

# Amazon Fine Food Reviews Analysis

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

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

The Amazon Fine Food Reviews dataset consists of reviews of fine foods from Amazon.

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

1. Id
2. ProductId - unique identifier for the product
3. UserId - unique identifier for the user
4. ProfileName
5. HelpfulnessNumerator - number of users who found the review helpful
6. HelpfulnessDenominator - number of users who indicated whether they found the review helpful or not
7. Score - rating between 1 and 5
8. Time - timestamp for the review
9. Summary - brief summary of the review
10. Text - text of the review

## Objective:

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

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

[Ans] We could use 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
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_val_score
from collections import Counter
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_validate

```

```

/anaconda3/lib/python3.6/site-packages/smart_open/ssh.py:34: UserWarning:
g: paramiko missing, opening SSH/SCP/SFTP paths will be disabled. `pip
install paramiko` to suppress
  warnings.warn('paramiko missing, opening SSH/SCP/SFTP paths will be d
isabled. `pip install paramiko` to suppress')

```

```

In [2]: # using SQLite Table to read data.
con = sqlite3.connect('/Users/puravshah/Downloads/amazon-fine-food-reviews/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 Sco
re != 3 LIMIT 500000""", con)
# for tsne assignment you can take 5k data points

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

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

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

```

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

Out[2]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpfulnes
--	----	-----------	--------	-------------	----------------------	------------

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

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

In [5]: `display[display['UserId']=='AZY10LLTJ71NX']`

Out[5]:

	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

	<b>Id</b>	<b>ProductId</b>	<b>UserId</b>	<b>ProfileName</b>	<b>HelpfulnessNumerator</b>	<b>Helpfuln</b>
<b>0</b>	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2
<b>1</b>	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2
<b>2</b>	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2
<b>3</b>	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2
<b>4</b>	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]: (87775, 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]: 87.775
```

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

```
In [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()
```

```
(87773, 10)
```

```
Out[13]: 1    73592  
         0    14181  
         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 observed 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)

```

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

=====

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

=====

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

=====

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

=====

```

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)
```

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

```
In [16]: # 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)
```

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

=====

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

=====

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

=====

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

```
In [17]: # 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]: sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

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

```
=====
In [19]: #remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

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

```
In [20]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

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

```
In [21]: # 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', 'yoursele
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 [22]: # 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%|██████████| 87773/87773 [00:32<00:00, 2734.29it/s]

```

```

In [23]: preprocessed_reviews[1500]

```

```

Out[23]: 'way hot blood took bite jig lol'

```



## [3.2] Preprocessing Review Summary

In [24]: `## Similarly you can do preprocessing for review summary also.`

## [4] Featurization

### [4.1] BAG OF WORDS

```
In [25]: #BoW
from sklearn.preprocessing import StandardScaler
scaler=StandardScaler(with_mean=False)
X_1, X_test, y_1, y_test = train_test_split(preprocessed_reviews, final
['Score'], test_size=0.3, random_state=0)
X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.3)
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(X_tr)

print("some feature names ", count_vect.get_feature_names()[:10])
print('='*50)

final_counts = count_vect.transform(X_tr)
final_counts=scaler.fit_transform(final_counts)
X_cv_bow=count_vect.transform(X_cv)
X_cv_bow=scaler.fit_transform(X_cv_bow)
X_test_bow=count_vect.transform(X_test)
X_test_bow=scaler.fit_transform(X_test_bow)
print("the type of count vectorizer ",type(final_counts))
print("the shape of out text BOW vectorizer ",final_counts.get_shape())
print("the number of unique words ", final_counts.get_shape()[1])
print(y_test.shape)

some feature names  ['aa', 'aaa', 'aaaa', 'aaaaa', 'aaaaaaaaaaaa', 'aaa
aaaaaaaaaaaa', 'aaaaaaahhhhhh', 'aaaaaaaarrrrrggghh', 'aaaaaawwwwwwww
w', 'aaaaah']
=====
```

```
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer (43008, 38907)
the number of unique words 38907
(26332,)
```

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

```
In [26]: #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

count_vect_gram = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
count_vect_gram.fit_transform(X_tr)
final_bigram_counts = count_vect_gram.transform(X_tr)
final_bigram_counts = scaler.fit_transform(final_bigram_counts)
X_cv_ngram = count_vect_gram.transform(X_cv)
X_cv_ngram = scaler.fit_transform(X_cv_ngram)
X_test_ngram = count_vect_gram.transform(X_test)
X_test_ngram = scaler.fit_transform(X_test_ngram)
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_shape()[1])
```

```
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer (43008, 5000)
the number of unique words including both unigrams and bigrams 5000
```

## [4.3] TF-IDF

```
In [27]: tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(X_tr)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names()[0:10])
print('='*50)

final_tf_idf = tf_idf_vect.transform(X_tr)
final_tf_idf=scaler.fit_transform(final_tf_idf)
X_cv_tfidf=tf_idf_vect.transform(X_cv)
X_cv_tfidf=scaler.fit_transform(X_cv_tfidf)
X_test_tfidf=tf_idf_vect.transform(X_test)
X_test_tfidf=scaler.fit_transform(X_test_tfidf)
print("the type of count vectorizer ",type(final_tf_idf))
print("the shape of out text TFIDF vectorizer ",final_tf_idf.get_shape())
print("the number of unique words including both unigrams and bigrams ", final_tf_idf.get_shape()[1])

some sample features(unique words in the corpus) ['abdominal', 'ability', 'able', 'able add', 'able buy', 'able chew', 'able drink', 'able eat', 'able enjoy', 'able feed']
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (43008, 25544)
the number of unique words including both unigrams and bigrams 25544
```

## [4.4] Word2Vec

```
In [28]: # Train your own Word2Vec model using your own text corpus
i=0
list_of_sentence_train=[]
for sentence in X_tr:
    list_of_sentence_train.append(sentence.split())
```

```
In [29]: i=0
list_of_sentence_cv=[]
```

```
for sentence in X_cv:  
    list_of_sentence_cv.append(sentence.split())
```

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

```
In [31]: # 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 val  
# ues  
# To use this code-snippet, download "GoogleNews-vectors-negative300.bi  
# n"  
# 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 variable 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_train=Word2Vec(list_of_sentence_train,min_count=5,size=50  
    , workers=4)  
    print(w2v_model_train.wv.most_similar('great'))
```

```

        print('='*50)
        print(w2v_model_train.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 ")
#w2v_model_train=scaler.fit_transform(w2v_model_train)

[('fantastic', 0.8307128548622131), ('awesome', 0.8240286111831665),
('good', 0.8003333806991577), ('excellent', 0.8002113103866577), ('wonderful', 0.7988190650939941), ('amazing', 0.7927408814430237), ('terrific', 0.7263221144676208), ('perfect', 0.7161662578582764), ('decent', 0.6933251619338989), ('incredible', 0.682373046875)]
=====
[('greatest', 0.8318613767623901), ('best', 0.7703869342803955), ('tastiest', 0.677887499332428), ('nastiest', 0.676857590675354), ('superior', 0.6608847379684448), ('closest', 0.6362959742546082), ('experience', 0.6348606944084167), ('hottest', 0.6183363795280457), ('coolest', 0.6037110090255737), ('ive', 0.5886297821998596)]

```

```

In [133]: '''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 occurred at least 5 times
    w2v_model_cv=Word2Vec(list_of_sentence_cv,min_count=5,size=50,workers=4)
    print(w2v_model_cv.wv.most_similar('great'))
    print('='*50)
    print(w2v_model_cv.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 ")
#w2v_model_cv=scaler.fit_transform(w2v_model_cv)'''

```

```

Out[133]: '\is_your_ram_gt_16g=False\nwant_to_use_google_w2v = False\nwant_to_train_w2v = True\n\nif want_to_train_w2v:\n    # min_count = 5 considers only words that occurred at least 5 times\n    w2v_model_cv=Word2Vec(list_of_sentence_cv,min_count=5,size=50, workers=4)\n    print(w2v_model_cv.wv.most_similar(\'great\'))\n    print(\'=\'*50)\n    print(w2v_model_cv.wv.most_similar(\'worst\'))\n    \nelif want_to_use_google_w2v and is_your_ram_gt_16g:\n    if os.path.isfile(\'GoogleNews-vectors-negative300.bin\'):\n        w2v_model=KeyedVectors.load_word2vec_format(\'GoogleNews-vectors-negative300.bin\', binary=True)\n        print(w2v_model.wv.most_similar(\'great\'))\n        print(w2v_model.wv.most_similar(\'worst\'))\n    else:\n        print("you don't have gogole's word2vec file, keep want_to_train_w2v = True, to train your own w2v ") \n#w2v_model_cv=scaler.fit_transform(w2v_model_cv)'

```

```

In [134]: '''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 occurred at least 5 times
    w2v_model_test=Word2Vec(list_of_sentence_test,min_count=5,size=50,
workers=4)
    print(w2v_model_test.wv.most_similar('great'))
    print('='*50)
    print(w2v_model_test.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 ")
#w2v_model_test=scaler.fit_transform(w2v_model_test)'''
```

```
Out[134]: '\\is_your_ram_gt_16g=False\\nwant_to_use_google_w2v = False\\nwant_to_train_w2v = True\\n\\nif want_to_train_w2v:\\n    # min_count = 5 considers only words that occurred at least 5 times\\n    w2v_model_test=Word2Vec(list_of_sentence_test,min_count=5,size=50, workers=4)\\n    print(w2v_model_test.wv.most_similar('\\great\\'))\\n    print('\\'=\\'*50)\\n    print(w2v_model_test.wv.most_similar('\\worst\\'))\\n    \\nelif want_to_use_google_w2v and is_your_ram_gt_16g:\\n    if os.path.isfile('\\GoogleNews-vectors-negative300.bin\\'):\\n        w2v_model=KeyedVectors.load_word2vec_format('\\GoogleNews-vectors-negative300.bin\\', binary=True)\\n        print(w2v_model.wv.most_similar('\\great\\'))\\n        print(w2v_model.wv.most_similar('\\worst\\'))\\n    else:\\n        print("you don't have gogole's word2vec file, keep want_to_train_w2v = True, to train your own w2v ")\\n#w2v_model_test=scaler.fit_transform(w2v_model_test)'
```

```
In [32]: w2v_words_train = list(w2v_model_train.wv.vocab)
print("number of words that occurred minimum 5 times ",len(w2v_words_train))
print("sample words ", w2v_words_train[0:50])
```

```
number of words that occurred minimum 5 times 12529
sample words ['chocolate', 'found', 'dean', 'twice', 'price', 'packaging', 'plain', 'plastic', 'bag', 'quality', 'good', 'put', 'coffee', 'low', 'calorie', 'treat', 'calories', 'serving', 'seems', 'like', 'indulgence', 'prompt', 'shipping', 'adds', 'value', 'not', 'speak', 'personal', 'experience', 'four', 'cats', 'prefer', 'fancy', 'feast', 'brands', 'cat', 'food', 'grilled', 'classic', 'shredded', 'elegant', 'medleys', 'part', 'healthy', 'eating', 'plan', 'gave', 'anything', 'made', 'processed']
```

```
In [136]: '''w2v_words_cv = list(w2v_model_cv.wv.vocab)
print("number of words that occurred minimum 5 times ",len(w2v_words_cv))
```

```
v))  
print("sample words ", w2v_words_cv[0:50])'''
```

```
Out[136]: '\\w2v_words_cv = list(w2v_model_cv.wv.vocab)\nprint("number of words t  
hat occured minimum 5 times ",len(w2v_words_cv))\nprint("sample words  
", w2v_words_cv[0:50])'
```

```
In [137]: '''w2v_words_test = list(w2v_model_test.wv.vocab)  
print("number of words that occured minimum 5 times ",len(w2v_words_tes  
t))  
print("sample words ", w2v_words_test[0:50])'''
```

```
Out[137]: '\\w2v_words_test = list(w2v_model_test.wv.vocab)\nprint("number of wor  
ds that occured minimum 5 times ",len(w2v_words_test))\nprint("sample w  
ords ", w2v_words_test[0:50])'
```

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

### [4.4.1.1] Avg W2v

```
In [33]: sent_vectors_train = []; # the avg-w2v for each sentence/review is stor  
ed in this list  
for sent in tqdm(list_of_sentence_train): # for each review/sentence  
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo  
u 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/re  
view  
    for word in sent: # for each word in a review/sentence  
        if word in w2v_words_train:  
            vec = w2v_model_train.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]))
sent_vectors_train=scaler.fit_transform(sent_vectors_train)
```

```
100%|██████████| 43008/43008 [01:05<00:00, 654.18it/s]
```

```
43008
```

```
50
```

```
In [34]: sent_vectors_cv = []; # the avg-w2v for each sentence/review is stored
         in this list
         for sent in tqdm(list_of_sentence_cv): # for each review/sentence
             sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
             u 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/re
             view
             for word in sent: # for each word in a review/sentence
                 if word in w2v_words_train:
                     vec = w2v_model_train.wv[word]
                     sent_vec += vec
                     cnt_words += 1
             if cnt_words != 0:
                 sent_vec /= cnt_words
             sent_vectors_cv.append(sent_vec)
         print(len(sent_vectors_cv))
         print(len(sent_vectors_cv[0]))
         sent_vectors_cv=scaler.fit_transform(sent_vectors_cv)
```

```
100%|██████████| 18433/18433 [00:30<00:00, 595.58it/s]
```

```
18433
```

```
50
```

```
In [35]: sent_vectors_test = []; # the avg-w2v for each sentence/review is store
         d in this list
         for sent in tqdm(list_of_sentence_test): # for each review/sentence
             sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
             u 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/re
```

```

view
    for word in sent: # for each word in a review/sentence
        if word in w2v_words_train:
            vec = w2v_model_train.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]))
sent_vectors_test=scaler.fit_transform(sent_vectors_test)

```

```

100%|██████████| 26332/26332 [00:35<00:00, 746.77it/s]

```

```

26332

```

```

50

```

#### [4.4.1.2] TFIDF weighted W2v

```

In [36]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model_train= TfidfVectorizer()
tf_idf_matrix_train = model_train.fit_transform(X_tr)
# we are converting a dictionary with word as a key, and the idf as a v
# value
dictionary_train = dict(zip(model_train.get_feature_names(), list(model
_train.idf_)))

```

```

In [144]: '''model_cv= TfidfVectorizer()
tf_idf_matrix_cv = model_cv.fit_transform(X_cv)
# we are converting a dictionary with word as a key, and the idf as a v
# value
dictionary_cv= dict(zip(model_cv.get_feature_names(), list(model_cv.idf
_)))'''

```

```

Out[144]: 'model_cv= TfidfVectorizer()\ntf_idf_matrix_cv = model_cv.fit_transform
(X_cv)\n# we are converting a dictionary with word as a key, and the id
f as a value\ndictionary_cv= dict(zip(model_cv.get_feature_names(), lis
t(model_cv.idf_)))'

```

```
In [146]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
'''model_test= TfidfVectorizer()
tf_idf_matrix_test = model_test.fit_transform(X_test)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary_test= dict(zip(model_test.get_feature_names(), list(model_test.idf_)))'''
```

```
Out[146]: 'model_test= TfidfVectorizer()\ntf_idf_matrix_test = model_test.fit_transform(X_test)\n# we are converting a dictionary with word as a key, and the idf as a value\ndictionary_test= dict(zip(model_test.get_feature_names(), list(model_test.idf_)))'
```

```
In [37]: # TF-IDF weighted Word2Vec
tfidf_feat_train = model_train.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_sentence_train): # for each review/sentence
    sent_vec_train = np.zeros(50) # as word vectors are of zero length
    weight_sum_train =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_train and word in tfidf_feat_train:
            vec_train = w2v_model_train.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_train[word]*(sent.count(word)/len(sent))
            sent_vec_train += (vec * tf_idf)
            weight_sum_train += tf_idf
    if weight_sum_train != 0:
        sent_vec_train /= weight_sum_train
```

```
tfidf_sent_vectors_train.append(sent_vec_train)
row += 1
tfidf_sent_vectors_train = scaler.fit_transform(tfidf_sent_vectors_train
)
```

```
100%|██████████| 43008/43008 [16:36:27<00:00, 1.39s/it]
```

In [38]: `tfidf_feat_cv = model_train.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_cv = []; # the tfidf-w2v for each sentence/review is
    stored in this list
row=0;
for sent in tqdm(list_of_sentence_cv): # for each review/sentence
    sent_vec_cv = np.zeros(50) # as word vectors are of zero length
    weight_sum_cv = 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_train and word in tfidf_feat_cv:
            vec_cv = w2v_model_train.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_train[word]*(sent.count(word)/len(sent
    ))
            sent_vec_cv += (vec * tf_idf)
            weight_sum_cv += tf_idf
        if weight_sum_cv != 0:
            sent_vec_cv /= weight_sum_cv
            tfidf_sent_vectors_cv.append(sent_vec_cv)
    row += 1
tfidf_sent_vectors_cv = scaler.fit_transform(tfidf_sent_vectors_cv)
```

```
100%|██████████| 18433/18433 [03:43<00:00, 82.39it/s]
```

In [39]: `tfidf_feat_test = model_train.get_feature_names() # tfidf words/col-names`

```

# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf

tfidf_sent_vectors_test = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentence_test): # for each review/sentence
    sent_vec_test = np.zeros(50) # as word vectors are of zero length
    weight_sum_test =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_train and word in tfidf_feat_test:
            vec_test = w2v_model_train.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_train[word]*(sent.count(word)/len(sent))
        sent_vec_test += (vec * tf_idf)
        weight_sum_test += tf_idf
    if weight_sum_test != 0:
        sent_vec_test /= weight_sum_test
    tfidf_sent_vectors_test.append(sent_vec_test)
    row += 1
tfidf_sent_vectors_test = scaler.fit_transform(tfidf_sent_vectors_test)
100%|██████████| 26332/26332 [08:50<00:00, 49.62it/s]

```

## [5] Assignment 5: Apply Logistic Regression

### 1. Apply Logistic Regression 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. Hyper parameter tuning (find best hyper parameters corresponding the algorithm that you choose)

- 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

## 3. Pertubation Test

- Get the weights  $W$  after fit your model with the data  $X$  i.e Train data.
- Add a noise to the  $X$  ( $X' = X + e$ ) and get the new data set  $X'$  (if  $X$  is a sparse matrix,  $X.data += e$ )
- Fit the model again on data  $X'$  and get the weights  $W'$
- Add a small eps value(to eliminate the divisible by zero error) to  $W$  and  $W'$  i.e  $W = W + 10^{-6}$  and  $W' = W' + 10^{-6}$
- Now find the % change between  $W$  and  $W'$  ( $| (W - W') / (W) | * 100$ )
- Calculate the 0th, 10th, 20th, 30th, ...100th percentiles, and observe any sudden rise in the values of `percentage_change_vector`
- Ex: consider your 99th percentile is 1.3 and your 100th percentiles are 34.6, there is sudden rise from 1.3 to 34.6, now calculate the 99.1, 99.2, 99.3,..., 100th percentile values and get the proper value after which there is sudden rise the values, assume it is 2.5
- Print the feature names whose % change is more than a threshold  $x$  (in our example it's 2.5)

## 4. Sparsity

- Calculate sparsity on weight vector obtained after using L1 regularization

NOTE: Do sparsity and multicollinearity for any one of the vectorizers. Bow or tf-idf is recommended.



## 5. Feature importance

- Get top 10 important features for both positive and negative classes separately.

## 6. Feature engineering

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

## 7. 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](#).



## 8. [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 Logistic Regression

### [5.1] Logistic Regression on BOW, SET 1

#### [5.1.1] Applying Logistic Regression with L1 regularization on BOW, SET 1

```
In [40]: # Please write all the code with proper documentation
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

tuned_parameters=[{'C':[10**-4,10**-3,10**-2,10**-1,1,10**1,10**2,10**3,10**4]}]
model=GridSearchCV(LogisticRegression(),tuned_parameters,scoring='roc_auc',cv=10)
model.fit(final_counts,y_tr)
print(model.best_estimator_)
optimal_c=model.best_estimator_.C
print(model.score(X_cv_bow,y_cv))
```

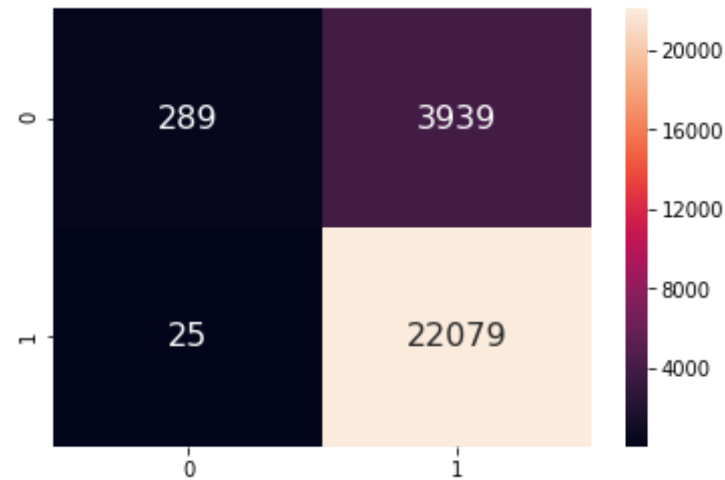
```
LogisticRegression(C=0.001, class_weight=None, dual=False, fit_intercept=True,
```

```
    intercept_scaling=1, l1_ratio=None, max_iter=100,
    multi_class='warn', n_jobs=None, penalty='l2',
    random_state=None, solver='warn', tol=0.0001, verbos
```

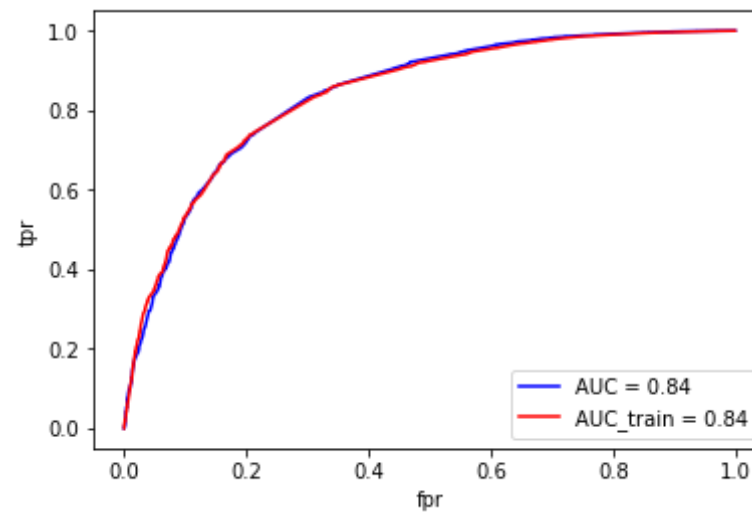


```
e=0,  
  
warm_start=False)  
0.9162728051362579
```

```
In [41]: from sklearn.metrics import confusion_matrix as cm  
from sklearn.metrics import roc_auc_score  
from sklearn.metrics import roc_curve  
  
model_new=LogisticRegression(C=optimal_c,penalty='l1')  
model_new.fit(final_counts,y_tr)  
pred=model_new.predict(X_test_bow)  
probab=model_new.predict_proba(X_test_bow)  
probab=probab[:,1]  
fpr, tpr, thresholds=metrics.roc_curve(y_test, probab)  
auc=roc_auc_score(y_test, probab)  
acc=accuracy_score(y_test, pred, normalize=True)*float(100)  
pred2=model_new.predict(final_counts)  
probab2=model_new.predict_proba(final_counts)  
probab2=probab2[:,1]  
fpr_train, tpr_train, thresholds=metrics.roc_curve(y_tr, probab2)  
auc2=roc_auc_score(y_test, probab)  
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), range(2), range(2))  
    #heatmap for visualization of matrix  
sns.heatmap(df_cm, annot=True, annot_kws={"size": 16}, fmt='g')  
plt.show()  
plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' %auc)  
plt.plot(fpr_train, tpr_train, 'r', label = 'AUC_train = %0.2f' %auc2)  
plt.legend(loc='lower right')  
plt.xlabel('fpr')  
plt.ylabel('tpr')
```



Out[41]: Text(0,0.5,'tpr')



**[5.1.1.1] Calculating sparsity on weight vector obtained using L1 regularization on BOW,**  
**SET 1**

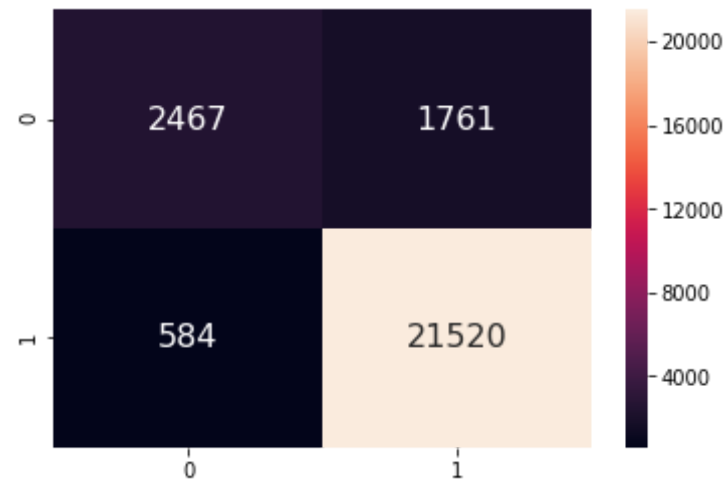
```
In [42]: # Please write all the code with proper documentation
```

```
clf=LogisticRegression(C=0.001,penalty='l1')
clf.fit(final_counts,y_tr)
w=clf.coef_
print(np.count_nonzero(w))
```

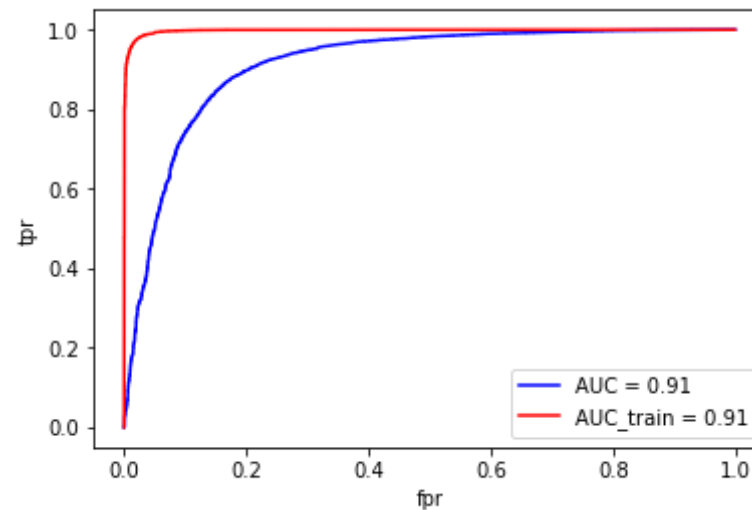
35

### [5.1.2] Applying Logistic Regression with L2 regularization on BOW, SET 1

```
In [43]: # Please write all the code with proper documentation
model_new=LogisticRegression(C=0.001)
model_new.fit(final_counts,y_tr)
pred=model_new.predict(X_test_bow)
probab=model_new.predict_proba(X_test_bow)
probab=probab[:,1]
fpr,tpr,thresholds=metrics.roc_curve(y_test,probab)
auc=roc_auc_score(y_test,probab)
acc=accuracy_score(y_test,pred,normalize=True)*float(100)
pred2=model_new.predict(final_counts)
probab2=model_new.predict_proba(final_counts)
probab2=probab2[:,1]
fpr_train,tpr_train,thresholds=metrics.roc_curve(y_tr,probab2)
auc2=roc_auc_score(y_test,probab)
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), range(2),range(2))
    #heatmap for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
plt.show()
plt.plot(fpr,tpr,'b', label = 'AUC = %0.2f' %auc)
plt.plot(fpr_train,tpr_train,'r', label = 'AUC_train = %0.2f' %auc2)
plt.legend(loc='lower right')
plt.xlabel('fpr')
plt.ylabel('tpr')
```



Out[43]: Text(0,0.5,'tpr')



#### [5.1.2.1] Performing pertubation test (multicollinearity check) on BOW, SET 1

```
In [44]: # Please write all the code with proper documentation
from scipy.sparse import csr_matrix
```

```

import numpy as np
w=model_new.coef_
print(w.shape)
print(final_counts.shape)
noise=np.random.normal(0,1)
#final_counts=csr_matrix.toarray(final_counts)
final_counts.data=final_counts.data+noise
model_new.fit(final_counts,y_tr)
w1=model_new.coef_
print(w1.shape)
w=w+10** -6
#print(w[0:10])
w1+=10** -6
#print(w1[0:10])
percentage_change=[]
percentage_change=np.absolute(np.divide(np.subtract(w[i],w1[i]),w[i]))*
100
print(percentage_change.T[0:20])
#From the below result it is clear that there is a high percentage difference between pairs of many features

```

```

(1, 38907)
(43008, 38907)
(1, 38907)
[ 4.56852039  8.93112335  6.2812383   8.91295642 15.69234769 26.5980902
9
 3.09536305  1.98147193  3.09536305  9.64945562 27.06394786 12.4320760
9
 1.91365457 26.62157442 23.70661379  4.62796287  7.51376343  3.4817694
6
 0.17037011 24.52926493]

```

### [5.1.3] Feature Importance on BOW, SET 1

#### [5.1.3.1] Top 10 important features of positive class from SET 1

In [45]: *# Please write all the code with proper documentation*

```

clf=LogisticRegression(C=0.001,penalty='l2')
clf.fit(final_counts,y_tr)
w=clf.coef_
#print(np.count_nonzero(w))
bow=count_vect.get_feature_names()
print(len(bow))
df=pd.DataFrame(w,columns=bow)
df=df.T

#df=df[0].sort_values(ascending=False)

#print(df.head(10))
print('The top ten positive features are-\n',df[0].sort_values(ascending=False)[0:10])

```

```

38907
The top ten positive features are-
great      0.368207
best       0.251581
love       0.249832
good       0.245001
delicious  0.212514
loves      0.184336
perfect    0.162284
favorite   0.161535
nice       0.156869
excellent  0.155096
Name: 0, dtype: float64

```

#### [5.1.3.2] Top 10 important features of negative class from SET 1

In [46]: *# Please write all the code with proper documentation*

```

print('The top ten negative features are-\n',df[0].sort_values(ascending=True)[0:10])

```

```

The top ten negative features are-
not      -0.229475
disappointed -0.144971
worst     -0.142356

```

```
terrible      -0.137295
awful         -0.121692
horrible      -0.116585
disappointing -0.113989
threw         -0.110944
disappointment -0.109563
weak          -0.104748
Name: 0, dtype: float64
```

## [5.2] Logistic Regression on TFIDF, SET 2

### [5.2.1] Applying Logistic Regression with L1 regularization on TFIDF, SET 2

```
In [47]: # Please write all the code with proper documentation
tuned_parameters=[{'C':[10**-4,10**-3,10**-2,10**-1,1,10**1,10**2,10**3,10**4]}]
model=GridSearchCV(LogisticRegression(),tuned_parameters,scoring='roc_auc',cv=10)
model.fit(final_tf_idf,y_tr)
print(model.best_estimator_)
optimal_c=model.best_estimator_.C
print(model.score(X_cv_tfidf,y_cv))
```

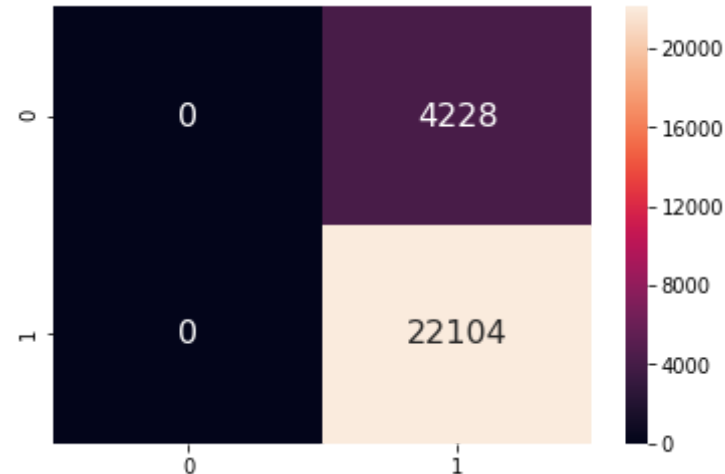
```
LogisticRegression(C=0.0001, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, l1_ratio=None, max_iter=100,
                    multi_class='warn', n_jobs=None, penalty='l2',
                    random_state=None, solver='warn', tol=0.0001, verbose=0,
                    warm_start=False)
0.9512659110154398
```

```
In [48]: model_new=LogisticRegression(C=optimal_c,penalty='l1')
model_new.fit(final_tf_idf,y_tr)
```

```

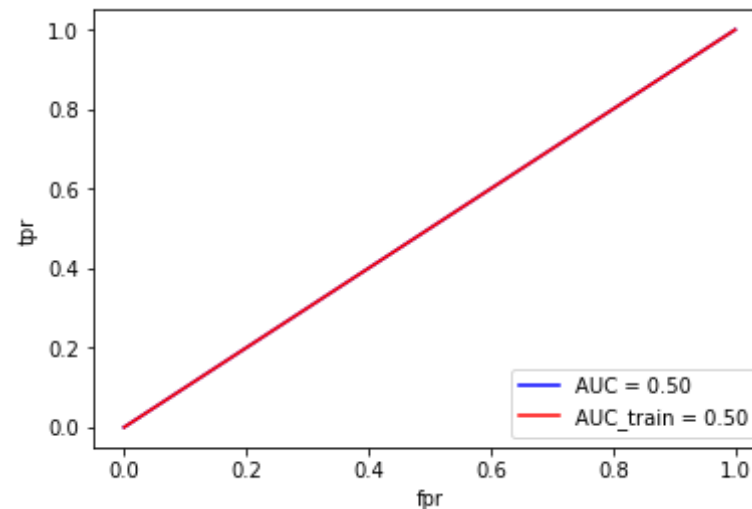
pred=model_new.predict(X_test_tfidf)
probab=model_new.predict_proba(X_test_tfidf)
probab=probab[:,1]
fpr,tpr,thresholds=metrics.roc_curve(y_test,probab)
auc=roc_auc_score(y_test,probab)
acc=accuracy_score(y_test,pred,normalize=True)*float(100)
pred2=model_new.predict(final_tf_idf)
probab2=model_new.predict_proba(final_tf_idf)
probab2=probab2[:,1]
fpr_train,tpr_train,thresholds=metrics.roc_curve(y_tr,probab2)
auc2=roc_auc_score(y_test,probab)
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), range(2),range(2))
    #heatmap for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
plt.show()
plt.plot(fpr,tpr,'b', label = 'AUC = %0.2f' %auc)
plt.plot(fpr_train,tpr_train,'r', label = 'AUC_train = %0.2f' %auc2)
plt.legend(loc='lower right')
plt.xlabel('fpr')
plt.ylabel('tpr')

```



Out[48]: Text(0,0.5,'tpr')

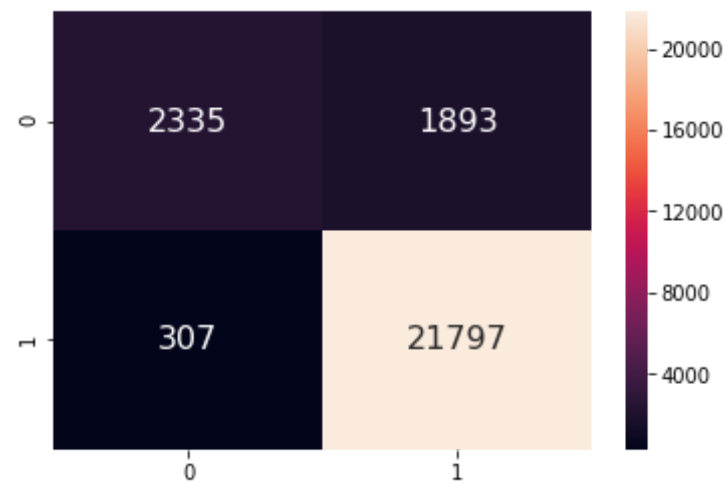




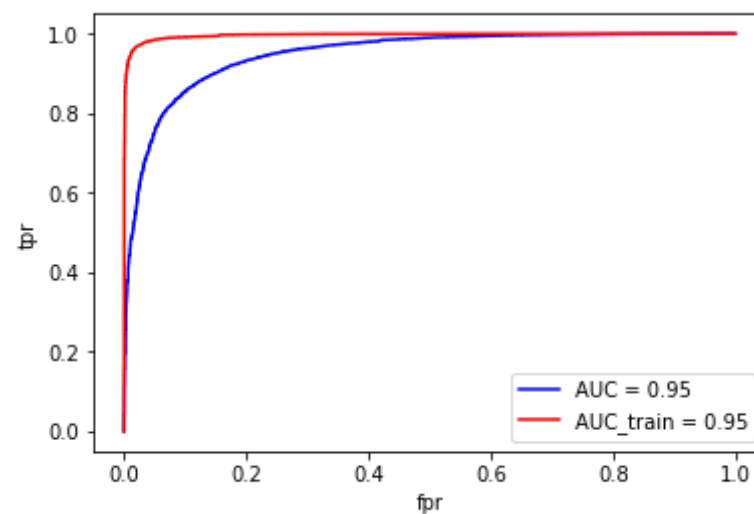
### [5.2.2] Applying Logistic Regression with L2 regularization on TFIDF, SET 2

```
In [49]: # Please write all the code with proper documentation
model_new=LogisticRegression(C=optimal_c)
model_new.fit(final_tf_idf,y_tr)
pred=model_new.predict(X_test_tfidf)
probab=model_new.predict_proba(X_test_tfidf)
probab=probab[:,1]
fpr,tpr,thresholds=metrics.roc_curve(y_test,probab)
auc=roc_auc_score(y_test,probab)
acc=accuracy_score(y_test,pred,normalize=True)*float(100)
pred2=model_new.predict(final_tf_idf)
probab2=model_new.predict_proba(final_tf_idf)
probab2=probab2[:,1]
fpr_train,tpr_train,thresholds=metrics.roc_curve(y_tr,probab2)
auc2=roc_auc_score(y_test,probab)
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), range(2),range(2))
#heatmap for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

```
plt.show()
plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' %auc)
plt.plot(fpr_train, tpr_train, 'r', label = 'AUC_train = %0.2f' %auc2)
plt.legend(loc='lower right')
plt.xlabel('fpr')
plt.ylabel('tpr')
```



Out[49]: Text(0,0.5,'tpr')



### [5.2.3] Feature Importance on TFIDF, SET 2

#### [5.2.3.1] Top 10 important features of positive class from SET 2

```
In [50]: # Please write all the code with proper documentation
# Please write all the code with proper documentation
clf=LogisticRegression(C=optimal_c,penalty='l2')
clf.fit(final_tf_idf,y_tr)
w=clf.coef_
tfidf=tf_idf_vect.get_feature_names()
df=pd.DataFrame(w,columns=tfidf)
df=df.T
print('The top ten positive features are-\n',df[0].sort_values(ascending=False)[0:10])
```

The top ten positive features are-

great	0.099722
love	0.076701
good	0.074882
best	0.069983
delicious	0.062689
loves	0.051904
perfect	0.048291
favorite	0.047066
nice	0.046251
excellent	0.044000

Name: 0, dtype: float64

#### [5.2.3.2] Top 10 important features of negative class from SET 2

```
In [51]: # Please write all the code with proper documentation
print('The top ten negative features are-\n',df[0].sort_values(ascending=True)[0:10])
```

The top ten negative features are-

not worth 0.040000

```

not worth      -0.048039
worst          -0.047183
not buy        -0.046487
disappointed   -0.044877
not good       -0.044158
awful          -0.042917
not recommend  -0.042789
terrible       -0.042111
disappointing  -0.039513
horrible       -0.036940
Name: 0, dtype: float64

```

## [5.3] Logistic Regression on AVG W2V, SET 3

### [5.3.1] Applying Logistic Regression with L1 regularization on AVG W2V SET 3

```

In [52]: # Please write all the code with proper documentation
tuned_parameters=[{'C':[10**-4,10**-3,10**-2,10**-1,1,10**1,10**2,10**3,10**4]}]
model=GridSearchCV(LogisticRegression(),tuned_parameters,scoring='roc_auc',cv=10)
model.fit(sent_vectors_train,y_tr)
print(model.best_estimator_)
optimal_c=model.best_estimator_.C
print(model.score(sent_vectors_cv,y_cv))

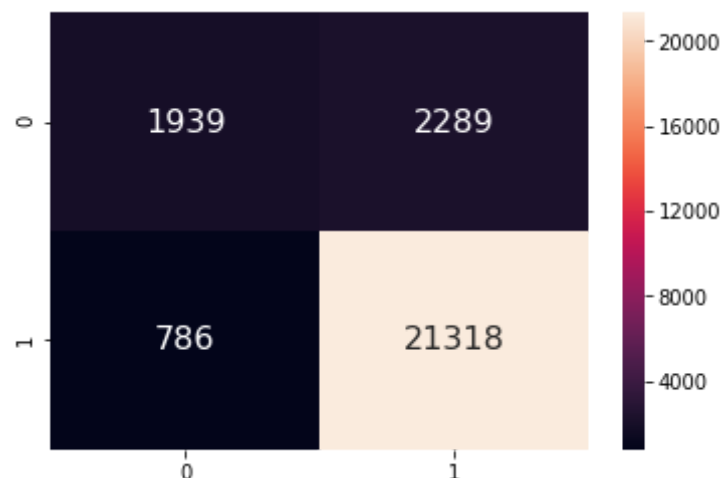
LogisticRegression(C=10, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, l1_ratio=None, max_iter=100,
                    multi_class='warn', n_jobs=None, penalty='l2',
                    random_state=None, solver='warn', tol=0.0001, verbose=0,
                    warm_start=False)
0.903650354808448

```

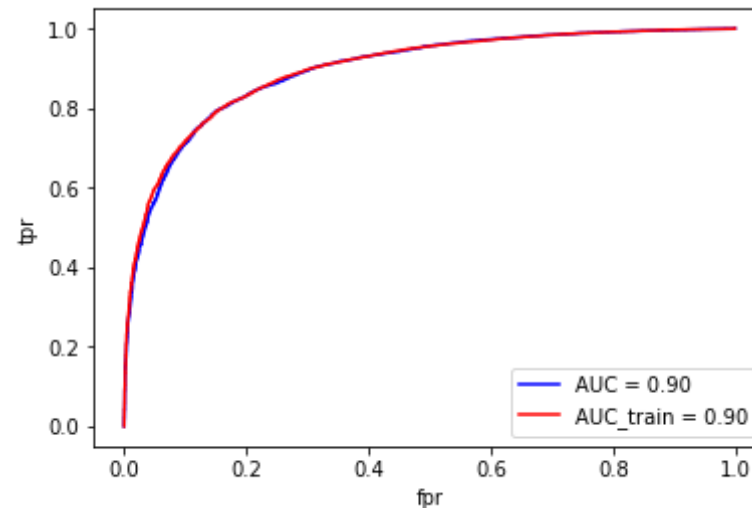
```

In [53]: model_new=LogisticRegression(C=optimal_c,penalty='l1')
model_new.fit(sent_vectors_train,y_tr)
pred=model_new.predict(sent_vectors_test)
probab=model_new.predict_proba(sent_vectors_test)
probab=probab[:,1]
fpr,tpr,thresholds=metrics.roc_curve(y_test,probab)
auc=roc_auc_score(y_test,probab)
acc=accuracy_score(y_test,pred,normalize=True)*float(100)
pred2=model_new.predict(sent_vectors_train)
probab2=model_new.predict_proba(sent_vectors_train)
probab2=probab2[:,1]
fpr_train,tpr_train,thresholds=metrics.roc_curve(y_tr,probab2)
auc2=roc_auc_score(y_test,probab)
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), range(2),range(2))
#heatmap for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
plt.show()
plt.plot(fpr,tpr,'b', label = 'AUC = %0.2f' %auc)
plt.plot(fpr_train,tpr_train,'r', label = 'AUC_train = %0.2f' %auc2)
plt.legend(loc='lower right')
plt.xlabel('fpr')
plt.ylabel('tpr')

```



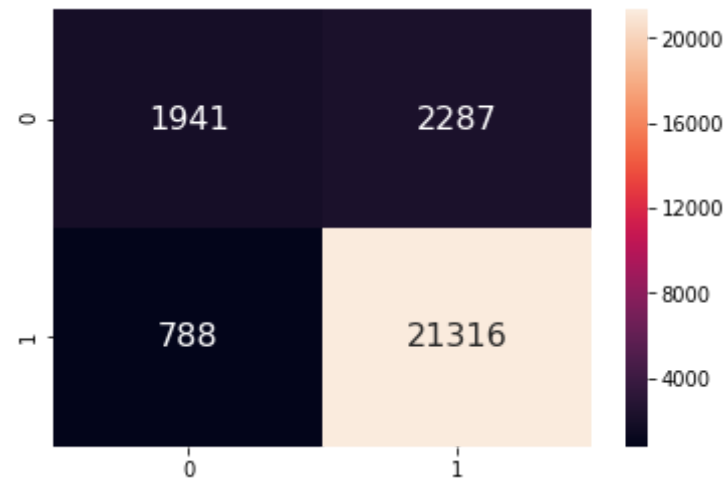
```
Out[53]: Text(0,0.5,'tpr')
```



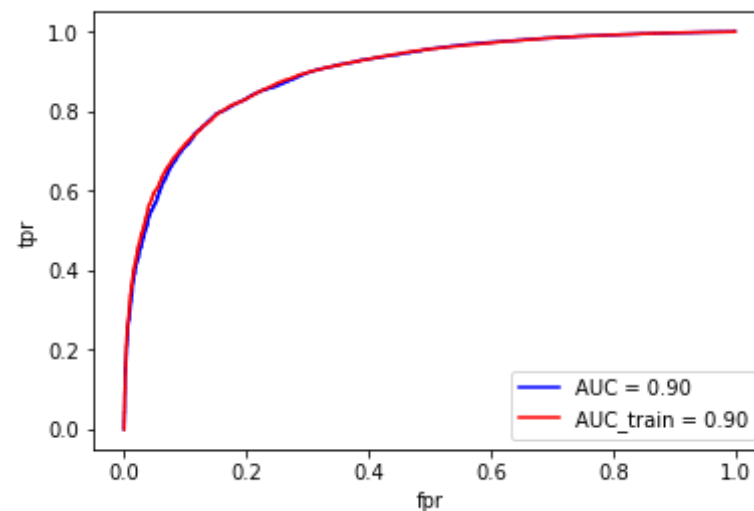
### [5.3.2] Applying Logistic Regression with L2 regularization on AVG W2V, SET 3

```
In [54]: # Please write all the code with proper documentation
model_new=LogisticRegression(C=optimal_c,penalty='l2')
model_new.fit(sent_vectors_train,y_tr)
pred=model_new.predict(sent_vectors_test)
probab=model_new.predict_proba(sent_vectors_test)
probab=probab[:,1]
fpr,tpr,thresholds=metrics.roc_curve(y_test,probab)
auc=roc_auc_score(y_test,probab)
acc=accuracy_score(y_test,pred,normalize=True)*float(100)
pred2=model_new.predict(sent_vectors_train)
probab2=model_new.predict_proba(sent_vectors_train)
probab2=probab2[:,1]
fpr_train,tpr_train,thresholds=metrics.roc_curve(y_tr,probab2)
auc2=roc_auc_score(y_test,probab)
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), range(2),range(2))
#heatmap for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
plt.show()
```

```
plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % auc)  
plt.plot(fpr_train, tpr_train, 'r', label = 'AUC_train = %0.2f' % auc2)  
plt.legend(loc='lower right')  
plt.xlabel('fpr')  
plt.ylabel('tpr')
```



Out[54]: Text(0,0.5,'tpr')



## [5.4] Logistic Regression on TFIDF W2V, SET 4

### [5.4.1] Applying Logistic Regression with L1 regularization on TFIDF W2V, SET 4

```
In [55]: # Please write all the code with proper documentation
tuned_parameters=[{'C':[10**-4,10**-3,10**-2,10**-1,1,10**1,10**2,10**3,10**4]}]
model=GridSearchCV(LogisticRegression(),tuned_parameters,scoring='roc_auc',cv=10)
model.fit( tfidf_sent_vectors_train,y_tr)
print(model.best_estimator_)
optimal_c=model.best_estimator_.C
print(model.score(tfidf_sent_vectors_cv,y_cv))
```

```
LogisticRegression(C=0.1, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, l1_ratio=None, max_iter=100,
                    multi_class='warn', n_jobs=None, penalty='l2',
                    random_state=None, solver='warn', tol=0.0001, verbose=0,
                    warm_start=False)
0.49600304219210767
```

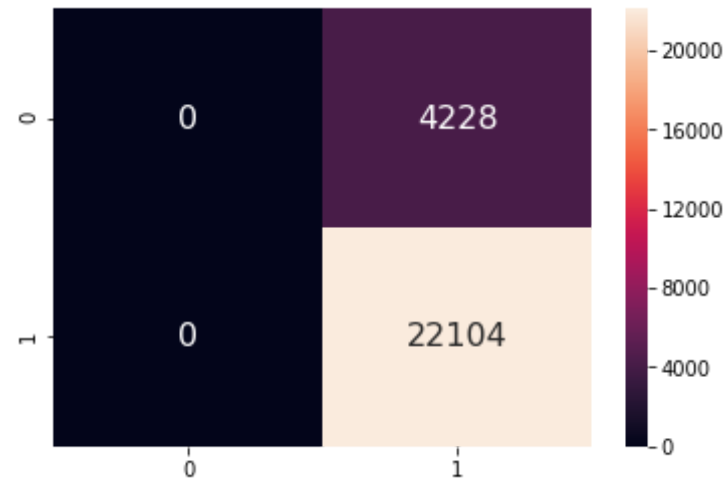
```
In [56]: model_new=LogisticRegression(C=optimal_c,penalty='l1')
model_new.fit(tfidf_sent_vectors_train,y_tr)
pred=model_new.predict(tfidf_sent_vectors_test)
probab=model_new.predict_proba(tfidf_sent_vectors_test)
probab=probab[:,1]
fpr,tpr,thresholds=metrics.roc_curve(y_test,probab)
auc=roc_auc_score(y_test,probab)
acc=accuracy_score(y_test,pred,normalize=True)*float(100)
pred2=model_new.predict(tfidf_sent_vectors_train)
probab2=model_new.predict_proba(tfidf_sent_vectors_train)
probab2=probab2[:,1]
fpr_train,tpr_train,thresholds=metrics.roc_curve(y_tr,probab2)
```



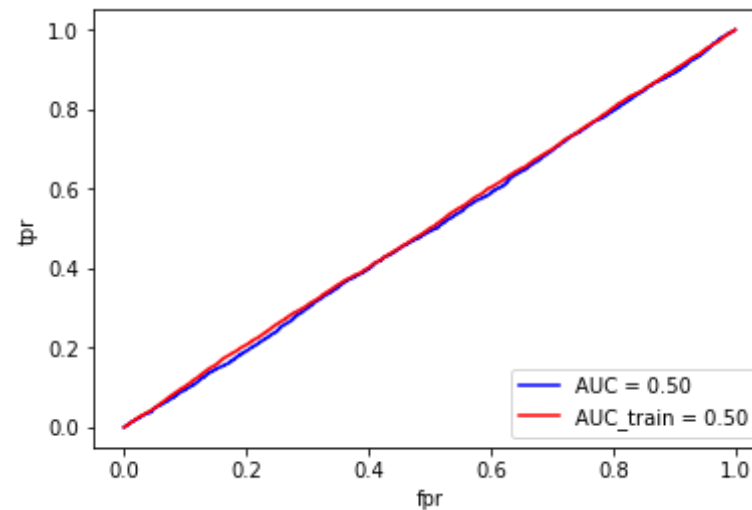
```

auc2=roc_auc_score(y_test,probab)
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), range(2),range(2))
#heatmap for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
plt.show()
plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' %auc)
plt.plot(fpr_train, tpr_train, 'r', label = 'AUC_train = %0.2f' %auc2)
plt.legend(loc='lower right')
plt.xlabel('fpr')
plt.ylabel('tpr')

```



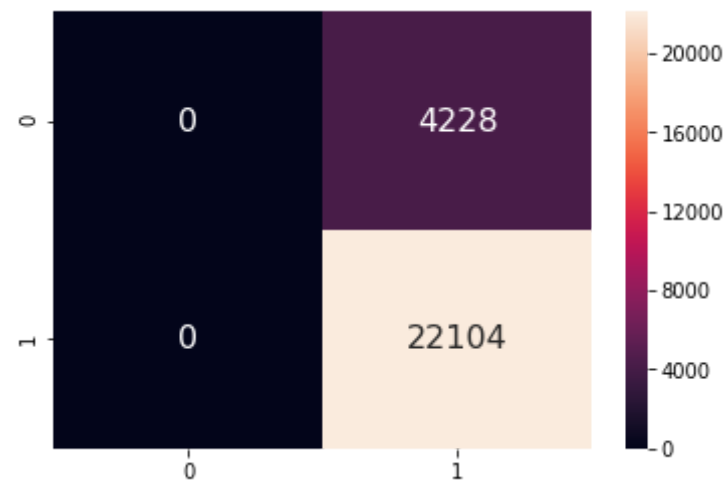
Out[56]: Text(0,0.5,'tpr')



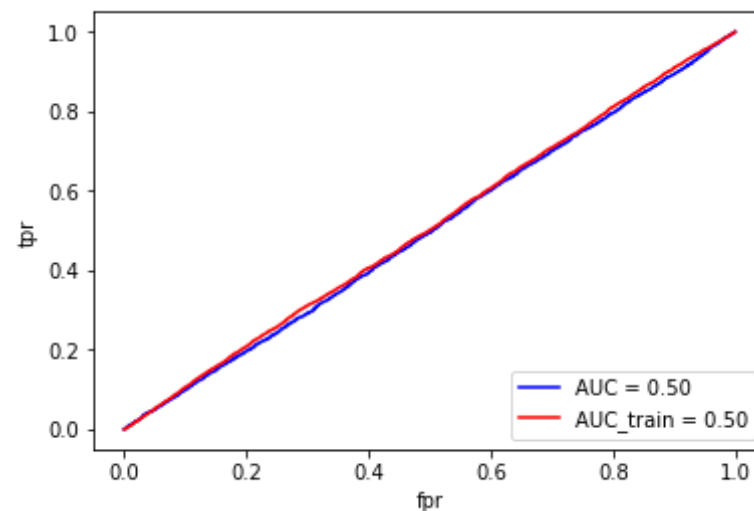
#### [5.4.2] Applying Logistic Regression with L2 regularization on TFIDF W2V, SET 4

```
In [57]: # Please write all the code with proper documentation
model_new=LogisticRegression(C=optimal_c,penalty='l2')
model_new.fit(tfidf_sent_vectors_train,y_tr)
pred=model_new.predict(tfidf_sent_vectors_test)
probab=model_new.predict_proba(tfidf_sent_vectors_test)
probab=probab[:,1]
fpr,tpr,thresholds=metrics.roc_curve(y_test,probab)
auc=roc_auc_score(y_test,probab)
acc=accuracy_score(y_test,pred,normalize=True)*float(100)
pred2=model_new.predict(tfidf_sent_vectors_train)
probab2=model_new.predict_proba(tfidf_sent_vectors_train)
probab2=probab2[:,1]
fpr_train,tpr_train,thresholds=metrics.roc_curve(y_tr,probab2)
auc2=roc_auc_score(y_test,probab)
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), range(2),range(2))
#heatmap for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
plt.show()
```

```
plt.plot(fpr,tpr,'b', label = 'AUC = %0.2f' %auc)
plt.plot(fpr_train,tpr_train,'r', label = 'AUC_train = %0.2f' %auc2)
plt.legend(loc='lower right')
plt.xlabel('fpr')
plt.ylabel('tpr')
```



Out[57]: Text(0,0.5,'tpr')



## [6] Conclusions

```
In [58]: # Please compare all your models using Prettytable library
from prettytable import PrettyTable
x=PrettyTable()
x.field_names=(['vectorizer','regularization','hyperparameter','AUC'])
x.add_row(['bag of words','l1',0.001,0.84])
x.add_row(['bag of words','l2',0.001,0.91])
x.add_row(['TFIDF','l1',0.0001,0.50])
x.add_row(['TFIDF','l2',0.0001,0.95])
x.add_row(['Avg W2V','l1',10,0.90])
x.add_row(['AVG W2V','l2',10,0.90])
x.add_row(['Weighted TFIDF','l1',0.1,0.50])
x.add_row(['Weighted TFIDF','l2',0.1,0.50])
print(x)
```

vectorizer	regularization	hyperparameter	AUC
bag of words	l1	0.001	0.84
bag of words	l2	0.001	0.91
TFIDF	l1	0.0001	0.5
TFIDF	l2	0.0001	0.95
Avg W2V	l1	10	0.9
AVG W2V	l2	10	0.9
Weighted TFIDF	l1	0.1	0.5
Weighted TFIDF	l2	0.1	0.5