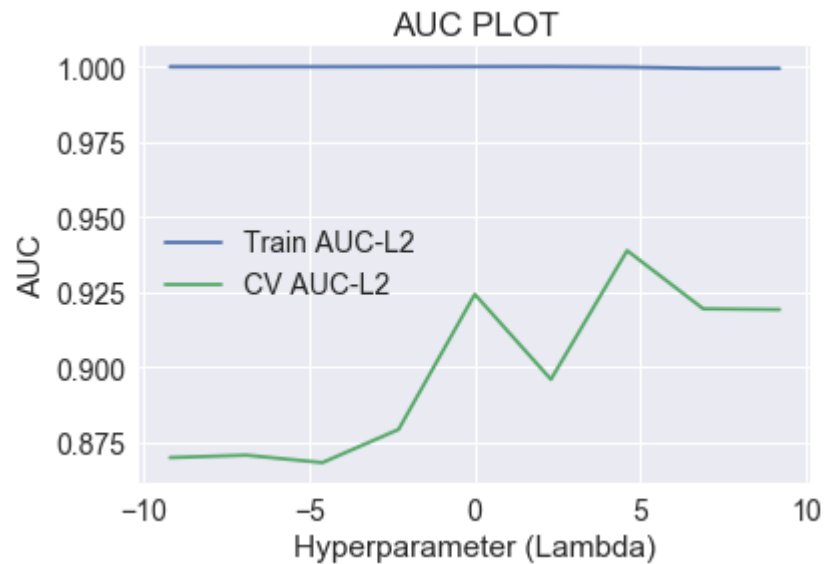
```
100%|███████████| 9/9 [00:28<00:00, 3.32s/it]
```



```
CV AUS Scores with Penalty=? Cv auc scores with penalty L1
[0.831609221053483, 0.7997954681454404, 0.7246442854099527, 0.62108846161426
4, 0.5570247894753875, 0.5010185390513413, 0.5000226264820345, 0.5, 0.5]
Maximun AUC value : 0.831609221053483
Index 0
```

```
In [112]: Optimal_Lamda_L2(X_Train_Bow, Y_Train, X_CV_Bow,Y_cv)
```

```
100%|███████████| 9/9 [00:32<00:00, 3.81s/it]
```



```
CV AUC Scores with Penalty=? Cv auc scores with penalty L2
[0.8699834246323106, 0.8707944295717452, 0.8683134256630607, 0.87934106212062
7, 0.924285089740454, 0.895990551044296, 0.9388157908518728, 0.91944666719080
58, 0.9191993897754295]
Maximun AUC value : 0.9388157908518728
Index 6
```

### **<font color = blue>[4.1.3] ROC Curve of SVM </font>**

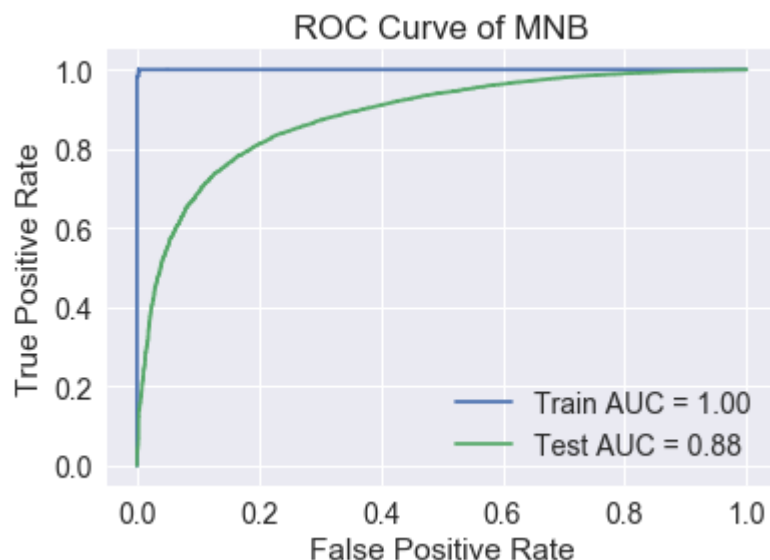
```
In [114]: #Testing with test data
clf = SGDClassifier(alpha=0.1, class_weight='balanced', loss='hinge')
calibrated_clf = CalibratedClassifierCV(clf, cv=5)
calibrated_clf.fit(X_Train_Bow,Y_Train)
prediction = calibrated_clf.predict_proba(X_Test_Bow)[: ,1]
print(prediction)
print(clf)

Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, calibrated_clf.predict_proba(X_Train_Bow)[: ,1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, calibrated_clf.predict_proba(X_Test_Bow)[: ,1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of MNB')
plt.show()
```

[0.90236714 0.90391762 0.7174323 ... 0.86114097 0.88722453 0.32113633]

```
SGDClassifier(alpha=0.1, average=False, class_weight='balanced', epsilon=0.1,
eta0=0.0, fit_intercept=True, l1_ratio=0.15,
learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
n_jobs=1, penalty='l2', power_t=0.5, random_state=None,
shuffle=True, tol=None, verbose=0, warm_start=False)
```



#### <font color = blue>[4.1.4]Train and Test Accuracy</font>



```
In [115]: Training_Accuracy_Bow = calibrated_clf.score(X_Train_Bow, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy_Bow)
Training_Error_Bow = 1 - Training_Accuracy_Bow
print('Training_Error=%0.3f'%Training_Error_Bow)

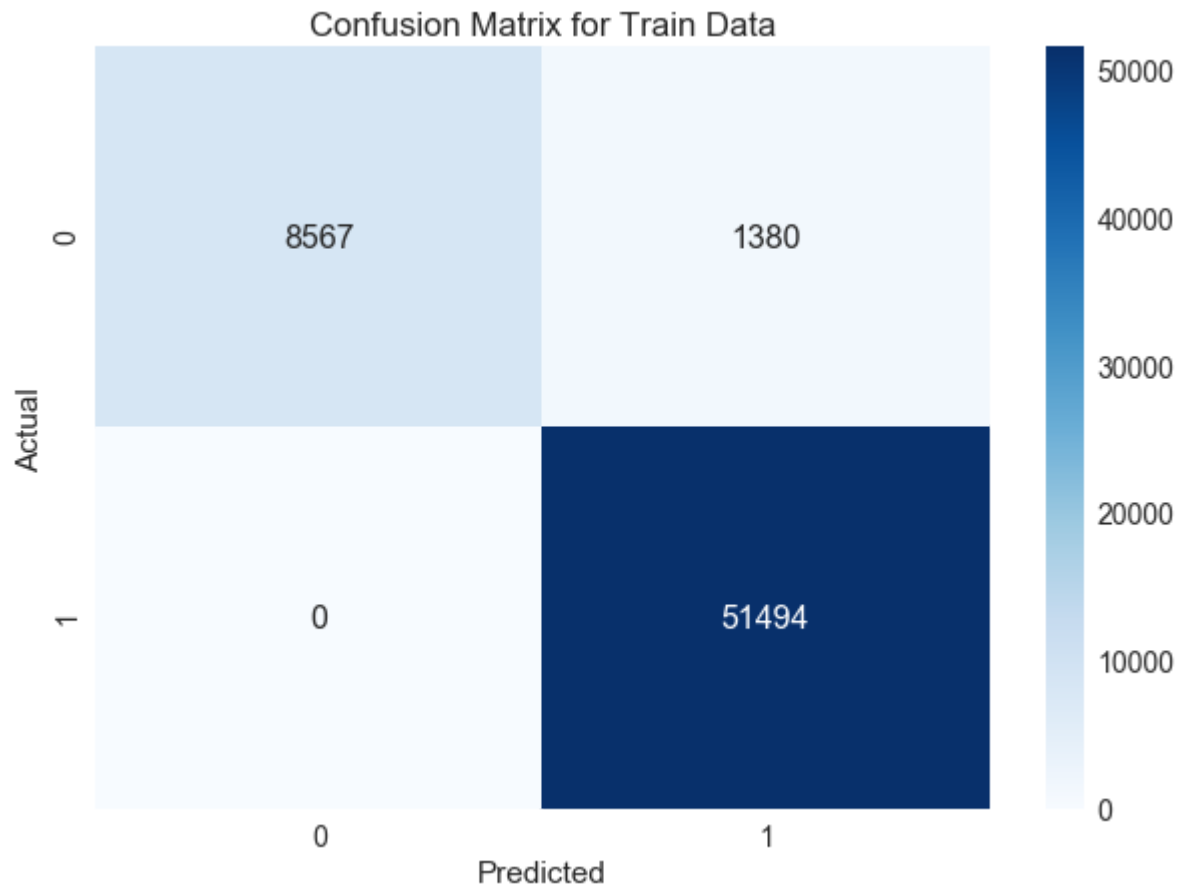
Test_Accuracy_Bow = accuracy_score(Y_Test, prediction.round())
print('Test_Accuracy=%0.3f'%Test_Accuracy_Bow)
Test_Error_Bow = 1 - Test_Accuracy_Bow
print('Test_Error=%0.3f'%Test_Error_Bow)
#print('\nThe accuracy of the MNB classifier for k = %d is %f%%' % (optimal_alpha_bow, Test_Accuracy_Bow))

Training_Accuracy=0.978
Training_Error=0.022
Test_Accuracy=0.873
Test_Error=0.127
```

## <font color = blue>[4.1.5] Confusion Matrix </font>

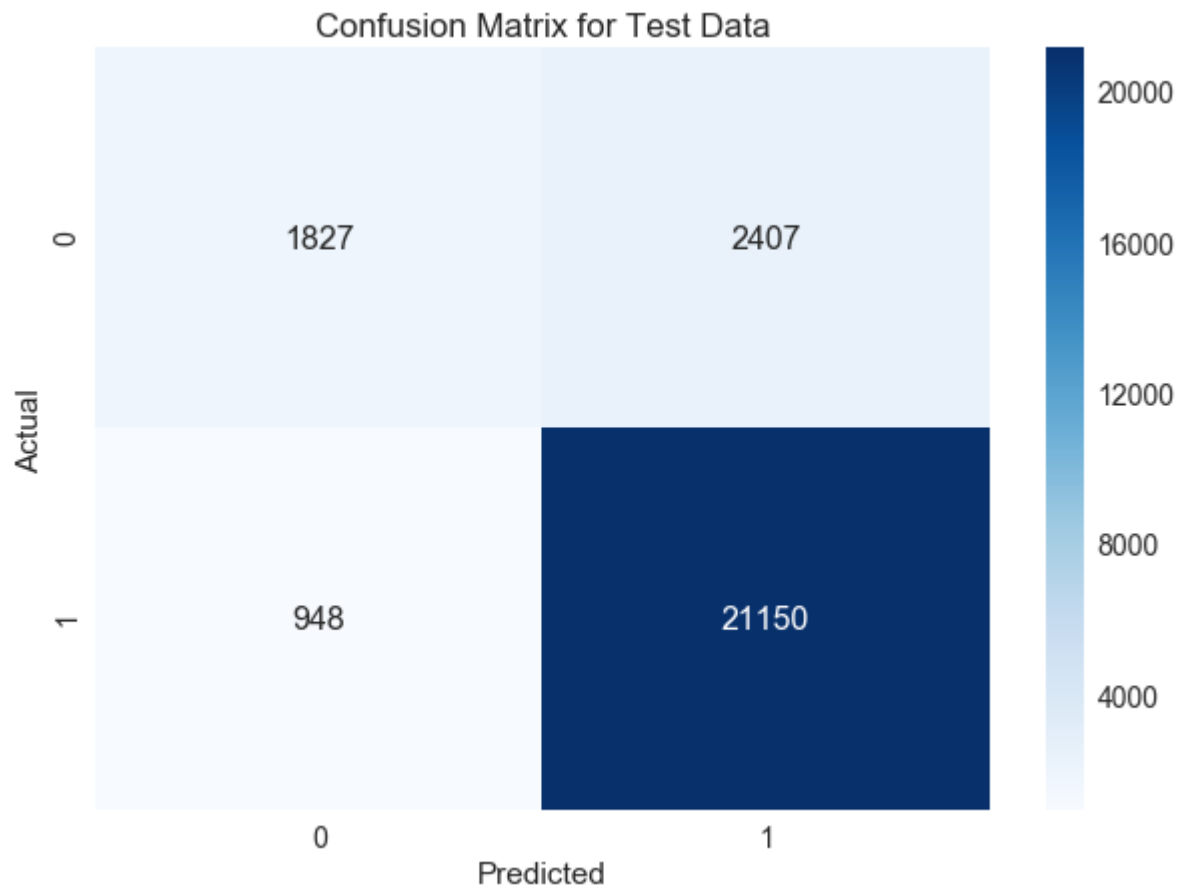
```
In [116]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Train, calibrated_clf.predict(X_Train_Bow))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Train), index=np.unique(Y_Train))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Train Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[116]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d131b2c50>



```
In [117]: #With the reference of below link:
#https://www.kaggle.com/agungor2/various-confusion-matrix-plots
from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, calibrated_clf.predict(X_Test_Bow))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Test Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[117]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d493c18d0>



## <font color = blue>[4.1.6] Classification Report</font>

```
In [118]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction.round()))
```

	precision	recall	f1-score	support
0	0.66	0.43	0.52	4234
1	0.90	0.96	0.93	22098
avg / total	0.86	0.87	0.86	26332

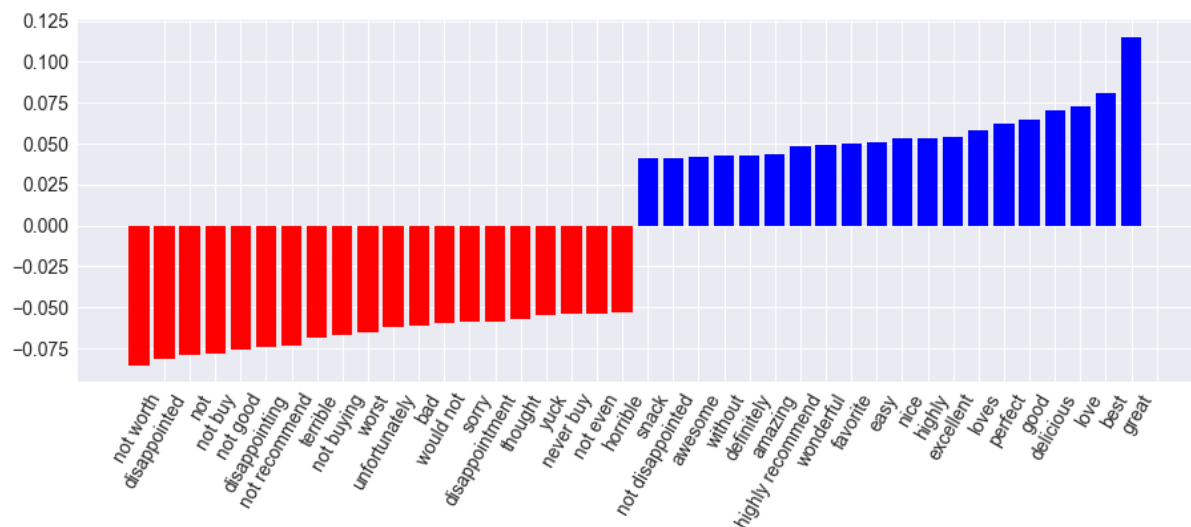
## <font color = Blue>[4.1.7] Feature Importance </font>

```
In [119]: #https://medium.com/@aneesha/visualising-top-features-in-linear-svm-with-scikit-learn-and-matplotlib-3454ab18a14d
def plot_coefficients(classifier, feature_names, top_features=20):
    coef = classifier.coef_.ravel()
    top_positive_coefficients = np.argsort(coef)[-top_features:]
    top_negative_coefficients = np.argsort(coef)[:top_features]
    top_coefficients = np.hstack([top_negative_coefficients, top_positive_coefficients])

    # create plot
    plt.figure(figsize=(15, 5))
    colors = ['red' if c < 0 else 'blue' for c in coef[top_coefficients]]
    plt.bar(np.arange(2 * top_features), coef[top_coefficients], color=colors)
    feature_names = np.array(feature_names)
    plt.xticks(np.arange(1, 1 + 2 * top_features), feature_names[top_coefficients], rotation=60, ha='right')
    plt.show()
```

## <font color = blue>Feature Importance for Positive and Negative Class</font>

```
In [121]: clf = SGDClassifier(alpha=0.1, class_weight='balanced')
clf.fit(X_Train_Bow, Y_Train)
plot_coefficients(clf, count_vect.get_feature_names())
```



## [4.2] TF-IDF

```
In [122]: #TF-IDF
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=5)
tf_idf_vect.fit_transform(X_Train)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names()[0:10])
print('='*50)

X_Train_TfIdf = tf_idf_vect.transform(X_Train)
X_Test_TfIdf = tf_idf_vect.transform(X_Test)
X_CV_TfIdf = tf_idf_vect.transform(X_cv)

#final_tf_idf = tf_idf_vect.transform(X_Test)
print("the type of count vectorizer ",type(X_Train_TfIdf))
print("the shape of out text TFIDF vectorizer ",X_Train_TfIdf.get_shape())
print("the shape of out text TFIDF vectorizer ",X_Test_TfIdf.get_shape())
print("the shape of out text TFIDF vectorizer ",X_CV_TfIdf.get_shape())
#print("the number of unique words including both unigrams and bigrams ", final_tf_idf.get_shape()[1])

some sample features(unique words in the corpus) ['aa', 'aaa', 'aafco', 'aback', 'abandon', 'abandoned', 'abdominal', 'abilities', 'ability', 'ability make']
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (61441, 80521)
the shape of out text TFIDF vectorizer (26332, 80521)
the shape of out text TFIDF vectorizer (26332, 80521)
```

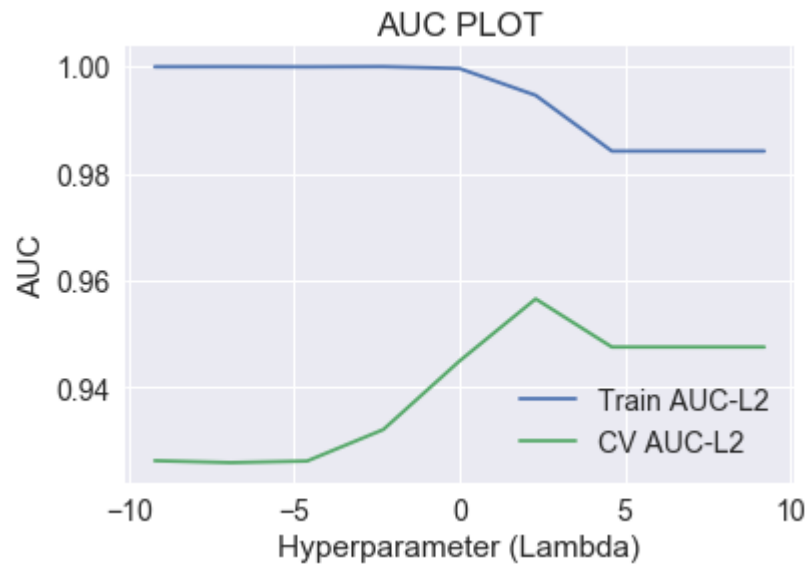
```
In [123]: scalar = StandardScaler(with_mean=False)
X_Train_TfIdf = scalar.fit_transform(X_Train_TfIdf)
X_Test_TfIdf = scalar.transform(X_Test_TfIdf)
X_CV_TfIdf = scalar.transform(X_CV_TfIdf)

print("the type of count vectorizer ",type(X_Train_TfIdf))
print("the shape of out text TFIDF vectorizer ",X_Train_TfIdf.get_shape())
print("the shape of out text TFIDF vectorizer ",X_Test_TfIdf.get_shape())
print("the shape of out text TFIDF vectorizer ",X_CV_TfIdf.get_shape())

the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (61441, 80521)
the shape of out text TFIDF vectorizer (26332, 80521)
the shape of out text TFIDF vectorizer (26332, 80521)
```

```
In [125]: Optimal_Lamda_L2(X_Train_TfIdf, Y_Train, X_CV_TfIdf,Y_cv)
```

```
100%|███████████| 9/9 [00:13<00:00, 1.37s/it]
```



```
CV AUS Scores with Penalty=? Cv auc scores with penalty L2
[0.9262989535214652, 0.9259391529115398, 0.9262368669677858, 0.93211943165697
29, 0.9449295154623841, 0.9565408553036794, 0.9475433070011101, 0.94754360626
49256, 0.9475434031930507]
Maximun AUC value : 0.9565408553036794
Index 5
```

### **<font color = blue>[4.2.3] ROC Curve of SVM</font>**

```

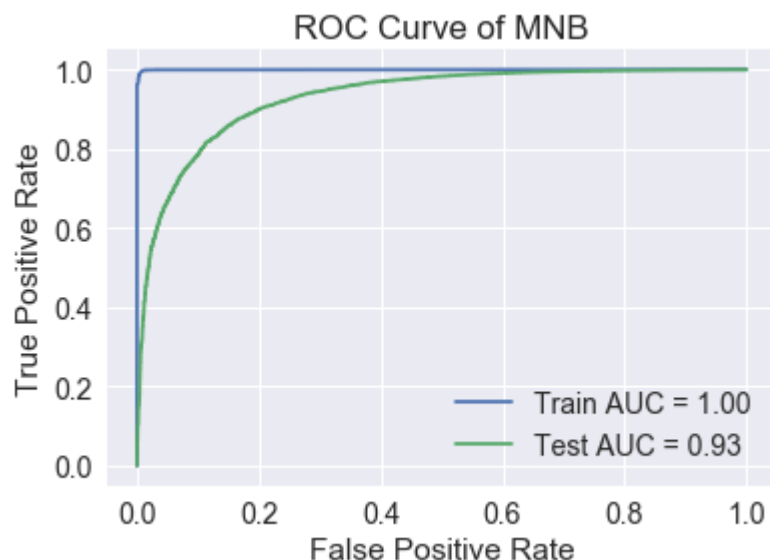
In [126]: #Testing with test data
clf = SGDClassifier(alpha=0.1, class_weight='balanced', loss='hinge')
calibrated_clf = CalibratedClassifierCV(clf, cv=3)
calibrated_clf.fit(X_Train_TfIdf, Y_Train)
prediction = calibrated_clf.predict_proba(X_Test_TfIdf)[: ,1]
print(prediction)
print(clf)

Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, calibrated_clf.predict_proba(X_Train_TfIdf)[: ,1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, calibrated_clf.predict_proba(X_Test_TfIdf)[: ,1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of MNB')
plt.show()

[0.92766173 0.99322665 0.13431146 ... 0.25847924 0.23463799 0.35375614]
SGDClassifier(alpha=0.1, average=False, class_weight='balanced', epsilon=0.1,
eta0=0.0, fit_intercept=True, l1_ratio=0.15,
learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
n_jobs=1, penalty='l2', power_t=0.5, random_state=None,
shuffle=True, tol=None, verbose=0, warm_start=False)

```



#### <font color = blue>[4.2.4]Train and Test Accuracy</font>



```
In [129]: Training_Accuracy_Tfidf = calibrated_clf.score(X_Train_Tfidf, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy_Tfidf)
Training_Error_Tfidf = 1 - Training_Accuracy_Tfidf
print('Training_Error=%0.3f'%Training_Error_Tfidf)

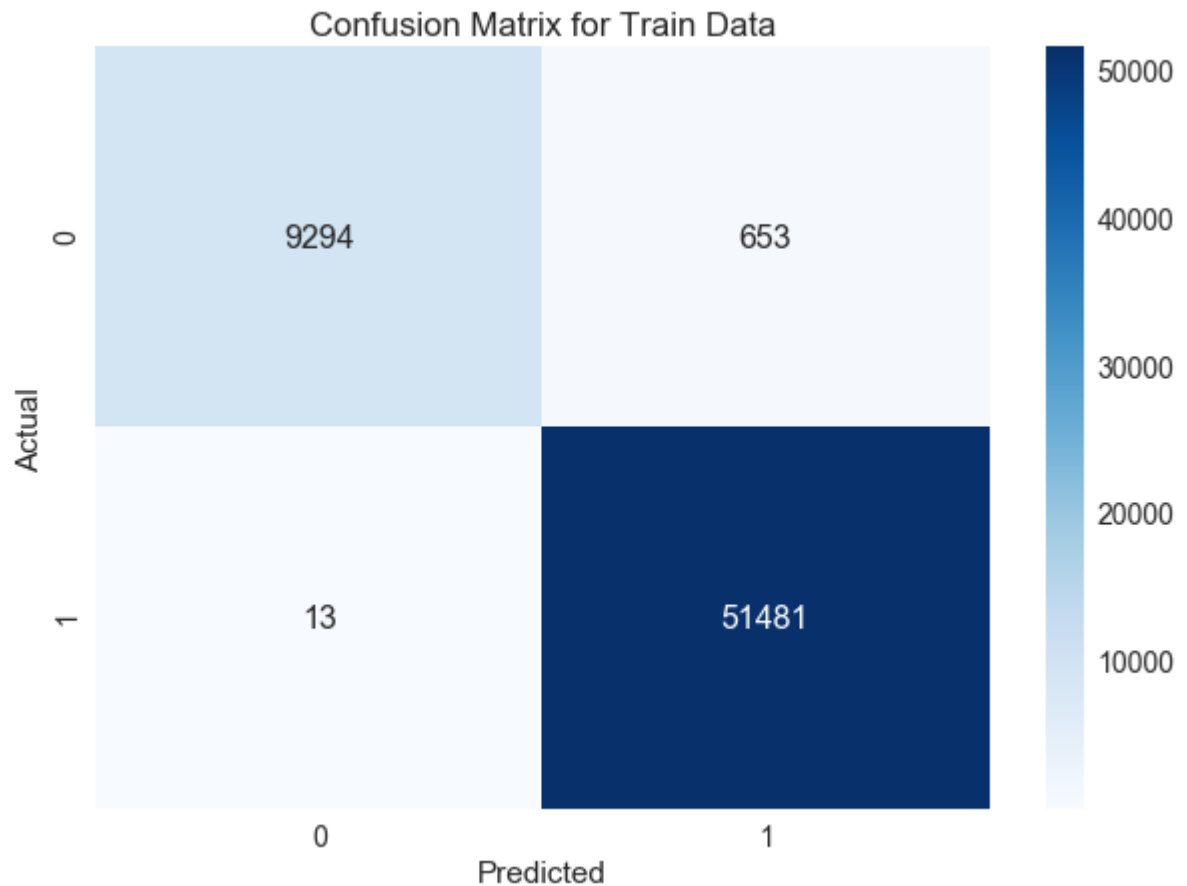
Test_Accuracy_Tfidf = accuracy_score(Y_Test, prediction.round())
print('Test_Accuracy=%0.3f'%Test_Accuracy_Tfidf)
Test_Error_Tfidf = 1 - Test_Accuracy_Tfidf
print('Test_Error=%0.3f'%Test_Error_Tfidf)
#print('\nThe accuracy of the MNB classifier for k = %d is %f%%' % (optimal_alpha_bow, Test_Accuracy_Bow))

Training_Accuracy=0.989
Training_Error=0.011
Test_Accuracy=0.908
Test_Error=0.092
```

## <font color = blue>[4.2.5] Confusion Matrix </font>

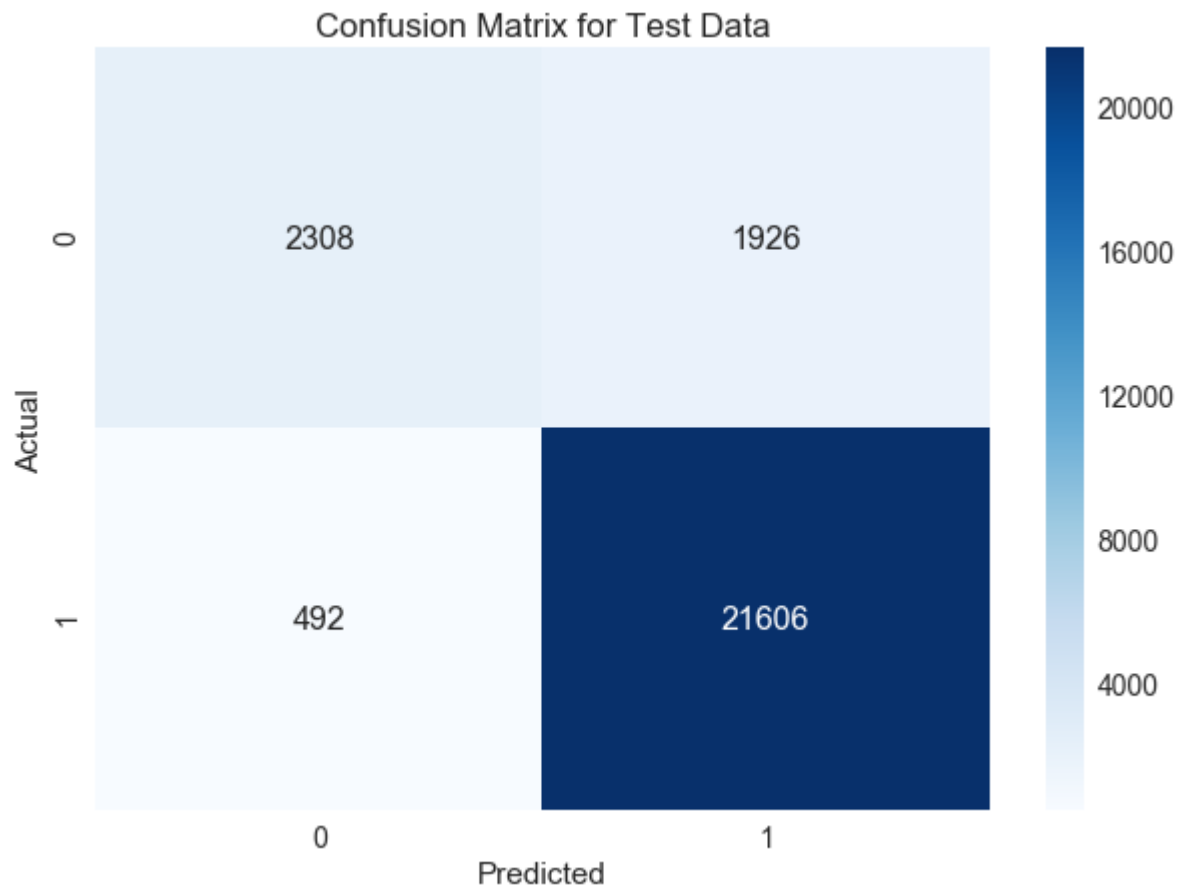
```
In [130]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Train, calibrated_clf.predict(X_Train_Tfidf))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Train), index=np.unique(Y_Train))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Train Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[130]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d48f79550>



```
In [131]: #With the reference of below link:
#https://www.kaggle.com/agungor2/various-confusion-matrix-plots
from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, calibrated_clf.predict(X_Test_Tfidf))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Test Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[131]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d37c3ee80>



## <font color = blue>[4.2.6] Classification Report</font>

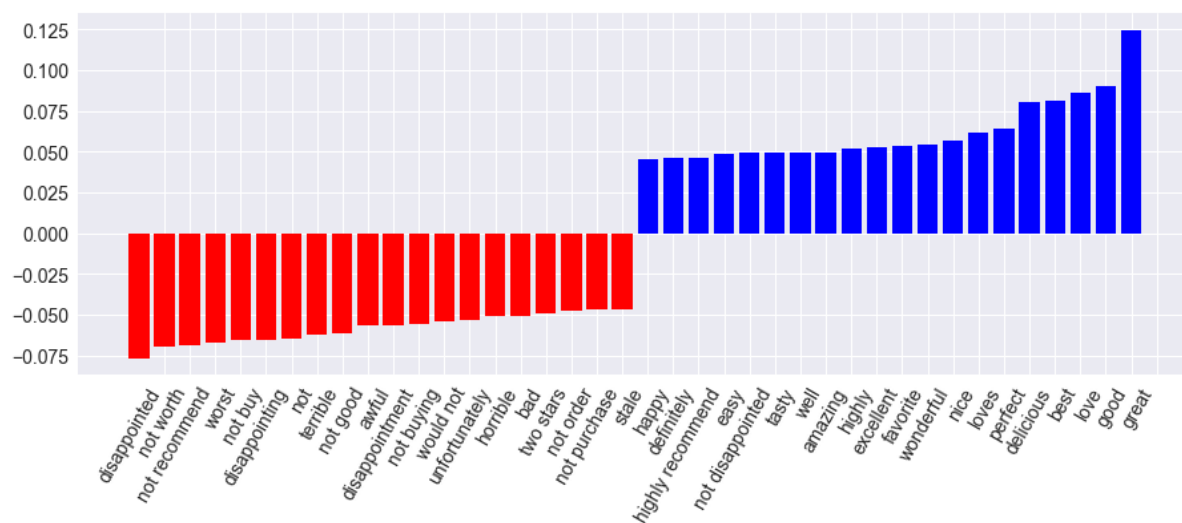
```
In [132]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction.round()))
```

	precision	recall	f1-score	support
0	0.82	0.55	0.66	4234
1	0.92	0.98	0.95	22098
avg / total	0.90	0.91	0.90	26332

## [4.2.7] Feature Importance

### Feature Importance for Positive and Negative Class

```
In [133]: clf = SGDClassifier(alpha=0.1, class_weight='balanced')
clf.fit(X_Train_TfIdf, Y_Train)
plot_coefficients(clf, tf_idf_vect.get_feature_names())
```



## [4.3] Word2Vec

```
In [134]: i=0
list_of_sentence_train=[]
for sentence in X_Train:
    list_of_sentence_train.append(sentence.split())

w2v_model=Word2Vec(list_of_sentence_train,min_count=5,size=50, workers=4)
w2v_words = list(w2v_model.wv.vocab)
print("number of words that occurred minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])

number of words that occurred minimum 5 times 14786
sample words ['weekend', 'week', 'long', 'fast', 'using', 'rice', 'green',
'tea', 'works', 'wonders', 'one', 'energy', 'level', 'tasty', 'even', 'bit',
'salt', 'makes', 'much', 'pleasant', 'family', 'favorite', 'flavor', 'hans
n', 'diet', 'sodas', 'clean', 'crisp', 'taste', 'enjoyable', 'meals', 'calm
s', 'upset', 'tummy', 'love', 'compared', 'ones', 'used', 'eat', 'like', 'nis
sin', 'maruchan', 'really', 'tell', 'difference', 'big', 'tub', 'spice', 'dro
ps', 'better']
```

### <font color = blue>[4.3.1] Computing avg w2v for train, test, and CV</font>



```
In [136]: %%time
i=0
list_of_sentence_cv=[]
for sentence in X_cv:
    list_of_sentence_cv.append(sentence.split())
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentence_cv): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might
    need to change this to 300 if you use google's w2v
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_cv.append(sent_vec)
sent_vectors_cv = np.array(sent_vectors_cv)
print(sent_vectors_cv.shape)
print(sent_vectors_cv[0])
```

```
100% |██████████████████████████████████████████████████████████████████████████|
██████████ 26332/26332 [01:48<00:00, 243.27it/s]
```

```
(26332, 50)
[-0.11573441  0.46256466 -1.50561294 -0.50534047  0.67303949 -0.03247621
  0.9521336   -1.07461442  0.59043062 -0.51276795  0.16654758  0.27134059
  0.38075584 -1.17755664 -0.36282212  0.80457097  0.41690589  0.79902533
 -0.41350211 -0.65742333  1.34620721 -0.2116666  -0.1831447  -0.21734889
 -0.09684955 -1.16451201  0.29688882 -0.48283021  1.6556564  -0.6978559
  1.13678383 -0.2415488  0.22276458 -0.93023639  0.3729921  0.26260851
 -0.32164465  0.67837226  0.28677943 -0.04293923  0.17995327  0.27606643
 -0.53558971  0.3690836  -0.25666084  1.2172374  0.93561432 -0.05929136
 -0.88247013 -0.71623034]
Wall time: 1min 48s
```

```

In [137]: %%time

i=0
list_of_santance_test=[]
for sentence in X_Test:
    list_of_santance_test.append(sentence.split())

# average Word2Vec
# compute average word2vec for each review.
sent_vectors_test = []; # the avg-w2v for each sentence/review is stored in th
is list
for sent in tqdm(list_of_santance_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might
need to change this to 300 if you use google's w2v
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_test.append(sent_vec)
sent_vectors_test = np.array(sent_vectors_test)
print(sent_vectors_test.shape)
print(sent_vectors_test[0])

100%|████████████████████████████████████████████████████████████████████████████████|
████████████████████████████████████████████████████████████████████████████████| 26332/26332 [01:37<00:00, 270.13it/s]

(26332, 50)
[-0.11573441  0.46256466 -1.50561294 -0.50534047  0.67303949 -0.03247621
 0.9521336  -1.07461442  0.59043062 -0.51276795  0.16654758  0.27134059
 0.38075584 -1.17755664 -0.36282212  0.80457097  0.41690589  0.79902533
-0.41350211 -0.65742333  1.34620721 -0.2116666  -0.1831447  -0.21734889
-0.09684955 -1.16451201  0.29688882 -0.48283021  1.6556564  -0.6978559
 1.13678383 -0.2415488  0.22276458 -0.93023639  0.3729921  0.26260851
-0.32164465  0.67837226  0.28677943 -0.04293923  0.17995327  0.27606643
-0.53558971  0.3690836  -0.25666084  1.2172374  0.93561432 -0.05929136
-0.88247013 -0.71623034]
Wall time: 1min 37s

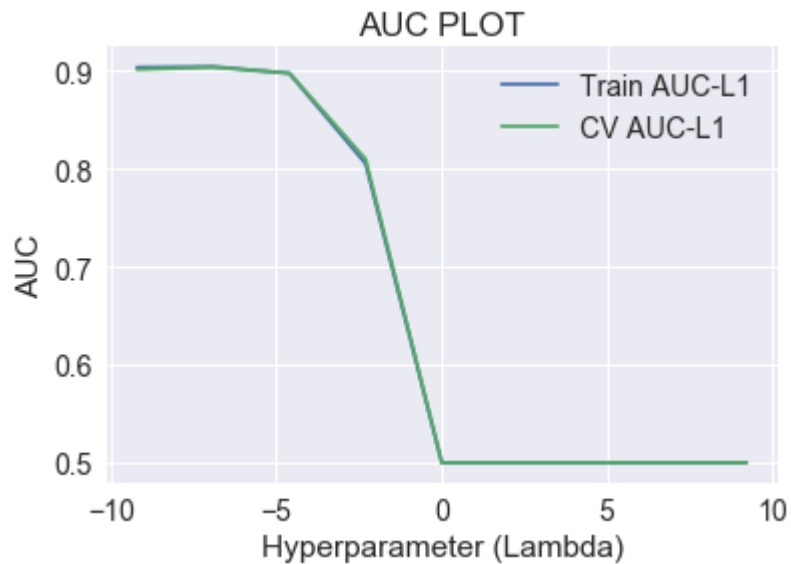
```

### [4.3.2] Hyperparameter tuning with L1 Regularizer and AUC Plot



```
In [138]: Optimal_Lamda_L1(sent_vectors_train, Y_Train, sent_vectors_cv, Y_cv)
```

```
100%|███████████ | 9/9 [00:14<00:00, 1.62s/it]
```

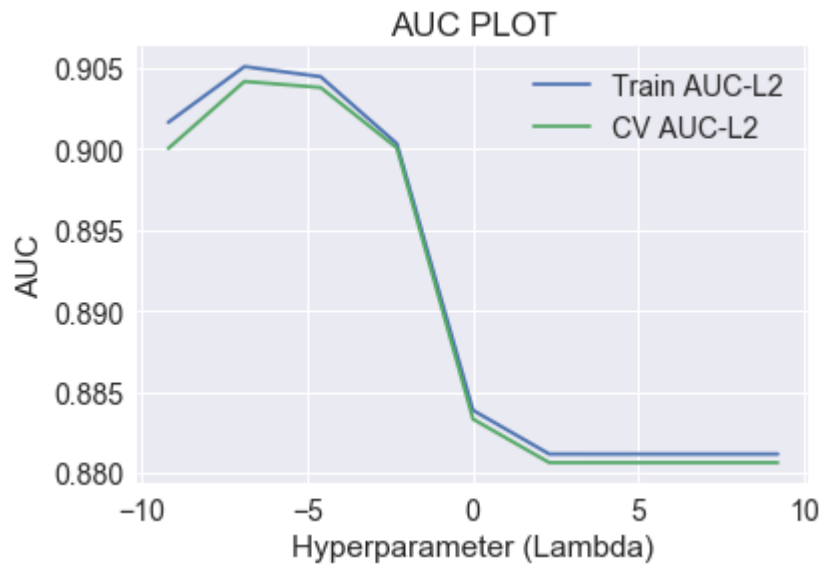


```
CV AUS Scores with Penalty=? Cv auc scores with penalty L1
[0.9016006360296618, 0.9038676342464342, 0.8979925190886493, 0.80996416401315
85, 0.5, 0.5, 0.5, 0.5, 0.5]
Maximun AUC value : 0.9038676342464342
Index 1
```

### [4.3.3] Hyperparameter tuning with L2 Regularizer and AUC Plot

```
In [139]: Optimal_Lamda_L2(sent_vectors_train, Y_Train, sent_vectors_cv, Y_cv)
```

```
100%|███████████| 9/9 [00:08<00:00, 1.08it/s]
```



```
CV AUS Scores with Penalty=? Cv auc scores with penalty L2
[0.9000430961270005, 0.9041814764847258, 0.9038064027322272, 0.90007921085670
98, 0.8833332521045835, 0.880647562434234, 0.8806469318426233, 0.880647081474
531, 0.880647081474531]
Maximun AUC value : 0.9041814764847258
Index 1
```

#### **<font color = blue>[4.3.4] ROC Curve of SVM</font>**

```

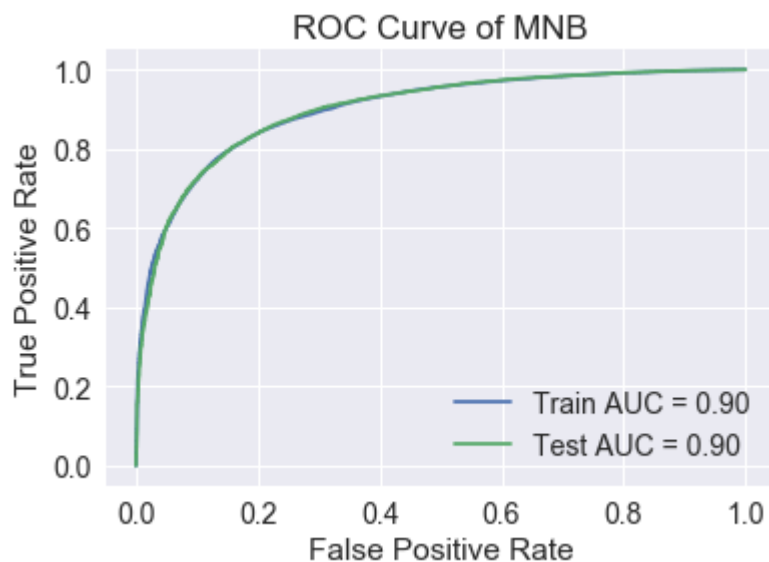
In [140]: #Testing with test data
clf = SGDClassifier(alpha=0.1, class_weight='balanced', loss='hinge')
calibrated_clf = CalibratedClassifierCV(clf, cv=5)
calibrated_clf.fit(sent_vectors_train,Y_Train)
prediction = calibrated_clf.predict_proba(sent_vectors_test)[: ,1]
print(prediction)
print(clf)

Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, calibrated_clf.predict_proba(sent_vectors_train)[: ,1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, calibrated_clf.predict_proba(sent_vectors_test)[: ,1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of MNB')
plt.show()

[0.97856529 0.96440308 0.06950247 ... 0.61533381 0.90356187 0.07763442]
SGDClassifier(alpha=0.1, average=False, class_weight='balanced', epsilon=0.1,
eta0=0.0, fit_intercept=True, l1_ratio=0.15,
learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
n_jobs=1, penalty='l2', power_t=0.5, random_state=None,
shuffle=True, tol=None, verbose=0, warm_start=False)

```



**<font color = blue>[4.3.5]Train and Test Accuracy </font>**

```
In [141]: Training_Accuracy_w2v = calibrated_clf.score(sent_vectors_train, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy_w2v)
Training_Error_w2v = 1 - Training_Accuracy_w2v
print('Training_Error=%0.3f'%Training_Error_w2v)

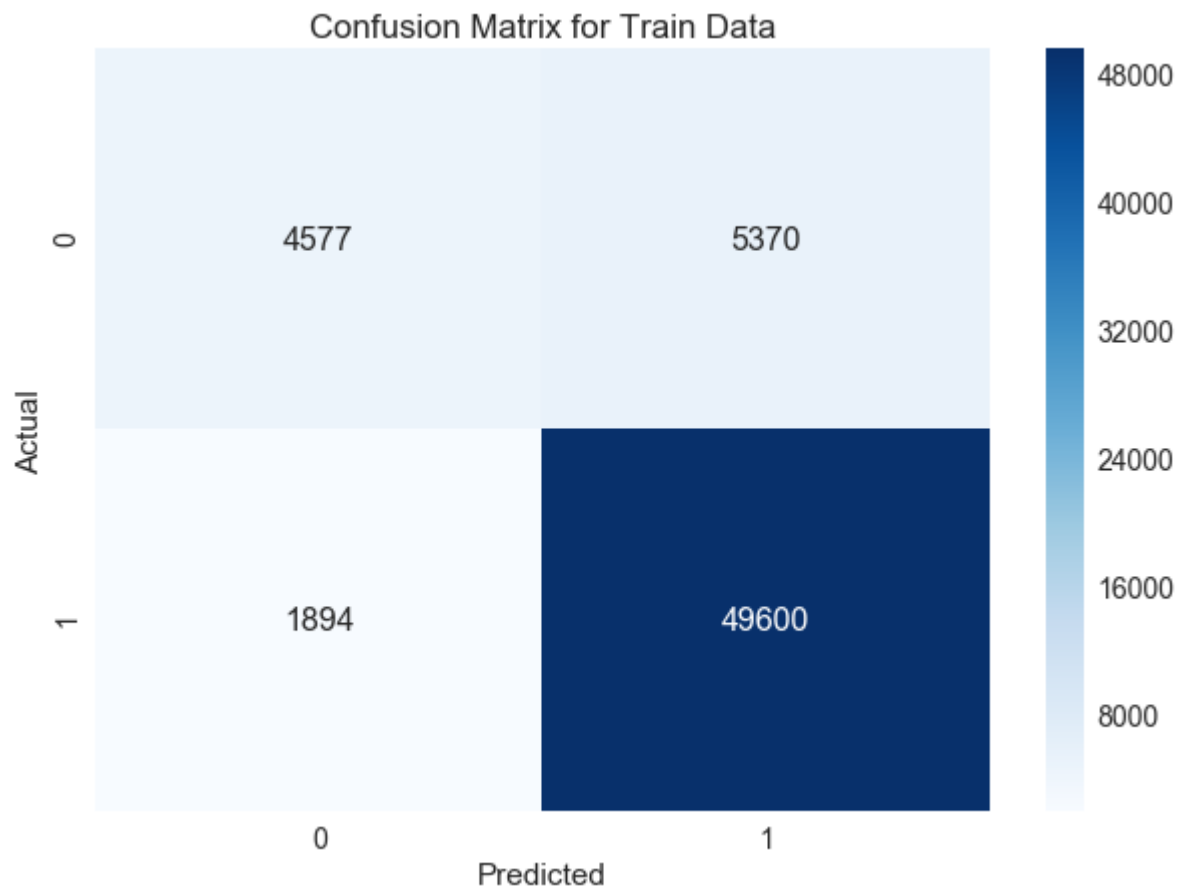
Test_Accuracy_w2v = accuracy_score(Y_Test, prediction.round())
print('Test_Accuracy=%0.3f'%Test_Accuracy_w2v)
Test_Error_w2v = 1 - Test_Accuracy_w2v
print('Test_Error=%0.3f'%Test_Error_w2v)
#print('\nThe accuracy of the MNB classifier for k = %d is %f%%' % (optimal_alpha_bow, Test_Accuracy_Bow)))

Training_Accuracy=0.882
Training_Error=0.118
Test_Accuracy=0.884
Test_Error=0.116
```

### <font color = blue>[4.3.6]Confusion Matrix </font>

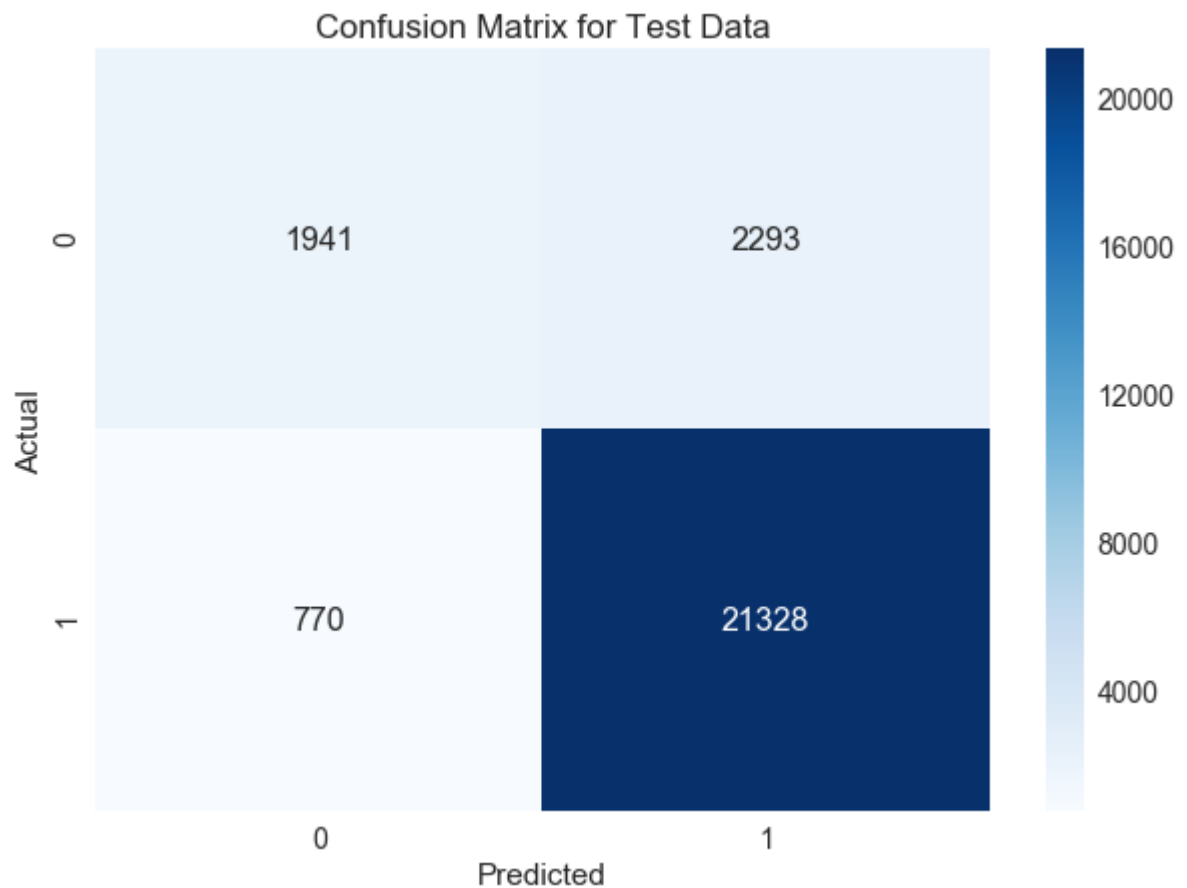
```
In [142]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Train, calibrated_clf.predict(sent_vectors_train))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Train), index=np.unique(Y_Train))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Train Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[142]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d2b08c8d0>



```
In [143]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, calibrated_clf.predict(sent_vectors_test))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Test Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[143]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d18755f60>



### <font color = blue>[4.3.7] Classification Report</font>

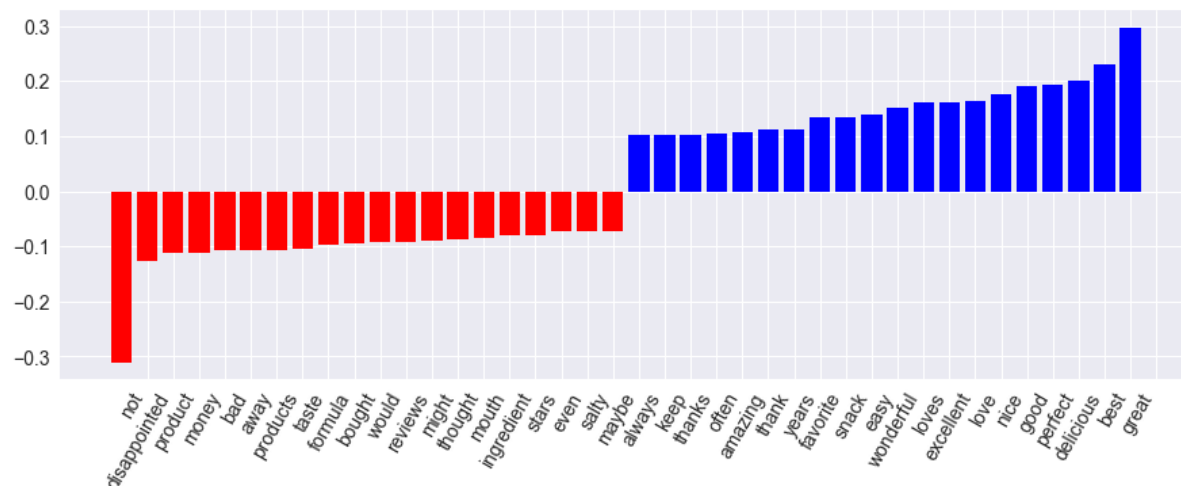
```
In [144]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction.round()))
```

	precision	recall	f1-score	support
0	0.72	0.46	0.56	4234
1	0.90	0.97	0.93	22098
avg / total	0.87	0.88	0.87	26332

## <font color = Blue>[4.2.7] Feature Importance </font>

### <font color = blue>Feature Importance for Positive and Negative Class</font>

```
In [235]: clf = SGDClassifier(alpha=0.1, class_weight='balanced')
          clf.fit(X_Train_Bow, Y_Train)
          plot_coefficients(clf, count_vect.get_feature_names())
```



## <font color = blue> [4.4] TFIDF weighted W2v </font>

```
In [145]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
          model = TfidfVectorizer()
          model.fit(X_Train)
          # we are converting a dictionary with word as a key, and the idf as a value
          dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

### <font color = blue>[4.4.1] Compute TF-IDF weighted Word2Vec for Train, Test, and CV </font>





```
In [147]: i=0
list_of_sentence_test=[]
for sentence in X_Test:
    list_of_sentence_test.append(sentence.split())
# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
= tfidf

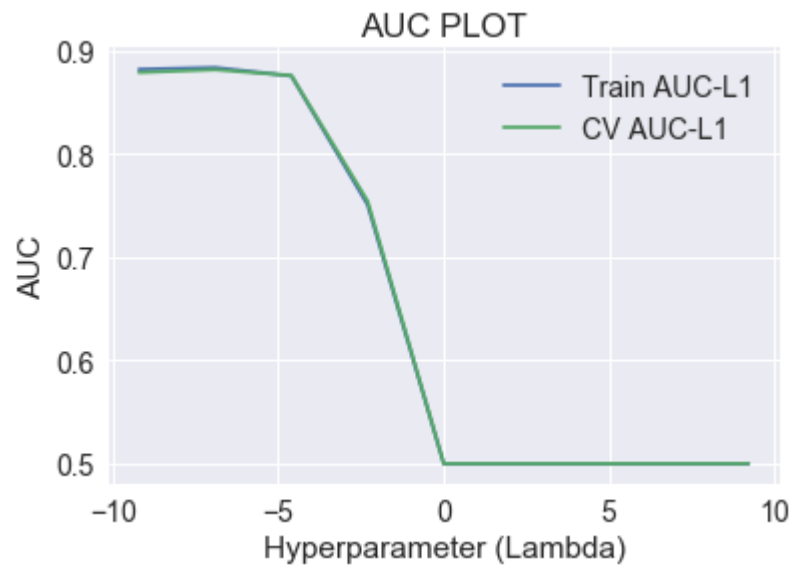
tfidf_sent_vectors_test = []; # the tfidf-w2v for each sentence/review is stor
ed in this list
row=0;
for sent in tqdm(list_of_sentence_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
#
            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors_test.append(sent_vec)
    row += 1
```

```
100%|███████████| 26332/26332 [20:23<00:00, 17.00it/s]
```



```
In [150]: Optimal_Lamda_L1(tfidf_sent_vectors_train, Y_Train, tfidf_sent_vectors_cv, Y_cv
)
```

```
100%|███████████| 9/9 [00:42<00:00, 3.80s/it]
```

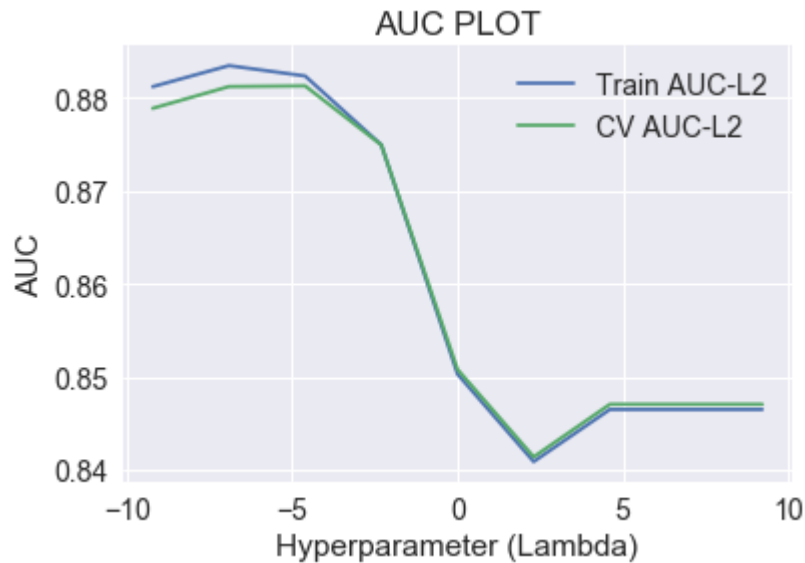


```
CV AUS Scores with Penalty=? Cv auc scores with penalty L1
[0.8787021445629771, 0.8813609004899504, 0.8756963815541822, 0.75417574558266
31, 0.5, 0.5, 0.5, 0.5, 0.5]
Maximun AUC value : 0.8813609004899504
Index 1
```

### [4.4.3] Hyperparameter tuning with L2 Regularizer and AUC Plot

```
In [151]: Optimal_Lamda_L2(tfidf_sent_vectors_train, Y_Train, tfidf_sent_vectors_cv, Y_cv
)
```

```
100%|██████████| 9/9 [00:10<00:00, 1.12s/it]
```



CV AUS Scores with Penalty=? Cv auc scores with penalty L2  
[0.8788924763495014, 0.8812327086970725, 0.8813056756280361, 0.87495876037745  
37, 0.8508129266406487, 0.8414037623361355, 0.8470622959955979, 0.84706205017  
17496, 0.8470620501717496]  
Maximun AUC value : 0.8813056756280361  
Index 2

#### **<font color = blue>[4.4.4] ROC Curve of SVM</font>**

```

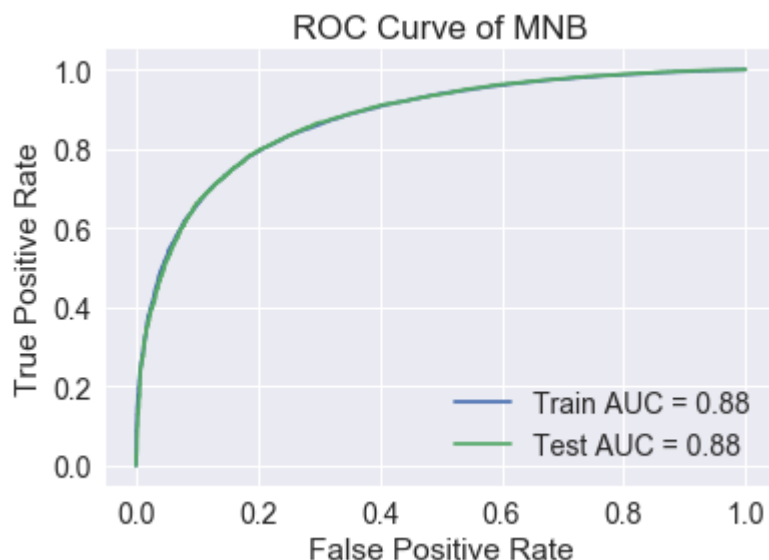
In [158]: #Testing with test data
clf = SGDClassifier(alpha=0.1, class_weight='balanced', loss='hinge')
calibrated_clf = CalibratedClassifierCV(clf, cv=5)
calibrated_clf.fit(tfidf_sent_vectors_train, Y_Train)
prediction = calibrated_clf.predict_proba(tfidf_sent_vectors_test)[: ,1]
print(prediction)
print(clf)

Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, calibrated_clf.predict_proba(tfidf_sent_vectors_train)[: ,1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, calibrated_clf.predict_proba(tfidf_sent_vectors_test)[: ,1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of MNB')
plt.show()

[0.93023877 0.9403652 0.10374825 ... 0.69068035 0.8233184 0.03266469]
SGDClassifier(alpha=0.1, average=False, class_weight='balanced', epsilon=0.1,
eta0=0.0, fit_intercept=True, l1_ratio=0.15,
learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
n_jobs=1, penalty='l2', power_t=0.5, random_state=None,
shuffle=True, tol=None, verbose=0, warm_start=False)

```



#### <font color = blue>[4.4.5]Train and Test Accuracy </font>

```
In [159]: Training_Accuracy_tfidf2v = calibrated_clf.score(tfidf_sent_vectors_train, Y_
Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy_tfidf2v)
Training_Error_tfidf2v = 1 - Training_Accuracy_tfidf2v
print('Training_Error=%0.3f'%Training_Error_tfidf2v)

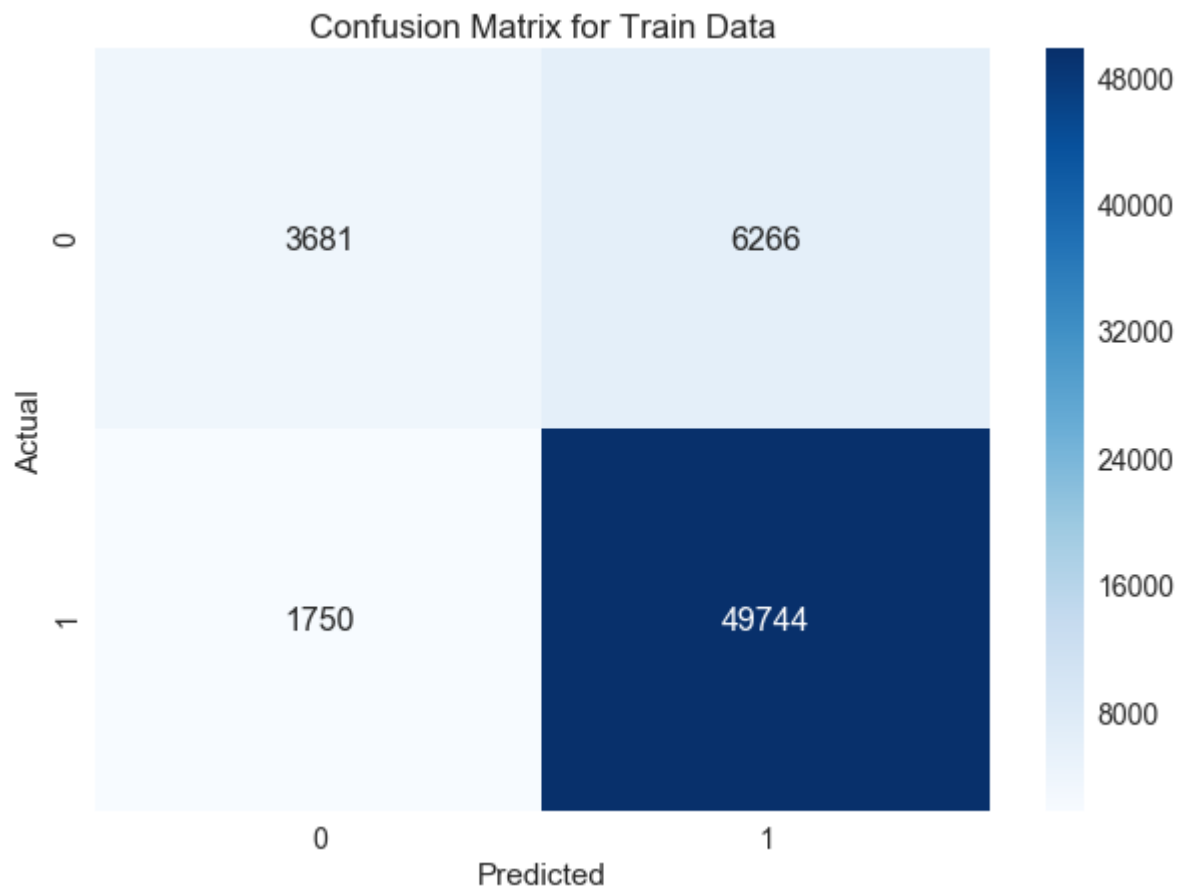
Test_Accuracy_tfidf2v = accuracy_score(Y_Test, prediction.round())
print('Test_Accuracy=%0.3f'%Test_Accuracy_tfidf2v)
Test_Error_tfidf2v = 1 - Test_Accuracy_tfidf2v
print('Test_Error=%0.3f'%Test_Error_tfidf2v)
#print('\nThe accuracy of the MNB classifier for k = %d is %f%%' % (optimal_al
pha_bow, Test_Accuracy_Bow))

Training_Accuracy=0.870
Training_Error=0.130
Test_Accuracy=0.871
Test_Error=0.129
```

#### <font color = blue>[4.4.6]Confusion Matrix </font>

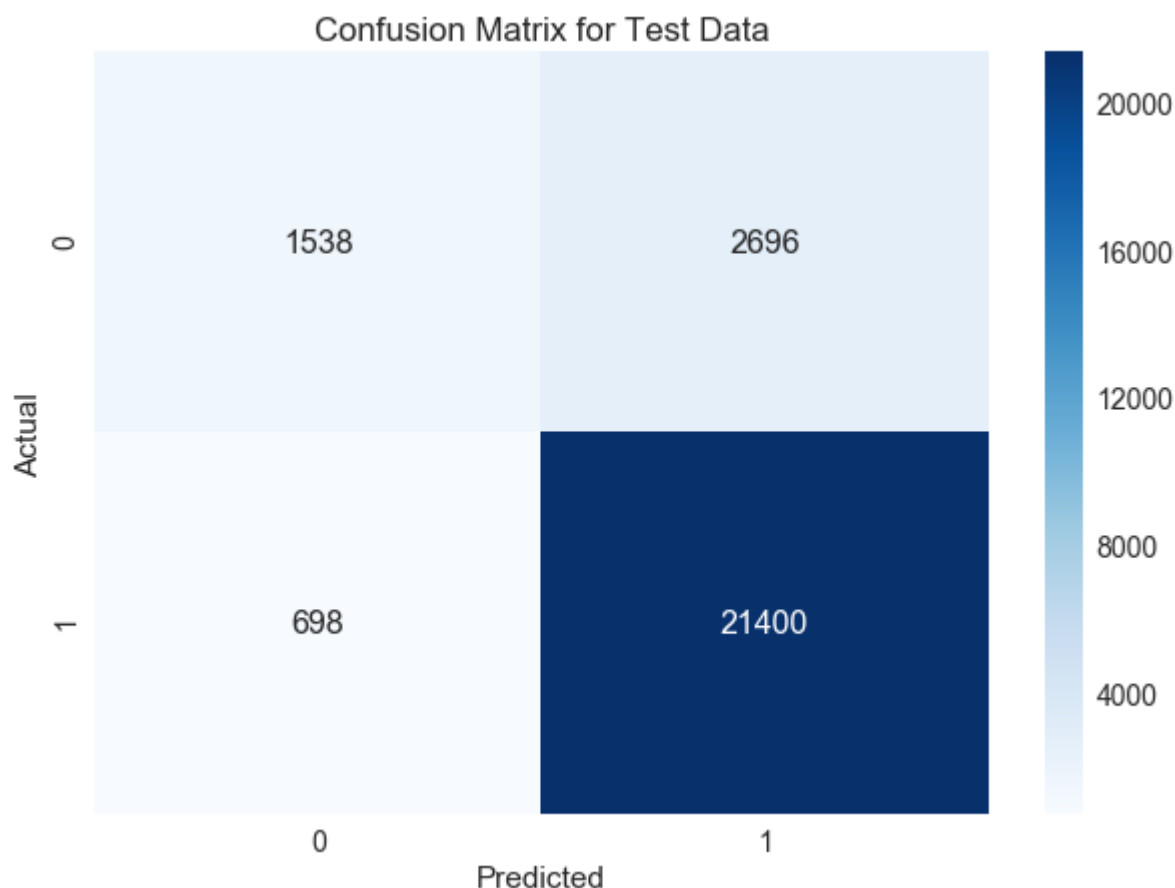
```
In [160]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Train, calibrated_clf.predict(tfidf_sent_vectors_train))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Train), index=np.unique(Y_Train))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Train Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[160]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d56f9ab00>



```
In [161]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, calibrated_clf.predict(tfidf_sent_vectors_test))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Test Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[161]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d3297d470>



### <font color = blue>[4.4.7] Classification Report</font>

```
In [162]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction.round()))
```

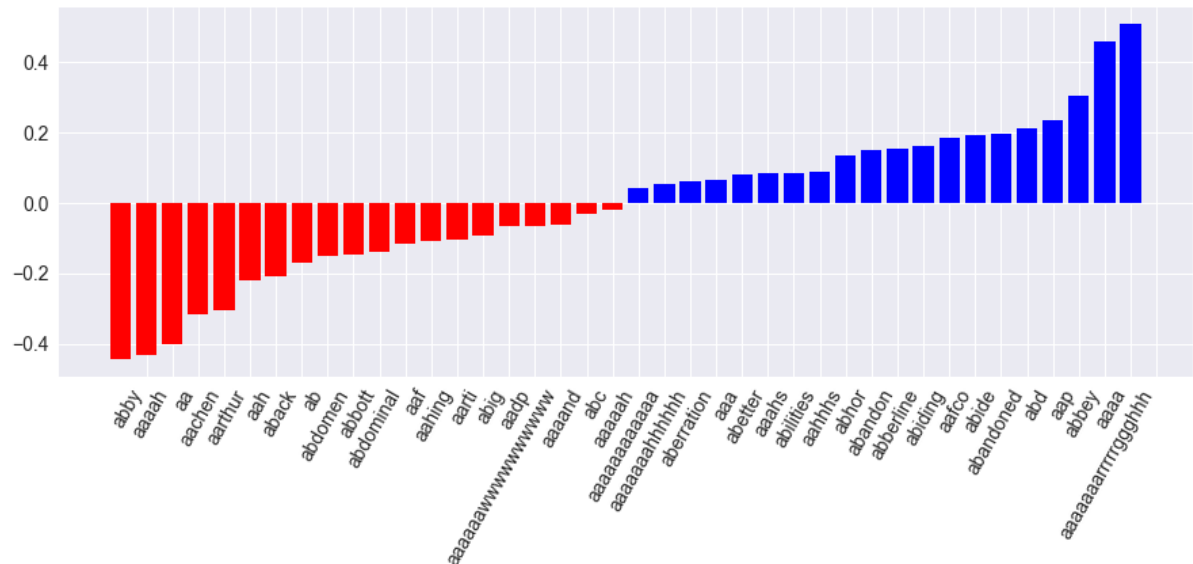
	precision	recall	f1-score	support
0	0.69	0.36	0.48	4234
1	0.89	0.97	0.93	22098
avg / total	0.86	0.87	0.85	26332



## <font color = Blue>[4.4.8] Feature Importance </font>

### <font color = blue>Feature Importance for Positive and Negative Class</font>

```
In [163]: clf = SGDClassifier(alpha=0.1, class_weight='balanced')
          clf.fit(tfidf_sent_vectors_train, Y_Train)
          plot_coefficients(clf, model.get_feature_names())
```



## [5] RBF Kernel

```
In [171]: Y = final['Score'].values
          X = final['cleaned_text'].values

          X=X[:20000]
          Y=Y[:20000]

          print(Y.shape)
          print(type(Y))
          print(X.shape)
          print(type(X))

          (20000,)
          <class 'numpy.ndarray'>
          (20000,)
          <class 'numpy.ndarray'>
```

```
In [173]: # split the data set into train and test
X_Train, X_Test, Y_Train, Y_Test = train_test_split(X,Y,test_size=0.3, random_
state=12, shuffle = False)

# split the train data set into cross validation train and cross validation te
st
X_tr, X_cv, Y_tr, Y_cv = train_test_split(X,Y, test_size=0.3, random_state=12,
shuffle = False)
```

```
In [174]: print('='*100)
print("After splitting")
print("X_Train Shape:",X_Train.shape, "Y_Train Shape:",Y_Train.shape)
print("X_cv Shape:",X_cv.shape, "Y_cv Shape",Y_cv.shape)
print("X_Test Shape",X_Test.shape, "Y_Test Shape",Y_Test.shape)
```

```
=====
=====
After splitting
X_Train Shape: (14000,) Y_Train Shape: (14000,)
X_cv Shape: (6000,) Y_cv Shape (6000,)
X_Test Shape (6000,) Y_Test Shape (6000,)
```

## [5] Featurization

### [5.1] BAG OF WORDS

```

In [238]: import math
def Optimal_Lamda(X_Train,Y_Train,X_CV,Y_CV):
    train_AUC = []
    CV_AUC = []
    tuned_parameters=[10**-4, 10**-3, 10**-2, 10**-1, 1,10**1, 10**2, 10**3, 1
0**4]
    for j in tqdm(tuned_parameters):
        clf = SVC(C=j, probability=True)
        clf.fit(X_Train, Y_Train)
        y_train_pred = clf.predict_proba(X_Train)[:,-1]
        y_cv_pred = clf.predict_proba(X_CV)[:,-1]
        train_AUC.append(roc_auc_score(Y_Train,y_train_pred))
        CV_AUC.append(roc_auc_score(Y_CV, y_cv_pred))

    #Error plots with penalty L1
    plt.plot(np.log(tuned_parameters), train_AUC, label='Train AUC')
    plt.plot(np.log(tuned_parameters), CV_AUC, label='CV AUC')
    plt.legend()
    plt.xlabel("Hyperparameter (Lambda)")
    plt.ylabel("AUC")
    plt.title("AUC PLOT")
    plt.show()

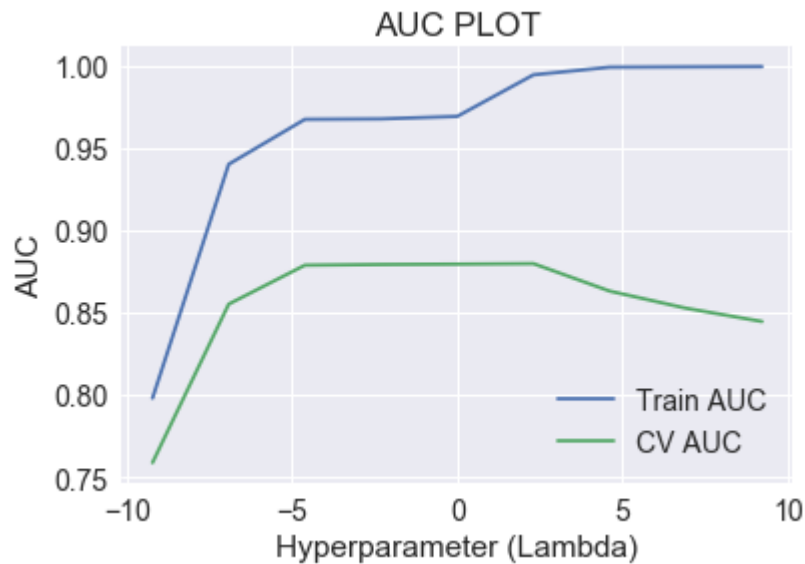
    #Cv auc scores
    print("CV AUS Scores with Penalty=? Cv auc scores")
    print(CV_AUC)
    print("Maximun AUC value :",max(CV_AUC))
    print("Index",CV_AUC.index(max(CV_AUC)))

```

### [5.1.2] Hyperameter tuning and AUC Plot

```
In [239]: Optimal_Lamda(X_Train_Bow, Y_Train, X_CV_Bow,Y_cv)
```

```
100%|██████████████████████████████████████████|  
██████████ | 9/9 [24:16<00:00, 164.61s/it]
```



CV AUS Scores with Penalty=? Cv auc scores

[0.7589056933521543, 0.8550491810624903, 0.878669198539132, 0.8790825883354619, 0.8793186350294638, 0.8796249243799152, 0.8629734937599427, 0.8525816136765925, 0.8445151352199145]

Maximun AUC value : 0.8796249243799152

Index 5

### **<font color = blue>[5.1.4] ROC Curve of SVM </font>**

```

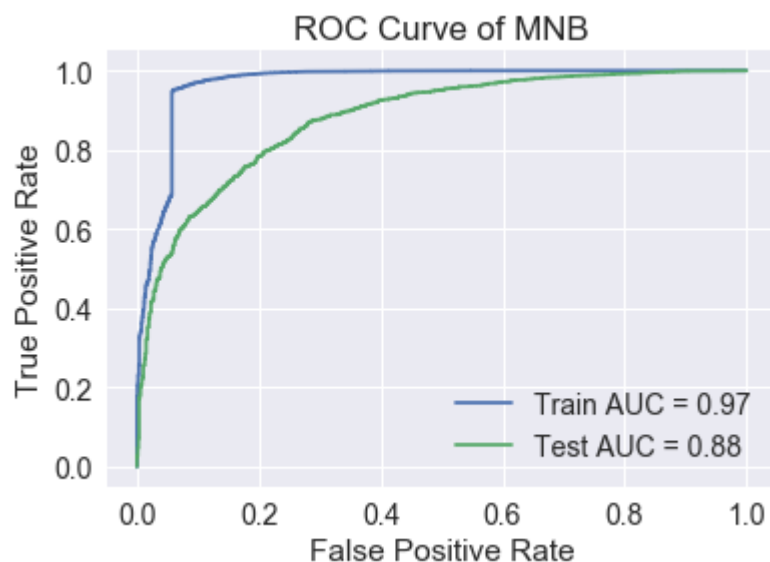
In [179]: #Testing with test data
clf = SVC(C=1, probability=True)
clf.fit(X_Train_Bow,Y_Train)
prediction = clf.predict_proba(X_Test_Bow)[: ,1]
print(prediction)
print(clf)

Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, clf.predict_proba(X_Train_Bow)[: ,1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, clf.predict_proba(X_Test_Bow)[: ,1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of MNB')
plt.show()

[0.10968408 0.84796084 0.99999183 ... 0.97645698 0.89394841 0.97582892]
SVC(C=1, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=True, random_state=None, shrinking=True,
    tol=0.001, verbose=False)

```



### <font color = blue>[5.1.5]Train and Test Accuracy </font>

```
In [180]: Training_Accuracy_Bow = clf.score(X_Train_Bow, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy_Bow)
Training_Error_Bow = 1 - Training_Accuracy_Bow
print('Training_Error=%0.3f'%Training_Error_Bow)

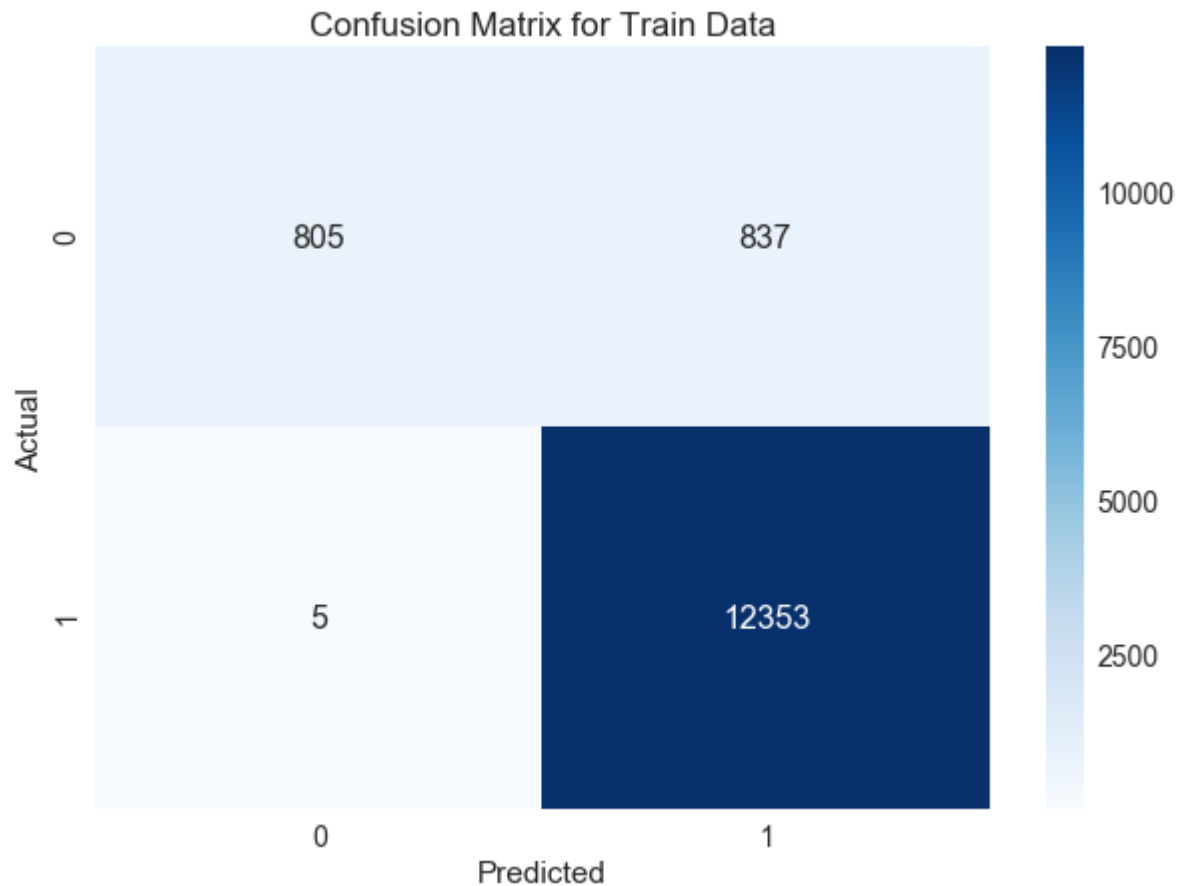
Test_Accuracy_Bow = accuracy_score(Y_Test, prediction.round())
print('Test_Accuracy=%0.3f'%Test_Accuracy_Bow)
Test_Error_Bow = 1 - Test_Accuracy_Bow
print('Test_Error=%0.3f'%Test_Error_Bow)
#print('\nThe accuracy of the MNB classifier for k = %d is %f%%' % (optimal_al
pha_bow, Test_Accuracy_Bow))

Training_Accuracy=0.940
Training_Error=0.060
Test_Accuracy=0.888
Test_Error=0.112
```

## <font color = blue>[5.1.6]Confusion Matrix </font>

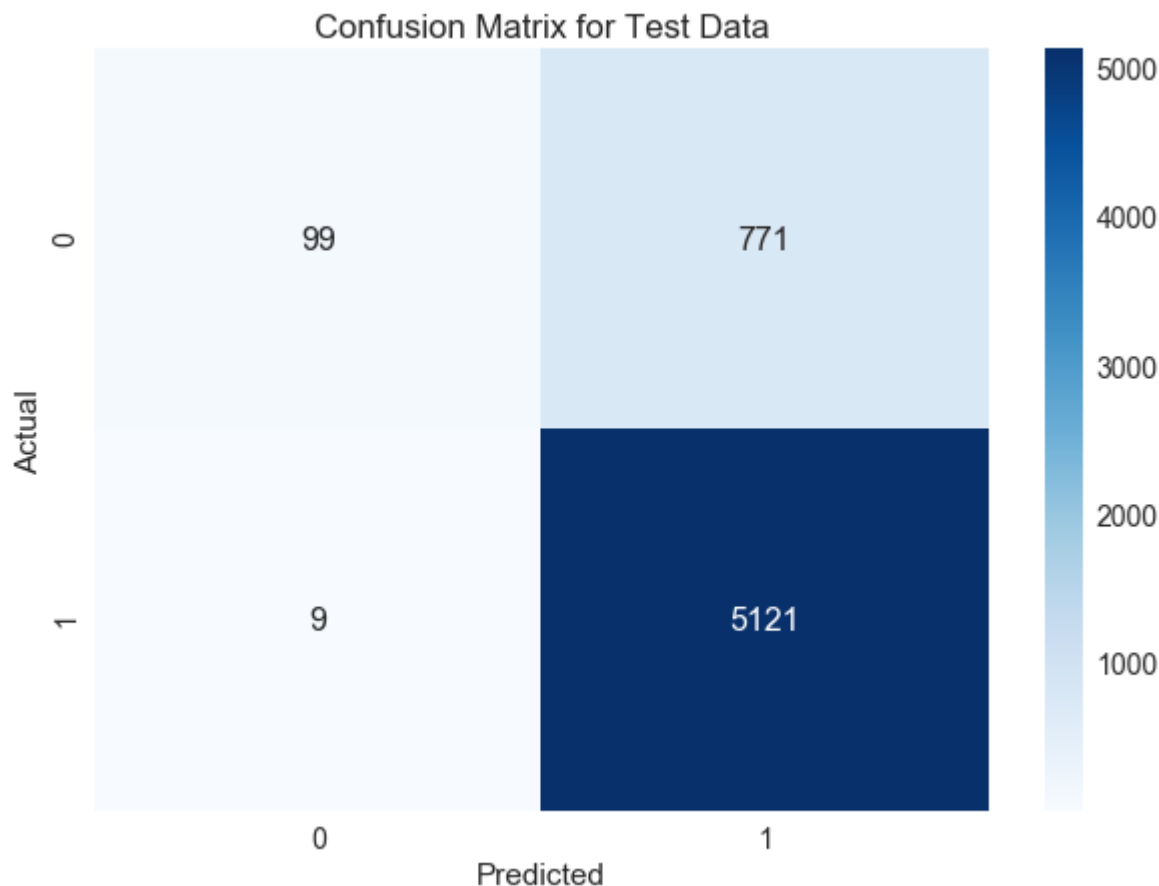
```
In [182]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Train, clf.predict(X_Train_Bow))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Train), index=np.unique(Y_Train))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Train Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[182]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d3280ea90>



```
In [183]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, clf.predict(X_Test_Bow))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Test Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[183]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d18114668>



### <font color = blue>[5.1.7] Classification Report</font>

```
In [184]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction.round()))
```

	precision	recall	f1-score	support
0	0.75	0.33	0.46	870
1	0.90	0.98	0.94	5130
avg / total	0.88	0.89	0.87	6000



## <font color = Blue>[5.1.8] Feature Importance </font>

### <font color = blue>Feature Importance for Positive and Negative Class</font>

```
In [234]: clf = SVC(C=1, class_weight='balanced', kernel = 'linear')
          clf.fit(X_Train_Bow, Y_Train)
          plot_coefficients(clf, count_vect.get_feature_names())
```

```
-----
AttributeError                                Traceback (most recent call last)
<ipython-input-234-145e40f58dcd> in <module>()
      1 clf = SVC(C=1, class_weight='balanced', kernel = 'linear')
      2 clf.fit(X_Train_Bow, Y_Train)
----> 3 plot_coefficients(clf, count_vect.get_feature_names())

<ipython-input-119-7beddd17e766> in plot_coefficients(classifier, feature_names, top_features)
      1 #https://medium.com/@aneesha/visualising-top-features-in-linear-svm-with-scikit-learn-and-matplotlib-3454ab18a14d
      2 def plot_coefficients(classifier, feature_names, top_features=20):
----> 3     coef = classifier.coef_.ravel()
      4     top_positive_coefficients = np.argsort(coef)[-top_features:]
      5     top_negative_coefficients = np.argsort(coef)[:top_features]

D:\Anaconda3\lib\site-packages\scipy\sparse\base.py in __getattr__(self, attr)
    684         return self.getnnz()
    685     else:
--> 686         raise AttributeError(attr + " not found")
    687
    688     def transpose(self, axes=None, copy=False):

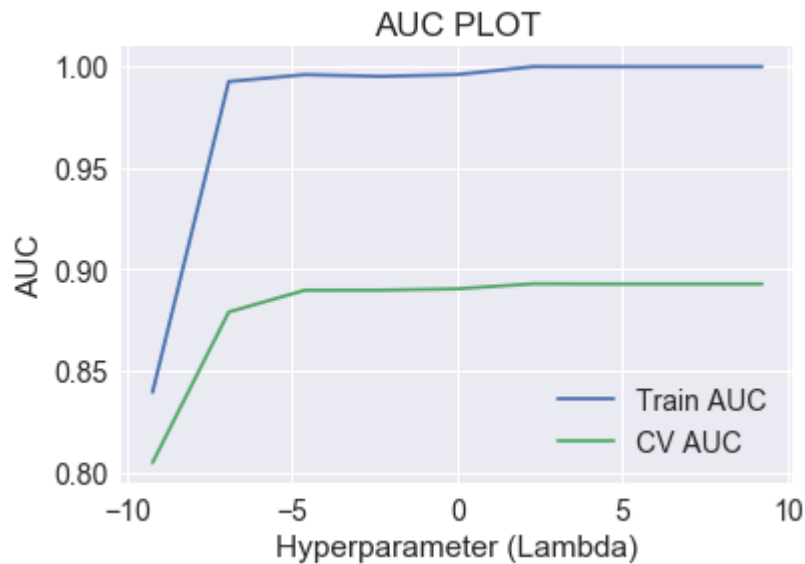
AttributeError: ravel not found
```

## [5.2] TF-IDF

### [5.2.1] Hyperparameter tuning with L1 Regularizer and AUC Plot

```
In [245]: Optimal_Lamda(X_Train_TfIdf, Y_Train, X_CV_TfIdf,Y_cv)
```

```
100%|███████████ | 9/9 [38:03<00:00, 338.87s/it]
```



```
CV AUS Scores with Penalty=? Cv auc scores
[0.8045293629988126, 0.8787419058501938, 0.8894954179830162, 0.88952387354081
24, 0.8901895543456343, 0.8926871457058994, 0.8925605520826331, 0.89256615357
03883, 0.8925666016894086]
Maximun AUC value : 0.8926871457058994
Index 5
```

### **<font color = blue>[5.2.2] ROC Curve of SVM</font>**

```

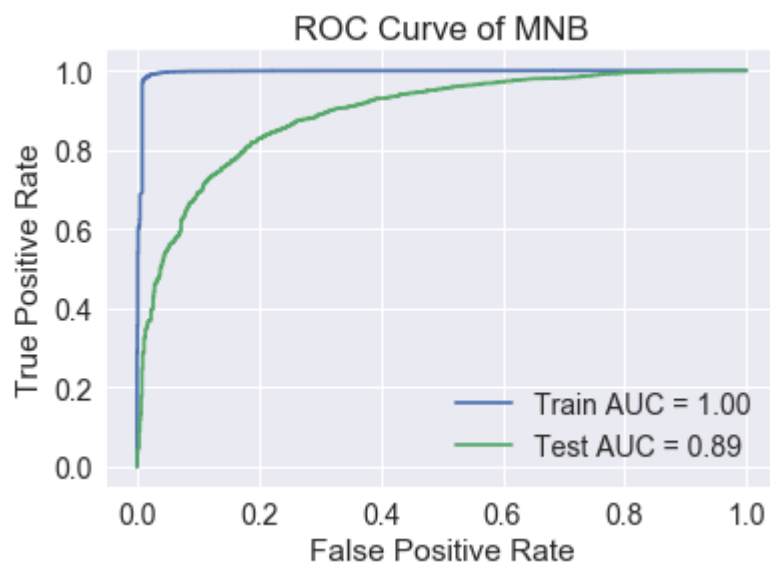
In [193]: #Testing with test data
clf = SVC(C=1, probability=True)
clf.fit(X_Train_TfIdf,Y_Train)
prediction = clf.predict_proba(X_Test_TfIdf)[:,:1]
print(prediction)
print(clf)

Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, clf.predict_proba(X_Train_TfIdf)[:,:1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, clf.predict_proba(X_Test_TfIdf)[:,:1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of MNB')
plt.show()

[0.18960522 0.89269833 0.99677773 ... 0.9488998 0.87401 0.96948229]
SVC(C=1, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=True, random_state=None, shrinking=True,
    tol=0.001, verbose=False)

```



### <font color = blue>[5.2.3]Train and Test Accuracy</font>

```
In [194]: Training_Accuracy_Tfidf = clf.score(X_Train_Tfidf, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy_Tfidf)
Training_Error_Tfidf = 1 - Training_Accuracy_Tfidf
print('Training_Error=%0.3f'%Training_Error_Tfidf)

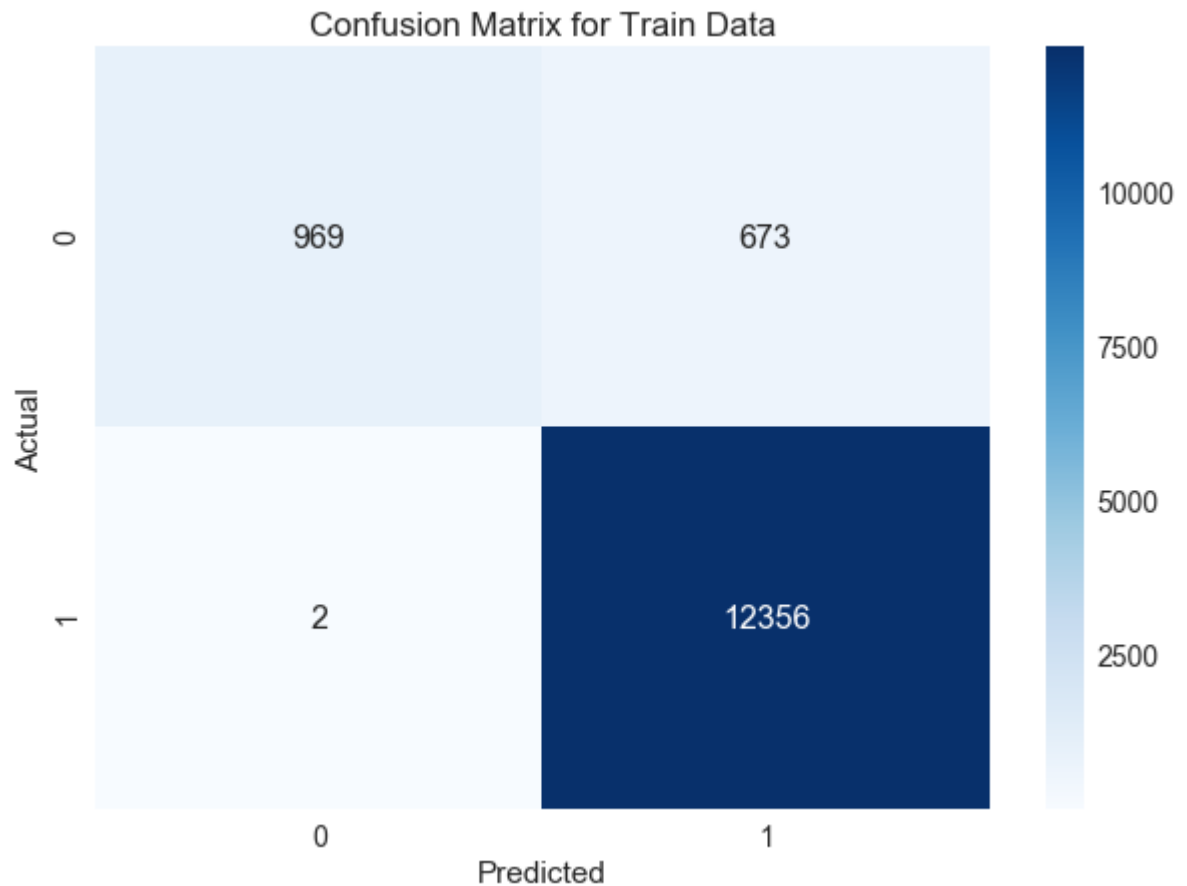
Test_Accuracy_Tfidf = accuracy_score(Y_Test, prediction.round())
print('Test_Accuracy=%0.3f'%Test_Accuracy_Tfidf)
Test_Error_Tfidf = 1 - Test_Accuracy_Tfidf
print('Test_Error=%0.3f'%Test_Error_Tfidf)
#print('\nThe accuracy of the MNB classifier for k = %d is %f%%' % (optimal_alpha_bow, Test_Accuracy_Bow))

Training_Accuracy=0.952
Training_Error=0.048
Test_Accuracy=0.888
Test_Error=0.112
```

## <font color = blue>[5.2.4] Confusion Matrix </font>

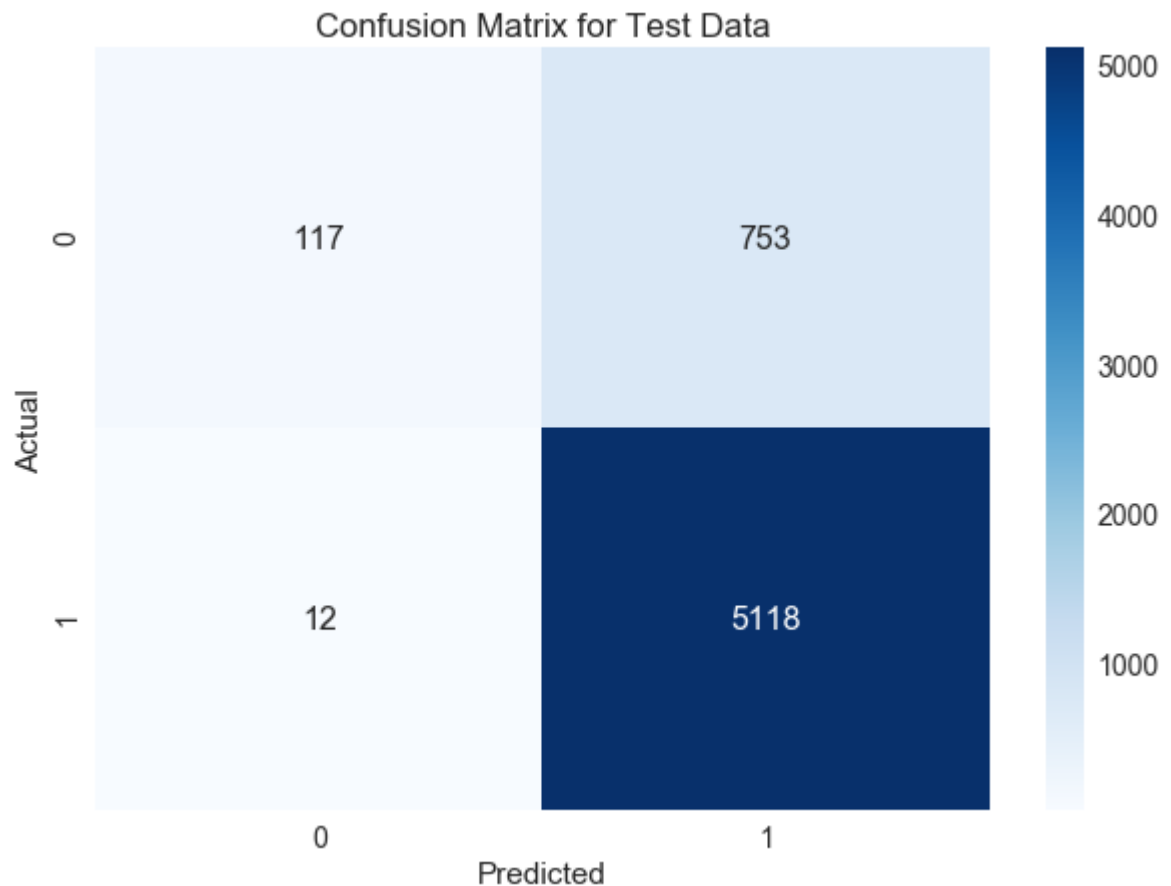
```
In [195]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Train, clf.predict(X_Train_Tfidf))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Train), index=np.unique(Y_Train))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Train Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[195]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d18a22be0>



```
In [196]: #With the reference of below link:
#https://www.kaggle.com/agungor2/various-confusion-matrix-plots
from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, clf.predict(X_Test_Tfidf))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Test Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[196]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d10368978>



### <font color = blue>[5.2.5] Classification Report</font>

```
In [197]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction.round()))
```

	precision	recall	f1-score	support
0	0.73	0.36	0.48	870
1	0.90	0.98	0.94	5130
avg / total	0.88	0.89	0.87	6000

## <font color = Blue>[5.2.6] Feature Importance </font>

### <font color = blue>Feature Importance for Positive and Negative Class</font>

```
In [233]: clf = SVC(C=1, class_weight='balanced', kernel='linear')
clf.fit(X_Train_TfIdf, Y_Train)
plot_coefficients(clf, tf_idf_vect.get_feature_names())
```

```
-----
AttributeError                                Traceback (most recent call last)
<ipython-input-233-2a1dd6b34b71> in <module>()
      1 clf = SVC(C=1, class_weight='balanced', kernel='linear')
      2 clf.fit(X_Train_TfIdf, Y_Train)
----> 3 plot_coefficients(clf, tf_idf_vect.get_feature_names())

<ipython-input-119-7beddd17e766> in plot_coefficients(classifier, feature_names, top_features)
      1 #https://medium.com/@aneesha/visualising-top-features-in-linear-svm-with-scikit-learn-and-matplotlib-3454ab18a14d
      2 def plot_coefficients(classifier, feature_names, top_features=20):
----> 3     coef = classifier.coef_.ravel()
      4     top_positive_coefficients = np.argsort(coef)[-top_features:]
      5     top_negative_coefficients = np.argsort(coef)[:top_features]

D:\Anaconda3\lib\site-packages\scipy\sparse\base.py in __getattr__(self, attr)
    684         return self.getnnz()
    685     else:
--> 686         raise AttributeError(attr + " not found")
    687
    688     def transpose(self, axes=None, copy=False):

AttributeError: ravel not found
```

## <font color = Blue>[5.3] Word2Vec </font>



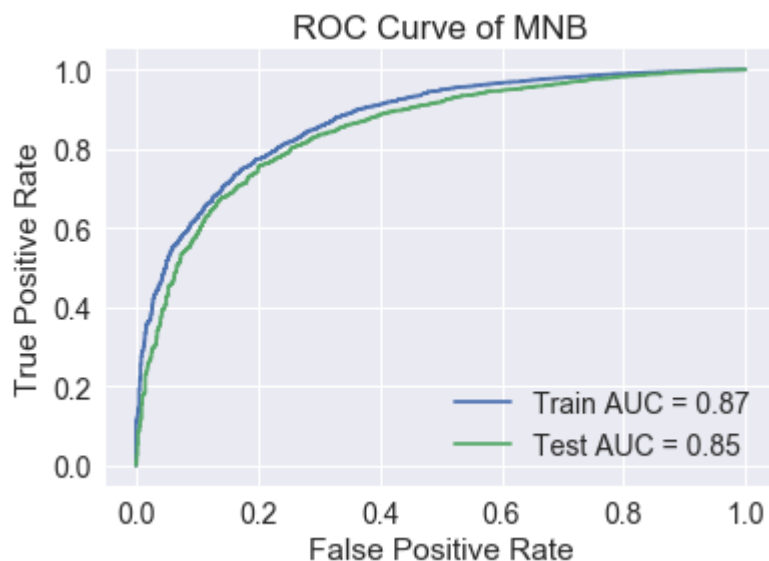


```
In [211]: #Testing with test data
clf = SVC(C=1, probability=True)
clf.fit(sent_vectors_train,Y_Train)
prediction = clf.predict_proba(sent_vectors_test)[:,-1]
print(prediction)
print(clf)

Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, clf.predict_proba(sent_vectors_train)[:,-1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, clf.predict_proba(sent_vectors_test)[:,-1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of MNB')
plt.show()
```

```
[0.65462583 0.64586442 0.9698258 ... 0.89672998 0.91763891 0.94664735]
SVC(C=1, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=True, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```



### <font color = blue>[5.3.4]Train and Test Accuracy </font>

```
In [212]: Training_Accuracy_w2v = clf.score(sent_vectors_train, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy_w2v)
Training_Error_w2v = 1 - Training_Accuracy_w2v
print('Training_Error=%0.3f'%Training_Error_w2v)

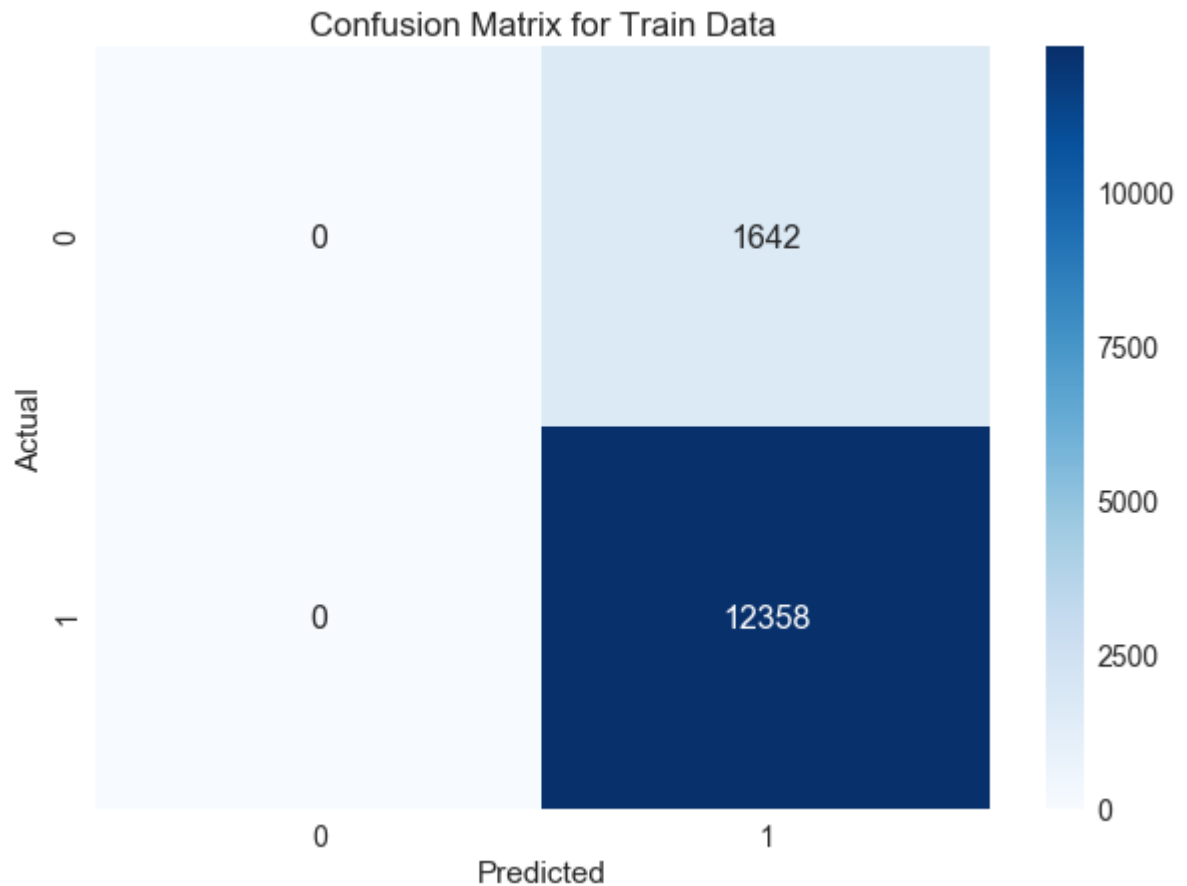
Test_Accuracy_w2v = accuracy_score(Y_Test, prediction.round())
print('Test_Accuracy=%0.3f'%Test_Accuracy_w2v)
Test_Error_w2v = 1 - Test_Accuracy_w2v
print('Test_Error=%0.3f'%Test_Error_w2v)
#print('\nThe accuracy of the MNB classifier for k = %d is %f%%' % (optimal_al
pha_bow, Test_Accuracy_Bow))

Training_Accuracy=0.883
Training_Error=0.117
Test_Accuracy=0.866
Test_Error=0.134
```

### <font color = blue>[5.3.5]Confusion Matrix </font>

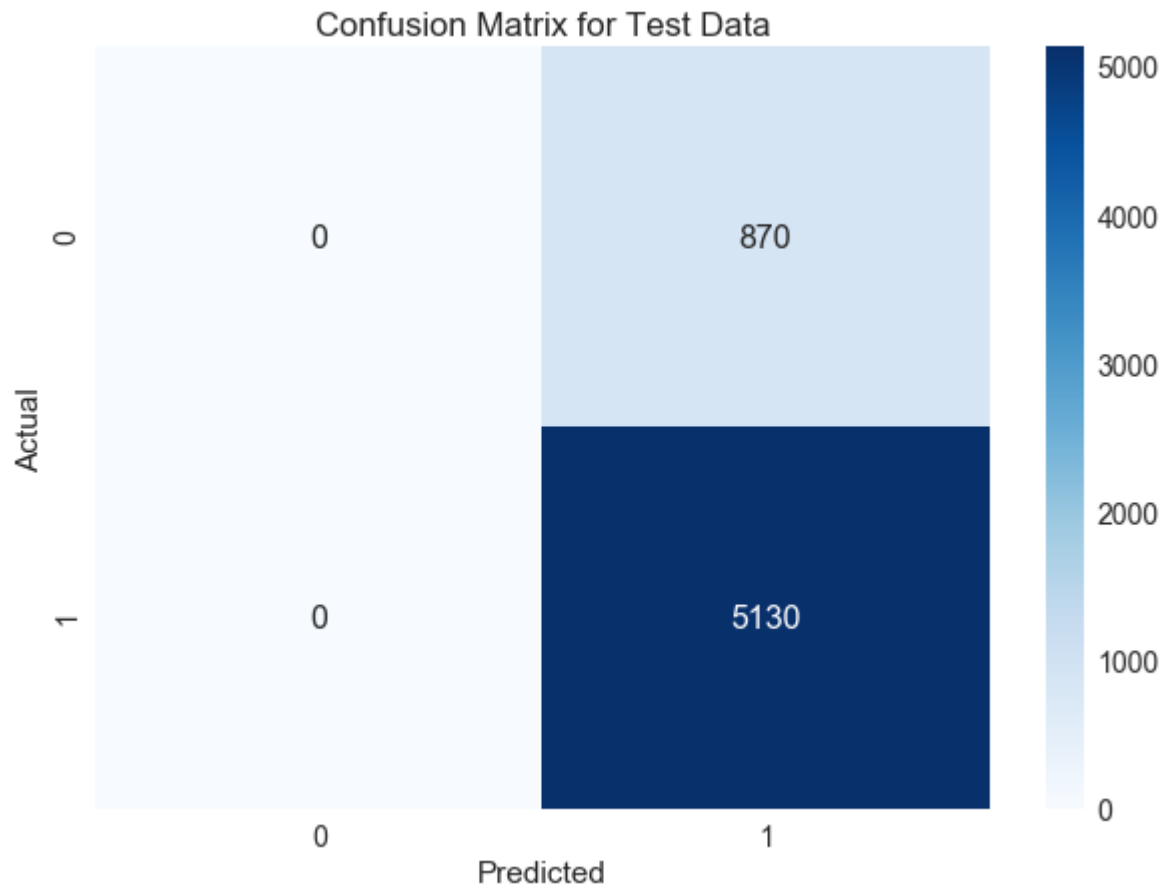
```
In [213]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Train, clf.predict(sent_vectors_train))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Train), index=np.unique(Y_Train))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Train Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[213]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d0d1a4748>



```
In [214]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, clf.predict(sent_vectors_test))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Test Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[214]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d2aea12e8>



### <font color = blue>[5.3.6] Classification Report</font>

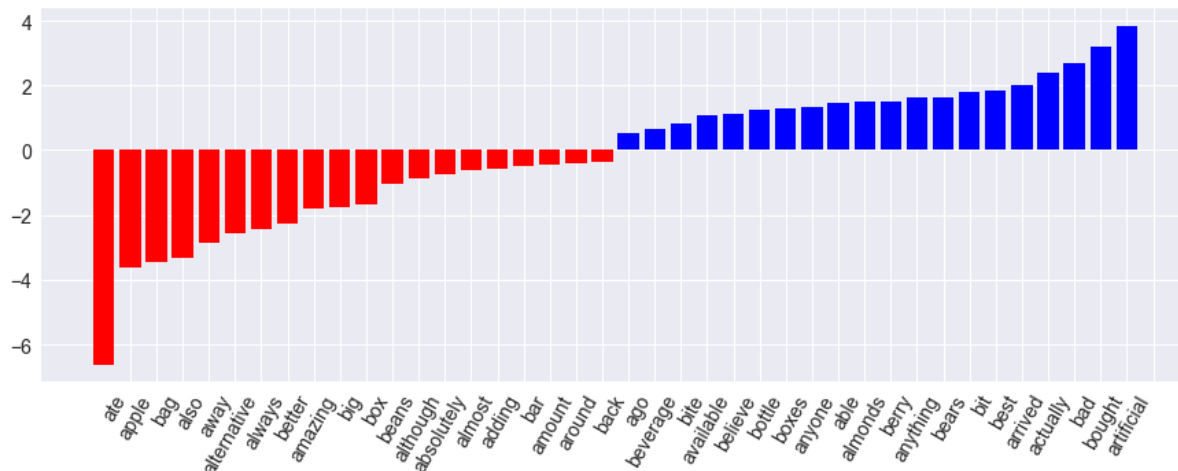
```
In [215]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction.round()))
```

	precision	recall	f1-score	support
0	0.68	0.14	0.23	870
1	0.87	0.99	0.93	5130
avg / total	0.84	0.87	0.83	6000

## <font color = Blue>[5.3.7] Feature Importance </font>

### <font color = blue>Feature Importance for Positive and Negative Class</font>

```
In [231]: clf = SVC(C=1, class_weight='balanced', kernel='linear')
          clf.fit(sent_vectors_train, Y_Train)
          plot_coefficients(clf, count_vect.get_feature_names())
```

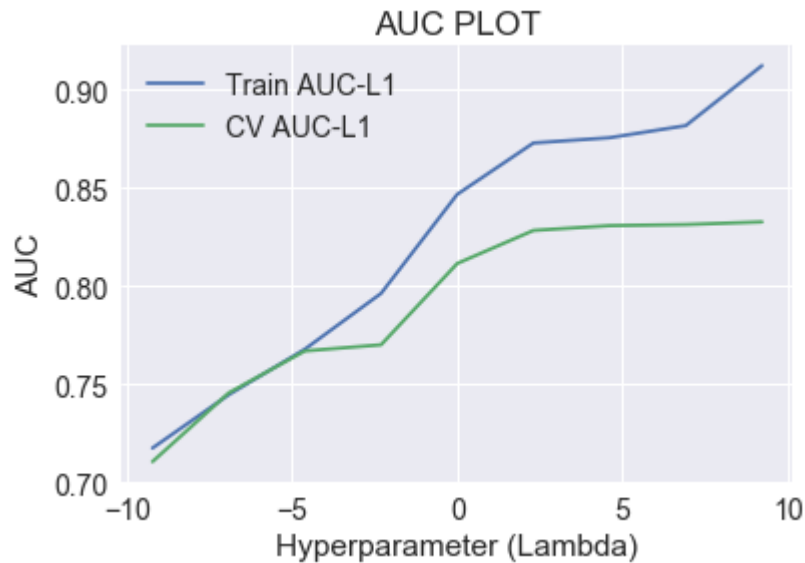


## <font color = blue> [5.4] TFIDF weighted W2v </font>

### [5.4.1] Hyperparameter tuning with L1 Regularizer and AUC Plot

```
In [251]: Optimal_Lamda_L1(tfidf_sent_vectors_train, Y_Train, tfidf_sent_vectors_cv, Y_cv
)
```

```
100%|███████████          | 9/9 [23:32<00:00, 301.39s/it]
```



```
CV AUS Scores with Penalty=? Cv auc scores with penalty L1
[0.7105182496471063, 0.7455152248437185, 0.7668969998431583, 0.77004369160449
02, 0.8114655732562569, 0.8283070063408842, 0.8307577692635164, 0.83130234590
3072, 0.832652976630593]
Maximun AUC value : 0.832652976630593
Index 8
```

### **<font color = blue>[5.4.2] ROC Curve of SVM</font>**

```

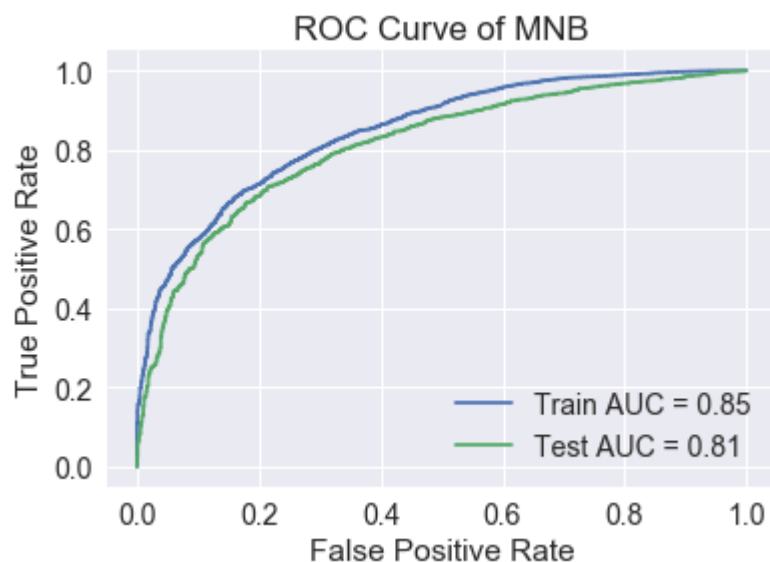
In [222]: #Testing with test data
clf = SVC(C=1, probability=True)
clf.fit(tfidf_sent_vectors_train,Y_Train)
prediction = clf.predict_proba(tfidf_sent_vectors_test)[: ,1]
print(prediction)
print(clf)

Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, clf.predict_proba(tfidf_s
ent_vectors_train)[: ,1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, clf.predict_proba(tfidf_sent
_vectors_test)[: ,1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of MNB')
plt.show()

[0.80022965 0.44810087 0.95012909 ... 0.90640622 0.93275641 0.92449651]
SVC(C=1, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=True, random_state=None, shrinking=True,
    tol=0.001, verbose=False)

```



### <font color = blue>[5.4.3]Train and Test Accuracy </font>

```
In [223]: Training_Accuracy_tfidf2v = clf.score(tfidf_sent_vectors_train, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy_tfidf2v)
Training_Error_tfidf2v = 1 - Training_Accuracy_tfidf2v
print('Training_Error=%0.3f'%Training_Error_tfidf2v)

Test_Accuracy_tfidf2v = accuracy_score(Y_Test, prediction.round())
print('Test_Accuracy=%0.3f'%Test_Accuracy_tfidf2v)
Test_Error_tfidf2v = 1 - Test_Accuracy_tfidf2v
print('Test_Error=%0.3f'%Test_Error_tfidf2v)
#print('\nThe accuracy of the MNB classifier for k = %d is %f%%' % (optimal_alpha_bow, Test_Accuracy_Bow))

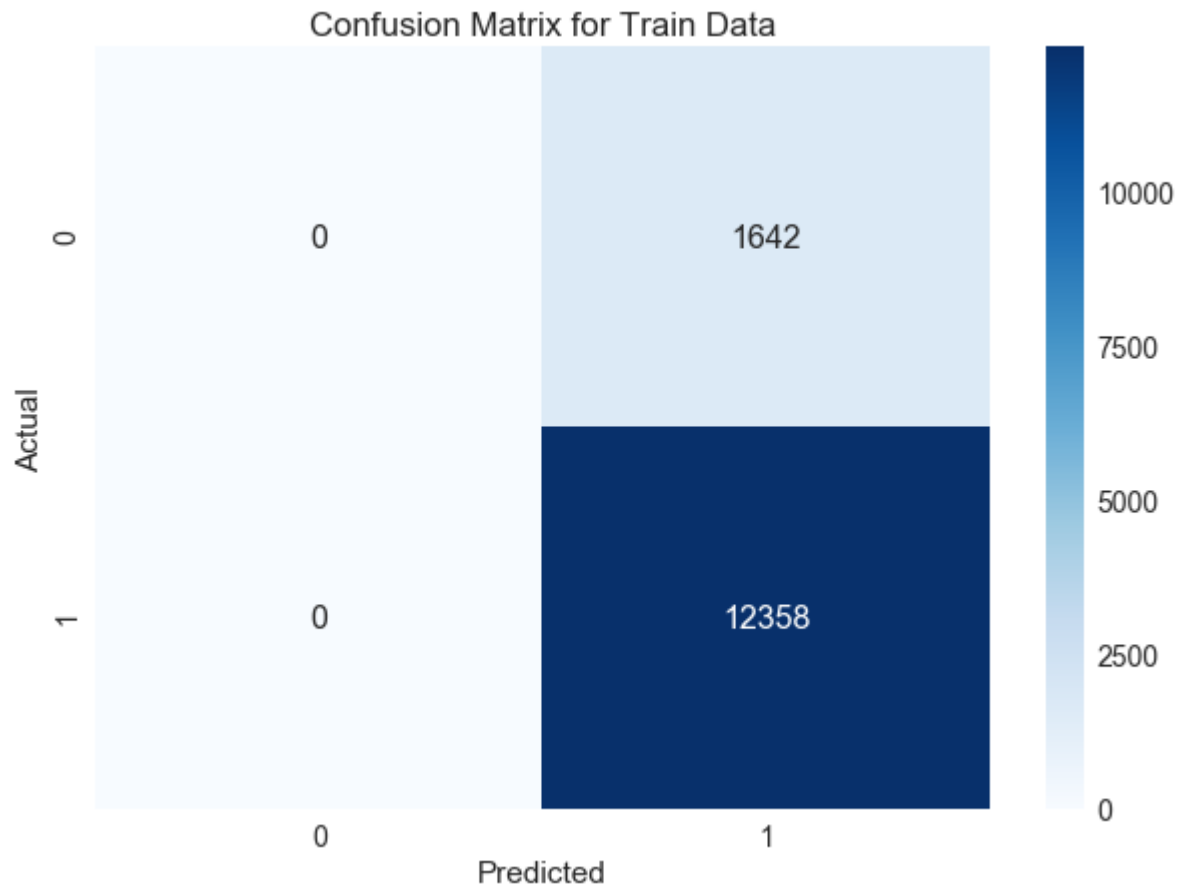
Training_Accuracy=0.883
Training_Error=0.117
Test_Accuracy=0.856
Test_Error=0.144
```

#### <font color = blue>[5.4.4]Confusion Matrix </font>



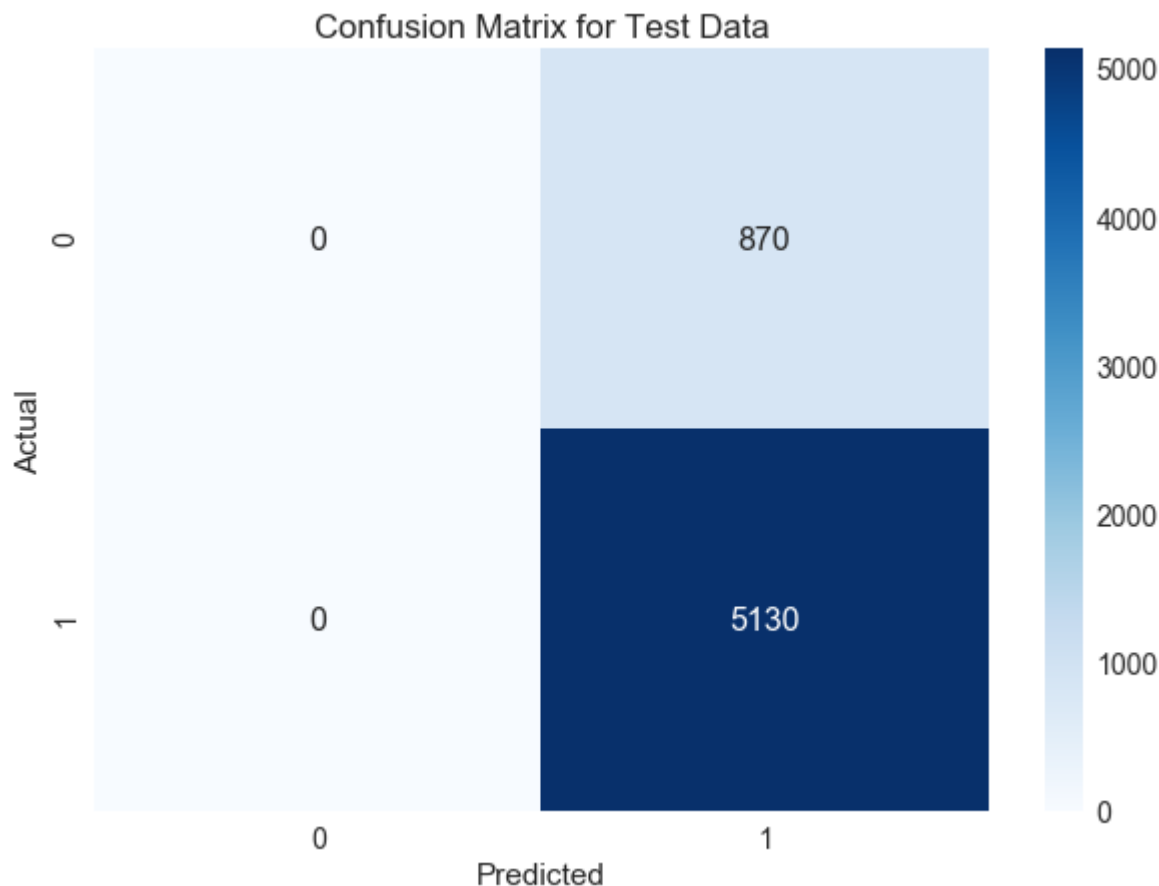
```
In [224]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Train, clf.predict(tfidf_sent_vectors_train))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Train), index=np.unique(Y_Train))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Train Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[224]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d2c91be48>



```
In [225]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, clf.predict(tfidf_sent_vectors_test))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
plt.title("Confusion Matrix for Test Data")
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[225]: <matplotlib.axes.\_subplots.AxesSubplot at 0x23d18e02080>



### <font color = blue>[5.4.5] Classification Report</font>

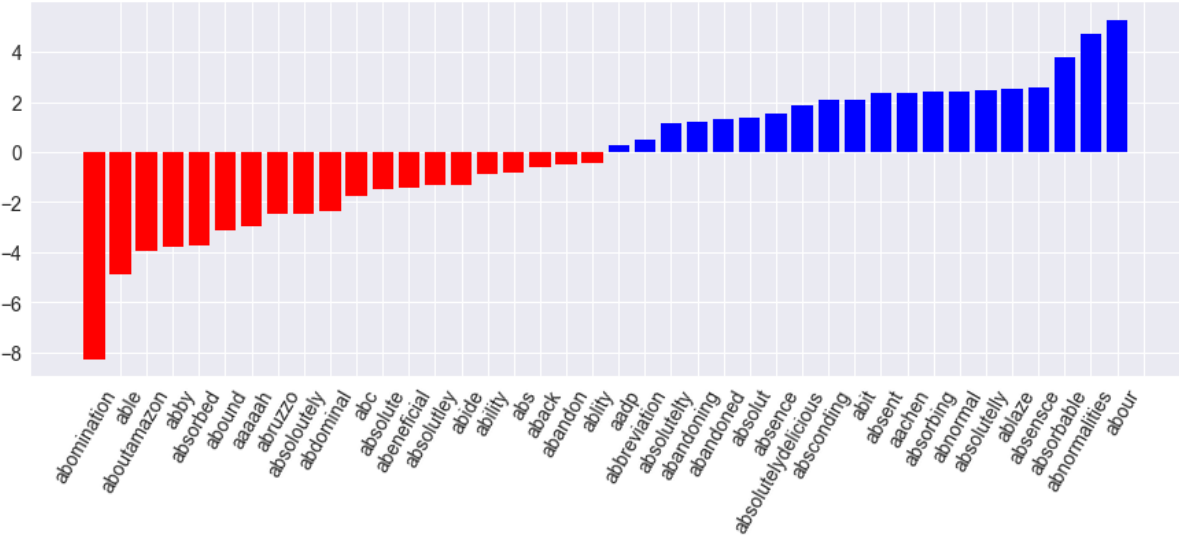
```
In [226]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction.round()))
```

	precision	recall	f1-score	support
0	0.54	0.06	0.11	870
1	0.86	0.99	0.92	5130
avg / total	0.82	0.86	0.80	6000

<font color = Blue>[5.4.6] Feature Importance </font>

<font color = blue>Feature Importance for Positive and Negative Class</font>

```
In [232]: clf = SVC(C=1, class_weight='balanced', kernel='linear')
          clf.fit(tfidf_sent_vectors_train, Y_Train)
          plot_coefficients(clf, model.get_feature_names())
```



<font color = Green>Pretty Table</font>

```
In [258]: from prettytable import PrettyTable
comparison = PrettyTable()
comparison.field_names = ["Vectorizer", "Kernel", "AUC - L1", "AUC-L2", "Training Error", "Test Error"]
comparison.add_row(["BOW", "Linear", 0.83, 0.93, 0.02, 0.12])
comparison.add_row(["TF-IDF", "Linear", 0.904, 0.95, 0.01, 0.09])
comparison.add_row(["Avg W2V", "Linear", 0.903, 0.904, 0.118, 0.116])
comparison.add_row(["TF-IDFWeighted W2V", "Linear", 0.88, 0.88, 0.130, 0.129])
comparison.add_row(["BoW", "RBF", "-", 0.87, 0.060, 0.888])
comparison.add_row(["TF-IDF", "RBF", "-", 0.89, 0.048, 0.888])
comparison.add_row(["Avg W2V", "RBF", "-", 0.85, 0.117, 0.134])
comparison.add_row(["TF-IDFWeighted W2V", "RBF", "-", 0.83, 0.117, 0.144])
print(comparison)
```

```
+-----+-----+-----+-----+-----+-----+
---+
| Vectorizer | Kernel | AUC - L1 | AUC-L2 | Training Error | Test Error |
+-----+-----+-----+-----+-----+-----+
---+
| BOW        | Linear | 0.83     | 0.93    | 0.02           | 0.12        |
| TF-IDF     | Linear | 0.904    | 0.95    | 0.01           | 0.09        |
| Avg W2V    | Linear | 0.903    | 0.904   | 0.118          | 0.116       |
| TF-IDFWeighted W2V | Linear | 0.88     | 0.88    | 0.13           | 0.129       |
| BoW        | RBF   | -        | 0.87    | 0.06           | 0.888       |
| TF-IDF     | RBF   | -        | 0.89    | 0.048          | 0.888       |
| Avg W2V    | RBF   | -        | 0.85    | 0.117          | 0.134       |
| TF-IDFWeighted W2V | RBF   | -        | 0.83    | 0.117          | 0.144       |
+-----+-----+-----+-----+-----+-----+
---+
```

**<font color = Green>Conclusion</font>**

1. Applied SVM on all the 4 vectorizers(BOW, TFIDF, AVG-W2V, TFIDF-AVG\_W2V).
2. Sorted the data based on Time and Considered 100 K data points for Training set 70K, Test set: 30K.  
Worked on 2 version of SVM 1-Linear Kernel and 2-RBF Kernel
3. While working with linear kernel used 'SGDClassifier' with 'hinge loss' and used CalibratedClassifierCV.  
and Class\_Weight has set to 'balanced'
4. While working with RBF Kernel, set min\_df=10 and max\_features=500 and considered 20K data points.
5. Used AUC as a metric for hyperparameter tuning. And took the range of lambda values between ( $10^{-4}$  to  $10^4$ ).
6. Found the top 20 features of positive and negative class for the featurizations Bow and TF-IDF, AvgW2V and TFIDF Weighted vector using Linear Kernel.
7. Found the top 20 features of positive and negative class for the featurizations AvgW2V and TF-IDFWeighted W2V using SVC Kernel. Got error for BoW and TF-IDF.
8. With reference to the pretty table, here is my understanding: a. Linear kernel SVM by using TF-IDF featurization is having the best AUC score: 0.95.
9. RBF Kernel SVM by using TF-IDF having best AUC score: 0.89.
10. Plotted ROC Curve and Confusion Matrix for train and test data for each vectorizer.