

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

[1]. Reading Data

[1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
```

```

import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

```

```

C:\ProgramData\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

```

```

In [2]: # using SQLite Table to read data.
con = sqlite3.connect('C:/Users/Excel/Desktop/vins/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 50
0000 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 20000""", 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 (20000, 10)

Out[2]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1

```
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(*)
--	--------	-----------	-------------	------	-------	------	----------

	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
--	--------	-----------	-------------	------	-------	------	-------

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

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

Out[7]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpfuln
--	----	-----------	--------	-------------	----------------------	----------

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpfuln
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 it can be seen above that same user has multiple reviews with same values for HelpfulnessNumerator, HelpfulnessDenominator, Score, Time, Summary and Text and on doing analysis it was found that

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

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

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

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

```
In [8]: #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 [9]: #Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time",
"Text"}, keep='first', inplace=False)
final.shape
```

```
Out[9]: (19354, 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]: 96.77
```

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 [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	Helpfuln
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2

```
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()
```

```
(19354, 10)
```

```
Out[13]: 1    16339  
        0     3015  
        Name: Score, dtype: int64
```

```
In [14]: final["Time"] = pd.to_datetime(final["Time"], unit = "s")  
        final = final.sort_values(by = "Time")
```

[3] Preprocessing

[3.1]. Preprocessing Review Text

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

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

1. Begin by removing the html tags
2. Remove any punctuations or limited set of special characters like , or . or # etc.
3. Check if the word is made up of english letters and is not alpha-numeric
4. Check to see if the length of the word is greater than 2 (as it was researched that there is no adjective in 2-letters)
5. Convert the word to lowercase
6. Remove Stopwords
7. Finally Snowball Stemming the word (it was observed to be better than Porter Stemming)

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

```
In [15]: # 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 was a really good idea and the final product is outstanding. I use the decals on my car window and everybody asks where i bought the decals i made. Two thumbs up!

=====

THIS COFFEE IS REALLY DELICIOUS.A COOL LATIN FLAVOR.EXCELLENT.5 STARS I GIVE TO THIS COFFEE.I HOPE AMAZON NEVER GET RID OF IT BECAUSE THIS COFFEE IS REALLY HARD TO FIND IN MY LOCAL SUPERMARKETS.AND I HPE ALWAYS THEY GOT THE SAME LOW PRICE FOR 6 PACKETS.THE LESS I GET THIS COFFEE BAGS IS \$2.49 PER PACKET.GOOD DEAL,AMAZON.HURRAY FOR YOU.

=====

I use this product frequently. Like most tofu, you need to press it (between paper towels with a heavy skillet on top works fine) to get the excess water out. The more water you can remove, the less you'll have to deal with sloppiness when cooking, something the previous reviewer did not like.

As for a blank slate on taste, that's exactly what makes tofu great. It will take on any flavor you want to impart. Our most common marinade is placing the pressed, cubed tofu in a Ziploc with a mix of soy, honey, and lemon or lime juice. A quick search online for a marinade will give you ideas (many also include minced garlic, ginger, etc). You can also purchase a premade marinade, like a teriyaki sauce. The more liquid the marinade (and the more water removed from pressing), the better it will penetrate the tofu. Even a half hour in the bag works fine.

We saute it right in the pan with the stir fry veggies (add tofu last--it just needs to get warm; or, for a fir

mer style, you can bake or "brown" it separately first) and pour in the remaining marinade as the final sauce (you can thicken with cornstarch if desired). Handle the tofu cubes gently as they are not firm like most meats, but broken up pieces taste just fine too.

Marinated tofu cubes also do well on kebabs with veggies on the grill.

Plain tofu "creams" well and is often good in a dish that requires thickness. We've made chocolate mousse with it as well as scrambling it like eggs.

We have great success breading pressed cubes of this tofu just like you would the chicken in a General Tso recipe. General Tso's Tofu is AWESOME. The flavor comes from the sauce so you don't need to marinate first.

Some people also freeze tofu--generally remove from packaging and press first--to give it a firmer consistency.

Tofu is all about how you prepare it. Most people will NOT like it plain, so don't expect to just dig a fork into it.

This particular product is great because of its all-purpose nature (firm is definitely better than soft for saute, etc) and its shelf life. I would have given it five stars but the must-be-refrigerated versions packaged in fluid do seem to be firmer. They don't last as long in the house, though, so this is a perfect choice to keep in the pantry (usually dated a few to several months out).

=====

I ordered three different kinds of Yogi tea based on the great reviews -- lemon ginger, stomach ease, and green tea super antioxidant. Unfortunately, they all taste exactly the same -- like licorice. If you like licorice flavor, you'll like this tea. If you're looking for something that actually tastes like lemon and ginger, look elsewhere.

=====

```
In [16]: # 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 was a really good idea and the final product is outstanding. I use the decals on my car window and everybody asks where i bought the decal

s i made. Two thumbs up!

```
In [17]: # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how
-to-remove-all-tags-from-an-element
from bs4 import BeautifulSoup

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 was a really good idea and the final product is outstanding. I use the decals on my car window and everybody asks where i bought the decals i made. Two thumbs up!

=====

THIS COFFEE IS REALLY DELICIOUS.A COOL LATIN FLAVOR.EXCELLENT.5 STARS I GIVE TO THIS COFFEE.I HOPE AMAZON NEVER GET RID OF IT BECAUSE THIS COFFEE IS REALLY HARD TO FIND IN MY LOCAL SUPERMARKETS.AND I HPE ALWAYS THE Y GOT THE SAME LOW PRICE FOR 6 PACKETS.THE LESS I GET THIS COFFEE BAGS IS \$2.49 PER PACKET.GOOD DEAL,AMAZON.HURRAY FOR YOU.

=====

I use this product frequently. Like most tofu, you need to press it (between paper towels with a heavy skillet on top works fine) to get the excess water out. The more water you can remove, the less you'll have to deal with sloppiness when cooking, something the previous reviewer d

id not like. As for a blank slate on taste, that's exactly what makes to fu great. It will take on any flavor you want to impart. Our most com mon marinade is placing the pressed, cubed tofu in a Ziploc with a mix of soy, honey, and lemon or lime juice. A quick search online for a ma rinade will give you ideas (many also include minced garlic, ginger, et c). You can also purchase a premade marinade, like a teriyaki sauce. The more liquid the marinade (and the more water removed from pressin g), the better it will penetrate the tofu. Even a half hour in the bag works fine. We saute it right in the pan with the stir fry veggies (add tofu last--it just needs to get warm; or, for a firmer style, you can b ake or "brown" it separately first) and pour in the remaining marinade as the final sauce (you can thicken with cornstarch if desired). Handl e the tofu cubes gently as they are not firm like most meats, but broke n up pieces taste just fine too. Marinated tofu cubes also do well on ke babs with veggies on the grill. Plain tofu "creams" well and is often go od in a dish that requires thickness. We've made chocolate mousse with it as well as scrambling it like eggs. We have great success breading pr essed cubes of this tofu just like you would the chicken in a General T so recipe. General Tso's Tofu is AWESOME. The flavor comes from the s auce so you don't need to marinate first. Some people also freeze tofu-- generally remove from packaging and press first--to give it a firmer co nsistency. Tofu is all about how you prepare it. Most people will NOT l ike it plain, so don't expect to just dig a fork into it. This particlua r product is great because of its all-purpose nature (firm is definitel y better than soft for saute, etc) and its shelf life. I would have gi ven it five stars but the must-be-refrigerated versions packaged in flu id do seem to be firmer. They don't last as long in the house, though, so this is a perfect choice to keep in the pantry (usually dated a few to several months out).

=====

I ordered three different kinds of Yogi tea based on the great reviews -- lemon ginger, stomach ease, and green tea super antioxidant. Unfortu nately, they all taste exactly the same -- like licorice. If you like l icorice flavor, you'll like this tea. If you're looking for something t hat actually tastes like lemon and ginger, look elsewhere.

In [18]: `# 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 [19]: sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

I use this product frequently. Like most tofu, you need to press it (between paper towels with a heavy skillet on top works fine) to get the excess water out. The more water you can remove, the less you will have to deal with sloppiness when cooking, something the previous reviewer did not like.

As for a blank slate on taste, that is exactly what makes tofu great. It will take on any flavor you want to impart. Our most common marinade is placing the pressed, cubed tofu in a Ziploc with a mix of soy, honey, and lemon or lime juice. A quick search online for a marinade will give you ideas (many also include minced garlic, ginger, etc). You can also purchase a premade marinade, like a teriyaki sauce. The more liquid the marinade (and the more water removed from pressing), the better it will penetrate the tofu. Even a half hour in the bag works fine.

We saute it right in the pan with the stir fry veggies (add tofu last--it just needs to get warm; or, for a firmer style, you can bake or "brown" it separately first) and pour in the remaining marinade as the final sauce (you can thicken with cornstarch if desired). Handle the tofu cubes gently as they are not firm like most meats, but broken up pieces taste just fine too.

Marinated tofu cubes also do well on kebabs with veggies on the grill.

/>Plain tofu "creams" well and is often good in a dish that requires thickiness. We have made chocolate mousse with it as well as scrambling it like eggs.

We have great success breading pressed cubes of this tofu just like you would the chicken in a General Tso recipe. General Tso is Tofu is AWESOME. The flavor comes from the sauce so you do not need to marinate first.

Some people also freeze tofu--generally remove from packaging and press first--to give it a firmer consistency.

Tofu is all about how you prepare it. Most people will NOT like it plain, so do not expect to just dig a fork into it.

This particular product is great because of its all-purpose nature (firm is definitely better than soft for saute, etc) and its shelf life. I would have given it five stars but the must-be-refrigerated versions packaged in fluid do seem to be firmer. They do not last as long in the house, though, so this is a perfect choice to keep in the pantry (usually dated a few to several months out).

=====

```
In [20]: #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 was a really good idea and the final product is outstanding. I use the decals on my car window and everybody asks where i bought the decals i made. Two thumbs up!

```
In [21]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

I use this product frequently Like most tofu you need to press it between paper towels with a heavy skillet on top works fine to get the excess water out The more water you can remove the less you will have to deal with sloppiness when cooking something the previous reviewer did not like br br As for a blank slate on taste that is exactly what makes tofu great It will take on any flavor you want to impart Our most common marinade is placing the pressed cubed tofu in a Ziploc with a mix of soy honey and lemon or lime juice A quick search online for a marinade will give you ideas many also include minced garlic ginger etc You can also

purchase a premade marinade like a teriyaki sauce The more liquid the m
arinade and the more water removed from pressing the better it will pen
etrate the tofu Even a half hour in the bag works fine br br We saute i
t right in the pan with the stir fry veggies add tofu last it just need
s to get warm or for a firmer style you can bake or brown it separately
first and pour in the remaining marinade as the final sauce you can thi
cken with cornstarch if desired Handle the tofu cubes gently as they ar
e not firm like most meats but broken up pieces taste just fine too br
br Marinated tofu cubes also do well on kebabs with veggies on the gril
l br br Plain tofu creams well and is often good in a dish that require
s thickness We have made chocolate mousse with it as well as scrambling
it like eggs br br We have great success breading pressed cubes of this
tofu just like you would the chicken in a General Tso recipe General Ts
o is Tofu is AWESOME The flavor comes from the sauce so you do not need
to marinate first br br Some people also freeze tofu generally remove f
rom packaging and press first to give it a firmer consistency br br Tof
u is all about how you prepare it Most people will NOT like it plain so
do not expect to just dig a fork into it br br This particluar product
is great because of its all purpose nature firm is definitely better th
an soft for saute etc and its shelf life I would have given it five sta
rs but the must be refrigerated versions packaged in fluid do seem to b
e firmer They do not last as long in the house though so this is a perf
ect choice to keep in the pantry usually dated a few to several months
out

```
In [22]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'no
t'
# <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', 'o
urs', 'ourselves', 'you', "you're", "you've",\
               "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselve
s', 'he', 'him', 'his', 'himself', \
               'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'it
s', 'itself', 'they', 'them', 'their',\
```

```

        'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'th
is', 'that', "that'll", 'these', 'those', \
        'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'h
ave', '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', 'h
ow', 'all', 'any', 'both', 'each', 'few', 'more', \
        'most', 'other', 'some', 'such', 'only', 'own', 'same', 's
o', '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', "doesn't", 'hadn', \
        "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "is
n'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 [23]: # Combining all the above students
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%|██| 19354/19354 [00:10<00:00, 1788.29it/s]

In [25]: `preprocessed_reviews[1500]`

Out[25]: 'use product frequently like tofu need press paper towels heavy skillet top works fine get excess water water remove less deal sloppiness cooking something previous reviewer not like blank slate taste exactly makes tofu great take flavor want impart common marinade placing pressed cube d tofu ziploc mix soy honey lemon lime juice quick search online marinade give ideas many also include minced garlic ginger etc also purchase premade marinade like teriyaki sauce liquid marinade water removed pressing better penetrate tofu even half hour bag works fine saute right pan stir fry veggies add tofu last needs get warm firmer style bake brown separately first pour remaining marinade final sauce thicken cornstarch desired handle tofu cubes gently not firm like meats broken pieces taste fine marinated tofu cubes also well kebabs veggies grill plain tofu creams well often good dish requires thickness made chocolate mousse well scrambling like eggs great success breading pressed cubes tofu like would chicken general tso recipe general tso tofu awesome flavor comes sauce not need marinate first people also freeze tofu generally remove packaging press first give firmer consistency tofu prepare people not like plain not expect dig fork particluar product great purpose nature firm definitely better soft saute etc shelf life would given five stars must refrigerated versions packaged fluid seem firmer not last long house though perfect choice keep pantry usually dated several months'

[3.2] Preprocessing Review Summary

In [26]: `## Similarly you can do preprocessing for review summary also.`

[4] Featurization

[4.1] BAG OF WORDS(RBF KERNEL)

```
In [28]: X=preprocessed_reviews  
Y=final["Score"]
```

```
In [29]: from sklearn.model_selection import train_test_split  
  
X_train,X_test,y_train,y_test=train_test_split(X,Y,test_size=0.33,shuffle=False)  
X_train,X_cv,y_train,y_cv=train_test_split(X_train,y_train,test_size=0.33,shuffle=False)  
print(np.shape(X_train),y_train.shape)  
print(np.shape(X_cv),y_cv.shape)  
print(np.shape(X_test),y_test.shape)  
  
(8687,) (8687,)  
(4280,) (4280,)  
(6387,) (6387,)
```

```
In [31]: from sklearn.feature_extraction.text import CountVectorizer  
vectorizer=CountVectorizer(min_df=10,max_features=500,ngram_range=(1,2))  
vectorizer.fit(X_train)  
  
X_train_bow=vectorizer.transform(X_train)  
X_cv_bow=vectorizer.transform(X_cv)  
X_test_bow=vectorizer.transform(X_test)  
  
print(np.shape(X_train_bow),np.shape(X_cv_bow),np.shape(X_test_bow))  
print(y_train.shape,y_cv.shape,y_test.shape)  
  
(8687, 500) (4280, 500) (6387, 500)  
(8687,) (4280,) (6387,)
```

```
In [33]: from sklearn.preprocessing import StandardScaler  
standardised=StandardScaler(with_mean=False)  
X_train_bow=standardised.fit_transform(X_train_bow)  
X_cv_bow=standardised.transform(X_cv_bow)  
X_test_bow=standardised.transform(X_test_bow)
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:
475: DataConversionWarning: Data with input dtype int64 was converted t
o float64 by StandardScaler.
warnings.warn(msg, DataConversionWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:
475: DataConversionWarning: Data with input dtype int64 was converted t
o float64 by StandardScaler.
warnings.warn(msg, DataConversionWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:
475: DataConversionWarning: Data with input dtype int64 was converted t
o float64 by StandardScaler.
warnings.warn(msg, DataConversionWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:
475: DataConversionWarning: Data with input dtype int64 was converted t
o float64 by StandardScaler.
warnings.warn(msg, DataConversionWarning)
```

```
In [37]: from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.linear_model import SGDClassifier
```

```
In [26]: def support(train,cv):
inv_lambda=[10** -4,10** -3,10** -2,10** -1,10**0,10**1,10**2,10**3,10*
*4]

parameter={"C":inv_lambda}
svm=GridSearchCV(SVC(),parameter,verbose=1,scoring="roc_auc")
svm.fit(train,y_train)

opt_c=svm.best_params_.get('C')
print("best optimized C:",opt_c)
train_score=svm.cv_results_.get("mean_train_score")
test_score=svm.cv_results_.get("mean_test_score")

plt.plot(np.log(inv_lambda),train_score,'r', label = 'Train AUC')
plt.plot(np.log(inv_lambda),test_score,'b', label = 'CV AUC')
plt.ylim(0,1)
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left', borderaxespa
```

```
d=0.)
plt.grid(True)
plt.title("AUC Values for Train and CV \n")
plt.xlabel("C:hyperparameter")
plt.ylabel("AUC Value")
plt.show()
```

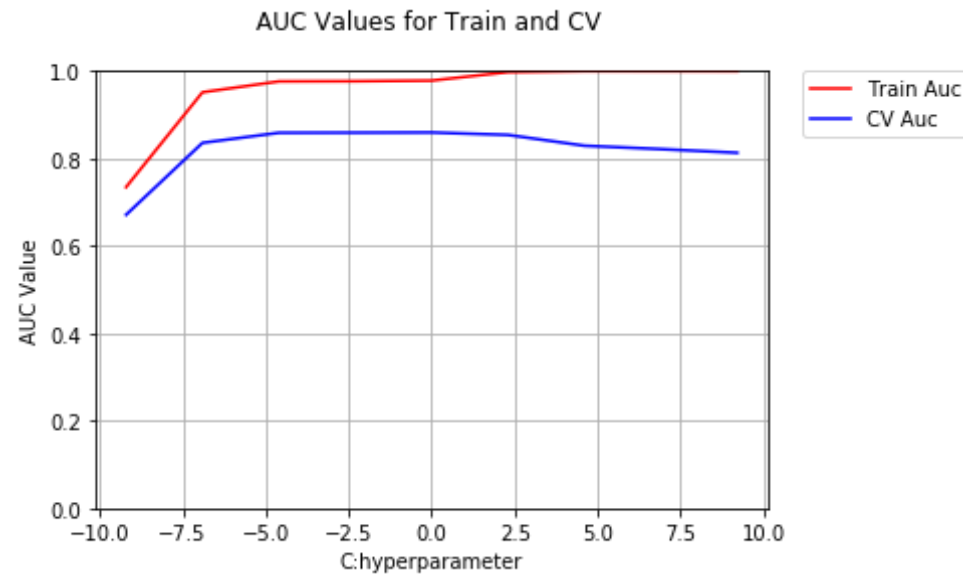
```
In [27]: def confusion_matrix(train,test):
        from sklearn.metrics import confusion_matrix
        from sklearn.metrics import confusion_matrix
        Y_test_pred=SGD.predict(test)
        Y_train_pred=SGD.predict(train)
        cm_train=confusion_matrix(y_train,Y_train_pred)
        cm_test=confusion_matrix(y_test,Y_test_pred)
        print(cm_train)
        print(cm_test)
        print(""*100)
        print("confusion matrix for test data")
        import seaborn as sns
        class_label=["0","1"]
        df_cm=pd.DataFrame(cm_test,index=class_label,columns=class_label)
        sns.heatmap(df_cm,annot=True,fmt="d")
        plt.title("confusion matrix")
        plt.xlabel("predicted label")
        plt.ylabel("true label")
        plt.show()
```

```
In [36]: support(X_train_bow,X_cv_bow)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 4.3min finished

best optimized C: 1



```
In [38]: SGD=SVC(probability=True,C=1)
from sklearn.metrics import roc_auc_score
SGD.fit(X_train_bow,y_train)

y_train_predict_proba=SGD.decision_function(X_train_bow)
y_test_predict_proba=SGD.decision_function(X_test_bow)
fpr,tpr,threshold=roc_curve(y_train,y_train_predict_proba[:])
fpr1,tpr1,threshold1=roc_curve(y_test,y_test_predict_proba[:])

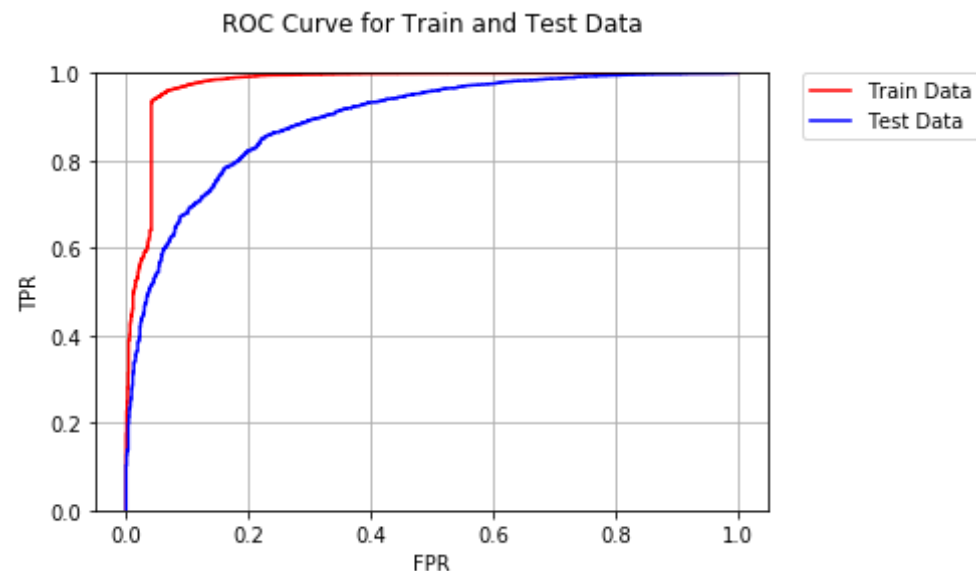
print("The AUC value for test data is ",roc_auc_score( y_test, y_test_p
redict_proba))

plt.plot(fpr,tpr,'r', label = 'Train Data')
plt.plot(fpr1,tpr1,'b', label = 'Test Data')
plt.ylim(0,1)
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left', borderaxespad=0.
)
plt.grid(True)
plt.title("ROC Curve for Train and Test Data\n")
plt.xlabel("FPR")
```



```
plt.ylabel("TPR")
plt.show()
```

The AUC value for test data is 0.8925903512751335



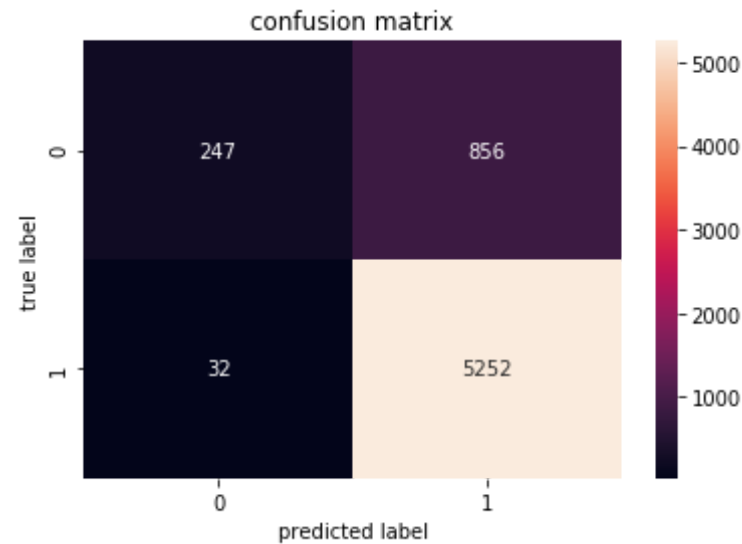
```
In [39]: confusion_matrix(X_train_bow,X_test_bow)
```

```
[[ 738  495]
 [   10 7444]]
[[ 247  856]
 [   32 5252]]
```

```
*****
```

```
*****
```

```
confusion matrix for test data
```



TF-IDF(RBF KERNEL)

```
In [41]: X=preprocessed_reviews  
Y=final["Score"]
```

```
In [42]: from sklearn.model_selection import train_test_split  
  
X_train,X_test,y_train,y_test=train_test_split(X,Y,test_size=0.33,shuffle=False)  
X_train,X_cv,y_train,y_cv=train_test_split(X_train,y_train,test_size=0.33,shuffle=False)
```

```
In [44]: from sklearn.feature_extraction.text import TfidfVectorizer  
vectorizer_TF = TfidfVectorizer(min_df=10,max_features=500,ngram_range=(1,2))  
vectorizer_TF.fit(X_train) # fit has to happen only on train data  
  
# we use the fitted CountVectorizer to convert the text to vector  
X_train_tf = vectorizer_TF.transform(X_train)
```

```
X_cv_tf = vectorizer_TF.transform(X_cv)
X_test_tf = vectorizer_TF.transform(X_test)

print("After TFIDF VEC")
print(X_train_tf.shape, y_train.shape)
print(X_cv_tf.shape, y_cv.shape)
print(X_test_tf.shape, y_test.shape)
```

```
After TFIDF VEC
(8687, 500) (8687,)
(4280, 500) (4280,)
(6387, 500) (6387,)
```

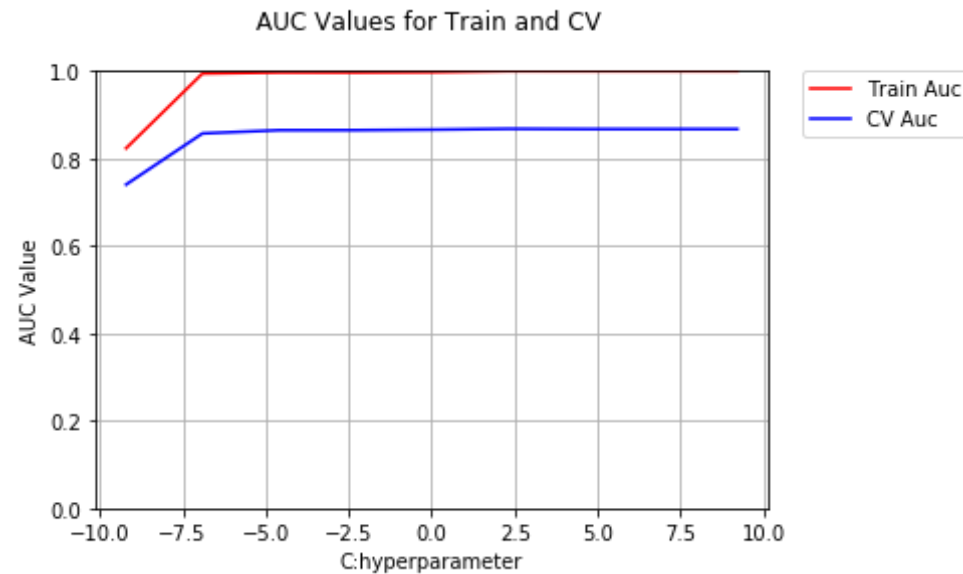
```
In [45]: from sklearn.preprocessing import StandardScaler
standardised=StandardScaler(with_mean=False)
X_train_tf=standardised.fit_transform(X_train_tf)
X_cv_tf=standardised.transform(X_cv_tf)
X_test_tf=standardised.transform(X_test_tf)
```

```
In [46]: support(X_train_tf,X_cv_tf)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 4.0min finished

best optimized C: 10



```
In [63]: SGD=SVC(probability=False,C=10)
from sklearn.metrics import roc_auc_score
SGD.fit(X_train_tf,y_train)

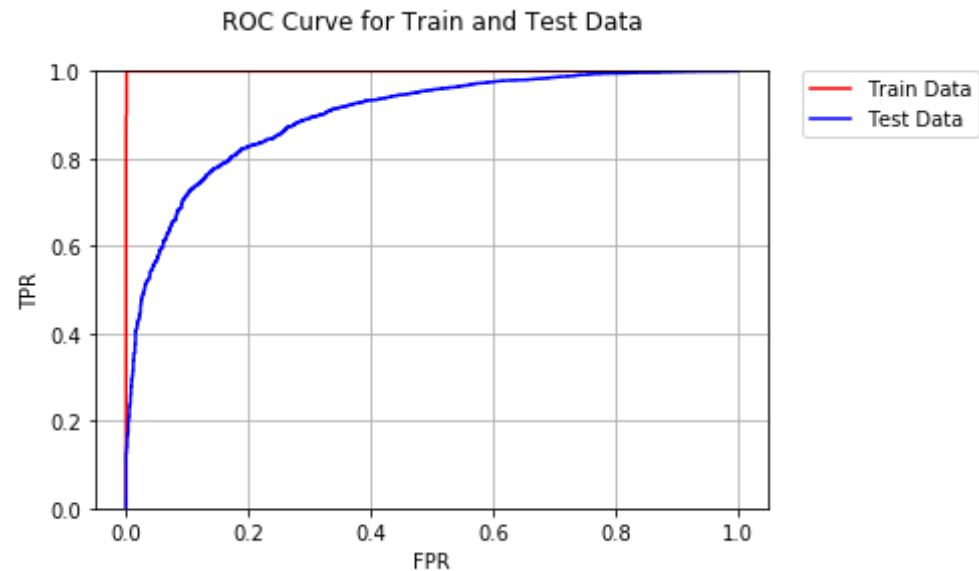
y_train_predict_proba=SGD.decision_function(X_train_tf)
y_test_predict_proba=SGD.decision_function(X_test_tf)
fpr,tpr,threshold=roc_curve(y_train,y_train_predict_proba[:])
fpr1,tpr1,threshold1=roc_curve(y_test,y_test_predict_proba[:])

print("The AUC value for test data is ",roc_auc_score( y_test, y_test_p
redict_proba))

plt.plot(fpr,tpr,'r', label = 'Train Data')
plt.plot(fpr1,tpr1,'b', label = 'Test Data')
plt.ylim(0,1)
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left', borderaxespad=0.
)
plt.grid(True)
plt.title("ROC Curve for Train and Test Data\n")
plt.xlabel("FPR")
```

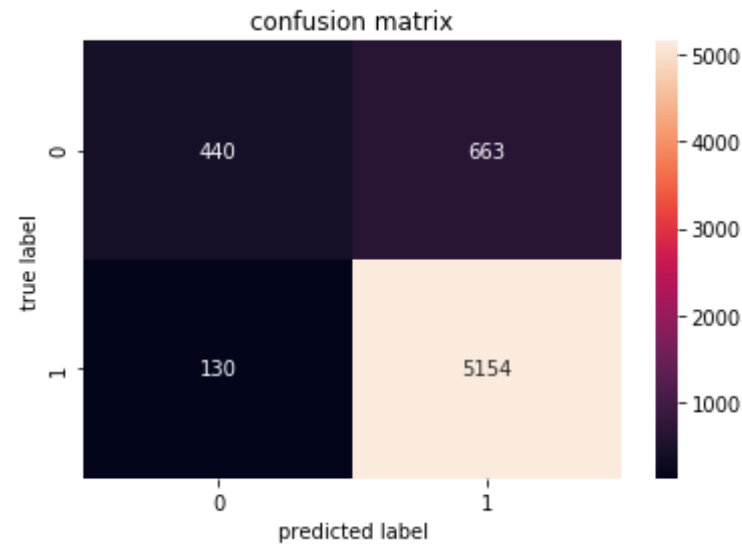
```
plt.ylabel("TPR")
plt.show()
```

The AUC value for test data is 0.8966680747503711



```
In [64]: confusion_matrix(X_train_tf,X_test_tf)
```

```
[[1232    1]
 [   1 7453]]
[[ 440   663]
 [ 130 5154]]
*****
*****
confusion matrix for test data
```



AVG-W2V(RBF KERNEL)

```
In [28]: X=preprocessed_reviews  
Y=final['Score']
```

```
In [29]: from sklearn.model_selection import train_test_split  
  
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.33,shuffle=False) # this is time based splitting  
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33,shuffle=False)
```

```
In [30]: i=0  
list_of_sentence=[]  
for sentence in preprocessed_reviews:  
    list_of_sentence.append(sentence.split())
```

```
In [31]: sent_of_train=[]
        for sent in X_train:
            sent_of_train.append(sent.split())
```

```
In [32]: sent_of_cv=[]
        for sent in X_cv:
            sent_of_cv.append(sent.split())

        sent_of_test=[]
        for sent in X_test:
            sent_of_test.append(sent.split())

        # Train your own Word2Vec model using your own train text corpus
        # min_count = 5 considers only words that occurred at least 5 times
        w2v_model=Word2Vec(sent_of_train,min_count=5,size=50, workers=4)

        w2v_words = list(w2v_model.wv.vocab)
```

```
In [33]: train_vectors = [];
        for sent in sent_of_train:
            sent_vec = np.zeros(50)
            cnt_words =0;
            for word in sent: #
                if word in w2v_words:
                    vec = w2v_model.wv[word]
                    sent_vec += vec
                    cnt_words += 1
            if cnt_words != 0:
                sent_vec /= cnt_words
            train_vectors.append(sent_vec)

        cv_vectors = [];
        for sent in sent_of_cv:
            sent_vec = np.zeros(50)
            cnt_words =0;
            for word in sent: #
                if word in w2v_words:
                    vec = w2v_model.wv[word]
```

```

        sent_vec += vec
        cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    cv_vectors.append(sent_vec)

```

```

# compute average word2vec for each review for X_test .
test_vectors = [];
for sent in sent_of_test:
    sent_vec = np.zeros(50)
    cnt_words = 0;
    for word in sent: #
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    test_vectors.append(sent_vec)

```

```

In [34]: X_train_wv=train_vectors
        X_cv_wv=cv_vectors
        X_test_wv=test_vectors

```

```

In [35]: from sklearn.preprocessing import StandardScaler
        standardised=StandardScaler(with_mean=False)
        X_train_wv=standardised.fit_transform(X_train_wv)
        X_cv_wv=standardised.transform(X_cv_wv)
        X_test_wv=standardised.transform(X_test_wv)

```

```

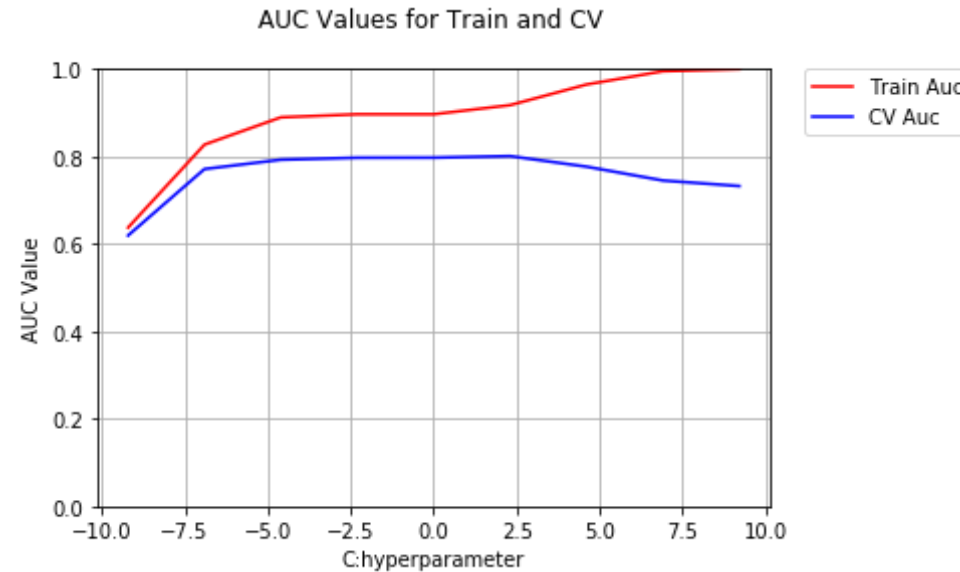
In [38]: support(X_train_wv,X_cv_wv)

```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 2.5min finished

best optimized C: 10



```
In [39]: SGD=SVC(C=10)
from sklearn.metrics import roc_auc_score
SGD.fit(X_train_wv,y_train)

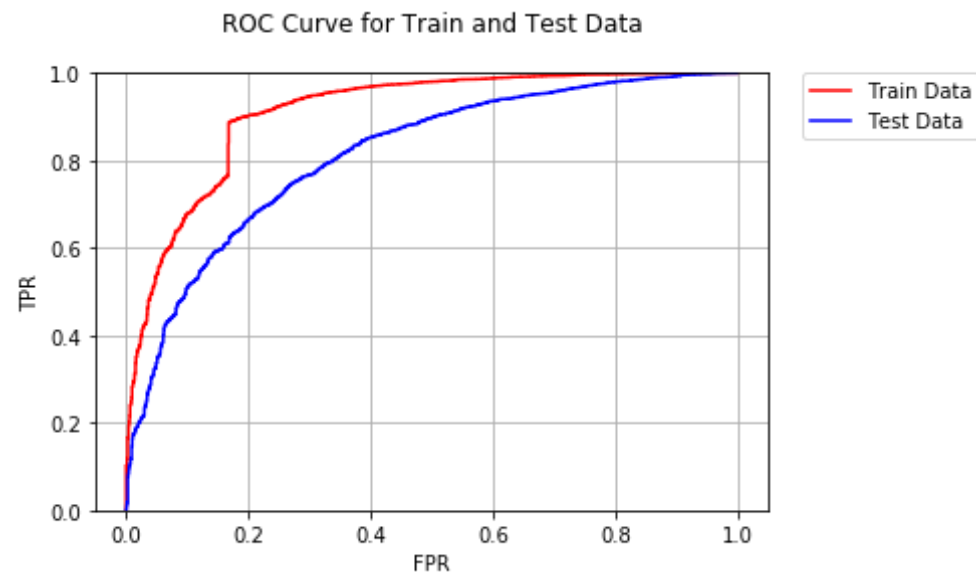
y_train_predict_proba=SGD.decision_function(X_train_wv)
y_test_predict_proba=SGD.decision_function(X_test_wv)
fpr, tpr, threshold=roc_curve(y_train,y_train_predict_proba[:])
fpr1, tpr1, threshold1=roc_curve(y_test,y_test_predict_proba[:])

print("The AUC value for test data is ",roc_auc_score( y_test, y_test_p
redict_proba))

plt.plot(fpr,tpr,'r', label = 'Train Data')
plt.plot(fpr1,tpr1,'b', label = 'Test Data')
plt.ylim(0,1)
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left', borderaxespad=0.
)
plt.grid(True)
plt.title("ROC Curve for Train and Test Data\n")
```

```
plt.xlabel("FPR")
plt.ylabel("TPR")
plt.show()
```

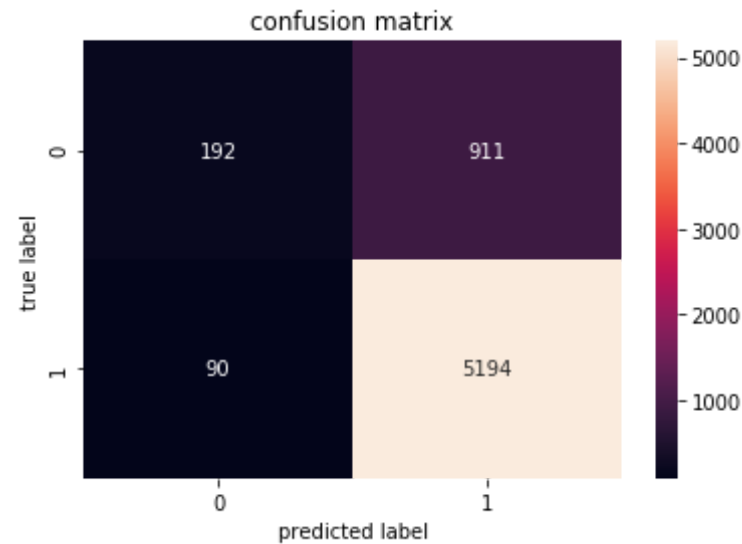
The AUC value for test data is 0.8135264226735563



```
In [40]: confusion_matrix(X_train_wv,X_test_wv)
```

```
[[ 372  861]
 [  45 7409]]
[[ 192  911]
 [  90 5194]]
```

```
*****
*****
confusion matrix for test data
```



TFIDF-W2Vec(RBF KERNEL)

```
In [45]: X=preprocessed_reviews  
Y=final["Score"]
```

```
In [46]: from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.33, shuffle=False) # this is random splitting  
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, shuffle=False)
```

```
In [47]: model = TfidfVectorizer()  
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 [48]: tfidf_feat = model.get_feature_names() # tfidf words/col-names
```

```
100%|██████████| 8687/8687 [01:22<00:00, 10  
5.82it/s]
```

```
In [49]: tfidf_cv_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(sent_of_cv): # 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]
            # tfidf = tfidf matrix[row, tfidf_feat.index(word)]
```

```

        # to reduce the computation we are
        # dictionary[word] = idf value of word in whole corpus
        # sent.count(word) = tf value 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_cv_vectors.append(sent_vec)
    row += 1

```

100%|██| 4280/4280 [00:39<00:00, 109.32it/s]

In [50]:

```

tfidf_test_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(sent_of_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 corpus
            # sent.count(word) = tf value 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_test_vectors.append(sent_vec)
    row += 1

```

100%|██| 6387/6387 [00:56<00:00, 112.64it/s]

```
In [51]: X_train_tw=tfidf_train_vectors
X_cv_tw=tfidf_cv_vectors
X_test_tw=tfidf_test_vectors
```

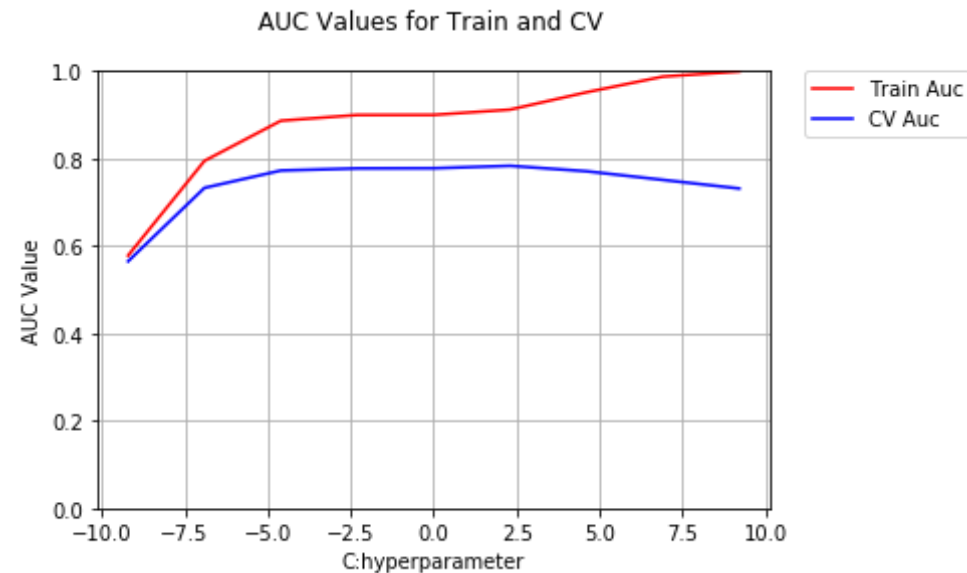
```
In [52]: from sklearn.preprocessing import StandardScaler
standardised=StandardScaler(with_mean=False)
X_train_tw=standardised.fit_transform(X_train_tw)
X_cv_tw=standardised.transform(X_cv_tw)
X_test_tw=standardised.transform(X_test_tw)
```

```
In [53]: support(X_train_tw,X_cv_tw)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 3.0min finished

best optimized C: 10



```
In [54]: SGD=SVC(C=10)
from sklearn.metrics import roc_auc_score
SGD.fit(X_train_tw,y_train)
```

```

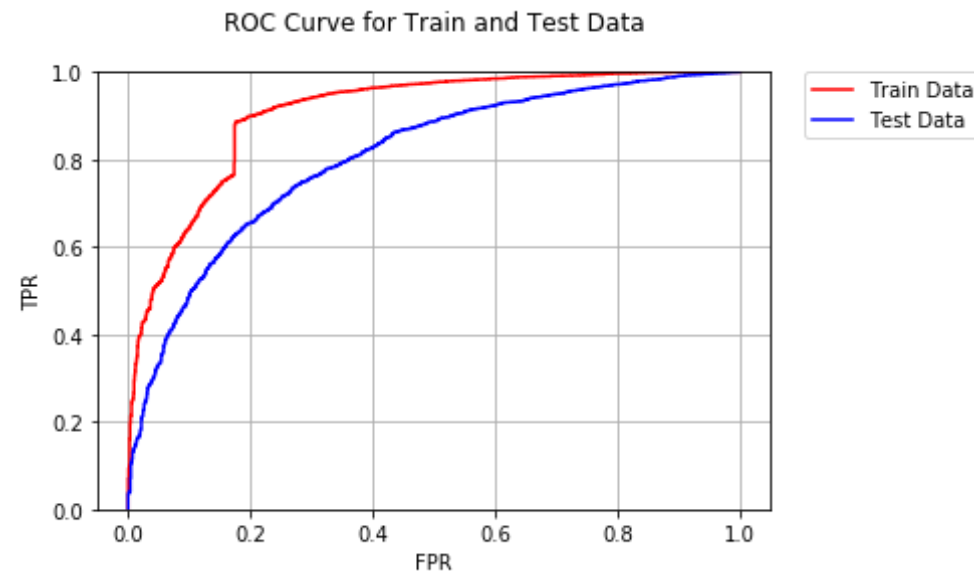
y_train_predict_proba=SGD.decision_function(X_train_tw)
y_test_predict_proba=SGD.decision_function(X_test_tw)
fpr,tpr,threshold=roc_curve(y_train,y_train_predict_proba[:])
fpr1,tpr1,threshold1=roc_curve(y_test,y_test_predict_proba[:])

print("The AUC value for test data is ",roc_auc_score( y_test, y_test_p
redict_proba))

plt.plot(fpr,tpr,'r', label = 'Train Data')
plt.plot(fpr1,tpr1,'b', label = 'Test Data')
plt.ylim(0,1)
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left', borderaxespad=0.
)
plt.grid(True)
plt.title("ROC Curve for Train and Test Data\n")
plt.xlabel("FPR")
plt.ylabel("TPR")
plt.show()

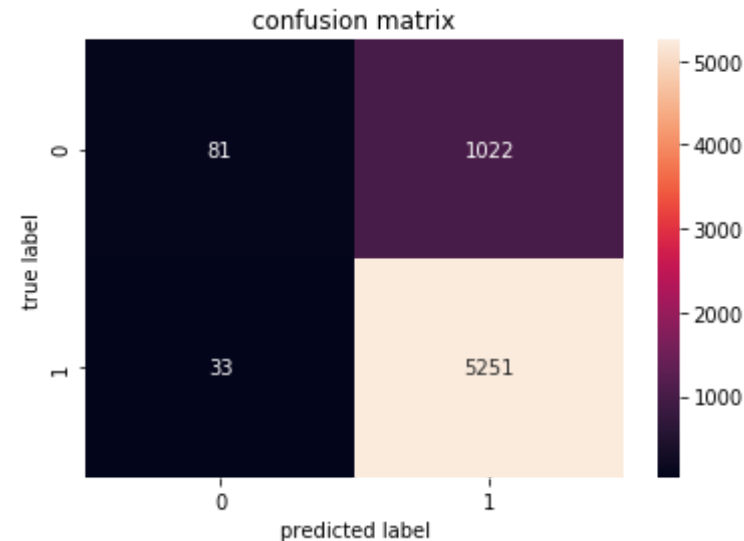
```

The AUC value for test data is 0.8044369049244954



```
In [55]: confusion_matrix(X_train_tw,X_test_tw)
```

```
[[ 188 1045]
 [  11 7443]]
[[  81 1022]
 [  33 5251]]
*****
*****
confusion matrix for test data
```



CONCLUSION

```
In [77]: print("RBF KERNEL")
from tabulate import tabulate
print(tabulate ([[ 'BOW', 1, 89], [ 'TF-IDF', 10, 89.66], [ 'AVG-W2V', 10, 81.94 ] , [ 'TFIDF-W2V', 10, 80.72]], headers=[ 'Vectorizer', 'best_C', 'AUC_test']))
```

```
RBF KERNEL
Vectorizer      best_C    AUC_test
-----
```


BOW	1	89
TF-IDF	10	89.66
AVG-W2V	10	81.94
TFIDF-W2V	10	80.72

1. Time computation is very high.
2. Here the good model in RBF kernel is Tf-idf.
3. Even we get more better results by taking more data points and featurization.

[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-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on your 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 [3]: `# Please write all the code with proper documentation`

[5.1.2] Applying Linear SVM on TFIDF, SET 2

In [3]: `# Please write all the code with proper documentation`

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

In [3]: *# Please write all the code with proper documentation*

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

In [3]: *# Please write all the code with proper documentation*

[5.2] RBF SVM

[5.2.1] Applying RBF SVM on BOW, SET 1

In [3]: *# Please write all the code with proper documentation*

[5.2.2] Applying RBF SVM on TFIDF, SET 2

In [3]: *# Please write all the code with proper documentation*

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

In [3]: *# Please write all the code with proper documentation*

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

In [3]: *# Please write all the code with proper documentation*

[6] Conclusions

```
In [4]: # Please compare all your models using Prettytable library
```