# Amazon Fine Food Reviews Analysis

Data Source: https://www.kaggle.com/snap/amazon-fine-food-reviews

EDA: 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 - unqiue 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 Score/Rating. A rating of 4 or 5 can be cosnidered as a positive review. A rating of 1 or 2 can be considered as negative one. A review of rating 3 is considered nuetral 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 wil be set to "positive". Otherwise, it will be set to "negative".

In [1]:

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")


import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```python
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
```

In [2]:

```python
# using SQLite Table to read data.
con = sqlite3.connect('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 """, 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 (525814, 10)

Out[2]:

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

| Id | ProductId | UserId | ProfileName"Cores" | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time | Summary |
|---|---|---|---|---|---|---|---|---|

In [3]:

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

In [4]:

```
print(display.shape)
display.head()
```

(80668, 7)

Out[4]:

| | 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 [5]:

```
display[display['UserId']=='AZY10LLTJ71NX']
```

Out[5]:

| | 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 [6]:

```
display['COUNT(*)'].sum()
```

Out[6]:

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 [7]:

```
display= pd.read_sql_query("""
SELECT *
FROM Reviews
```

```
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[7]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time | Summ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACK QUADRAT VANI WAFE |
| 1 | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACK QUADRAT VANI WAFE |
| 2 | 138277 | B000HDOPYM | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACK QUADRAT VANI WAFE |
| 3 | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACK QUADRAT VANI WAFE |
| 4 | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACK QUADRAT VANI WAFE |

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 delelte 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 [8]:

```
#Sorting data according to ProductId in ascending order
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False, kind='qui
cksort', na_position='last')
```

In [9]:

```
#Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first', inpl
ace=False)
final.shape
```

Out[9]:

```
(364173, 10)
```

In [10]:

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

Out[10]:

```
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 calcualtions

In [11]:

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

display.head()
```

Out[11]:

| | 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 [12]:

```
final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

In [13]:

```
#Before starting the next phase of preprocessing lets see the number of entries left
print(final.shape)

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

```
(364171, 10)
```

Out[13]:

```
1    307061
0     57110
Name: Score, dtype: int64
```

# [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 obsereved to be better than Porter Stemming)

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

In [14]:

```
# 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)
```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along a
nd he always can sing the refrain. he's learned about whales, India, drooping roses:  i love all t
he new words this book  introduces and the silliness of it all.  this is a classic book i am
willing to bet my son will STILL be able to recite from memory when he is  in college
==================================================
I was really looking forward to these pods based on the reviews.  Starbucks is good, but I prefer
bolder taste.... imagine my surprise when I ordered 2 boxes - both were expired! One expired back
in 2005 for gosh sakes.  I admit that Amazon agreed to credit me for cost plus part of shipping, b
ut geez, 2 years expired!!!  I'm hoping to find local San Diego area shoppe that carries pods so t
hat I can try something different than starbucks.
==================================================
Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing
I do not think belongs in it is Canola oil. Canola or rapeseed is not someting a dog would ever fi
nd in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's Food
industries have convinced the masses that Canola oil is a safe and even better oil than olive or v
irgin coconut, facts though say otherwise. Until the late 70's it was poisonous until they figured
out a way to fix that. I still like it but it could be better.
==================================================
Can't do sugar.  Have tried scores of SF Syrups.  NONE of them can touch the excellence of this
product.<br /><br />Thick, delicious.  Perfect.  3 ingredients: Water, Maltitol, Natural Maple
Flavor.  PERIOD.  No chemicals.  No garbage.<br /><br />Have numerous friends & family members
hooked on this stuff.  My husband & son, who do NOT like "sugar free" prefer this over major label
regular syrup.<br /><br />I use this as my SWEETENER in baking: cheesecakes, white brownies,
muffins, pumpkin pies, etc... Unbelievably delicious...<br /><br />Can you tell I like it? :)
==================================================


In [15]:

```
# remove urls from text python: https://stackoverflow.com/a/40823105/4084039
sent_0 = re.sub(r"http\S+", "", sent_0)
sent_1000 = re.sub(r"http\S+", "", sent_1000)
sent_150 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along a
nd he always can sing the refrain. he's learned about whales, India, drooping roses:  i love all t
he new words this book  introduces and the silliness of it all.  this is a classic book i am
willing to bet my son will STILL be able to recite from memory when he is  in college


In [16]:

```
# https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all-tags-from-an
-element
from bs4 import BeautifulSoup
```

```python
soup = BeautifulSoup(sent_0, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1000, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1500, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_4900, 'lxml')
text = soup.get_text()
print(text)
```

```
this witty little book makes my son laugh at loud. i recite it in the car as we're driving along a
nd he always can sing the refrain. he's learned about whales, India, drooping roses:  i love all t
he new words this book  introduces and the silliness of it all.  this is a classic book i am
willing to bet my son will STILL be able to recite from memory when he is  in college
==================================================
I was really looking forward to these pods based on the reviews.  Starbucks is good, but I prefer
bolder taste.... imagine my surprise when I ordered 2 boxes - both were expired! One expired back
in 2005 for gosh sakes.  I admit that Amazon agreed to credit me for cost plus part of shipping, b
ut geez, 2 years expired!!!  I'm hoping to find local San Diego area shoppe that carries pods so t
hat I can try something different than starbucks.
==================================================
Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing
I do not think belongs in it is Canola oil. Canola or rapeseed is not someting a dog would ever fi
nd in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's Food
industries have convinced the masses that Canola oil is a safe and even better oil than olive or v
irgin coconut, facts though say otherwise. Until the late 70's it was poisonous until they figured
out a way to fix that. I still like it but it could be better.
==================================================
Can't do sugar.  Have tried scores of SF Syrups.  NONE of them can touch the excellence of this
product.Thick, delicious.  Perfect.  3 ingredients: Water, Maltitol, Natural Maple Flavor.
PERIOD.  No chemicals.  No garbage.Have numerous friends & family members hooked on this stuff.  M
y husband & son, who do NOT like "sugar free" prefer this over major label regular syrup.I use thi
s as my SWEETENER in baking: cheesecakes, white brownies, muffins, pumpkin pies, etc...
Unbelievably delicious...Can you tell I like it? :)
```

In [17]:

```python
# 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 [18]:

```python
sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

```
Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing
I do not think belongs in it is Canola oil. Canola or rapeseed is not someting a dog would ever fi
```

I do not think belongs in it is Canola oil. Canola or rapeseed is not someting a dog would ever fi
nd in nature and if it did find rapeseed in nature and eat it, it would poison them. Today is Food
industries have convinced the masses that Canola oil is a safe and even better oil than olive or v
irgin coconut, facts though say otherwise. Until the late 70 is it was poisonous until they
figured out a way to fix that. I still like it but it could be better.
==================================================

In [19]:

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

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along a
nd he always can sing the refrain. he's learned about whales, India, drooping roses:  i love all t
he new words this book  introduces and the silliness of it all.  this is a classic book i am
willing to bet my son will STILL be able to recite from memory when he is  in college

In [20]:

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

Great ingredients although chicken should have been 1st rather than chicken broth the only thing I
do not think belongs in it is Canola oil Canola or rapeseed is not someting a dog would ever find
in nature and if it did find rapeseed in nature and eat it it would poison them Today is Food indu
stries have convinced the masses that Canola oil is a safe and even better oil than olive or virgi
n coconut facts though say otherwise Until the late 70 is it was poisonous until they figured out
a way to fix that I still like it but it could be better

In [21]:

```python
# 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', "y
ou'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', 'itself', 'they', 'them',
'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll",
'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', '
while', 'of', \
            'at', 'by', 'for', 'with', 'about', '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', 'all', 'any', 'both', 'e
ach', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll'
, 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "do
esn't", 'hadn',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn',
"mightn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',
"wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"])
```

In [22]:

```python
# Combining all the above stundents
from tqdm import tqdm
```

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

```
100%|████████| 364171/364171 [02:41<00:00, 2254.27it/s]
```

In [23]:

```python
preprocessed_reviews[1500]
```

Out[23]:

'great ingredients although chicken rather chicken broth thing not think belongs canola oil canola
rapeseed not someting dog would ever find nature find rapeseed nature eat would poison today food
industries convinced masses canola oil safe even better oil olive virgin coconut facts though say
otherwise late poisonous figured way fix still like could better'

## [3.2] Preprocessing Review Summary

In [24]:

```python
## Similartly you can do preprocessing for review summary also.
```

In [25]:

```python
## Summary preprocessing
from tqdm import tqdm
preprocessed_summary = []
# tqdm is for printing the status bar
for sentance in tqdm(final['Summary'].values):
    sentance = re.sub(r"http\S+", "", sentance)
    sentance = BeautifulSoup(sentance, 'lxml').get_text()
    sentance = decontracted(sentance)
    sentance = re.sub("\S*\d\S*", "", sentance).strip()
    sentance = re.sub('[^A-Za-z]+', ' ', sentance)
    # https://gist.github.com/sebleier/554280
    sentance = ' '.join(e.lower() for e in sentance.split() if e.lower() not in stopwords)
    preprocessed_summary.append(sentance.strip())
```

```
  0%|                                                      | 0/364
[00:00<?, ?it/s]C:\Users\Crypto Lab-1\Anaconda3\lib\site-packages\bs4\__init__.py:273:
UserWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pa
ss the filehandle into Beautiful Soup.
  ' Beautiful Soup.' % markup)
  6%|█                                                     | 22012/364171
[00:03<00:53, 6344.74it/s]C:\Users\Crypto Lab-1\Anaconda3\lib\site-packages\bs4\__init__.py:273: U
serWarning: "b'.'" looks like a filename, not markup. You should probably open this file and pass
the filehandle into Beautiful Soup.
  ' Beautiful Soup.' % markup)
 10%|██████                                                | 35647/364171
[00:05<00:52, 6231.94it/s]C:\Users\Crypto Lab-1\Anaconda3\lib\site-packages\bs4\__init__.py:273: U
serWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pas
s the filehandle into Beautiful Soup.
  ' Beautiful Soup.' % markup)
 10%|██████                                                | 36957/364171
[00:05<00:53, 6169.75it/s]C:\Users\Crypto Lab-1\Anaconda3\lib\site-packages\bs4\__init__.py:273: U
serWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pas
s the filehandle into Beautiful Soup.
  ' Beautiful Soup.' % markup)
 14%|████████                                              | 51402/364171
[00:08<00:50, 6234.81it/s]C:\Users\Crypto Lab-1\Anaconda3\lib\site-packages\bs4\__init__.py:273: U
serWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pas
s the filehandle into Beautiful Soup.
```

In [24]:

```python
final['Cleaned_text'] = preprocessed_reviews
```

```
### Sort data according to time series
final.sort_values('Time',inplace=True)
```

In [26]:

```
### Taking 100k samples
final_100k = final.sample(n=100000)
```

In [27]:

```
### Taking 20k samples
final_20k = final.sample(n=20000)
```

In [28]:

```
x = final_100k['Cleaned_text']
x.size
```

Out[28]:

```
100000
```

In [29]:

```
y = final_100k['Score']
y.size
```

Out[29]:

```
100000
```

In [30]:

```
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=42)
```

# [4] Featurization

## [4.1] BAG OF WORDS

In [31]:

```
#BoW
count_vect = CountVectorizer() #in scikit-learn
x_train_bow = count_vect.fit_transform(x_train)
print("some feature names ", count_vect.get_feature_names()[:10])
print('='*50)

x_test_bow = count_vect.transform(x_test)
print("the type of count vectorizer ",type(x_test_bow))
print("the shape of out text BOW vectorizer ",x_test_bow.get_shape())
print("the number of unique words ", x_test_bow.get_shape()[1])
```

```
some feature names  ['aa', 'aaa', 'aaaa', 'aaaaa', 'aaaaaaaaaaaaaaa', 'aaaahhhhhh', 'aaaallll', 'a
aah', 'aaahhhhhh', 'aabout']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (30000, 50725)
the number of unique words  50725
```

In [32]:

```
from sklearn import preprocessing
x_train_bow = preprocessing.normalize(x_train_bow)
x_test_bow = preprocessing.normalize(x_test_bow)
```

## [4.2] Bi-Grams and n-Grams.

In [33]:

```
#bi-gram, tri-gram and n-gram

#removing stop words like "not" should be avoided before building n-grams
# count_vect = CountVectorizer(ngram_range=(1,2))
# please do read the CountVectorizer documentation http://scikit-
learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html

# you can choose these numebrs min_df=10, max_features=5000, of your choice
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
final_bigram_counts = count_vect.fit_transform(x_train)
print("the type of count vectorizer ",type(final_bigram_counts))
print("the shape of out text BOW vectorizer ",final_bigram_counts.get_shape())
print("the number of unique words including both unigrams and bigrams ", final_bigram_counts.get_s
hape()[1])
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (70000, 5000)
the number of unique words including both unigrams and bigrams  5000
```

## [4.3] TF-IDF

In [33]:

```
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
x_train_tfidf = 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_test_tfidf = tf_idf_vect.transform(x_test)
print("the type of count vectorizer ",type(x_test_tfidf))
print("the shape of out text TFIDF vectorizer ",x_test_tfidf.get_shape())
print("the number of unique words including both unigrams and bigrams ", x_test_tfidf.get_shape()[
1])
```

```
some sample features(unique words in the corpus) ['abc', 'abdominal', 'ability', 'ability make', '
able', 'able add', 'able buy', 'able chew', 'able create', 'able drink']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer  (30000, 40705)
the number of unique words including both unigrams and bigrams  40705
```

In [34]:

```
from sklearn import preprocessing
x_train_tfidf = preprocessing.normalize(x_train_tfidf)
x_test_tfidf = preprocessing.normalize(x_test_tfidf)
```

## [4.4] Word2Vec

In [35]:

```
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance_train=[]
for sentance in x_train:
    list_of_sentance_train.append(sentance.split())
```

In [36]:

```
# Test your own Word2Vec model using your own text corpus
i=0
list_of_sentence_test=[]
for sentence in x_test:
    list_of_sentence_test.append(sentence.split())
```

In [37]:

```
# Using Google News Word2Vectors

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


# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17SRFAzZPY
# you can comment this whole cell
# or change these varible according to your need

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True

if want_to_train_w2v:
    # min_count = 5 considers only words that occured atleast 5 times
    w2v_model=Word2Vec(list_of_sentence_train,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))

elif want_to_use_google_w2v and is_your_ram_gt_16g:
    if os.path.isfile('GoogleNews-vectors-negative300.bin'):
        w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)
        print(w2v_model.wv.most_similar('great'))
        print(w2v_model.wv.most_similar('worst'))
    else:
        print("you don't have gogole's word2vec file, keep want_to_train_w2v = True, to train your
own w2v ")
```

```
[('excellent', 0.8485338687896729), ('fantastic', 0.8345503211021423), ('good',
0.8323901891708374), ('awesome', 0.8018529415130615), ('terrific', 0.8009918332099915),
('wonderful', 0.8004620671272278), ('perfect', 0.7499939799308777), ('incredible',
0.6992645263671875), ('nice', 0.6956198811531067), ('fabulous', 0.6833528876304626)]
==================================================
[('greatest', 0.8169127702713013), ('nastiest', 0.7796850800514221), ('best', 0.7342137098312378),
('tastiest', 0.6954528093338013), ('disgusting', 0.6851107478141785), ('nicest',
0.6435158848762512), ('closest', 0.6261910200119019), ('beats', 0.6145021915435791), ('softest', 0
.6076576113700867), ('horrible', 0.6008611917495728)]
```

In [38]:

```
w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])
```

```
number of words that occured minimum 5 times  16124
sample words  ['easy', 'make', 'meal', 'adding', 'meat', 'veggies', 'tastes', 'good', 'husband',
'starbucks', 'person', 'tasted', 'coffee', 'said', 'made', 'cups', 'right', 'nice', 'flavor', 'lov
e', 'coconut', 'pineapple', 'not', 'wild', 'one', 'water', 'splash', 'bad', 'mildly', 'coconutty',
'though', 'sweet', 'anything', 'special', 'rather', 'bland', 'hard', 'time', 'finishing',
'started', 'warm', 'since', 'competitive', 'athlete', 'great', 'need', 'replenish',
'electrolytes', 'would', 'drink']
```

## [4.4.1] Converting text into vectors using Avg W2V, TFIDF-W2V

**[4.4.1.1] Avg W2v**

```python
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentance_train): # 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_train.append(sent_vec)
print(len(sent_vectors_train))
print(len(sent_vectors_train[0]))
```

```
100%|██████████| 70000/70000 [02:59<00:00, 389.36it/s]
```

```
70000
50
```

```python
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentance_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)
print(len(sent_vectors_test))
print(len(sent_vectors_test[0]))
```

```
100%|██████████| 30000/30000 [01:19<00:00, 376.49it/s]
```

```
30000
50
```

```python
x_train_avgw2v=sent_vectors_train
x_test_avgw2v=sent_vectors_test
```

```python
from sklearn import preprocessing
x_train_avgw2v = preprocessing.normalize(x_train_avgw2v)
x_test_avgw2v = preprocessing.normalize(x_test_avgw2v)
```

**[4.4.1.2] TFIDF weighted W2v**

In [43]:

```python
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer(min_df=10, max_features=500)
tf_idf_matrix = model.fit_transform(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_)))
```

In [44]:

```python
# 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_train = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentance_train): # 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_train.append(sent_vec)
    row += 1
```

```
100%|██████████| 70000/70000 [03:22<00:00, 345.72it/s]
```

In [45]:

```python
tfidf_sent_vectors_test = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentance_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%|██████████| 30000/30000 [01:29<00:00, 336.48it/s]
```

In [46]:

```python
x_train_tfidfw2v = tfidf_sent_vectors_train
x_test_tfidfw2v = tfidf_sent_vectors_test
```

In [47]:

```python
from sklearn import preprocessing
x_train_tfidfw2v = preprocessing.normalize(x_train_tfidfw2v)
x_test_tfidfw2v = preprocessing.normalize(x_test_tfidfw2v)
```

# [5] Assignment 7: SVM

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
   - Similarly, like kdtree of 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 paramter 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 value
   - Find the best hyper paramter 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 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 with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps.

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

**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.

# Applying SVM

## [5.1] Linear SVM

### [5.1.1] Applying Linear SVM on BOW, SET 1

In [48]:

```python
from math import log
alpha_values = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]
x =[log(y,10) for y in alpha_values]
x
```

Out[48]:

```
[-3.99999999999999,
 -2.9999999999999996,
 -1.9999999999999996,
 -0.9999999999999998,
 0.0,
 1.0,
 2.0,
 4.0]
```

In [48]:

```python
## find hyperparameter using cross validation score and plot AUC
## find hyperparameter using cross validation score and plot AUC
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import TimeSeriesSplit
tscv = TimeSeriesSplit(n_splits=10)
from sklearn.model_selection import cross_val_score

alphas = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]

#empty lists that stores cv scores and training_scores
cv_scores = []
training_scores = []

#perform k fold cross validation
for alpha in alphas:
    clf = SGDClassifier(loss='hinge',alpha=alpha)
    scores = cross_val_score(clf, x_train_bow, y_train, cv=10, scoring='roc_auc')
    scores_training = clf.fit(x_train_bow, y_train).score(x_train_bow, y_train)

    cv_scores.append(scores.mean())
    training_scores.append(scores_training)

#plot cross-validated score, training score vs alpha
plt.plot(x, training_scores, label='Train AUC')
plt.plot(x, cv_scores, label='CV AUC')
plt.legend()
plt.xlabel("Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

In [53]:

```python
## find hyperparameter alpha using GridserachCV
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import SGDClassifier
tuned_parameters = [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l1','l2']}]
model = GridSearchCV(SGDClassifier(),tuned_parameters,cv=tscv,scoring='roc_auc')
model.fit(x_train_bow,y_train)
print(model.best_estimator_)
print(model.score(x_test_bow,y_test))
print("Best HyperParameter: ",model.best_params_)
print("Best Accuracy: %.2f%%"%(model.best_score_*100))
```

```
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=None,
       n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
       power_t=0.5, random_state=None, shuffle=True, tol=None,
       validation_fraction=0.1, verbose=0, warm_start=False)
0.9438214461598041
Best HyperParameter:  {'alpha': 0.0001, 'penalty': 'l2'}
Best Accuracy: 93.75%
```

Observation: BY observing the hyperparameter is 0.0001

In [55]:

```python
##Train the model with best hyperparameter
clf = SGDClassifier(loss='hinge',alpha=0.0001,penalty='l2')
clf.fit(x_train_bow,y_train)
y_pred = clf.predict(x_test_bow)
```

In [56]:

```python
### ROC Curve using false positive rate versus true positive rate
```

In [57]:

```python
from sklearn.metrics import roc_curve,auc
probas_ = clf.fit(x_train_bow, y_train).decision_function(x_test_bow)
# Compute ROC curve and area the curve
fpr, tpr, thresholds = roc_curve(y_test, probas_)
```

In [58]:

```python
probas_ = clf.fit(x_train_bow, y_train).decision_function(x_train_bow)
# Compute ROC curve and area the curve
fpr_, tpr_, thresholds = roc_curve(y_train, probas_)
```

In [59]:

```python
plt.title('Receiver Operating Characteristic')
plt.plot(fpr, tpr, 'b')
plt.plot(fpr_, tpr_, 'r')
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```
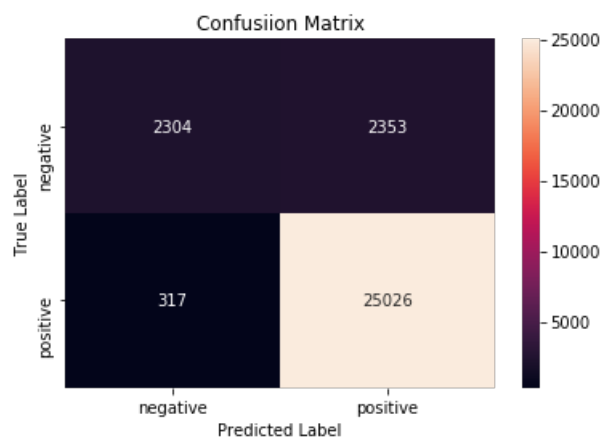
**In [60]:**

```python
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

**Out[60]:**

```
array([[ 2304,  2353],
       [  317, 25026]], dtype=int64)
```

**In [61]:**

```python
# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



# Top 10 important features of positive class

**In [62]:**

```python
#https://stackoverflow.com/questions/11116697/how-to-get-most-informative-features-for-scikit-lear
n-classifiers
def show_most_informative_features(vectorizer, clf, n=10):
    feature_names = vectorizer.get_feature_names()
    coefs_with_fns = sorted(zip(clf.coef_[0], feature_names),reverse=True)[:n]
    top = zip(coefs_with_fns[:n], coefs_with_fns[:-(n + 1):-1])
    print("Positive")

print("_____
")
    for (coef_1, fn_1),(coef_2, fn_2) in top:
        print("\t%.4f\t%-15s" % (coef_1, fn_1))

show_most_informative_features(count_vect,clf)
```

```
Positive

_____
 2.3833 start day
 1.9760 far superior
 1.8102 oily
 1.3500 rye
 1.2271 texture flavor
 1.1113 exercise
 0.9122 not product
```

```
 0.8373 bread mix
 0.8317 bring back
 0.7883 brewing
```

In [63]:

```python
##Top 10 important features of negative class
```

In [64]:

```python
#https://stackoverflow.com/questions/11116697/how-to-get-most-informative-features-for-scikit-lear
n-classifiers
def show_most_informative_features(vectorizer, clf, n=10):
    feature_names = vectorizer.get_feature_names()
    coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))[:n]
    top = zip(coefs_with_fns[:n], coefs_with_fns[:-(n + 1):-1])
    print("Negative")

print("_____
")
    for (coef_1, fn_1),(coef_2, fn_2) in top:
        print("\t%.4f\t%-15s" % (coef_2, fn_2))

show_most_informative_features(count_vect,clf)
```

```
Negative
_____
 -0.8277 overwhelming
 -0.8895 coffees
 -0.9992 love love
 -1.1380 noticeable
 -1.1938 promised
 -1.4863 offer
 -1.5285 topping
 -1.6625 paying
 -1.6654 story
 -3.2402 old daughter
```

## [5.1.2] Applying Linear SVM on TFIDF, SET 2

In [65]:

```python
# Please write all the code with proper documentation
```

In [49]:

```python
## find hyperparameter using cross validation score and plot AUC
## find hyperparameter using cross validation score and plot AUC
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import TimeSeriesSplit
tscv = TimeSeriesSplit(n_splits=10)
from sklearn.model_selection import cross_val_score

alphas = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]

#empty lists that stores cv scores and training_scores
cv_scores = []
training_scores = []

#perform k fold cross validation
for alpha in alphas:
    clf = SGDClassifier(loss='hinge',alpha=alpha)
    scores = cross_val_score(clf, x_train_tfidf, y_train, cv=10, scoring='roc_auc')
    scores_training = clf.fit(x_train_tfidf, y_train).score(x_train_tfidf, y_train)

    cv_scores.append(scores.mean())
    training_scores.append(scores_training)

#plot cross-validated score, training score vs alpha
plt.plot(x, training_scores, label='Train AUC')
plt.plot(x, cv_scores, label='CV AUC')
```

```
plt.legend()
plt.xlabel("Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

```
## find hyperparameter alpha using GridserachCV
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import SGDClassifier
tuned_parameters = [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l1','l2']}]
model = GridSearchCV(SGDClassifier(),tuned_parameters,cv=tscv,scoring='roc_auc')
model.fit(x_train_tfidf,y_train)
print(model.best_estimator_)
print(model.score(x_test_tfidf,y_test))
print("Best HyperParameter: ",model.best_params_)
print("Best Accuracy: %.2f%%"%(model.best_score_*100))
```

```
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=None,
       n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
       power_t=0.5, random_state=None, shuffle=True, tol=None,
       validation_fraction=0.1, verbose=0, warm_start=False)
0.9578276745224301
Best HyperParameter:  {'alpha': 0.0001, 'penalty': 'l2'}
Best Accuracy: 95.33%
```

Observation: BY observing the hyperparameter is 0.0001

In [70]:

```
##Train the model with best hyperparameter
clf = SGDClassifier(loss='hinge',alpha=0.0001,penalty='l2')
clf.fit(x_train_tfidf,y_train)
y_pred = clf.predict(x_test_tfidf)
```

In [71]:

```
### ROC Curve using false positive rate versus true positive rate
```

In [72]:

```
from sklearn.metrics import roc_curve,auc
probas_ = clf.fit(x_train_tfidf, y_train).decision_function(x_test_tfidf)
# Compute ROC curve and area the curve
fpr, tpr, thresholds = roc_curve(y_test, probas_)
probas_ = clf.fit(x_train_tfidf, y_train).decision_function(x_train_tfidf)
# Compute ROC curve and area the curve
fpr_, tpr_, thresholds = roc_curve(y_train, probas_)
plt.title('Receiver Operating Characteristic')
plt.plot(fpr, tpr, 'b')
plt.plot(fpr_, tpr_, 'r')
plt.ylabel('True Positive Rate')
```

```
plt.xlabel('False Positive Rate')
plt.show()
```



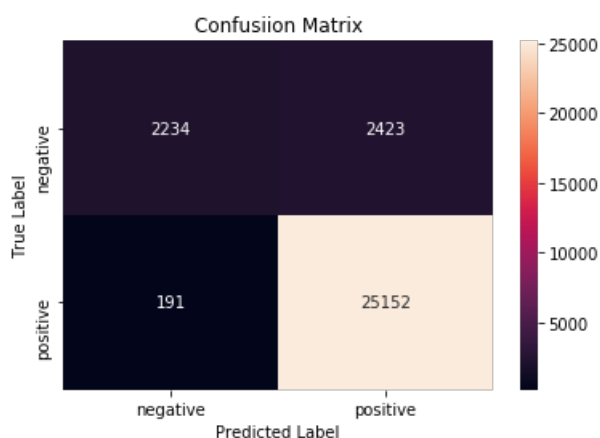Receiver Operating Characteristic

```
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

Out[73]:

```
array([[ 2234,  2423],
       [  191, 25152]], dtype=int64)
```

In [74]:

```
# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



Confusiion Matrix

In [75]:

```
# Top 10 important features of positive class
```

In [76]:

```
#https://stackoverflow.com/questions/11116697/how-to-get-most-informative-features-for-scikit-lear
n-classifiers
def show_most_informative_features(vectorizer, clf, n=10):
    feature_names = vectorizer.get_feature_names()
```

```
        coefs_with_fns = sorted(zip(clf.coef_[0], feature_names),reverse=True)[:n]
        top = zip(coefs_with_fns[:n], coefs_with_fns[:-(n + 1):-1])
        print("Positive")

print("_____
")
    for (coef_1, fn_1),(coef_2, fn_2) in top:
        print("\t%.4f\t%-15s" % (coef_1, fn_1))

show_most_informative_features(tf_idf_vect,clf)
```

```
Positive

_____
 3.6938 great
 2.7929 not disappointed
 2.6598 best
 2.5858 delicious
 2.3671 good
 2.2676 love
 2.1112 perfect
 1.9713 excellent
 1.8704 wonderful
 1.7743 loves
```

In [77]:

```
##Top 10 important features of negative class
```

In [78]:

```
#https://stackoverflow.com/questions/11116697/how-to-get-most-informative-features-for-scikit-lear
n-classifiers
def show_most_informative_features(vectorizer, clf, n=10):
    feature_names = vectorizer.get_feature_names()
    coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))[:n]
    top = zip(coefs_with_fns[:n], coefs_with_fns[:-(n + 1):-1])
    print("Negative")

print("_____
")
    for (coef_1, fn_1),(coef_2, fn_2) in top:
        print("\t%.4f\t%-15s" % (coef_2, fn_2))

show_most_informative_features(tf_idf_vect,clf)
```

```
Negative

_____
 -2.7094 return
 -2.7365 disappointing
 -3.0075 not
 -3.1061 not buy
 -3.1668 not recommend
 -3.2052 awful
 -3.2463 not worth
 -3.3000 terrible
 -3.7317 worst
 -4.4234 disappointed
```

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

In [79]:

```
# Please write all the code with proper documentation
```

In [49]:

```
## find hyperparameter using cross validation score and plot AUC
## find hyperparameter using cross validation score and plot AUC
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import TimeSeriesSplit
```

```
tscv = TimeSeriesSplit(n_splits=10)
from sklearn.model_selection import cross_val_score

alphas = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]

#empty lists that stores cv scores and training_scores
cv_scores = []
training_scores = []

#perform k fold cross validation
for alpha in alphas:
    clf = SGDClassifier(loss='hinge',alpha=alpha)
    scores = cross_val_score(clf, x_train_avgw2v, y_train, cv=10, scoring='roc_auc')
    scores_training = clf.fit(x_train_avgw2v, y_train).score(x_train_avgw2v, y_train)

    cv_scores.append(scores.mean())
    training_scores.append(scores_training)

#plot cross-validated score, training score vs alpha
plt.plot(x, training_scores, label='Train AUC')
plt.plot(x, cv_scores, label='CV AUC')
plt.legend()
plt.xlabel("Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```
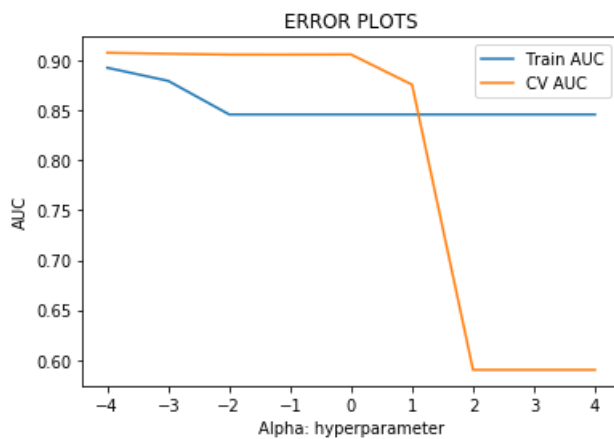


In [50]:

```
## find hyperparameter alpha using GridserachCV
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import SGDClassifier
tuned_parameters = [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l1','l2']}]
model = GridSearchCV(SGDClassifier(),tuned_parameters,cv=tscv,scoring='f1')
model.fit(x_train_avgw2v,y_train)
print(model.best_estimator_)
print(model.score(x_test_avgw2v,y_test))
print("Best HyperParameter: ",model.best_params_)
print("Best Accuracy: %.2f%%"%(model.best_score_*100))
```

```
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=None,
       n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
       power_t=0.5, random_state=None, shuffle=True, tol=None,
       validation_fraction=0.1, verbose=0, warm_start=False)
0.9338328508448279
Best HyperParameter:  {'alpha': 0.0001, 'penalty': 'l2'}
Best Accuracy: 93.75%
```

In [51]:

```
##Train the model with best hyperparameter
clf = SGDClassifier(loss='hinge',alpha=0.0001,penalty='l2')
clf.fit(x_train_avgw2v,y_train)
y_pred = clf.predict(x_test_avgw2v)
```
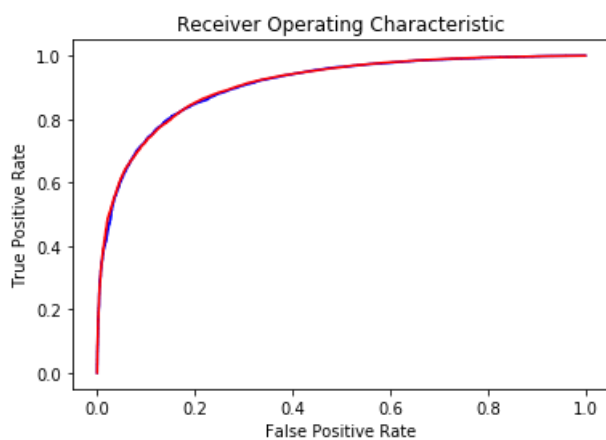
Observation: BY observing the hyperparameter is 0.0001

In [83]:

```
### ROC Curve using false positive rate versus true positive rate
```

In [52]:

```python
from sklearn.metrics import roc_curve,auc
probas_ = clf.fit(x_train_avgw2v, y_train).decision_function(x_test_avgw2v)
# Compute ROC curve and area the curve
fpr, tpr, thresholds = roc_curve(y_test, probas_)
probas_ = clf.fit(x_train_avgw2v, y_train).decision_function(x_train_avgw2v)
# Compute ROC curve and area the curve
fpr_, tpr_, thresholds = roc_curve(y_train, probas_)
plt.title('Receiver Operating Characteristic')
plt.plot(fpr, tpr, 'b')
plt.plot(fpr_, tpr_, 'r')
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```



In [53]:

```python
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

Out[53]:

```
array([[ 2197,  2639],
       [  702, 24462]])
```

In [54]:

```python
# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

positive    702         24462

negative    positive
Predicted Label

## [5.1.4] Applying Linear SVM on TFIDF W2V, <span style="color:red">SET 4</span>

In [87]:

```python
# Please write all the code with proper documentation
```

In [55]:

```python
## find hyperparameter using cross validation score and plot AUC
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import TimeSeriesSplit
tscv = TimeSeriesSplit(n_splits=10)
from sklearn.model_selection import cross_val_score

alphas = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]

#empty lists that stores cv scores and training_scores
cv_scores = []
training_scores = []

#perform k fold cross validation
for alpha in alphas:
    clf = SGDClassifier(loss='hinge',alpha=alpha)
    scores = cross_val_score(clf, x_train_tfidfw2v, y_train, cv=10, scoring='roc_auc')
    scores_training = clf.fit(x_train_tfidfw2v, y_train).score(x_train_tfidfw2v, y_train)

    cv_scores.append(scores.mean())
    training_scores.append(scores_training)

#plot cross-validated score, training score vs alpha
plt.plot(x, training_scores, label='Train AUC')
plt.plot(x, cv_scores, label='CV AUC')
plt.legend()
plt.xlabel("Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```
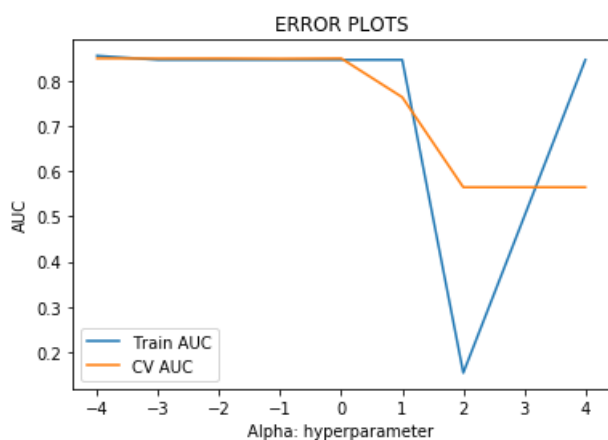


In [56]:

```python
## find hyperparameter alpha using GridserachCV
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import SGDClassifier
tuned_parameters = [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l1','l2']}]
model = GridSearchCV(SGDClassifier(),tuned_parameters,cv=tscv)
model.fit(x_train_tfidfw2v,y_train)
```

```
print(model.best_estimator_)
print(model.score(x_test_tfidfw2v,y_test))
print("Best HyperParameter: ",model.best_params_)
print("Best Accuracy: %.2f%%"%(model.best_score_*100))
```

```
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=None,
        n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l1',
        power_t=0.5, random_state=None, shuffle=True, tol=None,
        validation_fraction=0.1, verbose=0, warm_start=False)
0.8530333333333333
Best HyperParameter:  {'alpha': 0.0001, 'penalty': 'l1'}
Best Accuracy: 86.22%
```

Observation: BY observing the hyperparameter is 0.0001

In [57]:

```
##Train the model with best hyperparameter
clf = SGDClassifier(loss='hinge',alpha=0.0001,penalty='l1')
clf.fit(x_train_tfidfw2v,y_train)
y_pred = clf.predict(x_test_tfidfw2v)
```

In [96]:

```
### ROC Curve using false positive rate versus true positive rate
```

In [58]:

```
from sklearn.metrics import roc_curve,auc
probas_ = clf.fit(x_train_avgw2v, y_train).decision_function(x_test_avgw2v)
# Compute ROC curve and area the curve
fpr, tpr, thresholds = roc_curve(y_test, probas_)
probas_ = clf.fit(x_train_avgw2v, y_train).decision_function(x_train_avgw2v)
# Compute ROC curve and area the curve
fpr_, tpr_, thresholds = roc_curve(y_train, probas_)
plt.title('Receiver Operating Characteristic')
plt.plot(fpr, tpr, 'b')
plt.plot(fpr_, tpr_, 'r')
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```
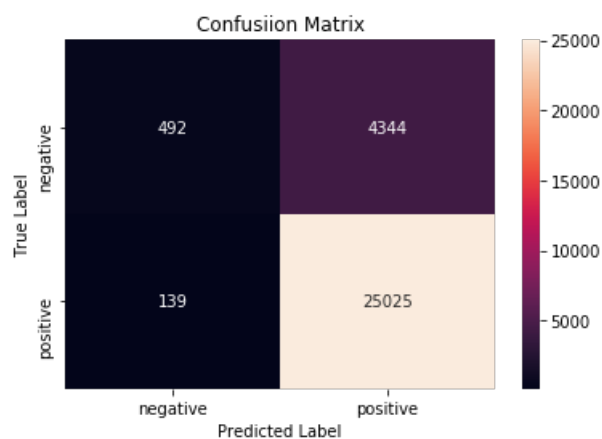


In [59]:

```
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

Out[59]:

```
array([[  492,  4344],
       [  139, 25025]])
```

```
# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



## [5.2] RBF SVM

In [87]:

```
x = final_20k['Cleaned_text']
x.size
```

Out[87]:

```
20000
```

In [88]:

```
y = final_20k['Score']
y.size
```

Out[88]:

```
20000
```

In [89]:

```
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=42)
```

### [5.2.1] Applying RBF SVM on BOW, SET 1

In [103]:

```
# Please write all the code with proper documentation
```

In [64]:

```
#BoW
count_vect = CountVectorizer(min_df = 10, max_features = 500) #in scikit-learn
x_train_bow = count_vect.fit_transform(x_train)
print("some feature names ", count_vect.get_feature_names()[:10])
```

```
print('='*50)

x_test_bow = count_vect.transform(x_test)
print("the type of count vectorizer ",type(x_test_bow))
print("the shape of out text BOW vectorizer ",x_test_bow.get_shape())
print("the number of unique words ", x_test_bow.get_shape()[1])
```

```
some feature names  ['able', 'absolutely', 'actually', 'add', 'added', 'aftertaste', 'ago',
'almost', 'also', 'alternative']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (6000, 500)
the number of unique words  500
```

In [65]:

```
from sklearn import preprocessing
x_train_bow = preprocessing.normalize(x_train_bow)
x_test_bow = preprocessing.normalize(x_test_bow)
```

In [66]:

```
from math import log
c_values = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]
x =[log(y,10) for y in c_values]
```

In [67]:

```
## find hyperparameter using cross validation score and plot AUC
from sklearn.svm import SVC
from sklearn.metrics import roc_auc_score
c_values = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]

#empty lists that stores cv scores and training_scores
cv_scores = []
training_scores = []

#perform k fold cross validation
for c in c_values:
    clf = SVC(C=c,kernel='rbf',probability=True)
    clf.fit(x_train_bow, y_train)
    y_train_pred = clf.predict_proba(x_train_bow)[:,1]
    y_test_pred =  clf.predict_proba(x_test_bow)[:,1]

    training_scores.append(roc_auc_score(y_train,y_train_pred))
    cv_scores.append(roc_auc_score(y_test, y_test_pred))


#plot cross-validated score, training score vs alpha
plt.plot(x, training_scores, label='Train AUC')
plt.plot(x, cv_scores, label='CV AUC')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```
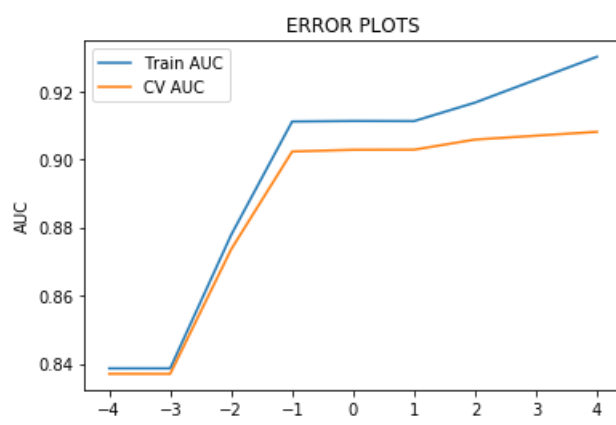
In [110]:

```python
## find hyperparameter alpha using GridserachCV
from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC
tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = GridSearchCV(SVC(),tuned_parameters,cv=tscv)
model.fit(x_train_bow,y_train)
print(model.best_estimator_)
print(model.score(x_test_bow,y_test))
print("Best HyperParameter: ",model.best_params_)
print("Best Accuracy: %.2f%%"%(model.best_score_*100))
```

```
SVC(C=10000, cache_size=200, class_weight=None, coef0=0.0,
  decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
  kernel='rbf', max_iter=-1, probability=False, random_state=None,
  shrinking=True, tol=0.001, verbose=False)
0.893
Best HyperParameter:  {'C': 10000}
Best Accuracy: 87.19%
```

In [138]:

```python
##Train the model with best hyperparameter
clf = SVC(C=10000,kernel='rbf',probability=True)
clf.fit(x_train_bow,y_train)
y_pred = clf.predict(x_test_bow)
```

Observation: BY observing the hyperparameter is 10000

In [139]:

```python
### ROC Curve using false positive rate versus true positive rate
```

In [68]:

```python
# https://scikit-
learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
clf = SVC(C=10000,kernel='rbf',probability=True)
clf.fit(x_train_bow, y_train)
y_pred=clf.predict(x_test_bow)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(x_train_bow)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(x_test_bow)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel('False positive rate')
plt.ylabel('True positive rate')
plt.title("ERROR PLOTS")
plt.show()

print("="*100)

from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, clf.predict(x_train_bow)))
print("Test confusion matrix")
print(confusion_matrix(y_test, clf.predict(x_test_bow)))
```
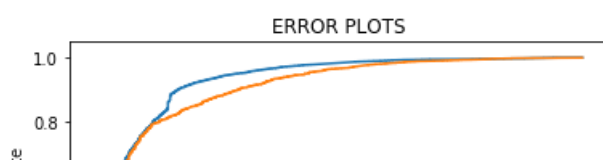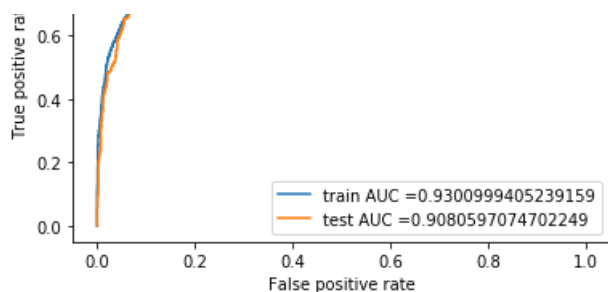
train AUC =0.9300999405239159
test AUC =0.9080597074702249

==============================================================================================

```
Train confusion matrix
[[ 1352   871]
 [  335 11442]]
Test confusion matrix
[[ 471   430]
 [ 189 4910]]
```
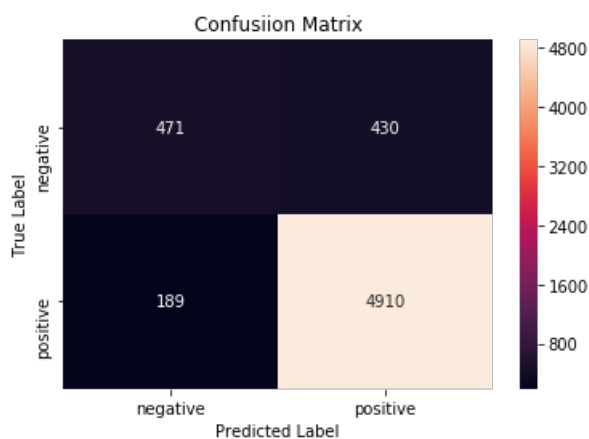
In [69]:

```python
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

Out[69]:

```
array([[ 471,  430],
       [ 189, 4910]])
```

In [70]:

```python
# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



## [5.2.2] Applying RBF SVM on TFIDF, SET 2

In [ ]:

```python
# Please write all the code with proper documentation
```

In [72]:

```
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df = 10, max_features = 500)
x_train_tfidf = 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_test_tfidf = tf_idf_vect.transform(x_test)
print("the type of count vectorizer ",type(x_test_tfidf))
print("the shape of out text TFIDF vectorizer ",x_test_tfidf.get_shape())
print("the number of unique words including both unigrams and bigrams ", x_test_tfidf.get_shape()[
1])
```

```
some sample features(unique words in the corpus) ['able', 'absolutely', 'actually', 'add',
'added', 'ago', 'almost', 'also', 'alternative', 'although']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer  (6000, 500)
the number of unique words including both unigrams and bigrams  500
```

In [73]:

```python
from sklearn import preprocessing
x_train_tfidf = preprocessing.normalize(x_train_tfidf)
x_test_tfidf = preprocessing.normalize(x_test_tfidf)
```
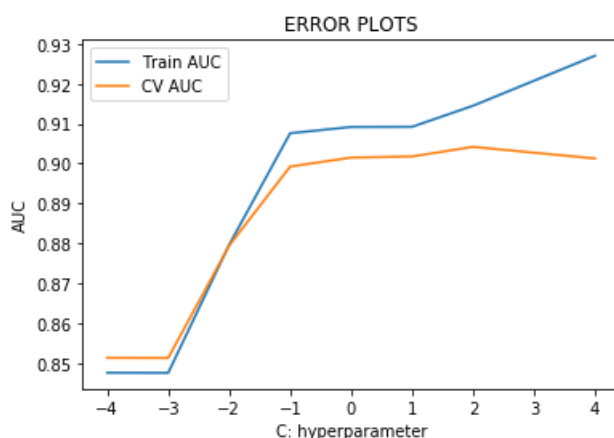
In [62]:

```python
## find hyperparameter using cross validation score and plot AUC
from sklearn.svm import SVC
from sklearn.metrics import roc_auc_score
c_values = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]

#empty lists that stores cv scores and training_scores
cv_scores = []
training_scores = []

#perform k fold cross validation
for c in c_values:
    clf = SVC(C=c,kernel='rbf',probability=True)
    clf.fit(x_train_tfidf, y_train)
    y_train_pred = clf.predict_proba(x_train_tfidf)[:,1]
    y_test_pred =  clf.predict_proba(x_test_tfidf)[:,1]

    training_scores.append(roc_auc_score(y_train,y_train_pred))
    cv_scores.append(roc_auc_score(y_test, y_test_pred))


#plot cross-validated score, training score vs alpha
plt.plot(x, training_scores, label='Train AUC')
plt.plot(x, cv_scores, label='CV AUC')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

```python
## find hyperparameter alpha using GridserachCV
from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC
tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = GridSearchCV(SVC(),tuned_parameters,cv=tscv)
model.fit(x_train_tfidf,y_train)
print(model.best_estimator_)
print(model.score(x_test_tfidf,y_test))
print("Best HyperParameter: ",model.best_params_)
print("Best Accuracy: %.2f%%"%(model.best_score_*100))
```

```
SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
  decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
  kernel='rbf', max_iter=-1, probability=False, random_state=None,
  shrinking=True, tol=0.001, verbose=False)
0.8905
Best HyperParameter:  {'C': 100}
Best Accuracy: 87.26%
```

Observation: BY observing the hyperparameter is 100
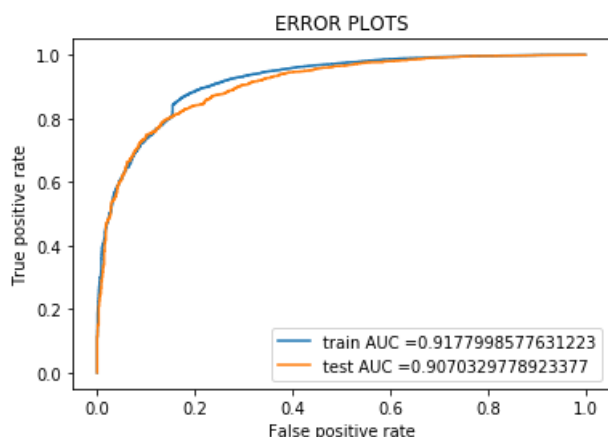
In [144]:

```python
## ROC Curve
```

In [74]:

```python
# https://scikit-
learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
clf = SVC(C=100,kernel='rbf',probability=True)
clf.fit(x_train_tfidf, y_train)
y_pred=clf.predict(x_test_tfidf)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(x_train_tfidf)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(x_test_tfidf)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel('False positive rate')
plt.ylabel('True positive rate')
plt.title("ERROR PLOTS")
plt.show()

print("="*100)

from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, clf.predict(x_train_tfidf)))
print("Test confusion matrix")
print(confusion_matrix(y_test, clf.predict(x_test_tfidf)))
```

```
==================================================================================================
Train confusion matrix
[[  915  1308]
 [  165 11612]]
Test confusion matrix
[[ 327  574]
 [  82 5017]]
```
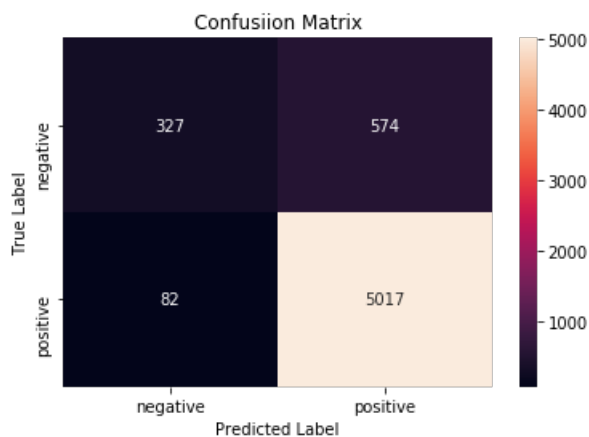
◀ ▊ ▶

In [75]:

```python
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

Out[75]:

```
array([[ 327,  574],
       [  82, 5017]])
```

In [76]:

```python
# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



## [5.2.3] Applying RBF SVM on AVG W2V, SET 3

In [100]:

```python
x = final_20k['Cleaned_text']
x.size
```

Out[100]:

```
20000
```

In [101]:

```python
y = final_20k['Score']
y.size
```

Out[101]:

```
20000
```

```python
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=42)
```

In [103]:

```python
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance_train=[]
for sentance in x_train:
    list_of_sentance_train.append(sentance.split())
```

In [104]:

```python
# Test your own Word2Vec model using your own text corpus
i=0
list_of_sentance_test=[]
for sentance in x_test:
    list_of_sentance_test.append(sentance.split())
```

In [105]:

```python
# Using Google News Word2Vectors

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


# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17SRFAzZPY
# you can comment this whole cell
# or change these varible according to your need

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True


if want_to_train_w2v:
    # min_count = 5 considers only words that occured atleast 5 times
    w2v_model=Word2Vec(list_of_sentance_train,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))

elif want_to_use_google_w2v and is_your_ram_gt_16g:
    if os.path.isfile('GoogleNews-vectors-negative300.bin'):
        w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=Tr
ue)
        print(w2v_model.wv.most_similar('great'))
        print(w2v_model.wv.most_similar('worst'))
    else:
        print("you don't have gogole's word2vec file, keep want_to_train_w2v = True, to train your
own w2v ")
```

```
[('good', 0.8624511361122131), ('excellent', 0.7976897358894348), ('wonderful',
0.7802258133888245), ('awesome', 0.7348822355270386), ('amazing', 0.710780143737793),
('fantastic', 0.7007759213447571), ('super', 0.6934992671012878), ('delicious',
0.6788982152938843), ('perfect', 0.6706241369247437), ('well', 0.6534604430198669)]
==================================================
[('varieties', 0.942334771156311), ('smoothest', 0.9324260950088501), ('absolute',
0.9312138557434082), ('ive', 0.9207708239555359), ('lundberg', 0.9155126214027405), ('zevia',
0.914898574352644), ('tastiest', 0.914873980293274), ('enjoyed', 0.9080132842063904),
('hottest', 0.9066132307052612), ('mum', 0.9039670825004578)]
```

```
w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])
```

```
number of words that occured minimum 5 times  7289
sample words  ['trio', 'bars', 'loved', 'disappointing', 'stale', 'sad', 'past', 'ordering', 'anot
her', 'company', 'treat', 'poodle', 'discerning', 'comes', 'treats', 'loves', 'product',
'apparently', 'taste', 'also', 'chew', 'factor', 'husky', 'gnaw', 'spends', 'hours', 'day', 'chewi
ng', 'bone', 'lasts', 'long', 'time', 'much', 'longer', 'others', 'tried', 'think', 'well', 'spent
', 'keeps', 'busy', 'happy', 'made', 'switch', 'iams', 'food', 'yrs', 'ago', 'actually', 'not']
```

```python
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentance_train): # 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_train.append(sent_vec)
print(len(sent_vectors_train))
print(len(sent_vectors_train[0]))
```

```
100%|██████████| 14000/14000 [00:23<00:00, 585.32it/s]
```

```
14000
50
```

```python
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentance_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)
print(len(sent_vectors_test))
print(len(sent_vectors_test[0]))
```

```
100%|██████████| 6000/6000 [00:10<00:00, 552.80it/s]
```

```
6000
50
```

```
x_train_avgw2v=sent_vectors_train
x_test_avgw2v=sent_vectors_test
```

In [114]:

```python
from sklearn import preprocessing
x_train_avgw2v = preprocessing.normalize(x_train_avgw2v)
x_test_avgw2v = preprocessing.normalize(x_test_avgw2v)
```

In [115]:

```python
from math import log
c_values = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]
x =[log(y,10) for y in c_values]
```
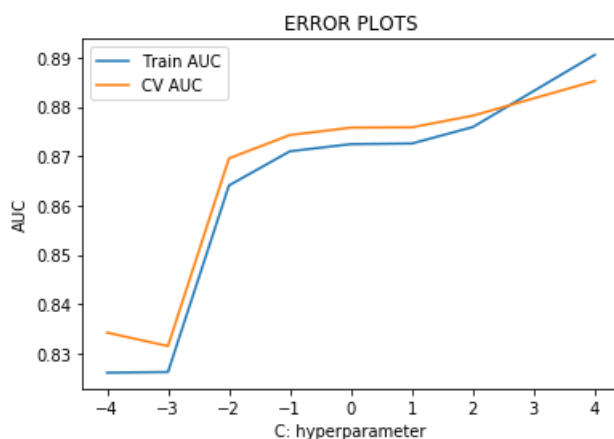
In [116]:

```python
## find hyperparameter using cross validation score and plot AUC
from sklearn.svm import SVC
from sklearn.metrics import roc_auc_score
c_values = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]

#empty lists that stores cv scores and training_scores
cv_scores = []
training_scores = []

#perform k fold cross validation
for c in c_values:
    clf = SVC(C=c,kernel='rbf',probability=True)
    clf.fit(x_train_avgw2v, y_train)
    y_train_pred = clf.predict_proba(x_train_avgw2v)[:,1]
    y_test_pred =  clf.predict_proba(x_test_avgw2v)[:,1]

    training_scores.append(roc_auc_score(y_train,y_train_pred))
    cv_scores.append(roc_auc_score(y_test, y_test_pred))


#plot cross-validated score, training score vs alpha
plt.plot(x, training_scores, label='Train AUC')
plt.plot(x, cv_scores, label='CV AUC')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [75]:

```python
## find hyperparameter alpha using GridserachCV
from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC
tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = GridSearchCV(SVC(),tuned_parameters,cv=tscv)
model.fit(x_train_avgw2v,y_train)
print(model.best_estimator_)
print(model.score(x_test_avgw2v,y_test))
print("Best HyperParameter: ",model.best_params_)
print("Best Accuracy: %.2f%%"%(model.best_score_*100))
```

```
SVC(C=10000, cache_size=200, class_weight=None, coef0=0.0,
  decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
  kernel='rbf', max_iter=-1, probability=False, random_state=None,
  shrinking=True, tol=0.001, verbose=False)
0.8838333333333334
Best HyperParameter:  {'C': 10000}
Best Accuracy: 87.32%
```

Observation: BY observing the hyperparameter is
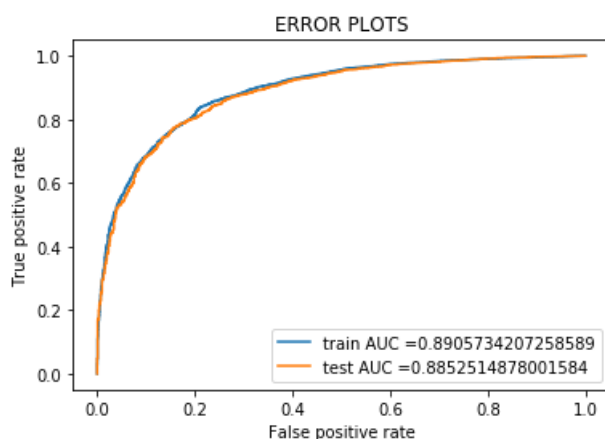
In [ ]:

```python
## ROC Curve
```

In [117]:

```python
# https://scikit-
learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
clf = SVC(C=10000,kernel='rbf',probability=True)
clf.fit(x_train_avgw2v, y_train)
y_pred=clf.predict(x_test_avgw2v)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(x_train_avgw2v)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(x_test_avgw2v)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel('False positive rate')
plt.ylabel('True positive rate')
plt.title("ERROR PLOTS")
plt.show()

print("="*100)

from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, clf.predict(x_train_avgw2v)))
print("Test confusion matrix")
print(confusion_matrix(y_test, clf.predict(x_test_avgw2v)))
```



```
====================================================================================================

Train confusion matrix
[[  783  1440]
 [  251 11526]]
Test confusion matrix
[[ 295  606]
 [ 107 4992]]
```
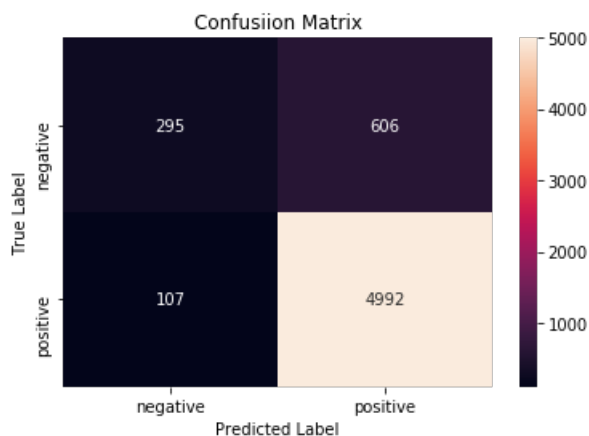
```
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

Out[118]:

```
array([[ 295,  606],
       [ 107, 4992]])
```

In [119]:

```
# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



## [5.2.4] Applying RBF SVM on TFIDF W2V, SET 4

In [120]:

```
x = final_20k['Cleaned_text']
x.size
```

Out[120]:

```
20000
```

In [121]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer(min_df = 10, max_features = 500)
tf_idf_matrix = model.fit_transform(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_)))
```

In [122]:

```
# 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_train = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentance_train): # 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_train.append(sent_vec)
        row += 1
```

```
100%|███████████| 14000/14000 [00:29<00:00, 470.35it/s]
```

In [123]:

```
tfidf_sent_vectors_test = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentance_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%|███████████| 6000/6000 [00:13<00:00, 456.99it/s]
```

In [124]:

```
x_train_tfidfw2v = tfidf_sent_vectors_train
x_test_tfidfw2v = tfidf_sent_vectors_test
```

In [125]:

```
from sklearn import preprocessing
x_train_tfidfw2v = preprocessing.normalize(x_train_tfidfw2v)
x_test_tfidfw2v = preprocessing.normalize(x_test_tfidfw2v)
```

In [126]:

```
from math import log
c_values = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]
x =[log(y,10) for y in c_values]
```

In [127]:

```
## find hyperparameter using cross validation score and plot AUC
from sklearn.svm import SVC
from sklearn.metrics import roc_auc_score
c_values = [10**-4,10**-3, 10**-2,10**-1, 10**0,10, 10**2, 10**4]

#empty lists that stores cv scores and training_scores
cv_scores = []
training_scores = []

#perform k fold cross validation
for c in c_values:
```
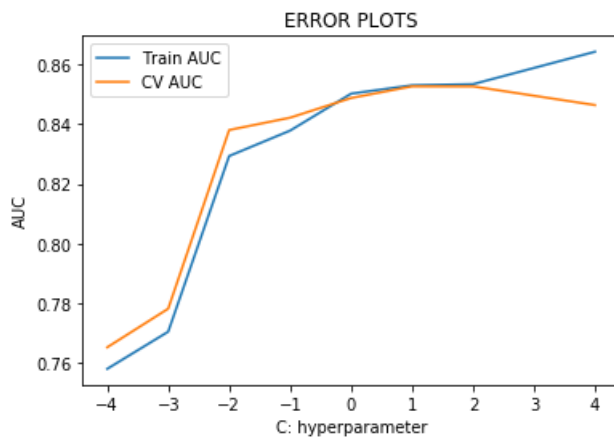
```
    clf = SVC(C=c,kernel='rbf',probability=True)
    clf.fit(x_train_tfidfw2v, y_train)
    y_train_pred = clf.predict_proba(x_train_tfidfw2v)[:,1]
    y_test_pred =  clf.predict_proba(x_test_tfidfw2v)[:,1]

    training_scores.append(roc_auc_score(y_train,y_train_pred))
    cv_scores.append(roc_auc_score(y_test, y_test_pred))


#plot cross-validated score, training score vs alpha
plt.plot(x, training_scores, label='Train AUC')
plt.plot(x, cv_scores, label='CV AUC')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [128]:

```
## find hyperparameter alpha using GridserachCV
from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC
tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = GridSearchCV(SVC(),tuned_parameters,cv=tscv)
model.fit(x_train_tfidfw2v,y_train)
print(model.best_estimator_)
print(model.score(x_test_tfidfw2v,y_test))
print("Best HyperParameter: ",model.best_params_)
print("Best Accuracy: %.2f%%"%(model.best_score_*100))
```

```
SVC(C=10000, cache_size=200, class_weight=None, coef0=0.0,
  decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
  kernel='rbf', max_iter=-1, probability=False, random_state=None,
  shrinking=True, tol=0.001, verbose=False)
0.8676666666666667
Best HyperParameter:  {'C': 10000}
Best Accuracy: 85.85%
```

Observation: BY observing the hyperparameter is 10000

In [ ]:

```
### ROC Curve using false positive rate versus true positive rate
```

In [129]:

```
# https://scikit-
learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
clf = SVC(C=10000,kernel='rbf',probability=True)
clf.fit(x_train_tfidfw2v, y_train)
y_pred=clf.predict(x_test_tfidfw2v)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
```
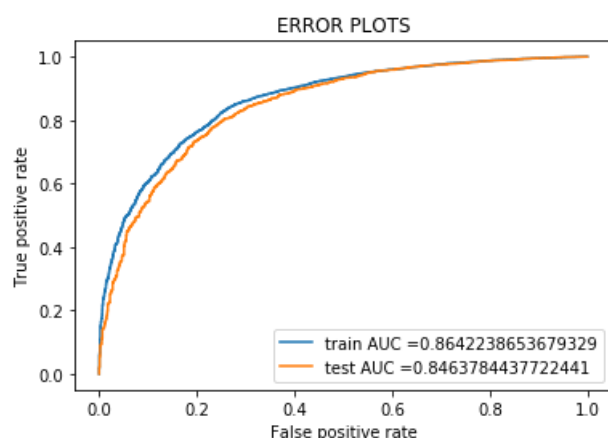
```
train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(x_train_tfidfw2v)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(x_test_tfidfw2v)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel('False positive rate')
plt.ylabel('True positive rate')
plt.title("ERROR PLOTS")
plt.show()

print("="*100)

from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, clf.predict(x_train_tfidfw2v)))
print("Test confusion matrix")
print(confusion_matrix(y_test, clf.predict(x_test_tfidfw2v)))
```



```
====================================================================================================

Train confusion matrix
[[  440  1783]
 [  141 11636]]
Test confusion matrix
[[ 165   736]
 [  58 5041]]
```

In [130]:

```
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

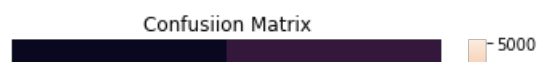Out[130]:

```
array([[ 165,   736],
       [  58, 5041]])
```
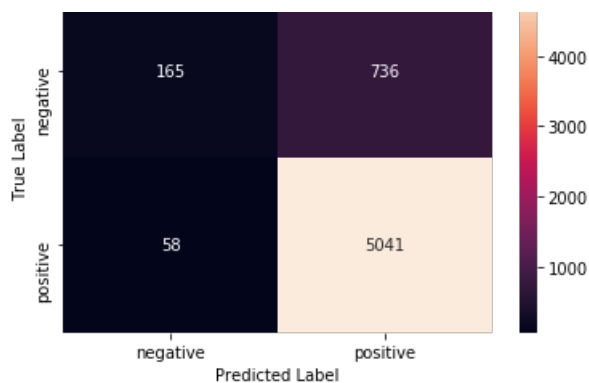
In [131]:

```
# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

# [6] Conclusions

## linear SVM

In [132]:

```python
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Vectorizer","Kernel","Hyperparameter","penalty","AUC"]

x.add_row(["BOW", "linear","0.0001","l2","93.75%"])
x.add_row(["TFIDF","linear","0.0001","l2","95.3%"])
x.add_row(["Avgw2v","linear", "0.0001","l2","93.75%"])
x.add_row(["TFIDF W2V", "linear","0.0001","l1","86.22%"])


print(x)
```

```
+------------+--------+----------------+---------+--------+
| Vectorizer | Kernel | Hyperparameter | penalty |  AUC   |
+------------+--------+----------------+---------+--------+
|    BOW     | linear |     0.0001     |    l2   | 93.75% |
|   TFIDF    | linear |     0.0001     |    l2   | 95.3%  |
|   Avgw2v   | linear |     0.0001     |    l2   | 93.75% |
| TFIDF W2V  | linear |     0.0001     |    l1   | 86.22% |
+------------+--------+----------------+---------+--------+
```

### rbf SVM

In [133]:

```python
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Vectorizer","Kernel","Hyperparameter","AUC"]


x.add_row(["BOW","rbf", "10000","90%"])
x.add_row(["TFIDF","rbf", "100","90%"])
x.add_row(["AVGW2V","rbf", "10000","88.5%"])
x.add_row(["TFIDF W2V", "rbf","10000","84.6%"])

print(x)
```

```
+------------+--------+----------------+-------+
| Vectorizer | Kernel | Hyperparameter |  AUC  |
+------------+--------+----------------+-------+
|    BOW     |  rbf   |     10000      |  90%  |
|   TFIDF    |  rbf   |      100       |  90%  |
|   AVGW2V   |  rbf   |     10000      | 88.5% |
| TFIDF W2V  |  rbf   |     10000      | 84.6% |
+------------+--------+----------------+-------+
```

1) SVM with rbf kernel requires more computation time. 2) SVM with linear kernel is faster. 3) TFIDF with linear SVM getting more

accuracy.