

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 os
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import nltk
import string
```

```

from sklearn.cross_validation import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.cross_validation import cross_val_score
from collections import Counter
from sklearn.metrics import accuracy_score
from sklearn import cross_validation
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
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 sklearn.naive_bayes import BernoulliNB
from sklearn.naive_bayes import GaussianNB
from sklearn.naive_bayes import MultinomialNB

from tqdm import tqdm
from sklearn.metrics import precision_recall_fscore_support
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.metrics import mean_squared_error
from sklearn.metrics import accuracy_score
os.getcwd()

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\cross_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.

"This module will be removed in 0.20.", DeprecationWarning)

C:\ProgramData\Anaconda3\lib\site-packages\gensim\utils.py:1209: UserWarning: detected Windows; aliasing chunkize to chunkize_serial

warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

Out[1]:

'C:\\Users\\admin\\Desktop\\AAIC_updated\\Assignments'

In [2]:

```

# using SQLite Table to read data.
con = sqlite3.connect('database.sqlite')

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

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

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

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

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

```

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

Out[2]:

		Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
0	1		B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	1	1303862400
1	2		B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	0	1346976000
2	3		B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	1	1219017600

In [3]:

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

In [4]:

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

(80668, 7)

Out[4]:

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

In [5]:

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

Out[5]:

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

In [6]:

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

Out[6]:

393063

[2] Exploratory Data Analysis

[2.1] Data Cleaning: Deduplication

It is observed (as shown in the table below) that the reviews data had many duplicate entries. Hence it was necessary to remove duplicates in order to get unbiased results for the analysis of the data. Following is an example:

In [7]:

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

Out[7]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Ti
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2	5	11995776
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2	5	11995776
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2	5	11995776
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2	5	11995776
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2	5	11995776

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	HelpfulnessDenominator	Score	Time
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5	12248926

1	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
	44737	B001EQ55RW	A2V0I904EH7ABY	Ram	3	2	4	12128832

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 ver y hard to find any chicken products made in the USA but they are out there, but this one isnt. It s too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

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

=====

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

=====

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

=====

In [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 because 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 because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

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

=====

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

=====

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

In [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 ver y hard to find any chicken products made in the USA but they are out there, but this one isnt. It s too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

In [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', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have revmoved in the 1st step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "y
ou're", "you've", \
    "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
'himself', \
    'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them',
'their', \
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll",
'these', 'those', \
    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
'do', 'does', \
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', '
while', 'of', \
    'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during',
'before', 'after', \
```


In [22]:

```
100%|██████████████████████████████████████████████████████████████████████████████| 87773/87773  
[00:39<00:00, 2241.37it/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.
```

In [25]:

```
preprocessed_reviews_fe = []
for i in preprocessed_reviews:
    count = 0;
    for j in i:
        count += 1;
    i = i + " " + str(count);
    preprocessed_reviews_fe.append(i)
```

[4] Featurization

In [26]:

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(preprocessed_reviews_fe, final['Score'], test_size=0.33) # this is random splitting
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33) # this is random splitting
```

[4.1] BAG OF WORDS

In [27]:

```
#BoW
count_vect = CountVectorizer(min_df=10, max_features=2000) #in scikit-learn
```

```
X_train_vect = count_vect.fit_transform(X_train)
X_train_vect = X_train_vect.toarray()
X_cv_vect = count_vect.transform(X_cv)
X_cv_vect = X_cv_vect.toarray()
X_test_vect = count_vect.transform(X_test)
X_test_vect = X_test_vect.toarray()
print("some feature names ", count_vect.get_feature_names()[:10])
print('='*50)
```

```
final_counts = count_vect.transform(preprocessed_reviews)
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])
```

```
some feature names  ['100', '101', '102', '103', '104', '105', '106', '107', '108', '109']
=====
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (87773, 2000)
the number of unique words  2000
```

[4.2] Bi-Grams and n-Grams.

In []:

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

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

```
# you can choose these numebtrs min_df=10, max_features=5000, of your choice
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
final_bigram_counts = count_vect.fit_transform(preprocessed_reviews)
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])
```

[4.3] TF-IDF

In [29]:

```
#tf_idf_vect.fit(preprocessed_reviews)
tf_idf = TfidfVectorizer(min_df=10, max_features=2000)
X_train_vect_tfidf = tf_idf.fit_transform(X_train)
X_train_vect_tfidf = X_train_vect_tfidf.toarray()
X_test_vect_tfidf = tf_idf.transform(X_test)
X_test_vect_tfidf = X_test_vect_tfidf.toarray()
X_cv_vect_tfidf = tf_idf.transform(X_cv)
X_cv_vect_tfidf = X_cv_vect_tfidf.toarray()
```

[4.4] Word2Vec

In []:

```
# Train your own Word2Vec model using your own text corpus
i=0
```

```
list_of_sentence=[]
for sentence in preprocessed_reviews:
    list_of_sentence.append(sentence.split())
```

In []:

```
# Using Google News Word2Vectors

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

# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.Wl7SRFazZPY
# 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 occurred atleast 5 times
    w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))

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

In []:

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

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

[4.4.1.1] Avg W2v

In []:

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

```
print(len(sent_vectors[0]))
```

[4.4.1.2] TFIDF weighted W2v

In []:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(preprocessed_reviews)
# 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 []:

```
# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf

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

[5] Assignment 4: Apply Naive Bayes

1. Apply Multinomial NaiveBayes 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)

2. The hyper parameter tuning (find best Alpha)

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Consider a wide range of alpha values for hyperparameter tuning, start as low as 0.00001
- Find the best hyper parameter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. Feature importance

- Find the top 10 features of positive class and top 10 features of negative class for both feature sets **Set 1** and **Set 2** using values of 'feature_log_prob_' parameter of [MultinomialNB](#) and print their corresponding feature names

4. 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.

5. 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. Here on X-axis you will have alpha values. since they have a wide range. just to represent

like shown in the figure. Here on X axis you will have alpha values, since they have a wide range, just to represent those alpha values on the graph, apply log function on those alpha values.

- 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](#).

6. 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 you train data, and apply the method transform() on cv/test data.
4. For more details please go through this [link](#).

Applying Multinomial Naive Bayes

[5.1] Applying Naive Bayes on BOW, SET 1

GridSearch CV

In [30]:

```
from sklearn.grid_search import GridSearchCV
from sklearn.linear_model import LogisticRegression
h_GSCV = [{'alpha': np.linspace(0.0001,5,100)}]
model = GridSearchCV(MultinomialNB(),h_GSCV, scoring = 'roc_auc')
model.fit(X_train_vect, y_train)
print(model.best_estimator_)
print(model.score(X_test_vect, y_test))
```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\grid_search.py:42: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. This module will be removed in 0.20.
DeprecationWarning)

MultinomialNB(alpha=3.3838707070707073, class_prior=None, fit_prior=True)
0.9222847348573702

Simple CV

In [31]:

```
from sklearn.metrics import accuracy_score
from sklearn.metrics import roc_auc_score

# Please write all the code with proper documentation
score = []
ind = []
dict_new = {}
metric_value_train = []
metric_value_cv = []
train_auc = []
cv_auc = []
h = np.linspace(0.0001,5,100)
for i in tqdm(h):
    nb = MultinomialNB(alpha=i)
    nb.fit(X_train_vect, y_train)
```

```

pred_tr = nb.predict(X_train_vect)
pred = nb.predict(X_cv_vect)
metric_value_train.append(accuracy_score(y_train,pred_tr))
metric_value_cv.append(accuracy_score(y_cv,pred))
score_measure = precision_recall_fscore_support(y_cv, pred, average = 'weighted')
score.append(score_measure[2])
ind.append(i)
y_train_pred = nb.predict_proba(X_train_vect)[:,-1]
y_cv_pred = nb.predict_proba(X_cv_vect)[:,-1]
train_auc.append(roc_auc_score(y_train,y_train_pred))
cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.figure(1)
plt.plot(metric_value_train,metric_value_cv,'o')
plt.xlabel('Train Accuracy')
plt.ylabel('CV Accuracy')

plt.figure(2)
plt.plot(h,score)
plt.xlabel('Hyperparameter alpha')
plt.ylabel('f-1 score')

optimal_h = ind[score.index(max(score))]
print('\nThe best hyperparameter is (according to f-1 score):', optimal_h)

optimal_h_auc = ind[cv_auc.index(max(cv_auc))]
print('\nThe optimal hyperparameter is (according to auc curve (max auc)): ', optimal_h_auc)

print("\n")

plt.figure(3)
plt.plot(h, train_auc, label='Train AUC')
plt.plot(h, cv_auc, label='CV AUC')
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()

plt.figure(4)
plt.scatter(h, train_auc, label='Train AUC')
plt.scatter(h, cv_auc, label='CV AUC')
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()

```

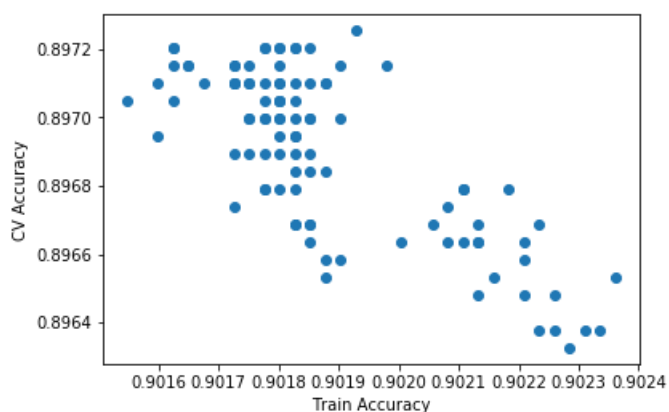
```

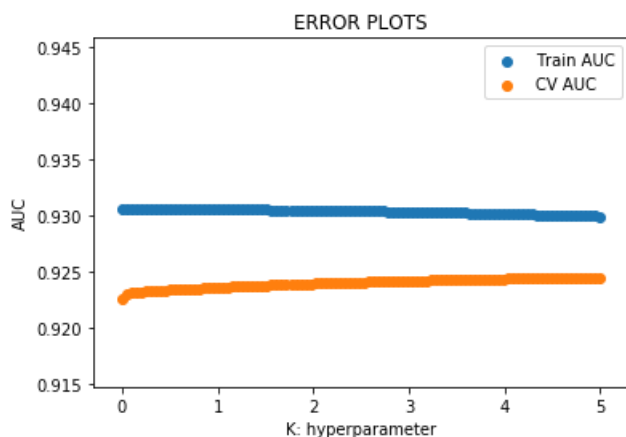
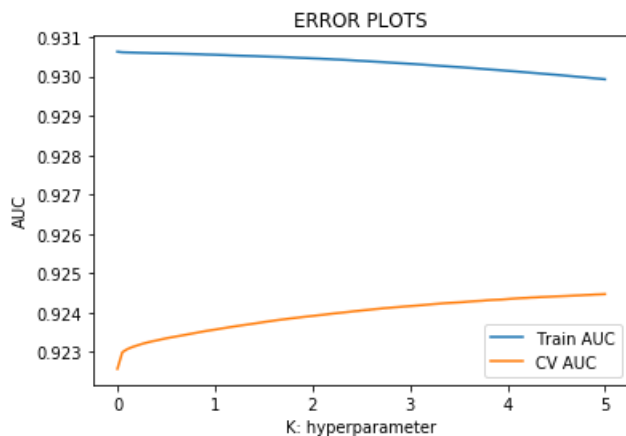
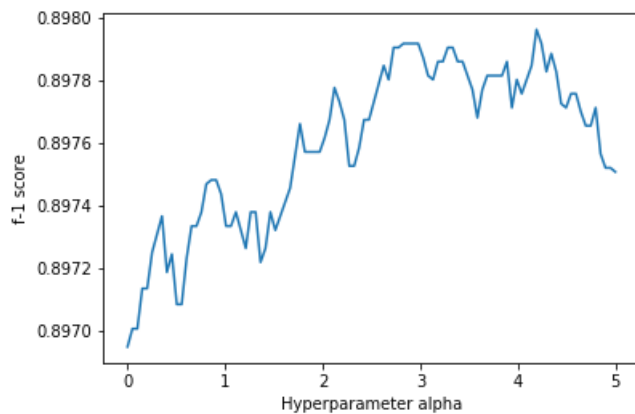
100% |████████████████████████████████████████████████████████████████████████████████| 100/100 [03
:54<00:00, 2.34s/it]

```

The best hyperparameter is (according to f-1 score): 4.191935353535354

The optimal hyperparameter is (according to auc curve (max auc)): 5.0





In [32]:

```
nb = MultinomialNB(alpha=optimal_h_auc)
nb.fit(X_train_vect,y_train)
pred_test = nb.predict(X_test_vect)
score_measure_test = precision_recall_fscore_support(y_test, pred_test, average = 'weighted')
acc = accuracy_score(y_test, pred_test, normalize=True) * float(100)
print('\n****Test accuracy for alpha = %f is %f%%' % (optimal_h_auc,acc))
print('\n****Test f-1 score for alpha = %f is %f%%' % (optimal_h_auc,score_measure_test[2]))
train_fpr, train_tpr, thresholds = roc_curve(y_train, nb.predict_proba(X_train_vect)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, nb.predict_proba(X_test_vect)[: ,1])
print("The Confusion Matrix is as following:")
print(confusion_matrix(list(y_test),list(pred_test),labels = [1,0]))
print("\n")
print("The classification Report can be presented as follows:")
print(classification_report(y_test, pred_test,labels = [0,1]))

plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
```

```
plt.show()

plt.scatter(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.scatter(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()

print("="*100)
```

****Test accuracy for alpha = 5.000000 is 89.370296%

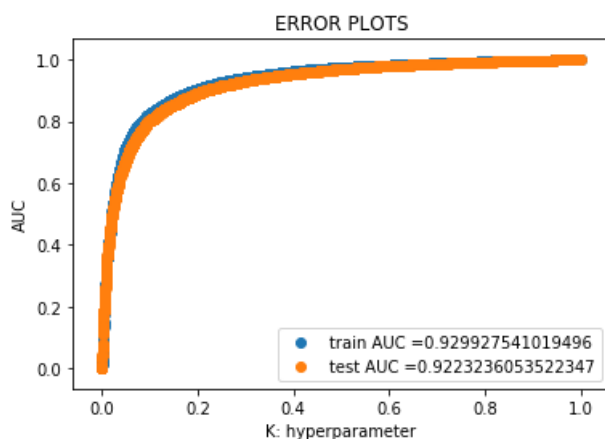
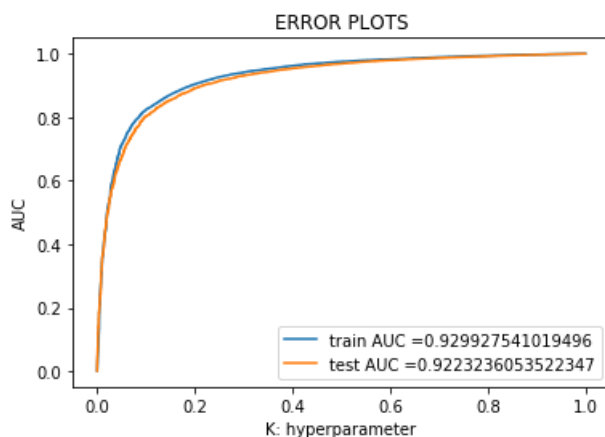
****Test f-1 score for alpha = 5.000000 is 0.893991%

The Confusion Matrix is as following:

```
[[22621  1572]
 [ 1507  3266]]
```

The classification Report can be presented as follows:

	precision	recall	f1-score	support
0	0.68	0.68	0.68	4773
1	0.94	0.94	0.94	24193
avg / total	0.89	0.89	0.89	28966



In [33]:

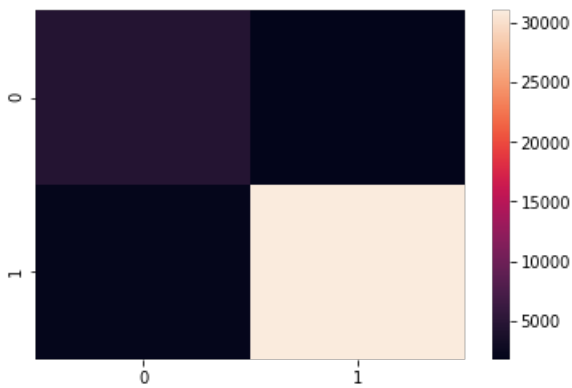
```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
sns.heatmap(confusion_matrix(y_train, nb.predict(X_train_vect)))
```

Train confusion matrix

Train Confusion Matrix

Out[33]:

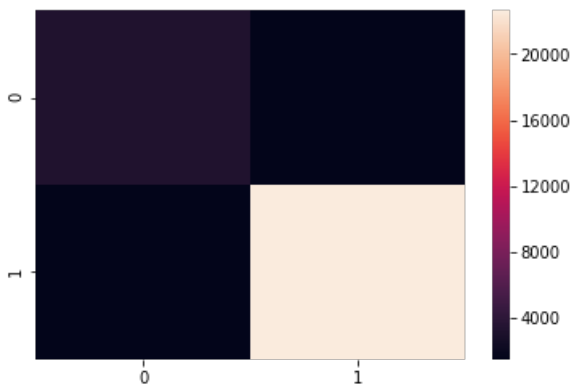
<matplotlib.axes._subplots.AxesSubplot at 0x1f5392be208>



In [34]:

```
print("Test confusion matrix")
print(sns.heatmap(confusion_matrix(y_test, nb.predict(X_test_vect))))
```

Test confusion matrix
AxesSubplot(0.125,0.125;0.62x0.755)



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

In [35]:

```
# Please write all the code with proper documentation
vocab = list(count_vect.get_feature_names())
len_of_vocab = len(vocab)
probs = nb.feature_log_prob_
dict_new = {'Words':vocab,'neg_prob':probs[0],'pos_prob' : probs[1]}
df_new = pd.DataFrame(dict_new)
df_neg_sorted =df_new.sort_values('neg_prob', axis=0, ascending=False, inplace=False, kind='quicksort', na_position='last')
df_pos_sorted =df_new.sort_values('pos_prob', axis=0, ascending=False, inplace=False, kind='quicksort', na_position='last')
neg_words = df_neg_sorted['Words']
pos_words = df_pos_sorted['Words']
print("Following are the best features of the positive class")
for i in pos_words[0:10]:
    print(i)
```

Following are the best features of the positive class
not
like
good
great
one

```
taste
coffee
flavor
love
would
```

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

In [36]:

```
# Please write all the code with proper documentation
print("Following are the best features of the negative class")
for i in neg_words[0:10]:
    print(i)
```

Following are the best features of the negative class

```
not
like
would
taste
product
one
coffee
good
flavor
no
```

[5.2] Applying Naive Bayes on TFIDF, SET 2

GridSearch CV

In [37]:

```
# Please write all the code with proper documentation
from sklearn.grid_search import GridSearchCV
from sklearn.linear_model import LogisticRegression
h_GSCV = [{'alpha': np.linspace(0.0001,5,100)}]

model = GridSearchCV(MultinomialNB(),h_GSCV, scoring = 'roc_auc')
model.fit(X_train_vect_tfidf, y_train)

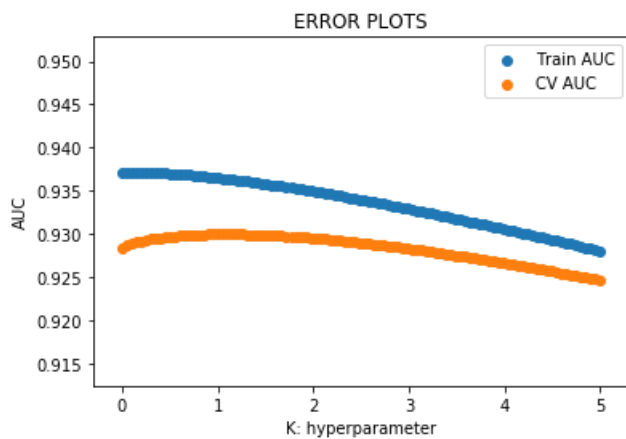
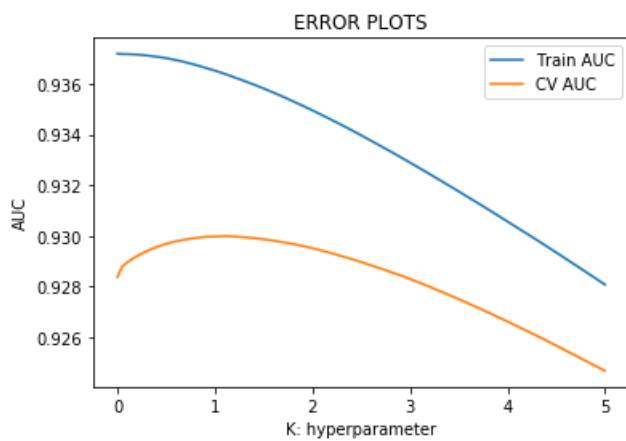
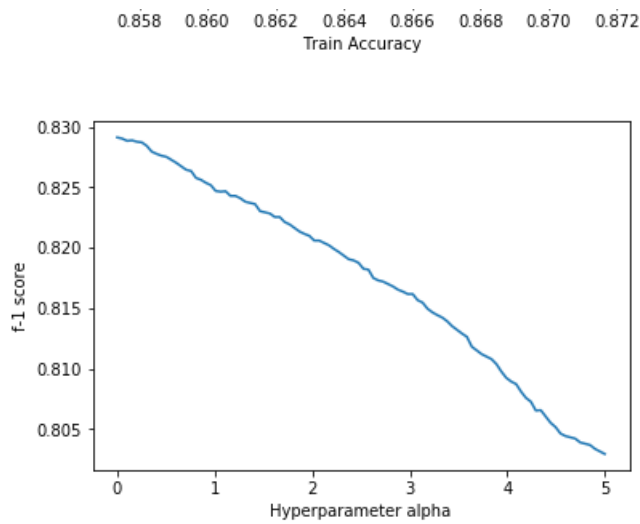
print(model.best_estimator_)
print(model.score(X_test_vect_tfidf, y_test))
```

```
MultinomialNB(alpha=0.7071565656565657, class_prior=None, fit_prior=True)
0.9263865484826957
```

In [38]:

```
from sklearn.metrics import accuracy_score
from sklearn.metrics import roc_auc_score

# Please write all the code with proper documentation
score = []
ind = []
dict_new = {}
metric_value_train = []
metric_value_cv = []
train_auc = []
cv_auc = []
h = np.linspace(0.0001,5,100)
for i in tqdm(h):
    nb = MultinomialNB(alpha=i)
    nb.fit(X_train_vect_tfidf, y_train)
    pred_tr = nb.predict(X_train_vect_tfidf)
    pred = nb.predict(X_cv_vect_tfidf)
    metric_value_train.append(accuracy_score(y_train,pred_tr))
    metric_value_cv.append(accuracy_score(y_train,pred))
```

The Confusion Matrix is as following:

```
[[ 3266  1507]
 [ 1572 22621]]
```

The classification Report can be presented as follows:

	precision	recall	f1-score	support
0	0.68	0.68	0.68	4773
1	0.94	0.94	0.94	24193
avg / total	0.89	0.89	0.89	28966

In [39]:

```
nb = MultinomialNB(alpha=optimal_h_auc)
nb.fit(X_train_vect_tfidf,y_train)
pred_test = nb.predict(X_test_vect_tfidf)
score_measure_test = precision_recall_fscore_support(y_test, pred_test, average = 'weighted')
```

```

score_measure_test = precision_recall_fscore_support(y_test, pred_test, average = 'weighted',
acc = accuracy_score(y_test, pred_test, normalize=True) * float(100)
print('\n****Test accuracy for alpha = %f is %f%%' % (optimal_h_auc, acc))
print('\n****Test f-1 score for alpha = %f is %f%%' % (optimal_h_auc, score_measure_test[2]))
train_fpr, train_tpr, thresholds = roc_curve(y_train, nb.predict_proba(X_train_vect_tfidf)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, nb.predict_proba(X_test_vect_tfidf)[: ,1])

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()

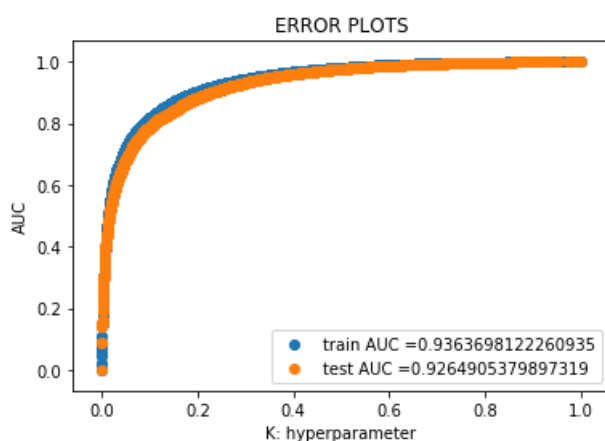
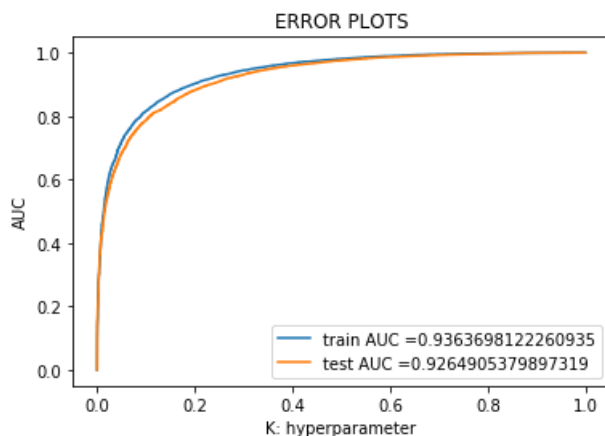
plt.scatter(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.scatter(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()

print("="*100)

```

****Test accuracy for alpha = 1.111189 is 86.083684%

****Test f-1 score for alpha = 1.111189 is 0.817550%



In [40]:

```

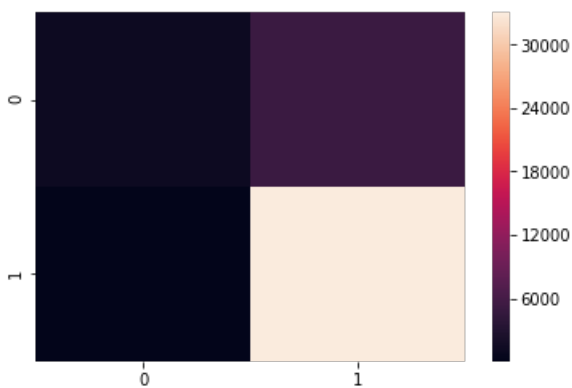
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
sns.heatmap(confusion_matrix(y_train, nb.predict(X_train_vect_tfidf)))

```

Train confusion matrix

Out[40]:

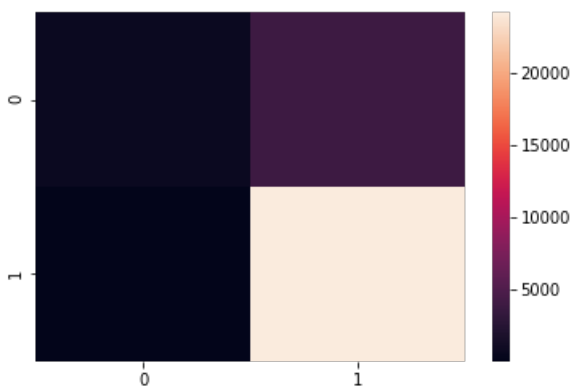
<matplotlib.axes._subplots.AxesSubplot at 0x1f56b743748>



In [41]:

```
print("Test confusion matrix")
print(sns.heatmap(confusion_matrix(y_test, nb.predict(X_test_vect_tfidf))))
```

Test confusion matrix
AxesSubplot(0.125,0.125;0.62x0.755)



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

In [42]:

```
# Please write all the code with proper documentation
# Please write all the code with proper documentation
vocab = list(tf_idf.get_feature_names())
len_of_vocab = len(vocab)
probs = nb.feature_log_prob_
dict_new = {'Words':vocab,'neg_prob':probs[0],'pos_prob' : probs[1]}
df_new = pd.DataFrame(dict_new)
df_neg_sorted =df_new.sort_values('neg_prob', axis=0, ascending=False, inplace=False, kind='quicksort', na_position='last')
df_pos_sorted =df_new.sort_values('pos_prob', axis=0, ascending=False, inplace=False, kind='quicksort', na_position='last')
neg_words = df_neg_sorted['Words']
pos_words = df_pos_sorted['Words']
print("Following are the best features of the positive class")
for i in pos_words[0:10]:
    print(i)
```

Following are the best features of the positive class

not
great
good
like
coffee

```
love
tea
one
taste
product
```

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

In [43]:

```
# Please write all the code with proper documentation
# Please write all the code with proper documentation
print("Following are the best features of the negative class")
for i in neg_words[0:10]:
    print(i)
print("\n")
```

Following are the best features of the negative class

```
not
like
product
taste
would
coffee
one
flavor
no
good
```

[6] Conclusions

In [44]:

```
# Please compare all your models using Prettytable library
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Vectorizer", "Type", "Hyperparameter (alpha)", "AUC"]
x.add_row(["BOW", "MultinomialNB", 5, 0.922])
x.add_row(["TFIDF", "MultinomialNB", 1.11, 0.926])

print(x)
```

Vectorizer	Type	Hyperparameter (alpha)	AUC
BOW	MultinomialNB	5	0.922
TFIDF	MultinomialNB	1.11	0.926

- f-1 score mainly decreases as the alpha increases.
- As the train accuracy increases, the CV accuracy also increases
- Here from the confusion matrix, it clearly seen that positive class is the dominant class.