

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 wil be set to "positive". Otherwise, it will be set to "negative".

```
In [3]: %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.metrics import accuracy_score
import seaborn as sn
from sklearn.metrics import classification_report

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

from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import validation_curve

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

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

from tqdm import tqdm
import os
```

```
In [4]: #loading train and test and validation dataset

file=open("x_train.pkl","rb")
x_train=pickle.load(file) # loading 'train' dataset

file=open("x_cv.pkl",'rb')
x_cv=pickle.load(file) # loading 'validation' dataset

file=open("x_test.pkl",'rb')
x_test=pickle.load(file) # loading 'test' dataset
```

```

file=open("y_train.pkl","rb")
y_train=pickle.load(file) # loading 'train' dataset

file=open("y_cv.pkl",'rb')
y_cv=pickle.load(file) # loading 'validation' dataset

file=open("y_test.pkl",'rb')
y_test=pickle.load(file) # loading 'test' dataset

#loading train_bow and test_bow
file=open('x_train_bow.pkl','rb')
x_train_bow=pickle.load(file)

file=open('x_test_bow.pkl','rb')
x_test_bow=pickle.load(file)

file=open('x_cv_bow.pkl','rb')
x_cv_bow=pickle.load(file)

#loading train_tf_idf and test_tf_idf
file=open('train_tf_idf.pkl','rb')
train_tf_idf=pickle.load(file)

file=open('cv_tf_idf.pkl','rb')
cv_tf_idf=pickle.load(file)

file=open('test_tf_idf.pkl','rb')
test_tf_idf=pickle.load(file)

#loading train_w2v and test_w2v
file=open('train_w2v.pkl','rb')
train_w2v=pickle.load(file)

file=open('cv_w2v.pkl','rb')
cv_w2v=pickle.load(file)

file=open('test_w2v.pkl','rb')
test_w2v=pickle.load(file)

#loading train_tf_idf_w2v and test_tf_idf_w2v
file=open('train_tf_idf_w2v.pkl','rb')
train_tf_idf_w2v=pickle.load(file)

file=open('cv_tf_idf_w2v.pkl','rb')
cv_tf_idf_w2v=pickle.load(file)

file=open('test_tf_idf_w2v.pkl','rb')
test_tf_idf_w2v=pickle.load(file)

```

In [7]: # using csv Table to read data.

```
dataset=pd.read_csv("Reviews.csv")
```

```
print(dataset.shape)
dataset.head(3)
```

```
(568454, 10)
```

Out[7]:

Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
-----------	------------------	---------------	--------------------	-----------------------------	-------------------------------

1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0
---	---	------------	----------------	--------	---	---

2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1
---	---	------------	---------------	--	---	---

```
In [8]: # 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)  
  
# taking reviews whose score is not equal to 3  
filtered_dataset=dataset[dataset['Score']!=3]  
filtered_dataset.shape  
  
#creating a function to filter the reviews (if score>3 --> positive , if score <3 --> negative)  
def partition(x):  
    if x>3:  
        return 1  
    return 0  
  
#changing reviews with score less than 3 to be positive and vice-versa  
actualScore = filtered_dataset['Score']  
positiveNegative = actualScore.map(partition)  
filtered_dataset['Score'] = positiveNegative  
print("Number of data points in our data", filtered_dataset.shape)  
filtered_dataset.head(3)
```

Number of data points in our data (525814, 10)

Out[8]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	\$
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	1
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	0

2	3	B000LQOCHO	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1
---	---	------------	---------------	--	---	---

[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:

observation:

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 [9]: # sorting the value
sorted_data=filtered_dataset.sort_values(by='Id',inplace=True )
#finding the duplicate values using 'df.duplicated'
filtered_dataset[filtered_dataset.duplicated(subset=['ProfileName','HelpfulnessNumerator','HelpfulnessDenominator','Score','Time'])].shape
#alternate way to drop duplicate values
dataset_no_dup=filtered_dataset.drop_duplicates(subset=['ProfileName','Score','Time','Summary'],keep='first')
print(f"before {dataset.shape}")
print(f"after removing duplicate values-->shape = {dataset_no_dup.shape}")
# %age of no. of review reamin in data set
print('percentage of data reamin after removing duplicate values and removing reviews with neutral scores % .2f'
      %((dataset_no_dup.size/dataset.size)*100))

before (568454, 10)
after removing duplicate values-->shape = (363255, 10)
```

percentage of data remain after removing duplicate values and removing reviews with neutral scores 63.90

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 [10]: # removing reviews where "HelpfulnessNumerator>HelpfulnessDenominator"
final=dataset_no_dup[dataset_no_dup['HelpfulnessNumerator']<=dataset_no_dup['HelpfulnessDenominator']]
```

```
In [11]: #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()
```

```
Out[11]: 1      306222  
          0      57031  
          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 [8]: # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all-tags-from-an-element
from bs4 import BeautifulSoup

# https://stackoverflow.com/a/47091490/4084039
import re

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

# 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 removed in the 1st
step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours',
'ourselves', 'you', "you're", "you've", \
               "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
               'him', 'his', 'himself', \
               'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself',
               'they', 'them', 'their', \
               'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that',
               "that'll", 'these', 'those', \
               'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have',
               'has', 'had', 'having', 'do', 'does', \
               'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because',
               'as', 'until', 'while', 'of', \
               'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
               'through', 'during', 'before', 'after', \
               'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on',
               'off', 'over', 'under', 'again', 'further', \
               'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how',
               'all', 'any', 'both', 'each', 'few', 'more', \
               'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than',
               'too', 'very', \
               's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've",
               'now', 'd', 'll', 'm', 'o', 're', \
               've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
               "didn't", 'doesn', "doesn't", 'hadn', \
               'hadn't', 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma',
               'mightn', "mightn't", 'mustn', \
               'mustn't', 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't",
               'wasn', "wasn't", 'weren', "weren't", \
               'won', "won't", 'wouldn', "wouldn't"])

```

In [9]: # Combining all the above students

```

from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentance in tqdm(final['Text'].values):
    sentance = re.sub(r"http\S+", "", sentance)
    sentance = BeautifulSoup(sentance, 'lxml').get_text()
    sentance = decontracted(sentance)
    sentance = re.sub("\S*\d\S*", "", sentance).strip()
    sentance = re.sub('[^A-Za-z]+', ' ', sentance)

```

```
# https://gist.github.com/sebleier/554280
sentance = ' '.join(e.lower() for e in sentance.split() if e.lower() not in stopwords)
preprocessed_reviews.append(sentance.strip())
```

100% [██████████] | 363253/363253 [02:56<00:00, 2063.41it/s]

In [10]: preprocessed_reviews[:5]

Out[10]: ['bought several vitality canned dog food products found good quality product looks like stew processed meat smells better labrador finicky appreciates product better', 'product arrived labeled jumbo salted peanuts peanuts actually small sized unsalted not sure error vendor intended represent product jumbo', 'confection around centuries light pillow citrus gelatin nuts case filberts cut tiny squares liberally coated powdered sugar tiny mouthful heaven not chewy flavorful highly recommend yummy treat familiar story c lewis lion which wardrobe treat seduces edmund selling brother sisters witch', 'looking secret ingredient robitussin believe found got addition root beer extract ordered good made cherry soda flavor medicinal', 'great taffy great price wide assortment yummy taffy delivery quick taffy 1 over deal']

In [11]: final['preprocessed_reviews']=preprocessed_reviews

In [12]: # splitting the dataset in train , cv and test

```
n=final.shape[0] #size of final dataset
train=final.iloc[:round(0.60*n),:]
cv=final.iloc[round(0.60*n):round(0.80*n),:]
test=final.iloc[round(0.80*n):round(1.0*n),:]
```

In [13]: from sklearn.model_selection import train_test_split
x_training,x_test,y_training,y_test=train_test_split(preprocessed_reviews,final["Score"], test_size=0.20, random_state=42)
x_train,x_cv,y_train,y_cv=train_test_split(x_training,y_training, test_size=0.25, random_state=42)

In [48]: len(x_train),len(y_train),len(x_test),len(y_test),len(x_cv),len(y_cv)

Out[48]: (217951, 217951, 72651, 72651, 72651, 72651)

In [14]: # saving train and test dataset using pickle for future use

```
'''file=open("x_train.pkl","wb")
pickle.dump(x_train,file)
file.close()

file=open('x_cv.pkl','wb')
pickle.dump(x_cv,file)
file.close

file=open("x_test.pkl",'wb')
pickle.dump(x_test,file)
file.close()'''
```

```

file=open("y_train.pkl","wb")
pickle.dump(y_train,file)
file.close()

file=open('y_cv.pkl','wb')
pickle.dump(y_cv,file)
file.close()

file=open("y_test.pkl",'wb')
pickle.dump(y_test,file)
file.close()

'''
```

In [2]:

```

#loading train and test and validation dataset

file=open("x_train.pkl","rb")
x_train=pickle.load(file) # loading 'train' dataset

file=open("x_cv.pkl",'rb')
x_cv=pickle.load(file) # loading 'validation' dataset

file=open("x_test.pkl",'rb')
x_test=pickle.load(file) # loading 'test' dataset

file=open("y_train.pkl","rb")
y_train=pickle.load(file) # loading 'train' dataset

file=open("y_cv.pkl",'rb')
y_cv=pickle.load(file) # loading 'validation' dataset

file=open("y_test.pkl",'rb')
y_test=pickle.load(file) # loading 'test' dataset
```

[4] Featurization

[4.1] BAG OF WORDS

In [71]:

```

#BoW

count_vect = CountVectorizer() #in scikit-learn

x_train_bow=count_vect.fit_transform(x_train)
print("some feature names ", count_vect.get_feature_names() [:10])
print('='*100)

# transform cv and test dataset
x_cv_bow=count_vect.transform(x_cv)
x_test_bow=count_vect.transform(x_test)

some feature names  ['aa', 'aaa', 'aaaa', 'aaaaaa', 'aaaaaaa', 'aaaaaaaaaaaaa',
'aaaaaaaaaaaaa', 'aaaaaaaaaaaaaa', 'aaaaaaaaaaaaaaa', 'aaaaaaaaaaaaaaaaaaaa
aaaaaaaaaaaaaaaaaaaaa']

=====
```

In [15]:

```
x_train_bow.shape,x_test_bow.shape,x_cv_bow.shape
```

```
Out[15]: ((217951, 90000), (72651, 90000), (72651, 90000))
```

```
In [17]: # saving train_bow and test_bow dataset using pickle for future use

'''file=open("x_train_bow.pkl","wb")
pickle.dump(x_train_bow,file)
file.close()

file=open("x_test_bow.pkl",'wb')
pickle.dump(x_test_bow,file)
file.close()

file=open("x_cv_bow.pkl",'wb')
pickle.dump(x_cv_bow,file)
file.close()
'''
```

```
In [15]: #loading train_bow and test_bow
file=open('x_train_bow.pkl','rb')
x_train_bow=pickle.load(file)

file=open('x_test_bow.pkl','rb')
x_test_bow=pickle.load(file)

file=open('x_cv_bow.pkl','rb')
x_cv_bow=pickle.load(file)
```

[4.3] TF-IDF

```
In [108]: # tf-idf "from sklearn.feature_extraction.text.TfidfVectorizer"
tf_idf=TfidfVectorizer()

train_tf_idf=tf_idf.fit_transform(x_train)
cv_tf_idf=tf_idf.transform(x_cv)
test_tf_idf=tf_idf.transform(x_test)
```

```
In [109]: from sklearn.preprocessing import StandardScaler

sc=StandardScaler(with_mean=False)

train_tf_idf=sc.fit_transform(train_tf_idf)
cv_tf_idf=sc.transform(cv_tf_idf)
test_tf_idf=sc.transform(test_tf_idf)
```

```
In [21]: # saving train_tf_idf and test_tf_idf dataset using pickle for future use

'''file=open("train_tf_idf.pkl","wb")
pickle.dump(train_tf_idf,file)
file.close()

file=open("cv_tf_idf.pkl",'wb')
pickle.dump(cv_tf_idf,file)
file.close()

file=open("test_tf_idf.pkl",'wb')
pickle.dump(test_tf_idf,file)
file.close()'''
```

```
'''
```

```
In [4]: #loading train_tf_idf and test_tf_idf
file=open('train_tf_idf.pkl','rb')
train_tf_idf=pickle.load(file)

file=open('cv_tf_idf.pkl','rb')
cv_tf_idf=pickle.load(file)

file=open('test_tf_idf.pkl','rb')
test_tf_idf=pickle.load(file)
```

[4.4] Word2Vec

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

[4.4.1.1] Avg W2v

```
In [ ]: # converting our text-->vector using w2v with 50-dim
# more the dimension of each word = better the semantic of word
# using lib from "gensim.models.Word2Vec"
# to run w2v we need list of list of the words as w2v covert each word into n
umber of dim

# for train_w2v
list_of_sent_train=[]
for sent in x_train:
    list_of_sent_train.append((str(sent)).split())
w2v_model=Word2Vec(list_of_sent_train,min_count=5,size=50)

# vocabulary of w2v model of amazon dataset
vocab=w2v_model.wv.vocab
len(vocab)

#-----
-----

# for test_w2v
list_of_sent_cv=[]
for sent in x_cv:
    list_of_sent_cv.append((str(sent)).split())

#-----
-----

# for test_w2v
list_of_sent_test=[]
for sent in x_test:
    list_of_sent_test.append((str(sent)).split())
```

```
C:\Users\Adarsh\Anaconda3\lib\site-packages\gensim\models\base_any2vec.py:74
```

```
3: UserWarning: C extension not loaded, training will be slow. Install a C c
ompiler and reinstall gensim for fast training.
```

"C extension not loaded, training will be slow. "

In [7]:

```
'''  
    -->procedure to make avg w2v of each reviews  
    1. find the w2v of each word  
    2. sum-up w2v of each word in a sentence  
    3. divide the total w2v of sentence by total no. of words in the sentence  
# average Word2Vec  
# compute average word2vec for each review.  
train_w2v = [] # the avg-w2v for each sentence/review in train dataset is stored in this list  
  
for sent in list_of_sent_train: # for each review/sentence  
    sent_vec = np.zeros(50) # as word vectors are of zero length  
    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 vocab:  
            vec = w2v_model.wv[word]  
            sent_vec += vec  
            cnt_words += 1  
        if cnt_words != 0:  
            sent_vec /= cnt_words  
    train_w2v.append(sent_vec)  
  
print(len(train_w2v))  
  
#-----  
-----  
  
cv_w2v = [] # the avg-w2v for each sentence/review in test dataset is stored in this list  
  
for sent in list_of_sent_cv: # for each review/sentence  
    sent_vec = np.zeros(50) # as word vectors are of zero length  
    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 vocab:  
            vec = w2v_model.wv[word]  
            sent_vec += vec  
            cnt_words += 1  
        if cnt_words != 0:  
            sent_vec /= cnt_words  
    cv_w2v.append(sent_vec)  
  
print(len(cv_w2v))  
  
#-----  
-----  
  
test_w2v = [] # the avg-w2v for each sentence/review in test dataset is stored in this list  
  
for sent in list_of_sent_test: # for each review/sentence  
    sent_vec = np.zeros(50) # as word vectors are of zero length  
    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 vocab:  
            vec = w2v_model.wv[word]
```

```

        sent_vec += vec
        cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    test_w2v.append(sent_vec)

print(len(test_w2v))

```

```

217951
72651
72651

```

```

In [9]: from sklearn.preprocessing import StandardScaler
sc=StandardScaler(with_mean=True)

train_w2v=sc.fit_transform(train_w2v)
cv_w2v=sc.transform(cv_w2v)
test_w2v=sc.transform(test_w2v)

```

```

In [10]: # saving train_w2v and test_w2v dataset using pickle for future use

'''file=open("train_w2v.pkl","wb")
pickle.dump(train_w2v,file)
file.close()

file=open("cv_w2v.pkl",'wb')
pickle.dump(cv_w2v,file)
file.close()

file=open("test_w2v.pkl",'wb')
pickle.dump(test_w2v,file)
file.close()
'''

```

```

In [5]: #loading train_w2v and test_w2v
file=open('train_w2v.pkl','rb')
train_w2v=pickle.load(file)

file=open('cv_w2v.pkl','rb')
cv_w2v=pickle.load(file)

file=open('test_w2v.pkl','rb')
test_w2v=pickle.load(file)

```

```
In [21]: train_w2v[0]
```

```

Out[21]: array([-0.0080363 , -0.55650226, -2.24174777,  1.02574886,  0.12327979,
   -0.19916425, -1.06522974,  0.89144715, -1.13231167,  2.54008377,
   0.8032532 ,  0.3404576 ,  1.6792167 , -0.98081078,  1.08851643,
   -0.72007858, -0.65714762, -0.56007184,  0.01985994,  2.12137305,
   -0.09203752, -0.23671867, -1.63326771,  1.04496922,  0.45004579,
   0.3219116 ,  0.78335079,  0.54301334, -2.4968575 ,  0.35478244,
   1.46397278, -0.01982212, -0.1817636 ,  1.35729521, -0.61338792,
   -1.68822842, -0.84256537, -0.59978494,  0.40587478, -0.49775708,
   0.31289323,  0.34938107, -0.18756661, -2.25982333,  0.01440547,
   -0.97699964, -0.10107761,  0.28043456,  1.88480264, -0.81507891])

```

```

In [12]: w2v_words = list(w2v_model.wv.vocab)
print("number of words that occurred minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])

```

```
number of words that occurred minimum 5 times 26749
sample words ['really', 'nice', 'seasoning', 'bought', 'product', 'sam', 'c
lub', 'happy', 'able', 'purchase', 'cannot', 'get', 'anymore', 'use', 'meat
s', 'spaghetti', 'coffee', 'not', 'bad', 'defiantly', 'anything', 'special',
'price', 'could', 'ethical', 'fair', 'trade', 'organic', 'shade', 'grown',
'etc', 'taste', 'ok', 'stick', 'dean', 'beans', 'smooth', 'flavorful', 'medi
um', 'roast', 'pleasantly', 'surprised', 'k', 'cup', 'would', 'deal', 'dre
w', 'glad', 'dogs', 'love']
```

[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(x_train)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

```
In [ ]: # TF-IDF weighted Word2Vec Train
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

train_tf_idf_w2v = []; # the tfidf-w2v for each sentence/review is stored in t
his list
row=0;

#list_of_sentence_train=list_of_sent_train[:30000] # reducing the size of train
n list due to computational constrain

for sent in tqdm(list_of_sent_train): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
            #
            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
        if weight_sum != 0:
            sent_vec /= weight_sum
    train_tf_idf_w2v.append(sent_vec)
    row += 1

len(train_tf_idf_w2v)
```

```
In [15]: # TF-IDF weighted Word2Vec cv
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

cv_tf_idf_w2v = []; # the tfidf-w2v for each sentence/review is stored in this
list
row=0;
```

```

#list_of_sentence_train=list_of_sent_train[:30000] # reducing the size of train list due to computational constrain

for sent in tqdm(list_of_sent_cv[:20000]): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
        #
            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    cv_tf_idf_w2v.append(sent_vec)
    row += 1

len(cv_tf_idf_w2v)

```

100%|██████████| 20000/20000 [17:19<00:00, 19.23it/s]

Out[15]: 20000

```

In [16]: # TF-IDF weighted Word2Vec Test
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

test_tf_idf_w2v = [] # the tfidf-w2v for each sentence/review is stored in this list
row=0;

#list_of_sentence_train=list_of_sent_train[:30000] # reducing the size of train list due to computational constrain

for sent in tqdm(list_of_sent_test[:20000]): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
        #
            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    test_tf_idf_w2v.append(sent_vec)
    row += 1

len(test_tf_idf_w2v)

```

100%|██████████| 20000/20000 [17:31<00:00, 19.03it/s]

Out[16]: 20000

```
In [15]: from sklearn.preprocessing import StandardScaler
sc=StandardScaler(with_mean=True)

train_tf_idf_w2v=sc.fit_transform(train_tf_idf_w2v)
cv_tf_idf_w2v=sc.transform(cv_tf_idf_w2v)
test_tf_idf_w2v=sc.transform(test_tf_idf_w2v)
```

```
In [16]: # saving train_tf_idf_w2v and test_tf_idf_w2v dataset using pickle for future use

'''file=open("train_tf_idf_w2v.pkl","wb")
pickle.dump(train_tf_idf_w2v,file)
file.close()

file=open("cv_tf_idf_w2v.pkl",'wb')
pickle.dump(cv_tf_idf_w2v,file)
file.close()

file=open("test_tf_idf_w2v.pkl",'wb')
pickle.dump(test_tf_idf_w2v,file)
file.close()

'''
```

```
In [6]: #loading train_tf_idf_w2v and test_tf_idf_w2v
file=open('train_tf_idf_w2v.pkl','rb')
train_tf_idf_w2v=pickle.load(file)

file=open('cv_tf_idf_w2v.pkl','rb')
cv_tf_idf_w2v=pickle.load(file)

file=open('test_tf_idf_w2v.pkl','rb')
test_tf_idf_w2v=pickle.load(file)
```

[5] Assignment 9: Random Forests

1. Apply Random Forests & GBDT 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. The hyper paramter tuning (Consider two hyperparameters: n_estimators & max_depth)

- 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. Feature importance

- Get top 20 important features and represent them in a word cloud. Do this for BOW & TFIDF.

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 with X-axis  as **n_estimators**, Y-axis as **max_depth**, and Z-axis as **AUC Score** , we have given the notebook which explains how to plot this 3d plot, you can find it in the same drive [3d_scatter_plot.ipynb](#)

(or)

- 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 [seaborn heat maps](#) with rows as **n_estimators**, columns as **max_depth**, and values inside the cell representing **AUC Score**
- You choose either of the plotting techniques out of 3d plot or heat map
- 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 your train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](#).

[5.1] Applying RF

USING ONLY 60k DATAPOINTS to train because of computational strain

[5.1.1] Applying Random Forests on BOW, SET 1

```
In [31]: # By using grid search
from sklearn.ensemble import RandomForestClassifier
```

```

from sklearn.model_selection import GridSearchCV

param={'n_estimators':[5, 10, 50, 100, 200, 500, 1000], 'max_depth': [2, 4, 6, 8, 10, 12, 14, 16, 18, 20]}

# taking decision tree as estimator for grid search
rf_clf=RandomForestClassifier(criterion='gini', bootstrap=True, n_jobs=1, oob_score=False, class_weight='balanced')

gsCV=GridSearchCV(estimator=rf_clf, param_grid=param, scoring='roc_auc', return_train_score=True, n_jobs=-2)
gsCV.fit(x_train_bow[:60000],y_train[:60000])

# storing results
result=gsCV.cv_results_

roc_auc_cv=result['mean_test_score'].reshape(10,7) #reshaping for heatmap visualisation
roc_auc_train=result['mean_train_score'].reshape(10,7) #reshaping for heatmap visualisation

```

In [33]: # visualising train_ROC_AUC and cv_ROC_AUC using heatmap

```

from matplotlib import gridspec

fig = plt.figure(figsize=(16,7))
gs = gridspec.GridSpec(1, 2)

ax0 = plt.subplot(gs[0])

# train heatmap
sn.heatmap(data=roc_auc_train, linewidths=0.01, annot=True, xticklabels=param['n_estimators'],
            yticklabels=param['max_depth'], linecolor='black', fmt='.2g', ax=ax0)

plt.xlabel('n_estimators')
plt.ylabel('depth of tree')
plt.title('heatmap of train ROC AUC\n')

#####
#####

ax1 = plt.subplot(gs[1])

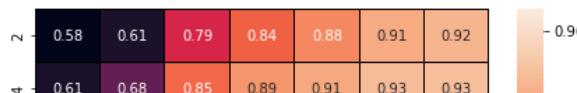
# cv heatmap
sn.heatmap(data=roc_auc_cv, linewidths=0.01, annot=True, xticklabels=param['n_estimators'],
            yticklabels=param['max_depth'], linecolor='black', fmt='.2g', ax=ax1)

plt.xlabel('n_estimators')
plt.ylabel('depth of tree')
plt.title('heatmap of cv ROC AUC\n')

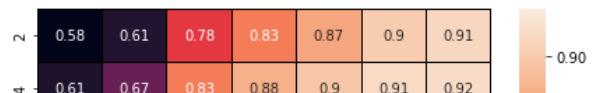
plt.show()

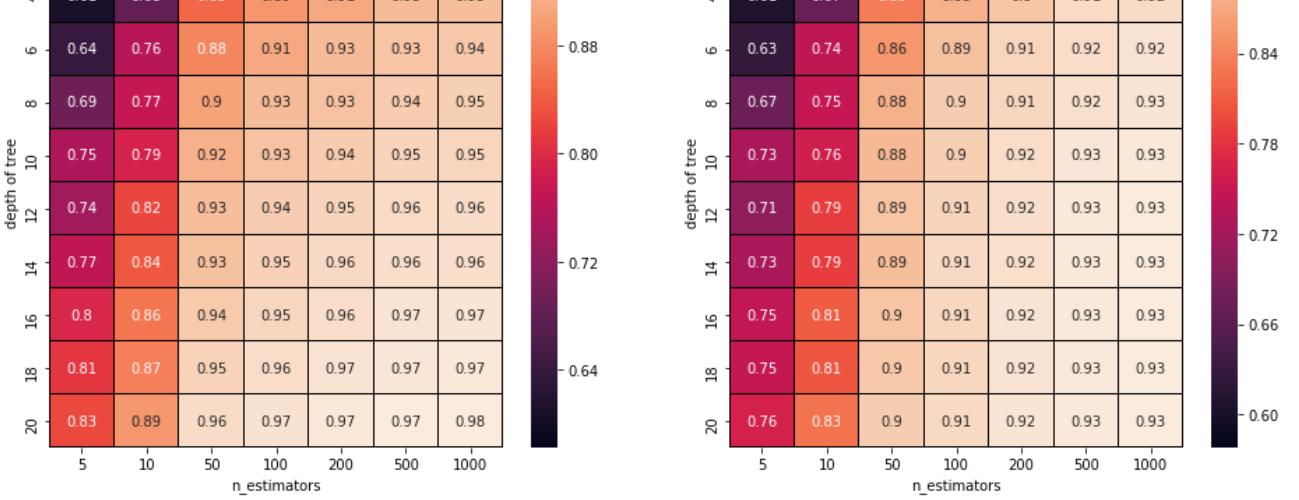
```

heatmap of train ROC AUC



heatmap of cv ROC AUC





best hyperparameter

```
In [35]: # best estimators
print(f"best hyperparameters are: {gsCV.best_params_}")
```

```
best hyperparameters are: {'max_depth': 16, 'n_estimators': 1000}
```

```
In [41]: #Testing
from sklearn.metrics import roc_auc_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve

#probabilty score for ROC_AUC score
y_train_pred_proba=gsCV.predict_proba(x_train_bow[:60000])[:,1]
y_test_pred_proba=gsCV.predict_proba(x_test_bow[:40000])[:,1]

train_roc_score=roc_auc_score(y_train[:60000],y_train_pred_proba)
test_roc_score=roc_auc_score(y_test[:40000],y_test_pred_proba)

#ploting confusion matrix
y_pred=gsCV.predict(x_test_bow[:40000])
sn.heatmap(confusion_matrix(y_test[:40000],y_pred), annot=True, fmt="d", linewidths=.5)
plt.title('Confusion Matrix')
plt.xlabel('Predicted values')
plt.ylabel('Actual values')
plt.show()
print("\n\nclassification report:\n",classification_report(y_test[:40000],y_pred))

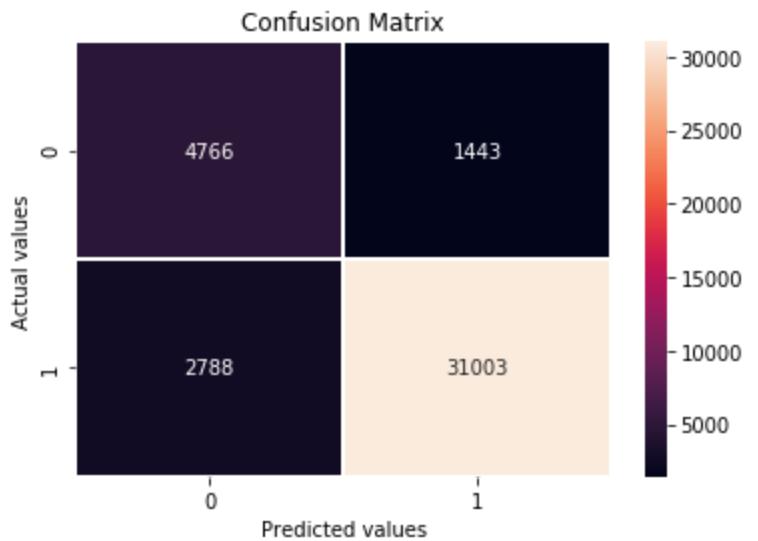
# ROC Curve (reference:stack overflow with little modification)
train_fpr, train_tpr, train_threshold =roc_curve(y_train[:60000], y_train_pred_proba)
test_fpr, test_tpr, test_threshold =roc_curve(y_test[:40000], y_test_pred_proba)

#ploting ROC curve
plt.figure(figsize=(7,7))
plt.title('Receiver Operating Characteristic')
plt.plot(train_fpr, train_tpr, 'b', label = 'train_AUC = %0.2f' % train_roc_score)
plt.plot(test_fpr, test_tpr, 'r', label = 'test_AUC = %0.2f' % test_roc_score)
```

```

plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([-0.02, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.grid(linestyle='--', linewidth=1)
plt.show()

```

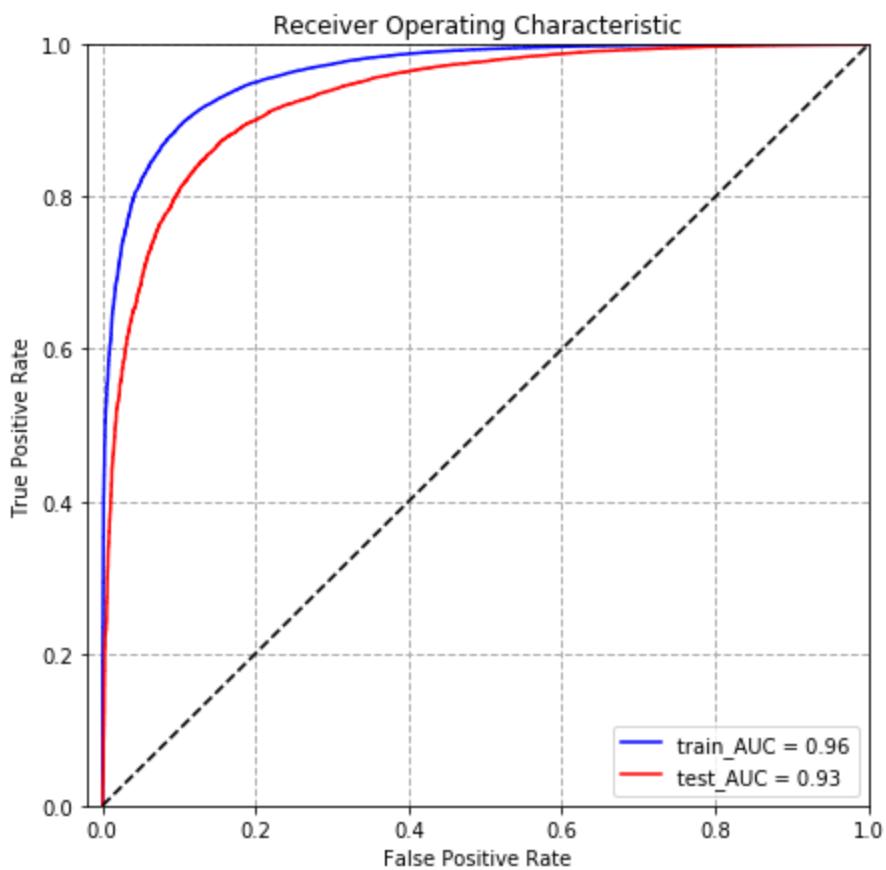


```

classification report:
      precision    recall    f1-score   support
      0          0.63     0.77     0.69     6209
      1          0.96     0.92     0.94    33791

      micro avg     0.89     0.89     0.89     40000
      macro avg     0.79     0.84     0.81     40000
  weighted avg     0.91     0.89     0.90     40000

```



[5.1.2] Wordcloud of top 20 important features from SET 1

```
In [46]: # feature importance scores from decision tree classifier
feature_importance_scores=(gscv.best_estimator_.feature_importances_)

# storing the index of top 20 features
top_features_index=feature_importance_scores.argsort()[-21:-1]

#storing the name of features from BOW vectorizer
bow_feat=count_vect.get_feature_names()
feature_names=pd.DataFrame(bow_feat)

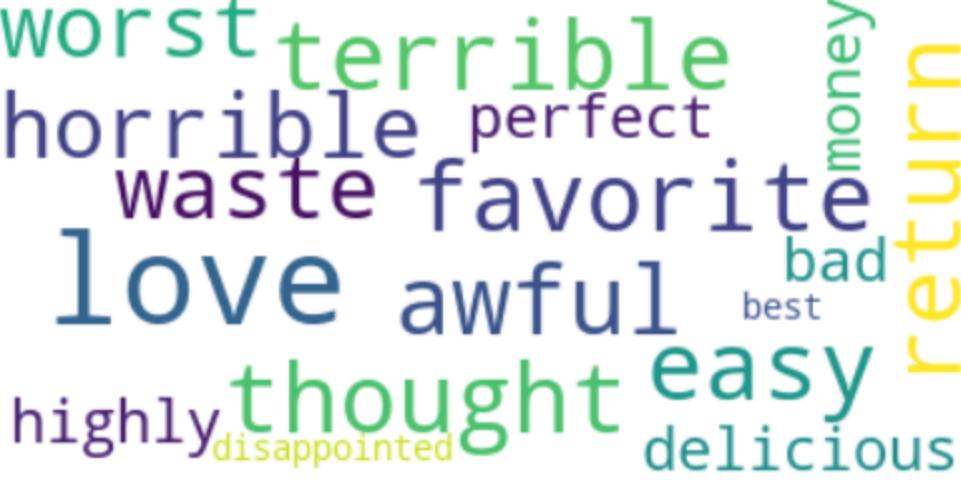
top_features=pd.DataFrame({'feature importance scores':feature_importance_scores[top_features_index],
                           'feature names':feature_names.iloc[top_features_index[0]])}
```

```
In [101]: # https://www.datacamp.com/community/tutorials/wordcloud-python

text = " ".join(list(top_features.iloc[:,1]))

# Create and generate a word cloud image:
wordcloud = WordCloud(max_font_size=50, max_words=100, background_color="white").generate(text)

# Display the generated image:
plt.figure(figsize=(9,7))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



[5.1.3] Applying Random Forests on TFIDF, SET 2

```
In [103]: # By using grid search
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV

param={'n_estimators':[5, 10, 50, 100, 200, 500, 1000], 'max_depth': [2, 4, 6,
8, 10, 12, 14, 16, 18,20]}

# taking decision tree as estimator for grid search
rf_clf=RandomForestClassifier(criterion='gini', bootstrap=True, n_jobs=1, oob_
score=False, class_weight='balanced')

gsCV=GridSearchCV(estimator=rf_clf, param_grid=param, scoring='roc_auc', return_
_train_score=True,n_jobs=-2)
gsCV.fit(train_tf_idf[:60000],y_train[:60000])

# storing results
result=gsCV.cv_results_

roc_auc_cv=result['mean_test_score'].reshape(10,7) #reshaping for heatmap visu
alisation
roc_auc_train=result['mean_train_score'].reshape(10,7) #reshaping for heatmap
visualisation
```

```
In [104]: # visualising train_ROC_AUC and cv_ROC_AUC using heatmap

from matplotlib import gridspec

fig = plt.figure(figsize=(16,7))
gs  = gridspec.GridSpec(1, 2)

ax0 = plt.subplot(gs[0])

# train heatmap
sn.heatmap(data=roc_auc_train,linewidths=0.01,annot=True,xticklabels=param['n_
estimators'],
            yticklabels=param['max_depth'],linecolor='black',fmt='%.2g',ax=ax0)

plt.xlabel('n_estimators')
plt.ylabel('depth of tree')
plt.title('heatmap of train ROC AUC\n')
```

```

*****  

*****  

*****  

*****  

*****  

*****  

*****  

ax1 = plt.subplot(gs[1])  

# cv heatmap  

sn.heatmap(data=roc_auc_cv, linewidths=0.01, annot=True, xticklabels=param['n_estimators'],  

           yticklabels=param['max_depth'], linecolor='black', fmt='%.2g', ax=ax1)  

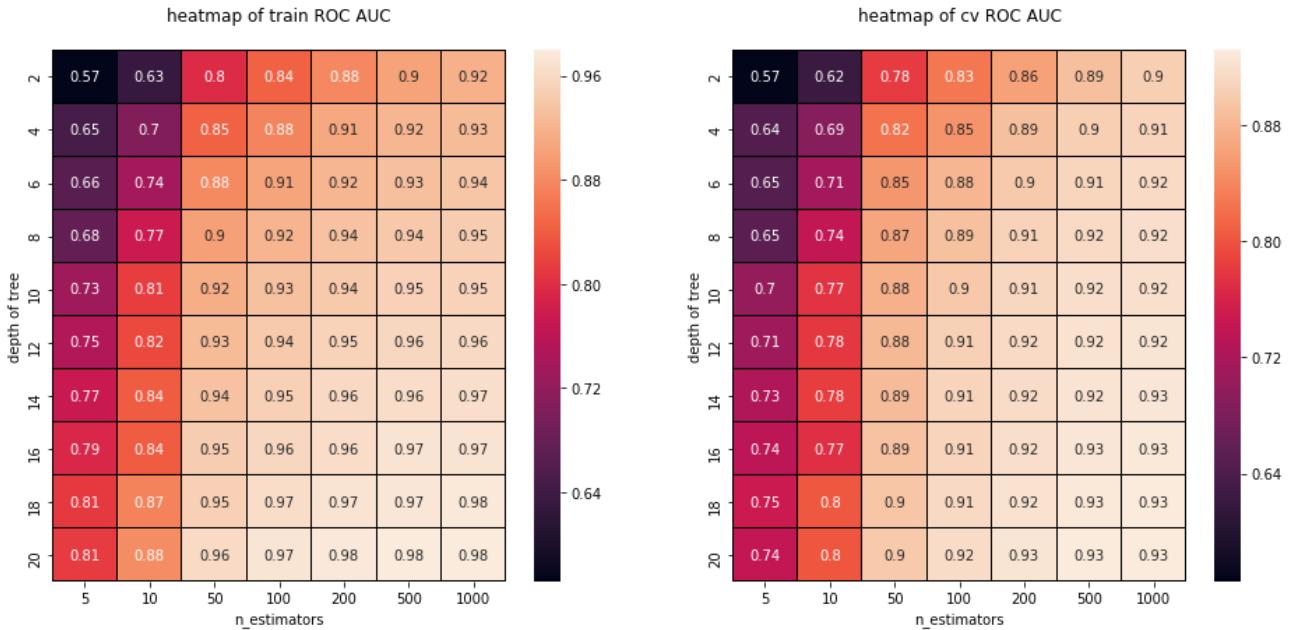
plt.xlabel('n_estimators')  

plt.ylabel('depth of tree')  

plt.title('heatmap of cv ROC AUC\n')  

plt.show()

```



best hyperparameter

```
In [105]: # best estimators
print(f"best hyperparameters are: {gsCV.best_params_} ")
```

```
best hyperparameters are: {'max_depth': 20, 'n_estimators': 1000}
```

```
In [110]: #Testing
from sklearn.metrics import roc_auc_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve

#probabiltiy score for ROC_AUC score
y_train_pred_proba=gsCV.predict_proba(train_tf_idf[:60000])[:,1]
y_test_pred_proba=gsCV.predict_proba(test_tf_idf[:40000])[:,1]

train_roc_score=roc_auc_score(y_train[:60000],y_train_pred_proba)
test_roc_score=roc_auc_score(y_test[:40000],y_test_pred_proba)

#ploting confusion matrix
y_pred=gsCV.predict(test_tf_idf[:40000])
sn.heatmap(confusion_matrix(y_test[:40000],y_pred), annot=True, fmt="d", linewidths=.5)
plt.title('Confusion Matrix')
```

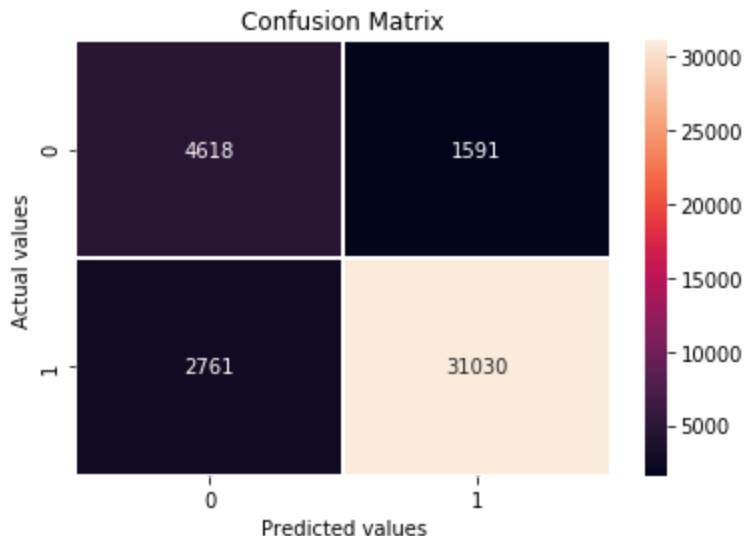
```

plt.xlabel('Predicted values')
plt.ylabel('Actual values')
plt.show()
print("\n\nclassification report:\n", classification_report(y_test[:40000], y_pred))

# ROC Curve (reference: stack overflow with little modification)
train_fpr, train_tpr, train_threshold =roc_curve(y_train[:60000], y_train_pred_proba)
test_fpr, test_tpr, test_threshold =roc_curve(y_test[:40000], y_test_pred_prob_a)

#ploting ROC curve
plt.figure(figsize=(7,7))
plt.title('Receiver Operating Characteristic')
plt.plot(train_fpr, train_tpr, 'b', label = 'train_AUC = %0.2f' % train_roc_score)
plt.plot(test_fpr, test_tpr, 'r', label = 'test_AUC = %0.2f' % test_roc_score)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([-0.02, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.grid(linestyle='--', linewidth=1)
plt.show()

```



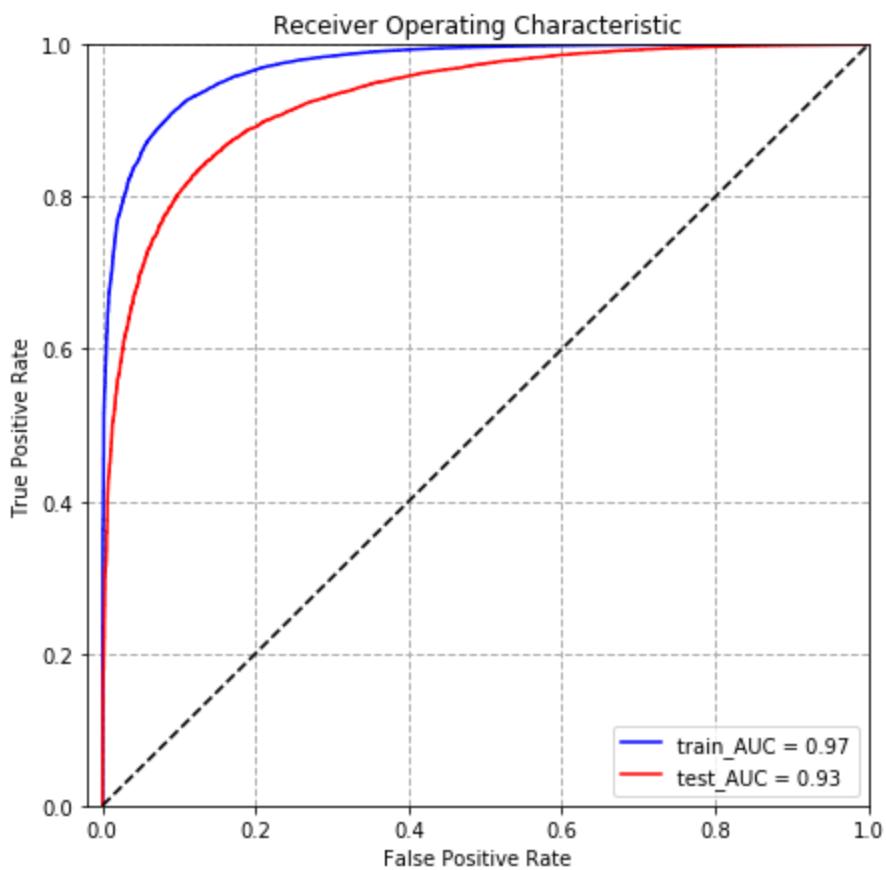
```

classification report:
      precision    recall  f1-score   support

          0       0.63      0.74      0.68      6209
          1       0.95      0.92      0.93     33791

      micro avg       0.89      0.89      0.89     40000
      macro avg       0.79      0.83      0.81     40000
  weighted avg       0.90      0.89      0.89     40000

```



[5.1.4] Wordcloud of top 20 important features from SET 2

```
In [111]: # feature importance scores from decision tree classifier
feature_importance_scores=(gsCV.best_estimator_.feature_importances_)

# storing the index of top 20 features
top_features_index=feature_importance_scores.argsort()[-21:-1]

#storing the name of features from BOW vectorizer
tf_idf_feat=tf_idf.get_feature_names()
feature_names=pd.DataFrame(tf_idf_feat)

top_features=pd.DataFrame({'feature importance scores':feature_importance_scores[top_features_index],
                           'feature names':feature_names.iloc[top_features_index][0]})
```

```
In [112]: # https://www.datacamp.com/community/tutorials/wordcloud-python

text = " ".join(list(top_features.iloc[:,1]))

# Create and generate a word cloud image:
wordcloud = WordCloud(max_font_size=50, max_words=100, background_color="white").generate(text)

# Display the generated image:
plt.figure(figsize=(9,7))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



[5.1.5] Applying Random Forests on AVG W2V, SET 3

```
In [116]: # By using grid search
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV

param={'n_estimators':[5, 10, 50, 100, 200, 500, 1000], 'max_depth': [2, 4, 6,
8, 10, 12, 14, 16, 18,20]}

# taking decision tree as estimator for grid search
rf_clf=RandomForestClassifier(criterion='gini', bootstrap=True, n_jobs=1, oob_
score=False, class_weight='balanced')

gsCV=GridSearchCV(estimator=rf_clf, param_grid=param, scoring='roc_auc', return_
_train_score=True,n_jobs=-2)
gsCV.fit(train_w2v[:60000],y_train[:60000])

# storing results
result=gsCV.cv_results_

roc_auc_cv=result['mean_test_score'].reshape(10,7) #reshaping for heatmap visu
alisation
roc_auc_train=result['mean_train_score'].reshape(10,7) #reshaping for heatmap
visualisation
```

```
In [117]: # visualising train_ROC_AUC and cv_ROC_AUC using heatmap
from matplotlib import gridspec

fig = plt.figure(figsize=(16,7))
gs  = gridspec.GridSpec(1, 2)

ax0 = plt.subplot(gs[0])

# train heatmap
sn.heatmap(data=roc_auc_train,linewidths=0.01,annot=True,xticklabels=param['n_
estimators'],
            yticklabels=param['max_depth'],linecolor='black',fmt='%.2g',ax=ax0)

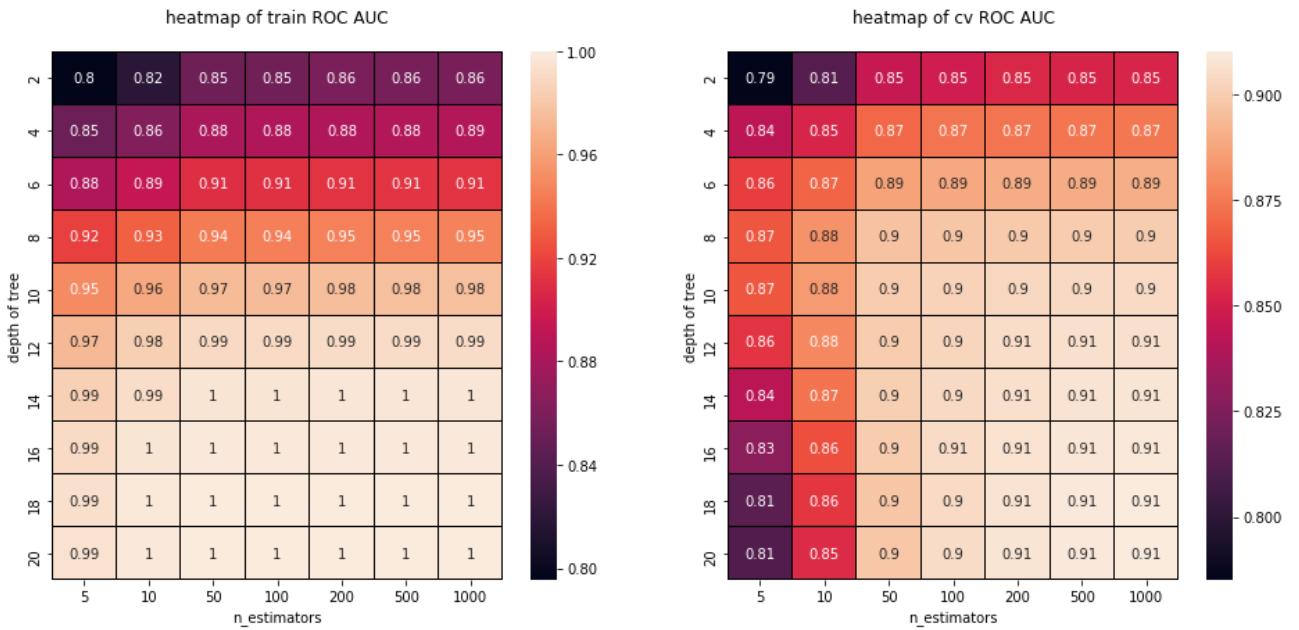
plt.xlabel('n_estimators')
plt.ylabel('depth of tree')
plt.title('heatmap of train ROC AUC\n')

*****
```

```
*****
ax1 = plt.subplot(gs[1])

# cv heatmap
sn.heatmap(data=roc_auc_cv, linewidths=0.01, annot=True, xticklabels=param['n_estimators'],
            yticklabels=param['max_depth'], linecolor='black', fmt='.2g', ax=ax1)

plt.xlabel('n_estimators')
plt.ylabel('depth of tree')
plt.title('heatmap of cv ROC AUC\n')
plt.show()
```



best hyperparameter

```
In [120]: # best estimators
print(f"best hyperparameters are: {gsCV.best_params_} ")
```

```
best hyperparameters are: {'max_depth': 20, 'n_estimators': 1000}
```

```
In [119]: #Testing
from sklearn.metrics import roc_auc_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve

#probability score for ROC_AUC score
y_train_pred_proba=gsCV.predict_proba(train_w2v[:60000])[:,1]
y_test_pred_proba=gsCV.predict_proba(test_w2v[:40000])[:,1]

train_roc_score=roc_auc_score(y_train[:60000],y_train_pred_proba)
test_roc_score=roc_auc_score(y_test[:40000],y_test_pred_proba)

#ploting confusion matrix
y_pred=gsCV.predict(test_w2v[:40000])
sn.heatmap(confusion_matrix(y_test[:40000],y_pred), annot=True, fmt="d", linewidths=.5)
plt.title('Confusion Matrix')
plt.xlabel('Predicted values')
```

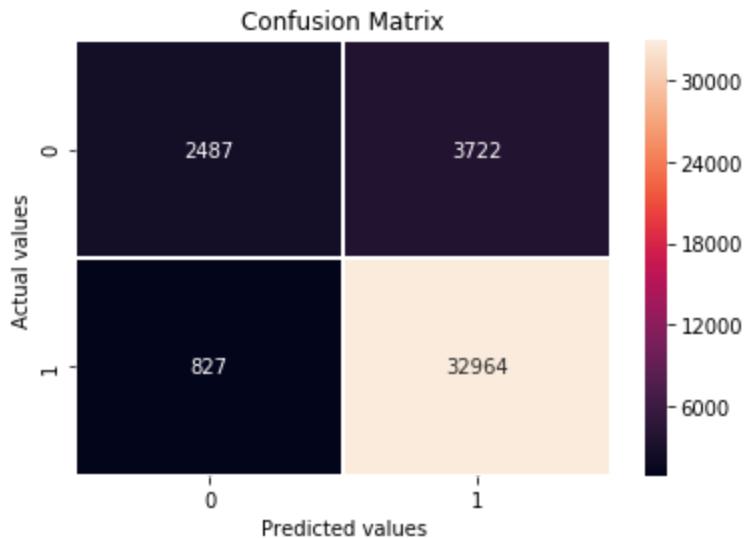
```

plt.ylabel('Actual values')
plt.show()
print("\n\nclassification report:\n", classification_report(y_test[:40000], y_pred))

# ROC Curve (reference:stack overflow with little modification)
train_fpr, train_tpr, train_threshold =roc_curve(y_train[:60000], y_train_pred_proba)
test_fpr, test_tpr, test_threshold =roc_curve(y_test[:40000], y_test_pred_proba)

#ploting ROC curve
plt.figure(figsize=(7,7))
plt.title('Receiver Operating Characteristic')
plt.plot(train_fpr, train_tpr, 'b', label = 'train_AUC = %0.2f' % train_roc_score)
plt.plot(test_fpr, test_tpr, 'r', label = 'test_AUC = %0.2f' % test_roc_score)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([-0.02, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.grid(linestyle='--', linewidth=1)
plt.show()

```



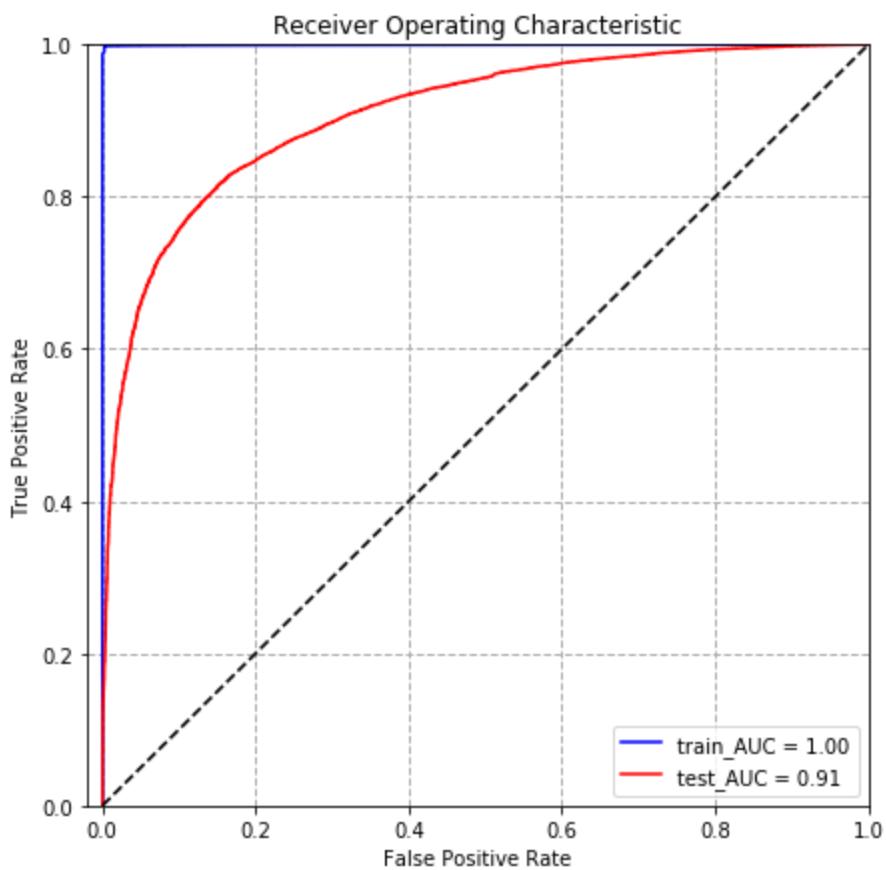
```

classification report:
      precision    recall  f1-score   support

          0       0.75     0.40      0.52      6209
          1       0.90     0.98      0.94     33791

      micro avg       0.89     0.89      0.89      40000
      macro avg       0.82     0.69      0.73      40000
  weighted avg       0.88     0.89      0.87      40000

```



[5.1.6] Applying Random Forests on TFIDF W2V, SET 4

```
In [121]: # By using grid search
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV

param={'n_estimators':[5, 10, 50, 100, 200, 500, 1000], 'max_depth': [2, 4, 6, 8, 10, 12, 14, 16, 18, 20]}

# taking decision tree as estimator for grid search
rf_clf=RandomForestClassifier(criterion='gini', bootstrap=True, n_jobs=1, oob_score=False, class_weight='balanced')

gsCV=GridSearchCV(estimator=rf_clf, param_grid=param, scoring='roc_auc', return_train_score=True, n_jobs=-2)
gsCV.fit(train_tf_idf_w2v[:60000], y_train[:60000])

# storing results
result=gsCV.cv_results_

roc_auc_cv=result['mean_test_score'].reshape(10,7) #reshaping for heatmap visualisation
roc_auc_train=result['mean_train_score'].reshape(10,7) #reshaping for heatmap visualisation
```

```
In [122]: # visualising train_ROC_AUC and cv_ROC_AUC using heatmap

from matplotlib import gridspec

fig = plt.figure(figsize=(16,7))
gs = gridspec.GridSpec(1, 2)
```

```

ax0 = plt.subplot(gs[0])

# train heatmap
sn.heatmap(data=roc_auc_train, linewidths=0.01, annot=True, xticklabels=param['n_estimators'],
            yticklabels=param['max_depth'], linecolor='black', fmt='.2g', ax=ax0)

plt.xlabel('n_estimators')
plt.ylabel('depth of tree')
plt.title('heatmap of train ROC AUC\n')

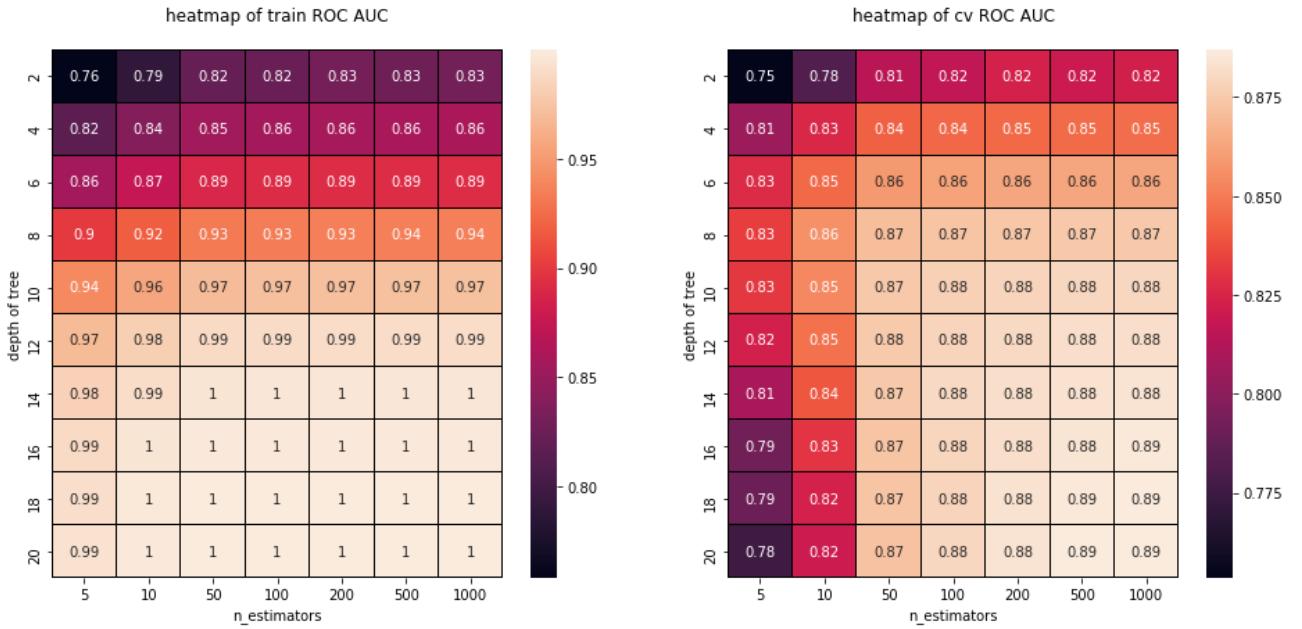
*****
*****
```

ax1 = plt.subplot(gs[1])

```

# cv heatmap
sn.heatmap(data=roc_auc_cv, linewidths=0.01, annot=True, xticklabels=param['n_estimators'],
            yticklabels=param['max_depth'], linecolor='black', fmt='.2g', ax=ax1)

plt.xlabel('n_estimators')
plt.ylabel('depth of tree')
plt.title('heatmap of cv ROC AUC\n')
plt.show()
```



best hyperparameter

```
In [123]: # best estimators
print(f"best hyperparameters are: {gsCV.best_params_} ")
```

best hyperparameters are: {'max_depth': 20, 'n_estimators': 1000}

```
In [124]: from sklearn.metrics import roc_auc_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve

#probability score for ROC_AUC score
y_train_pred_proba=gsCV.predict_proba(train_tf_idf_w2v[:60000])[:,1]
```

```

y_test_pred_proba=gCSV.predict_proba(test_tf_idf_w2v[:20000])[:,1]

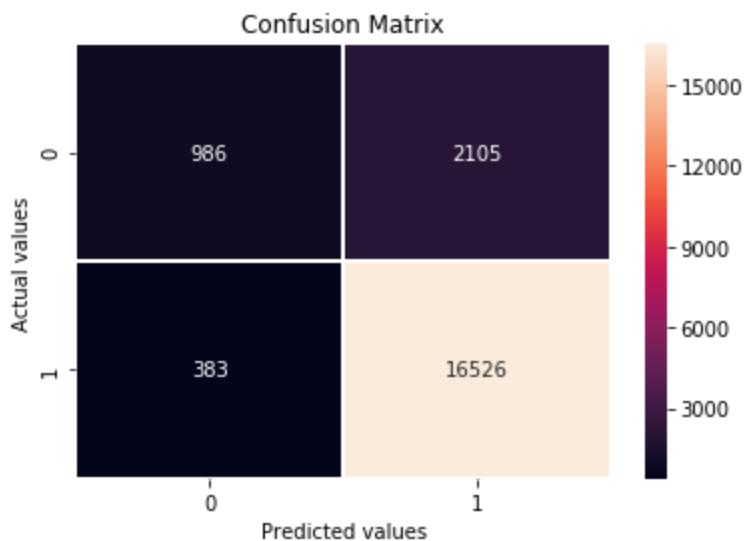
train_roc_score=roc_auc_score(y_train[:60000],y_train_pred_proba)
test_roc_score=roc_auc_score(y_test[:20000],y_test_pred_proba)

#ploting confusion matrix
y_pred=gCSV.predict(test_tf_idf_w2v[:20000])
sn.heatmap(confusion_matrix(y_test[:20000],y_pred),annot=True, fmt="d", linewidths=.5)
plt.title('Confusion Matrix')
plt.xlabel('Predicted values')
plt.ylabel('Actual values')
plt.show()
print("\n\nclassification report:\n",classification_report(y_test[:20000],y_pred))

# ROC Curve (reference:stack overflow with little modification)
train_fpr, train_tpr, train_threshold =roc_curve(y_train[:60000], y_train_pred_proba)
test_fpr, test_tpr, test_threshold =roc_curve(y_test[:20000], y_test_pred_proba)

#ploting ROC curve
plt.figure(figsize=(7,7))
plt.title('Receiver Operating Characteristic')
plt.plot(train_fpr, train_tpr, 'b', label = 'train_AUC = %0.2f' % train_roc_score)
plt.plot(test_fpr, test_tpr, 'r', label = 'test_AUC = %0.2f' % test_roc_score)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([-0.02, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.grid(linestyle='--', linewidth=1)
plt.show()

```



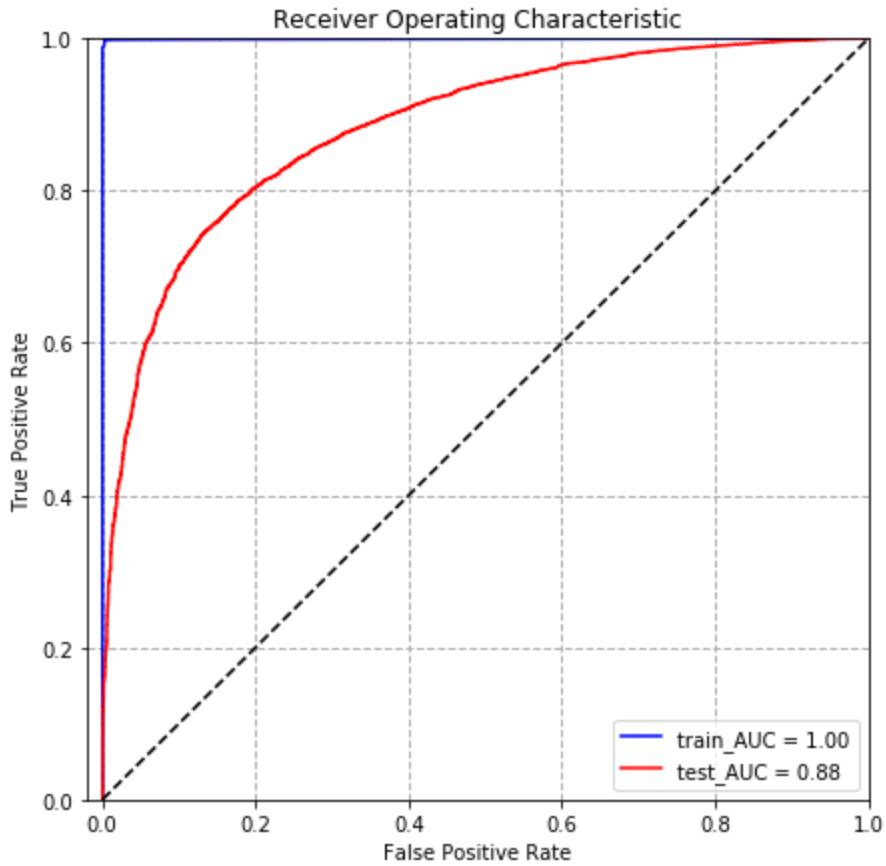
```

classification report:
      precision    recall  f1-score   support

          0       0.72      0.32      0.44     3091
          1       0.89      0.98      0.93    16909

```

micro avg	0.88	0.88	0.88	20000
macro avg	0.80	0.65	0.69	20000
weighted avg	0.86	0.88	0.85	20000



[5.2] Applying GBDT using XGBOOST

Using LightGBM for faster traing

[5.2.1] Applying LightGBM on BOW, SET 1

```
In [6]: from sklearn.preprocessing import StandardScaler
sc=StandardScaler(with_mean=False)
train_bow=sc.fit_transform(x_train_bow)
test_bow=sc.transform(x_test_bow)
```

```
In [8]: # By using grid search
from lightgbm import LGBMClassifier
from sklearn.model_selection import GridSearchCV

param={'num_leaves':[5, 10, 15, 20, 25, 30], 'learning_rate': [0.1,0.2,0.3,0.4, 0.5,0.6,0.7,0.8,0.9,1],
       'n_estimators':[5, 10, 50, 100, 200, 500, 1000]}

# taking decision tree as estimator for grid search
clf_lightgb=LGBMClassifier(boosting_type='gbdt',class_weight='balanced',random_state=1,n_jobs=-1)

gsCV=GridSearchCV(estimator=clf_lightgb, param_grid=param, scoring='roc_auc', r
```

```

    return_train_score=True, n_jobs=-2)
gsCV.fit(train_bow[:60000], y_train[:60000])

# storing results
result=gsCV.cv_results_

roc_auc_cv=result['mean_test_score'].reshape(10,7,6) #reshaping for heatmap visualisation
roc_auc_train=result['mean_train_score'].reshape(10,7,6) #reshaping for heatmap visualisation

```

In [11]: # visualising train_ROC_AUC and cv_ROC_AUC using heatmap

```

for index,lr in enumerate(param['learning_rate']):
    from matplotlib import gridspec

    fig = plt.figure(figsize=(15,5))
    gs = gridspec.GridSpec(1, 2)

    ax0 = plt.subplot(gs[0])

    # train heatmap
    sn.heatmap(data=roc_auc_train[index], linewidths=0.01, annot=True, xticklabels=param['num_leaves'],
                yticklabels=param['n_estimators'], linecolor='black', fmt='.2g', ax=ax0)

    plt.xlabel('num_leaves')
    plt.ylabel('n_estimators')
    plt.title(f"heatmap of train ROC AUC with learning rate: {lr} \n")
    ****
    ****

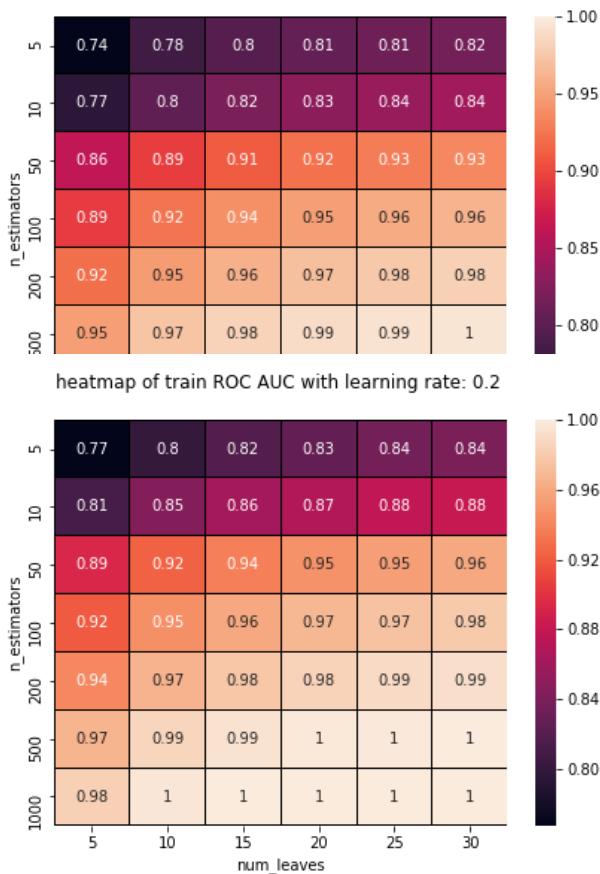
    ax1 = plt.subplot(gs[1])

    # cv heatmap
    sn.heatmap(data=roc_auc_cv[index], linewidths=0.01, annot=True, xticklabels=param['num_leaves'],
                yticklabels=param['n_estimators'], linecolor='black', fmt='.2g', ax=ax1)

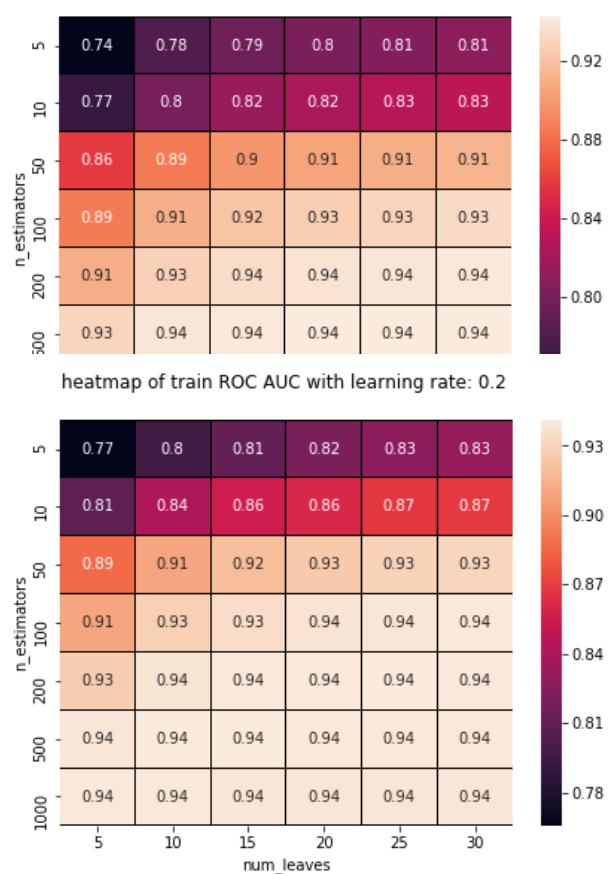
    plt.xlabel('num_leaves')
    plt.ylabel('n_estimators')
    plt.title(f"heatmap of train ROC AUC with learning rate: {lr} \n")
    plt.show()

```

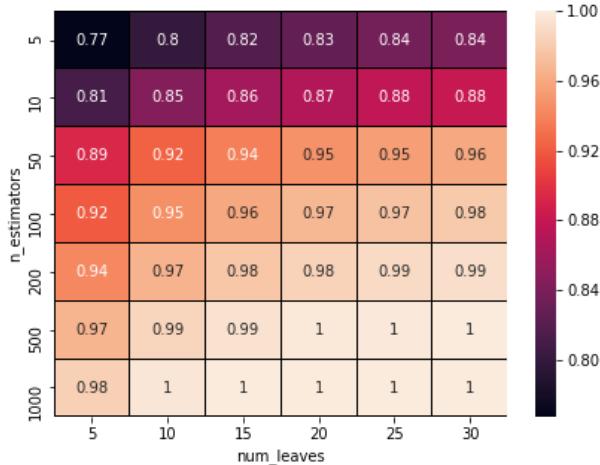
heatmap of train ROC AUC with learning rate: 0.1



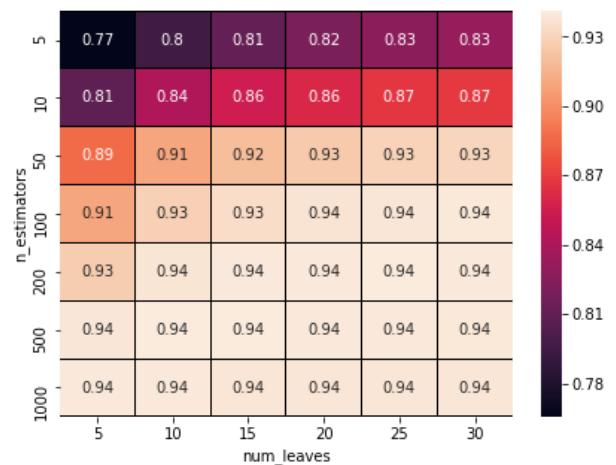
heatmap of train ROC AUC with learning rate: 0.1



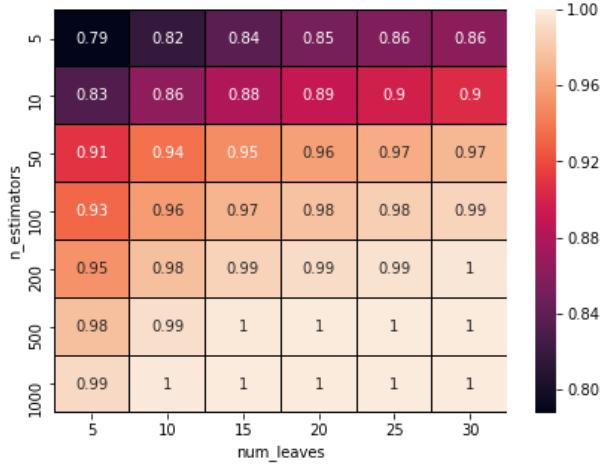
heatmap of train ROC AUC with learning rate: 0.2



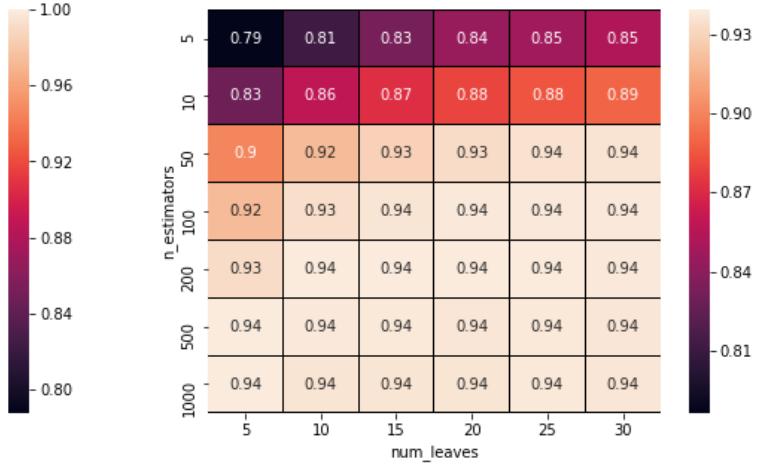
heatmap of train ROC AUC with learning rate: 0.2



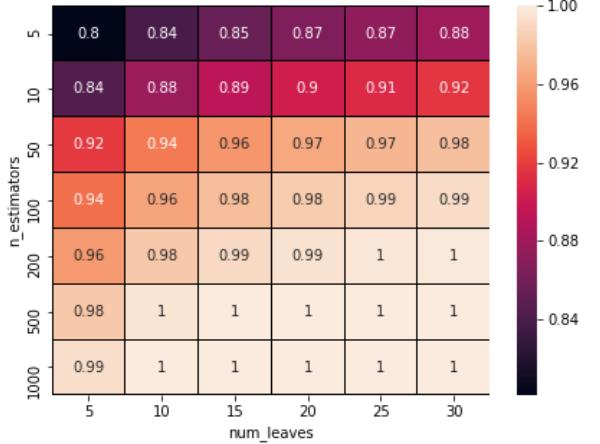
heatmap of train ROC AUC with learning rate: 0.3



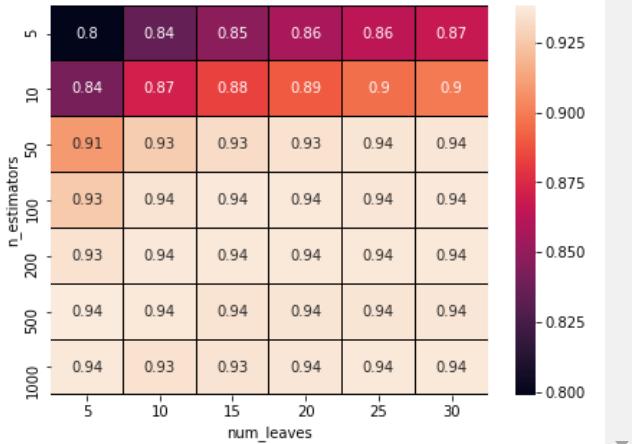
heatmap of train ROC AUC with learning rate: 0.3

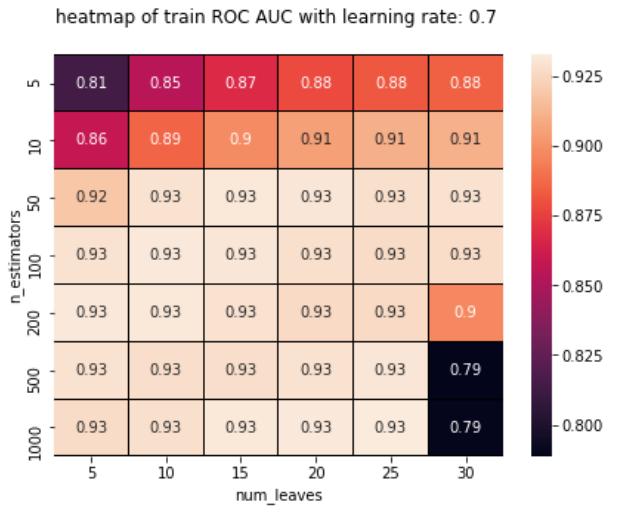
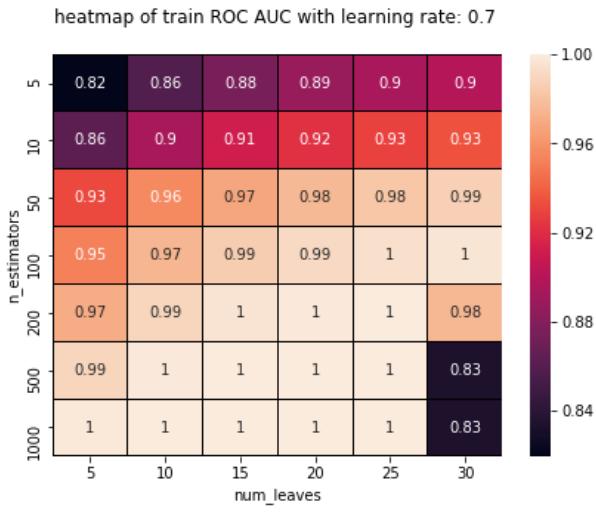
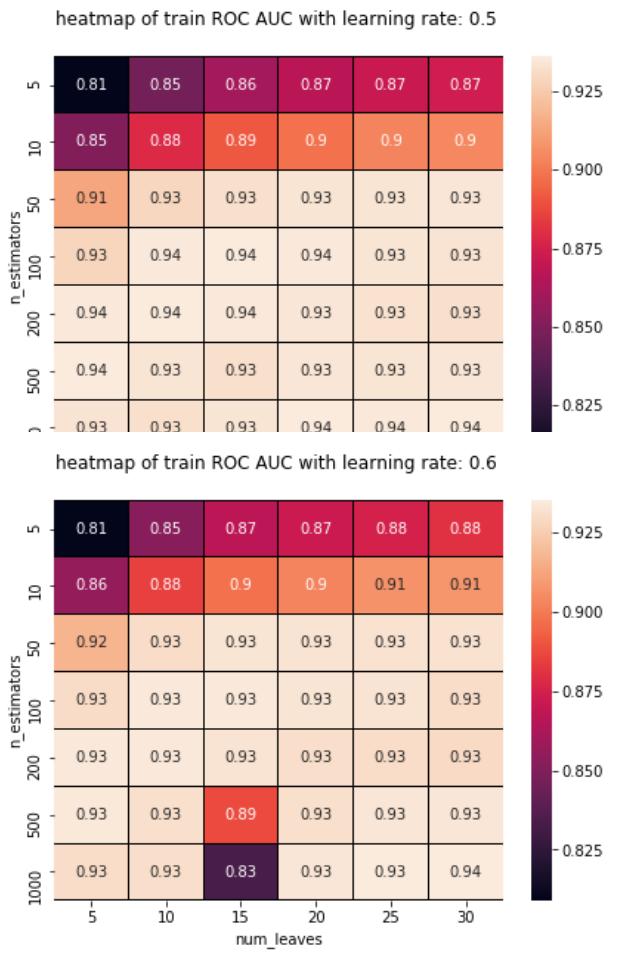
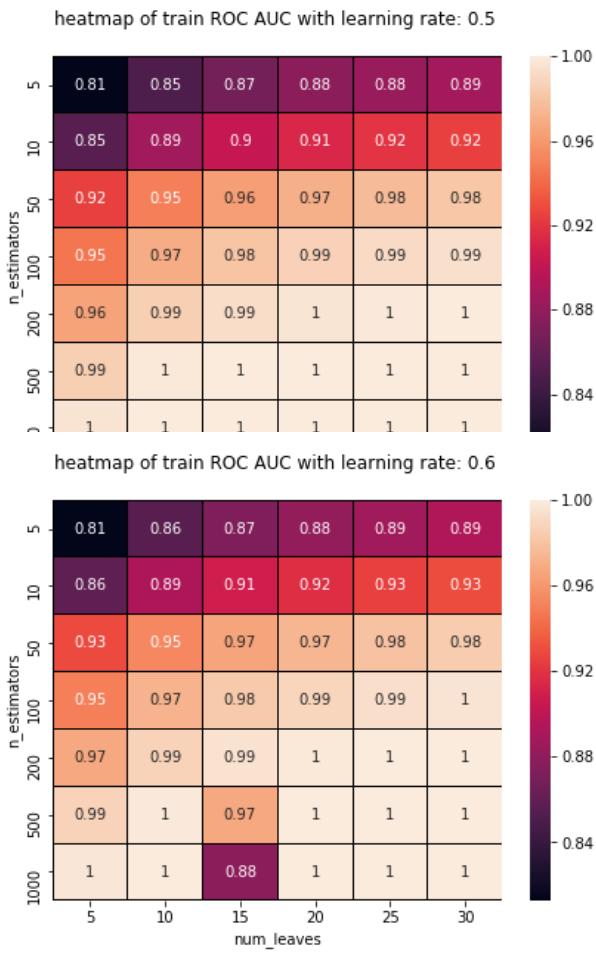


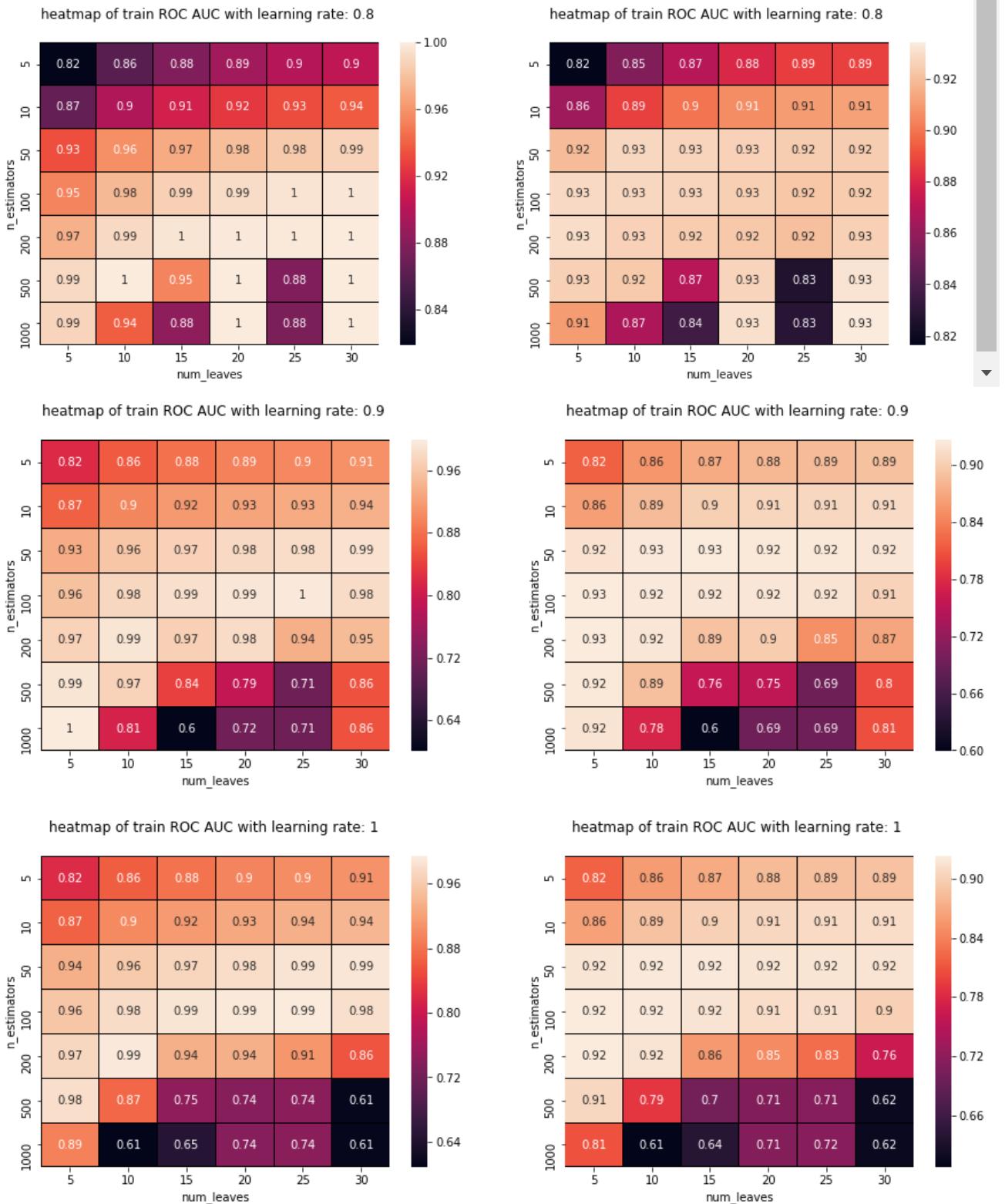
heatmap of train ROC AUC with learning rate: 0.4



heatmap of train ROC AUC with learning rate: 0.4







best hyperparameter

```
In [12]: # best estimators
print(f"best hyperparameters are: {gsCV.best_params_}")
```

```
best hyperparameters are: {'learning_rate': 0.1, 'n_estimators': 500, 'num_leaves': 30}
```

```
In [14]: #Testing
from sklearn.metrics import roc_auc_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve
```

```

#probability score for ROC_AUC score
y_train_pred_proba=gsCV.predict_proba(train_bow[:60000])[:,1]
y_test_pred_proba=gsCV.predict_proba(test_bow[:40000])[:,1]

train_roc_score=roc_auc_score(y_train[:60000],y_train_pred_proba)
test_roc_score=roc_auc_score(y_test[:40000],y_test_pred_proba)

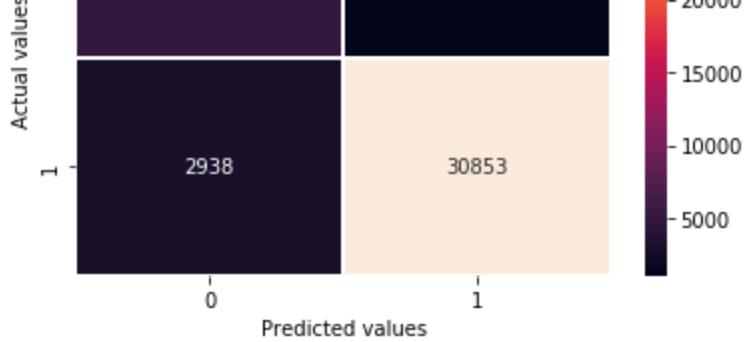
#ploting confusion matrix
y_pred=gsCV.predict(test_bow[:40000])
sn.heatmap(confusion_matrix(y_test[:40000],y_pred), annot=True, fmt="d", linewidths=.5)
plt.title('Confusion Matrix')
plt.xlabel('Predicted values')
plt.ylabel('Actual values')
plt.show()
print("\n\nclassification report:\n",classification_report(y_test[:40000],y_pred))

# ROC Curve (reference:stack overflow with little modification)
train_fpr, train_tpr, train_threshold =roc_curve(y_train[:60000], y_train_pred_proba)
test_fpr, test_tpr, test_threshold =roc_curve(y_test[:40000], y_test_pred_proba)

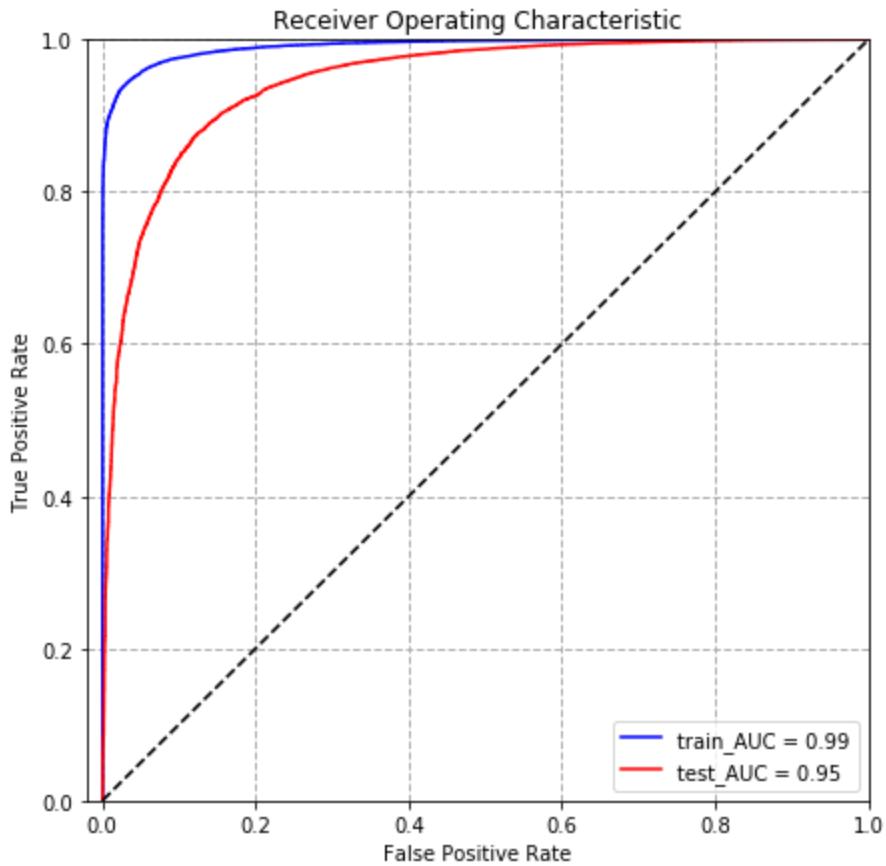
#ploting ROC curve
plt.figure(figsize=(7,7))
plt.title('Receiver Operating Characteristic')
plt.plot(train_fpr, train_tpr, 'b', label = 'train_AUC = %0.2f' % train_roc_score)
plt.plot(test_fpr, test_tpr, 'r', label = 'test_AUC = %0.2f' % test_roc_score)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([-0.02, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.grid(linestyle='--', linewidth=1)
plt.show()

```





```
classification report:
              precision    recall   f1-score   support
0            0.64      0.83      0.72      6209
1            0.97      0.91      0.94     33791
               micro avg       0.90      0.90      0.90      40000
               macro avg       0.80      0.87      0.83      40000
               weighted avg    0.91      0.90      0.90      40000
```



[5.2.2] Applying LightGBM on TFIDF, SET 2

```
In [24]: # By using grid search
from lightgbm import LGBMClassifier
from sklearn.model_selection import GridSearchCV

param={'num_leaves':[5, 10, 15, 20, 25, 30], 'learning_rate': [0.1,0.2,0.3,0.4, 0.5,0.6,0.7,0.8,0.9,1],
       'n_estimators':[5,15,25,35]}
```

```

# taking decision tree as estimator for grid search
clf_lightgb=LGBMClassifier(boosting_type='gbdt',class_weight='balanced',random_
_state=1,n_jobs=-1)

gsCV=GridSearchCV(estimator=clf_lightgb, param_grid=param, scoring='roc_auc',r
eturn_train_score=True,n_jobs=-2)
gsCV.fit(train_tf_idf[:60000],y_train[:60000])

# storing results
result=gsCV.cv_results_

roc_auc_cv=result['mean_test_score'].reshape(10,4,6) #reshaping for heatmap vi
sualisation
roc_auc_train=result['mean_train_score'].reshape(10,4,6) #reshaping for heatma
p visualisation

```

In [57]: # visualising train_ROC_AUC and cv_ROC_AUC using heatmap

```

for index,lr in enumerate(param['learning_rate']):
    from matplotlib import gridspec

    fig = plt.figure(figsize=(15,5))
    gs  = gridspec.GridSpec(1, 2)

    ax0 = plt.subplot(gs[0])

    # train heatmap
    sn.heatmap(data=roc_auc_train[index],lineweights=0.01,annot=True,xticklabel
s=param['num_leaves'],
                yticklabels=param['n_estimators'],linecolor='black',fmt='.2g',a
x=ax0)

    plt.xlabel('num_leaves')
    plt.ylabel('n_estimators')
    plt.title(f"heatmap of train ROC AUC with learning rate: {lr} \n")
    ****
    ****

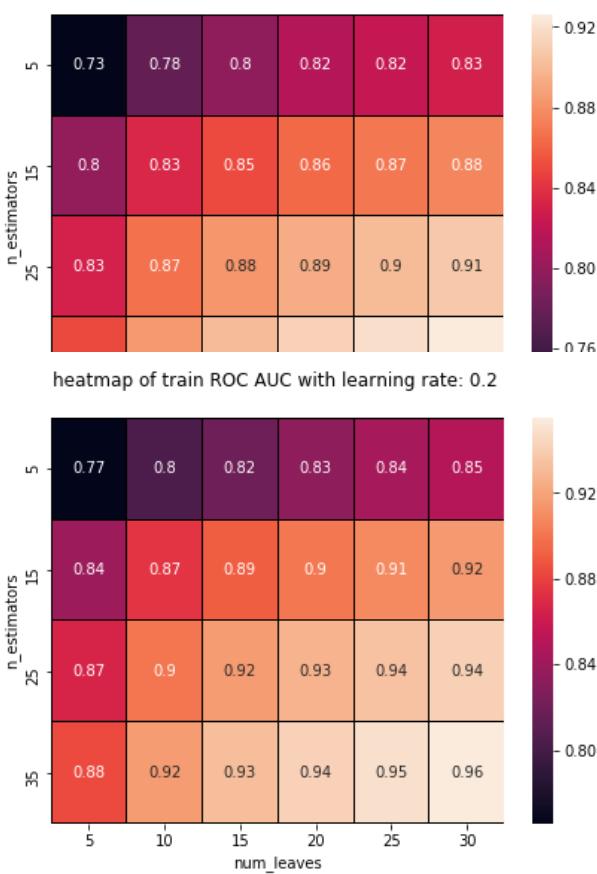
    ax1 = plt.subplot(gs[1])

    # cv heatmap
    sn.heatmap(data=roc_auc_cv[index],lineweights=0.01,annot=True,xticklabels=p
aram['num_leaves'],
                yticklabels=param['n_estimators'],linecolor='black',fmt='.2g',a
x=ax1)

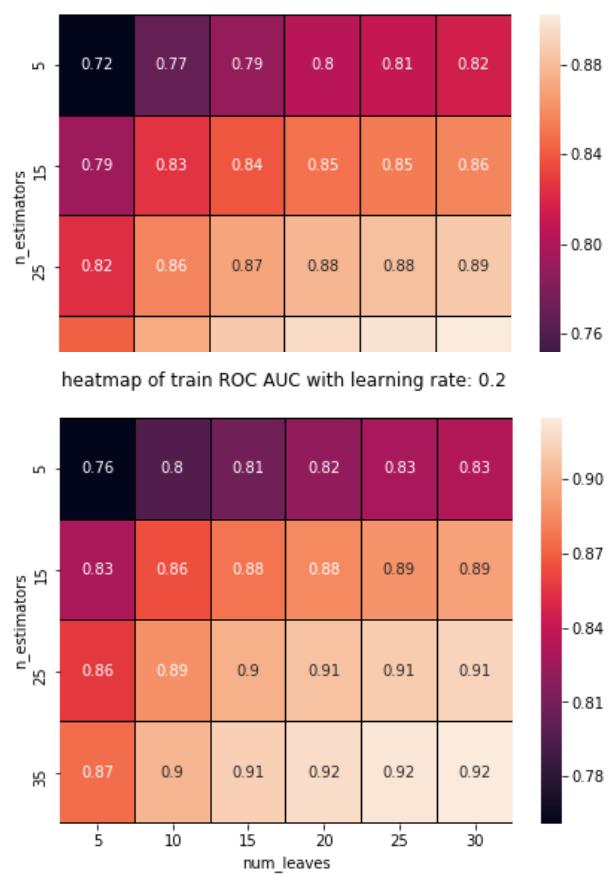
    plt.xlabel('num_leaves')
    plt.ylabel('n_estimators')
    plt.title(f"heatmap of train ROC AUC with learning rate: {lr} \n")
    plt.show()

```

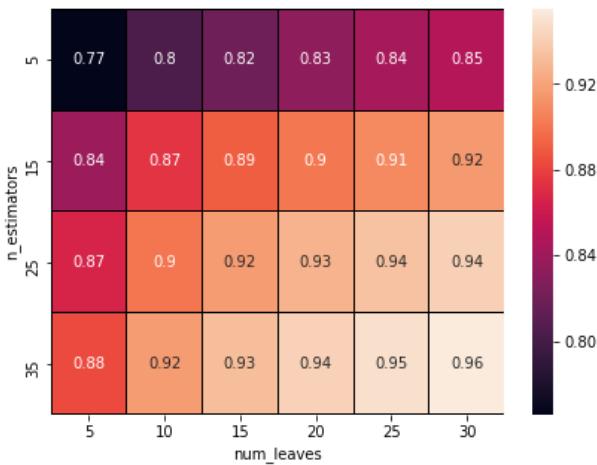
heatmap of train ROC AUC with learning rate: 0.1



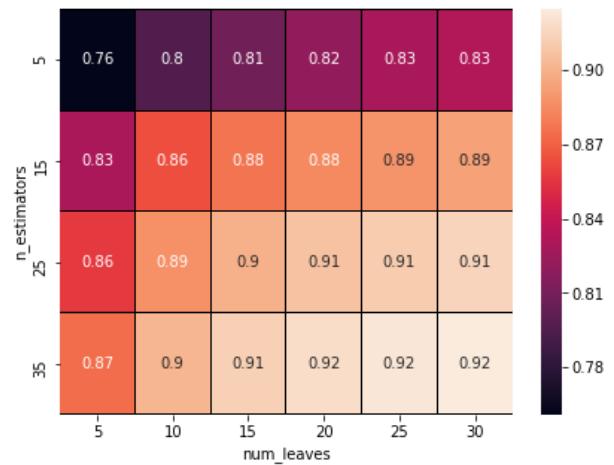
heatmap of train ROC AUC with learning rate: 0.1



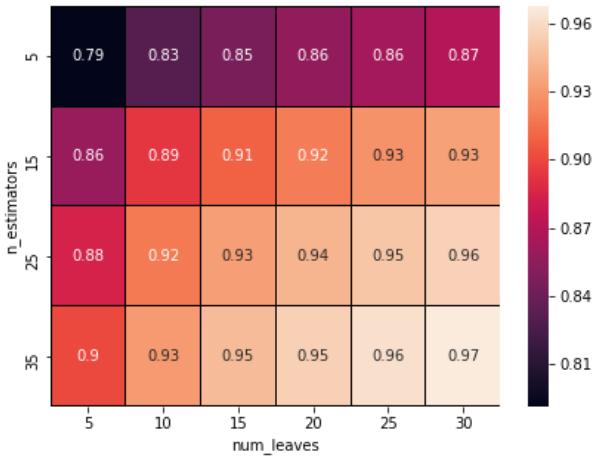
heatmap of train ROC AUC with learning rate: 0.2



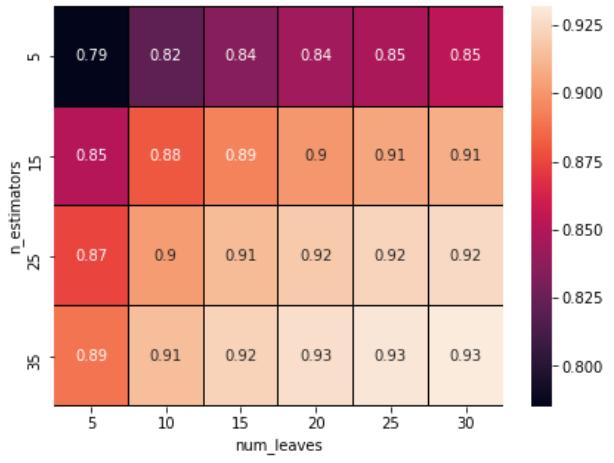
heatmap of train ROC AUC with learning rate: 0.2



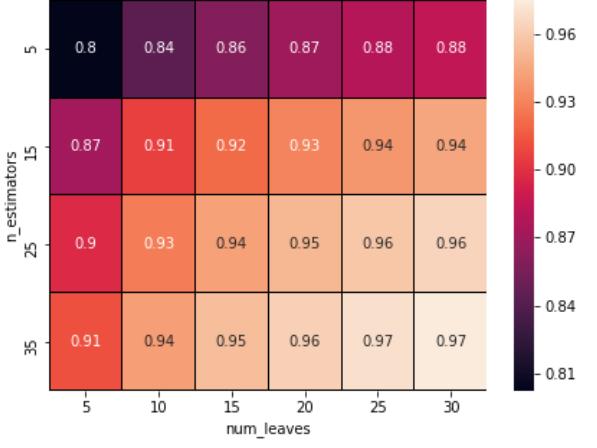
heatmap of train ROC AUC with learning rate: 0.3



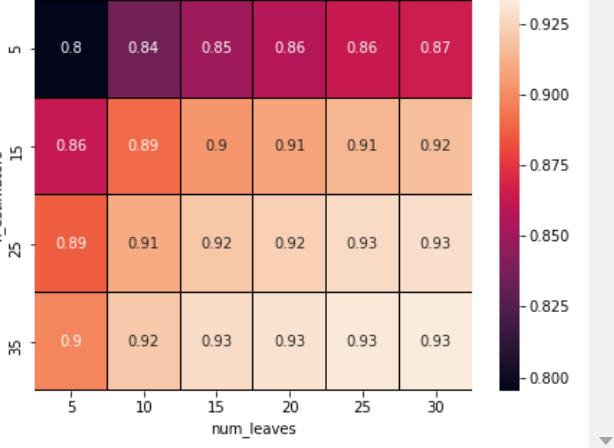
heatmap of train ROC AUC with learning rate: 0.3



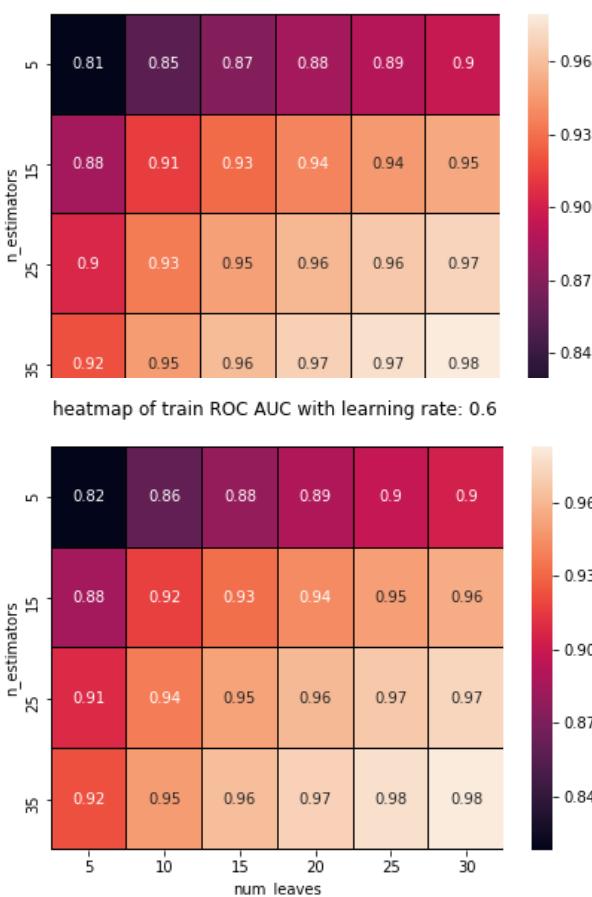
heatmap of train ROC AUC with learning rate: 0.4



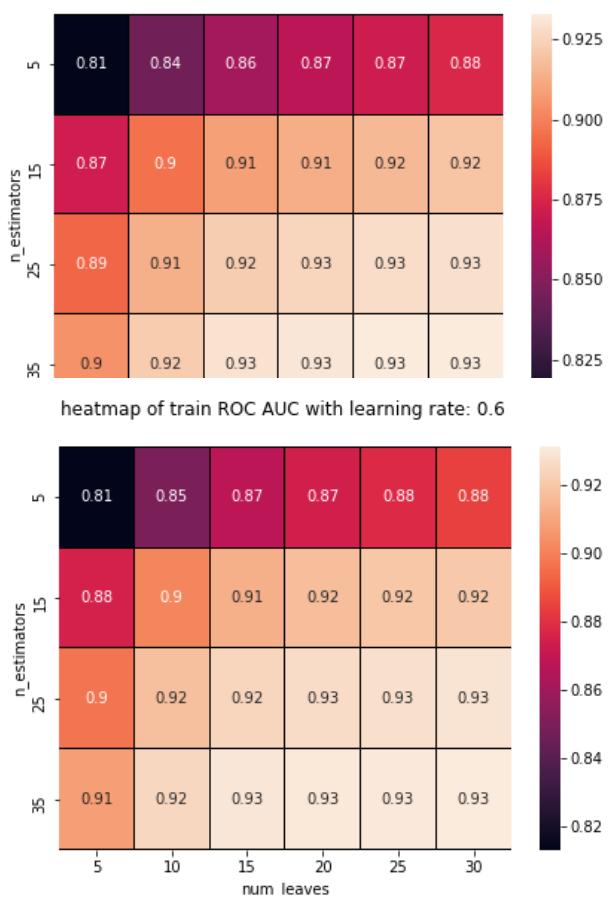
heatmap of train ROC AUC with learning rate: 0.4



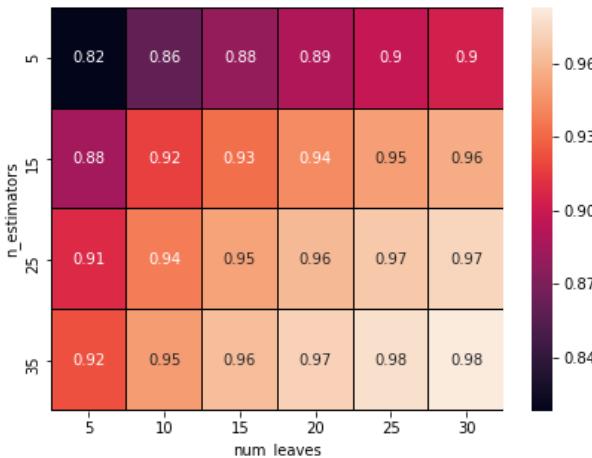
heatmap of train ROC AUC with learning rate: 0.5



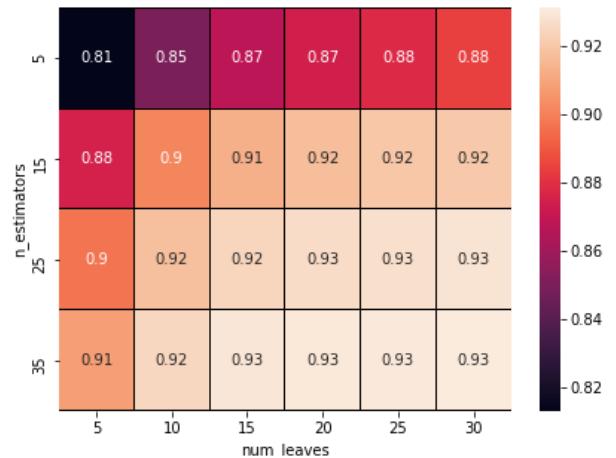
heatmap of train ROC AUC with learning rate: 0.5



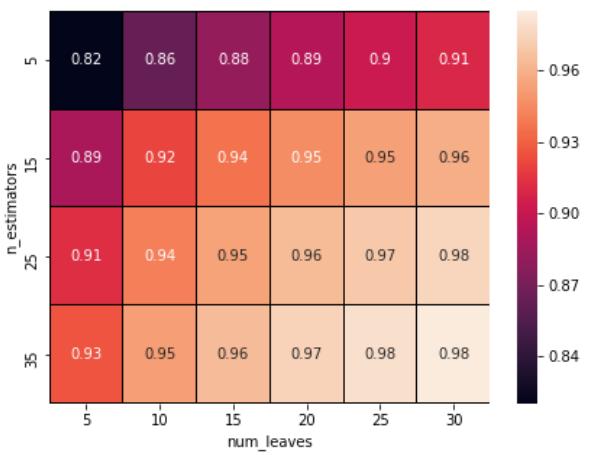
heatmap of train ROC AUC with learning rate: 0.6



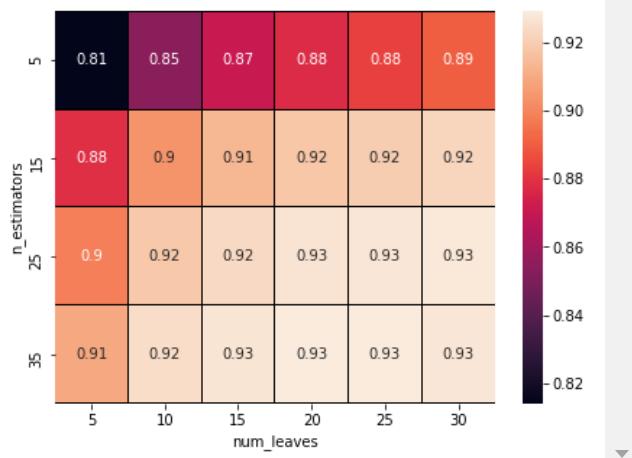
heatmap of train ROC AUC with learning rate: 0.6



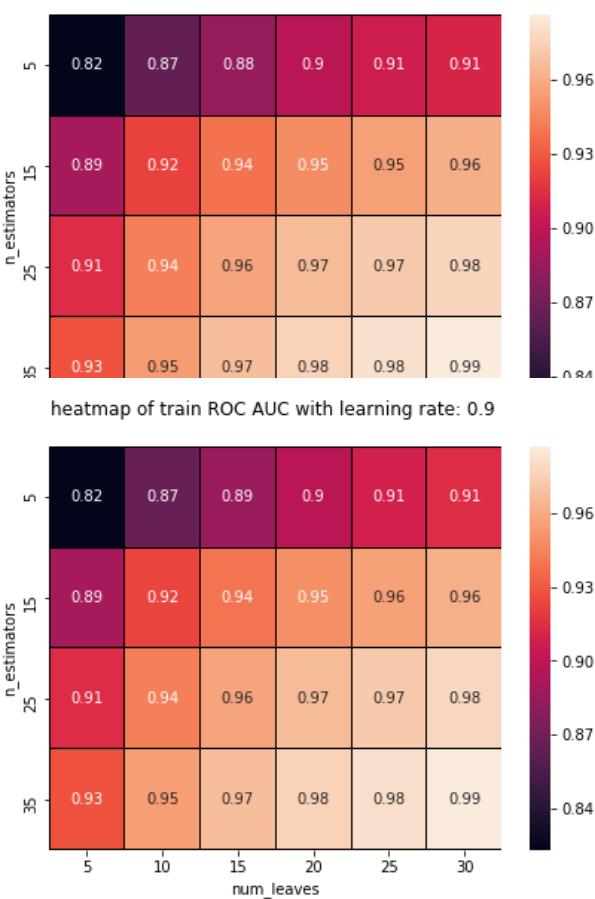
heatmap of train ROC AUC with learning rate: 0.7



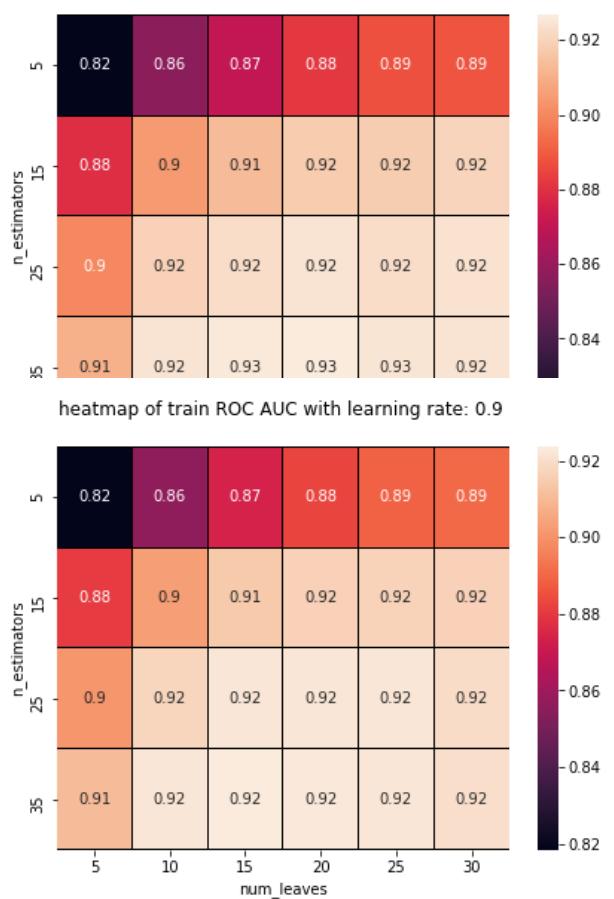
heatmap of train ROC AUC with learning rate: 0.7



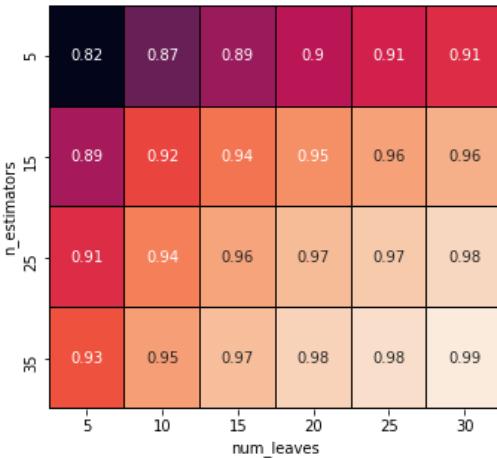
heatmap of train ROC AUC with learning rate: 0.8



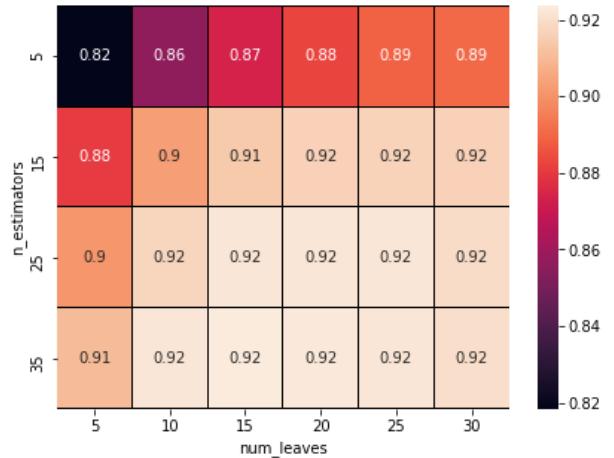
heatmap of train ROC AUC with learning rate: 0.8



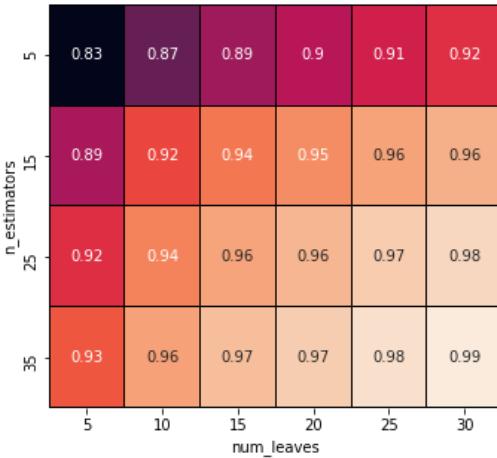
heatmap of train ROC AUC with learning rate: 0.9



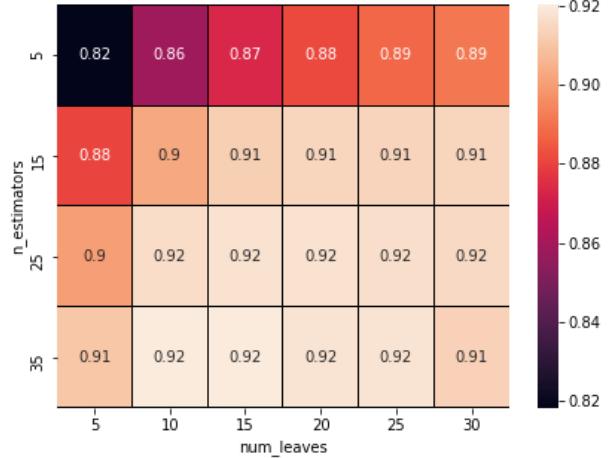
heatmap of train ROC AUC with learning rate: 0.9



heatmap of train ROC AUC with learning rate: 1



heatmap of train ROC AUC with learning rate: 1



best hyperparameter

```
In [58]: # best estimators
print(f"best hyperparameters are: {gsCV.best_params_}")
```

```
best hyperparameters are: {'learning_rate': 0.4, 'n_estimators': 35, 'num_leaves': 30}
```

```
In [59]: #Testing
from sklearn.metrics import roc_auc_score
```

```

from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve

#probability score for ROC_AUC score
y_train_pred_proba=gsCV.predict_proba(train_tf_idf[:60000])[:,1]
y_test_pred_proba=gsCV.predict_proba(test_tf_idf[:40000])[:,1]

train_roc_score=roc_auc_score(y_train[:60000],y_train_pred_proba)
test_roc_score=roc_auc_score(y_test[:40000],y_test_pred_proba)

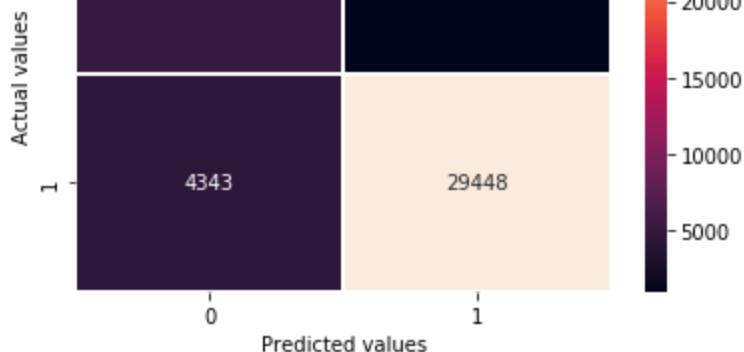
#ploting confusion matrix
y_pred=gsCV.predict(test_tf_idf[:40000])
sn.heatmap(confusion_matrix(y_test[:40000],y_pred), annot=True, fmt="d", linewidths=.5)
plt.title('Confusion Matrix')
plt.xlabel('Predicted values')
plt.ylabel('Actual values')
plt.show()
print("\n\nclassification report:\n",classification_report(y_test[:40000],y_pred))

# ROC Curve (reference:stack overflow with little modification)
train_fpr, train_tpr, train_threshold =roc_curve(y_train[:60000], y_train_pred_proba)
test_fpr, test_tpr, test_threshold =roc_curve(y_test[:40000], y_test_pred_proba)

#ploting ROC curve
plt.figure(figsize=(7,7))
plt.title('Receiver Operating Characteristic')
plt.plot(train_fpr, train_tpr, 'b', label = 'train_AUC = %0.2f' % train_roc_score)
plt.plot(test_fpr, test_tpr, 'r', label = 'test_AUC = %0.2f' % test_roc_score)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([-0.02, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.grid(linestyle='--', linewidth=1)
plt.show()

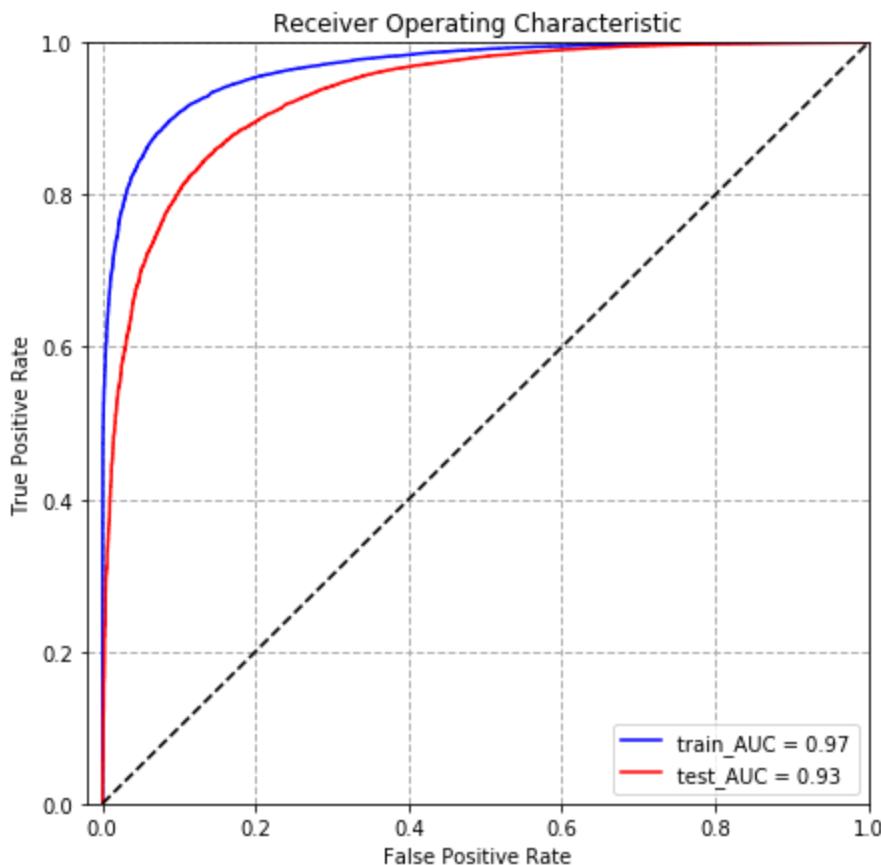
```





classification report:

	precision	recall	f1-score	support
0	0.54	0.84	0.66	6209
1	0.97	0.87	0.92	33791
micro avg	0.87	0.87	0.87	40000
macro avg	0.76	0.85	0.79	40000
weighted avg	0.90	0.87	0.88	40000



[5.2.3] Applying XGBOOST on AVG W2V, SET 3

In [60]:

```
# By using grid search
from lightgbm import LGBMClassifier
from sklearn.model_selection import GridSearchCV

param={'num_leaves':[5, 10, 15, 20, 25, 30], 'learning_rate': [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1],
       'n_estimators':[5, 15, 25, 35]}
```

```

# taking decision tree as estimator for grid search
clf_lightgb=LGBMClassifier(boosting_type='gbdt',class_weight='balanced',random_
_state=1,n_jobs=-1)

gsCV=GridSearchCV(estimator=clf_lightgb, param_grid=param, scoring='roc_auc',r
eturn_train_score=True,n_jobs=-2)
gsCV.fit(train_w2v[:60000],y_train[:60000])

# storing results
result=gsCV.cv_results_

roc_auc_cv=result['mean_test_score'].reshape(10,4,6) #reshaping for heatmap vi
sualisation
roc_auc_train=result['mean_train_score'].reshape(10,4,6) #reshaping for heatmap
p visualisation

```

In [62]: # visualising train_ROC_AUC and cv_ROC_AUC using heatmap

```

for index,lr in enumerate(param['learning_rate']):
    from matplotlib import gridspec

    fig = plt.figure(figsize=(15,5))
    gs  = gridspec.GridSpec(1, 2)

    ax0 = plt.subplot(gs[0])

    # train heatmap
    sn.heatmap(data=roc_auc_train[index],lineweights=0.01,annot=True,xticklabel
s=param['num_leaves'],
                yticklabels=param['n_estimators'],linecolor='black',fmt='%.2g',a
x=ax0)

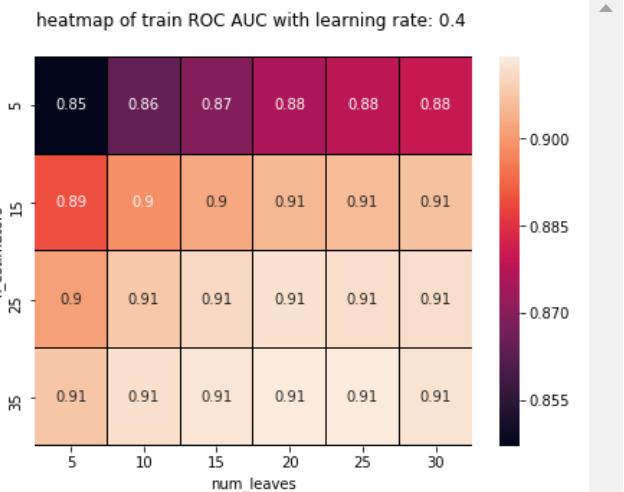
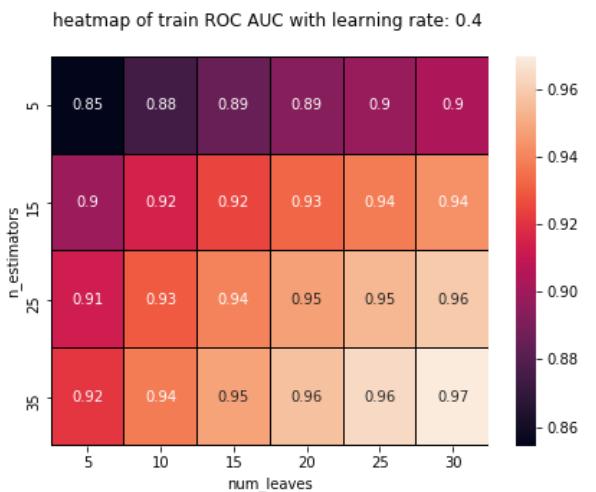
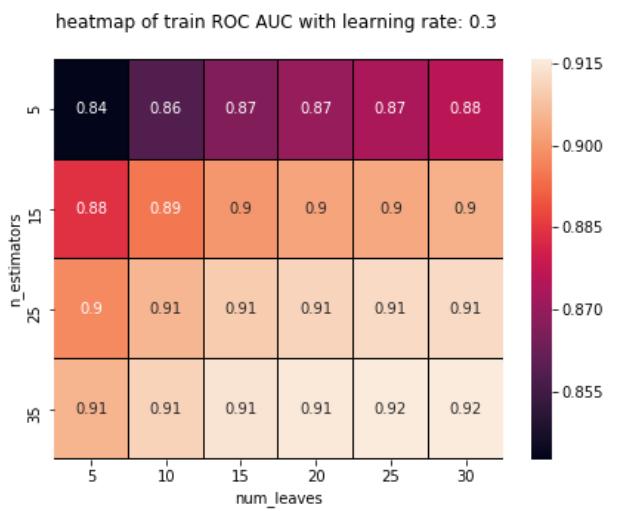
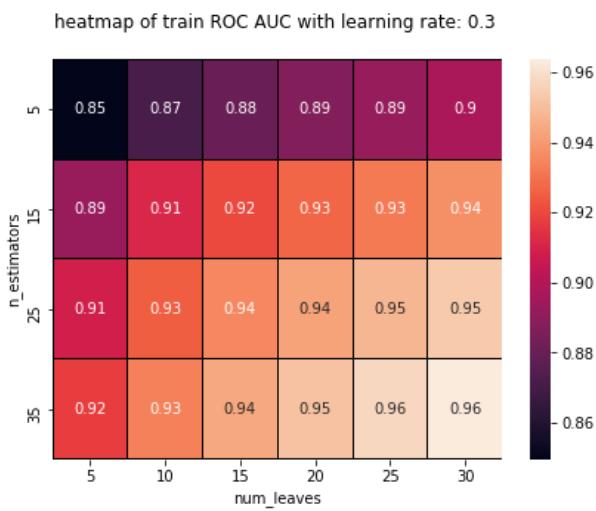
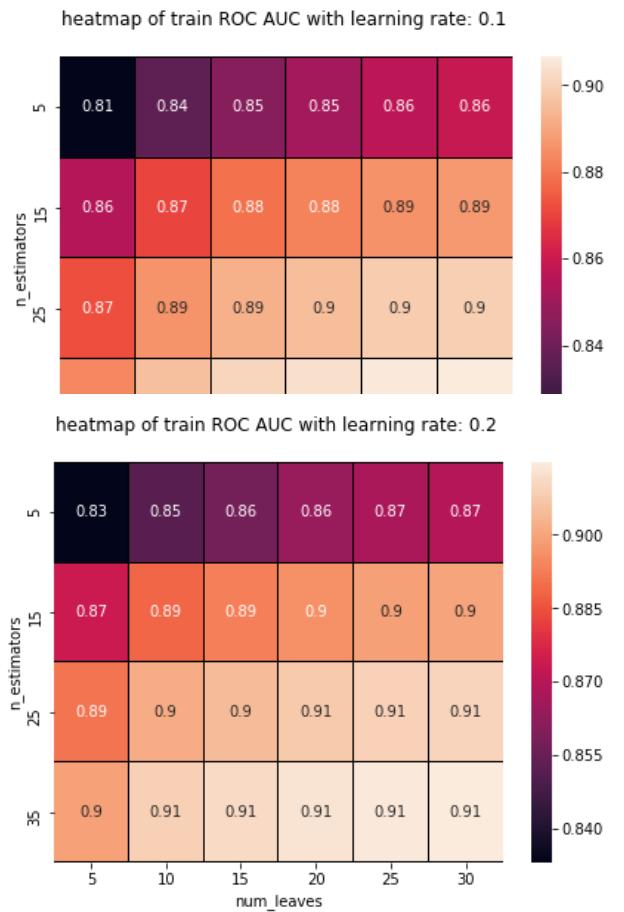
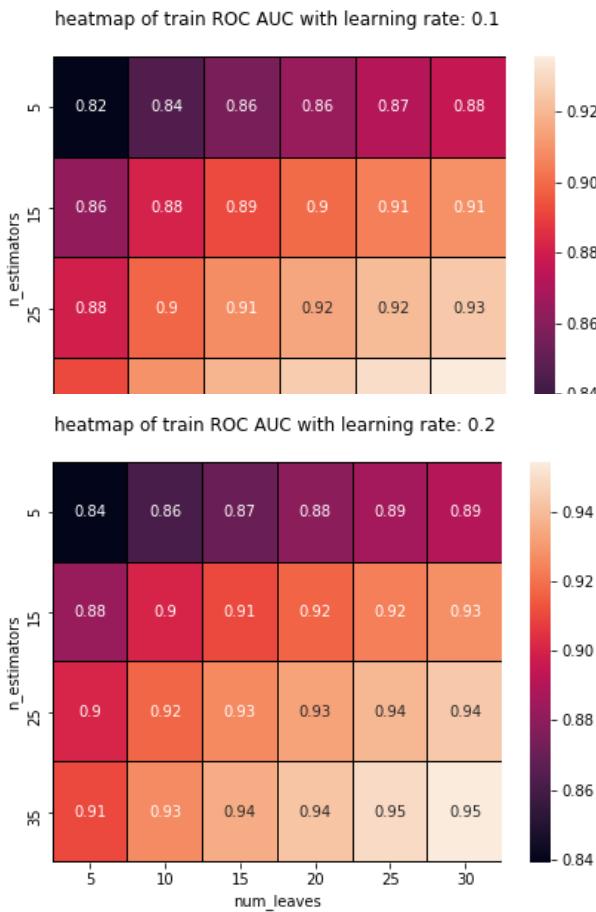
    plt.xlabel('num_leaves')
    plt.ylabel('n_estimators')
    plt.title(f"heatmap of train ROC AUC with learning rate: {lr} \n")
    ****
    ****

    ax1 = plt.subplot(gs[1])

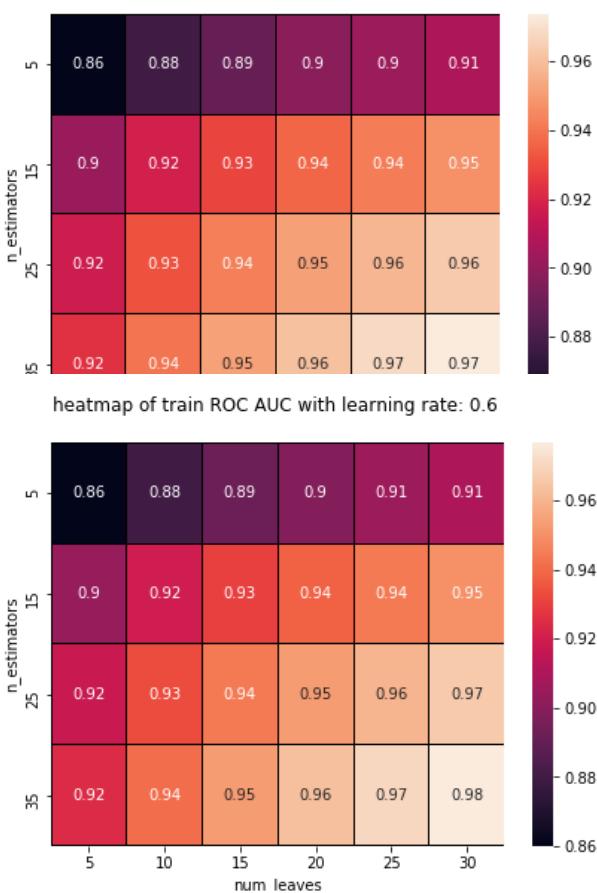
    # cv heatmap
    sn.heatmap(data=roc_auc_cv[index],lineweights=0.01,annot=True,xticklabels=p
aram['num_leaves'],
                yticklabels=param['n_estimators'],linecolor='black',fmt='%.2g',a
x=ax1)

    plt.xlabel('num_leaves')
    plt.ylabel('n_estimators')
    plt.title(f"heatmap of train ROC AUC with learning rate: {lr} \n")
    plt.show()

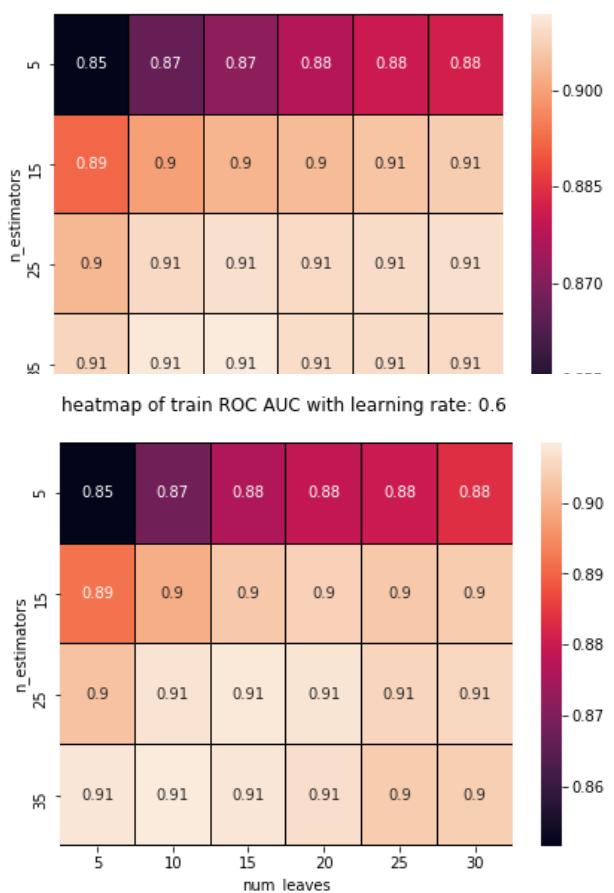
```



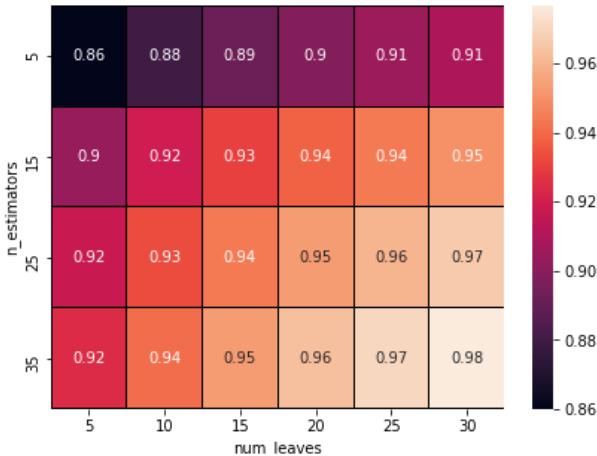
heatmap of train ROC AUC with learning rate: 0.5



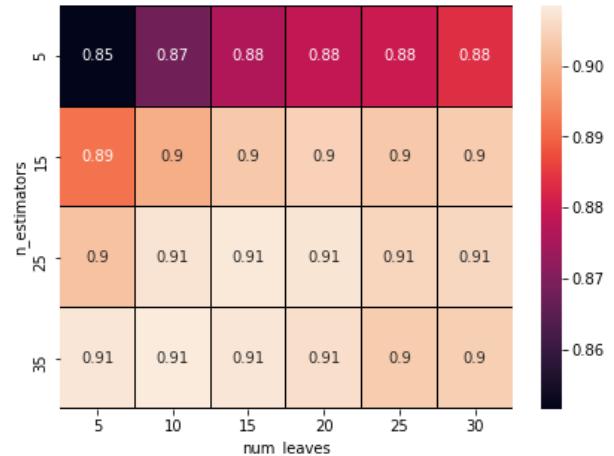
heatmap of train ROC AUC with learning rate: 0.5



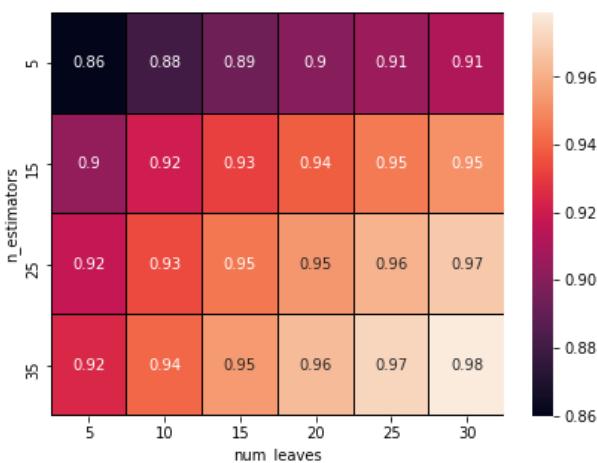
heatmap of train ROC AUC with learning rate: 0.6



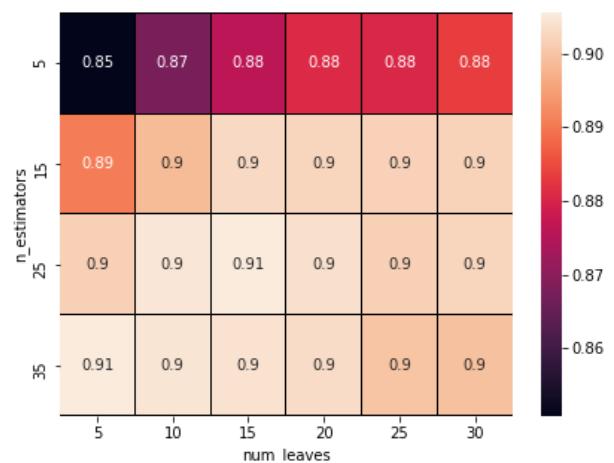
heatmap of train ROC AUC with learning rate: 0.6



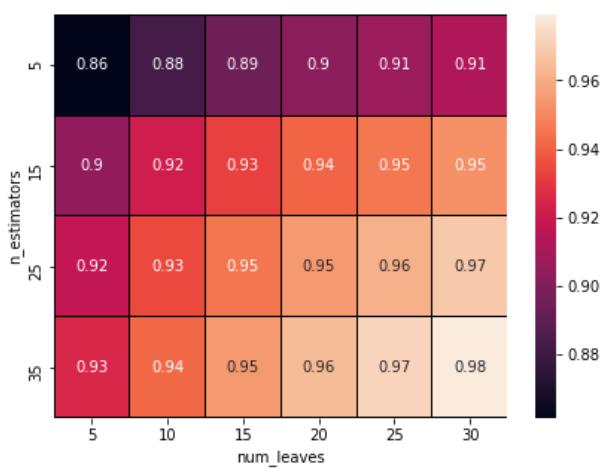
heatmap of train ROC AUC with learning rate: 0.7



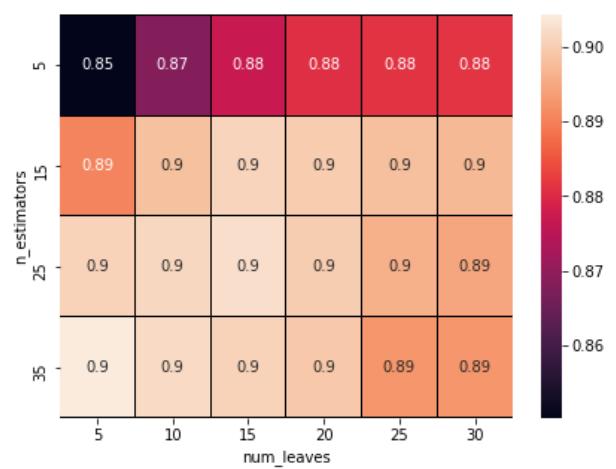
heatmap of train ROC AUC with learning rate: 0.7



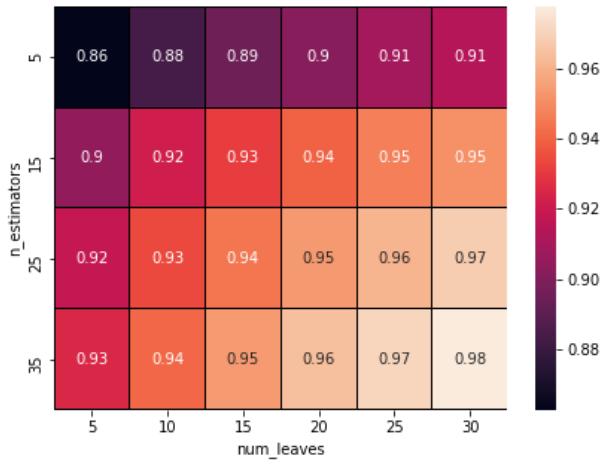
heatmap of train ROC AUC with learning rate: 0.8



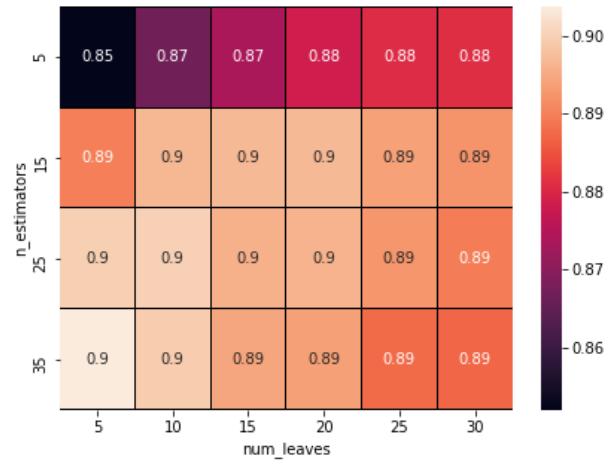
heatmap of train ROC AUC with learning rate: 0.8



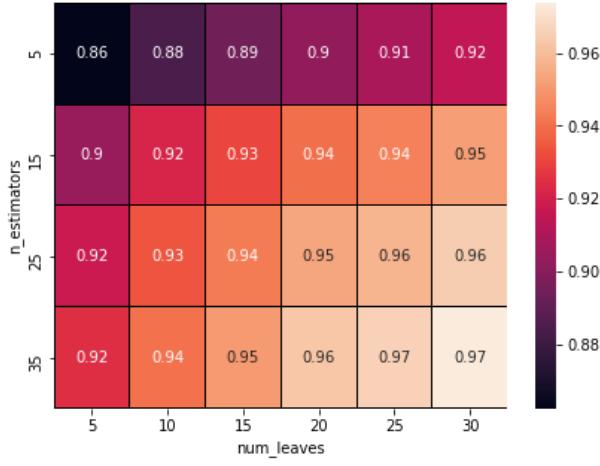
heatmap of train ROC AUC with learning rate: 0.9



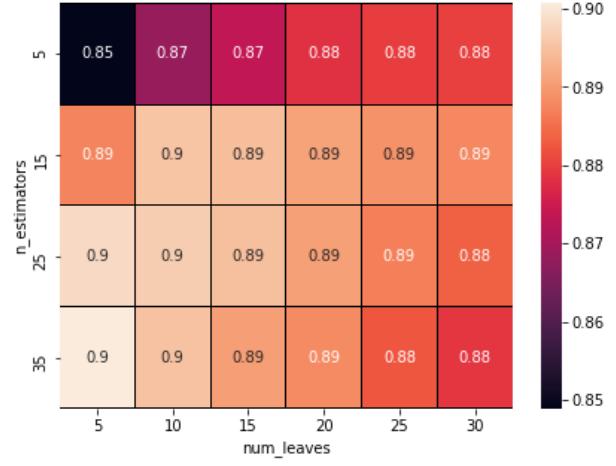
heatmap of train ROC AUC with learning rate: 0.9



heatmap of train ROC AUC with learning rate: 1



heatmap of train ROC AUC with learning rate: 1



best hyperparameter

```
In [63]: # best estimators
print(f"best hyperparameters are: {gsCV.best_params_}")
```

```
best hyperparameters are: {'learning_rate': 0.3, 'n_estimators': 35, 'num_leaves': 25}
```

```
In [64]: #Testing
from sklearn.metrics import roc_auc_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve
```

```

#probability score for ROC_AUC score
y_train_pred_proba=gsCV.predict_proba(train_w2v[:60000])[:,1]
y_test_pred_proba=gsCV.predict_proba(test_w2v[:40000])[:,1]

train_roc_score=roc_auc_score(y_train[:60000],y_train_pred_proba)
test_roc_score=roc_auc_score(y_test[:40000],y_test_pred_proba)

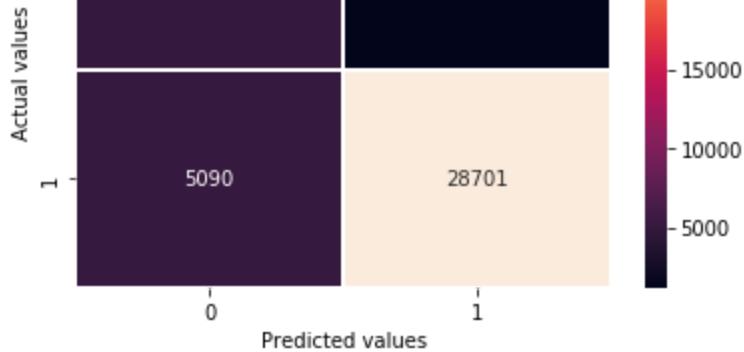
#ploting confusion matrix
y_pred=gsCV.predict(test_w2v[:40000])
sn.heatmap(confusion_matrix(y_test[:40000],y_pred), annot=True, fmt="d", linewidths=.5)
plt.title('Confusion Matrix')
plt.xlabel('Predicted values')
plt.ylabel('Actual values')
plt.show()
print("\n\nclassification report:\n",classification_report(y_test[:40000],y_pred))

# ROC Curve (reference:stack overflow with little modification)
train_fpr, train_tpr, train_threshold =roc_curve(y_train[:60000], y_train_pred_proba)
test_fpr, test_tpr, test_threshold =roc_curve(y_test[:40000], y_test_pred_proba)

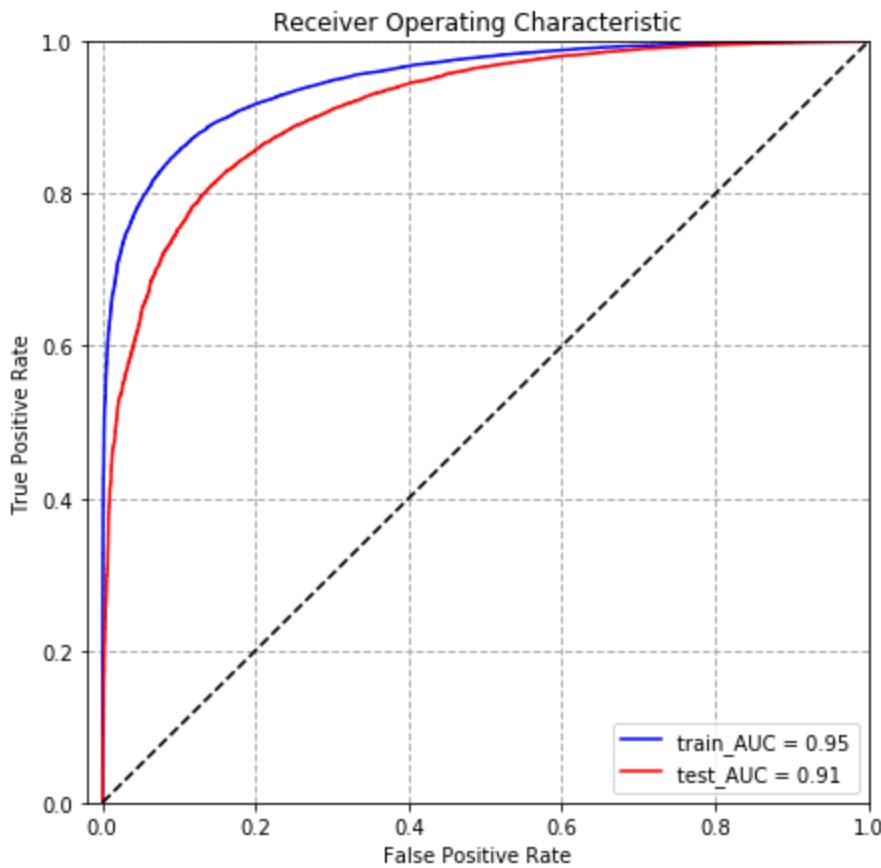
#ploting ROC curve
plt.figure(figsize=(7,7))
plt.title('Receiver Operating Characteristic')
plt.plot(train_fpr, train_tpr, 'b', label = 'train_AUC = %0.2f' % train_roc_score)
plt.plot(test_fpr, test_tpr, 'r', label = 'test_AUC = %0.2f' % test_roc_score)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([-0.02, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.grid(linestyle='--', linewidth=1)
plt.show()

```





```
classification report:
      precision    recall   f1-score   support
0       0.50      0.81      0.62     6209
1       0.96      0.85      0.90    33791
           micro avg   0.84      0.84      0.84     40000
           macro avg   0.73      0.83      0.76     40000
           weighted avg  0.89      0.84      0.86     40000
```



[5.2.4] Applying XGBOOST on TFIDF W2V, SET 4

```
In [65]: # By using grid search
from lightgbm import LGBMClassifier
from sklearn.model_selection import GridSearchCV

param={'num_leaves':[5, 10, 15, 20, 25, 30], 'learning_rate': [0.1,0.2,0.3,0.4, 0.5,0.6,0.7,0.8,0.9,1],
       'n_estimators':[5,15,25,35]}
```

```

# taking decision tree as estimator for grid search
clf_lightgb=LGBMClassifier(boosting_type='gbdt',class_weight='balanced',random_
_state=1,n_jobs=-1)

gsCV=GridSearchCV(estimator=clf_lightgb, param_grid=param, scoring='roc_auc',r
eturn_train_score=True,n_jobs=-2)
gsCV.fit(train_tf_idf_w2v[:60000],y_train[:60000])

# storing results
result=gsCV.cv_results_

roc_auc_cv=result['mean_test_score'].reshape(10,4,6) #reshaping for heatmap vi
sualisation
roc_auc_train=result['mean_train_score'].reshape(10,4,6) #reshaping for heatma
p visualisation

```

In [66]: # visualising train_ROC_AUC and cv_ROC_AUC using heatmap

```

for index,lr in enumerate(param['learning_rate']):
    from matplotlib import gridspec

    fig = plt.figure(figsize=(15,5))
    gs  = gridspec.GridSpec(1, 2)

    ax0 = plt.subplot(gs[0])

    # train heatmap
    sn.heatmap(data=roc_auc_train[index],lineweights=0.01,annot=True,xticklabel
s=param['num_leaves'],
                yticklabels=param['n_estimators'],linecolor='black',fmt='.2g',a
x=ax0)

    plt.xlabel('num_leaves')
    plt.ylabel('n_estimators')
    plt.title(f"heatmap of train ROC AUC with learning rate: {lr} \n")
    ****
    ****

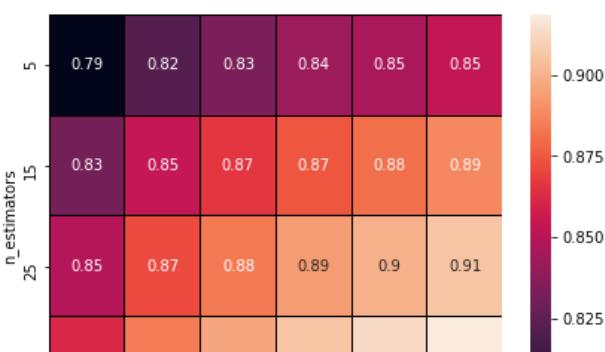
    ax1 = plt.subplot(gs[1])

    # cv heatmap
    sn.heatmap(data=roc_auc_cv[index],lineweights=0.01,annot=True,xticklabels=p
aram['num_leaves'],
                yticklabels=param['n_estimators'],linecolor='black',fmt='.2g',a
x=ax1)

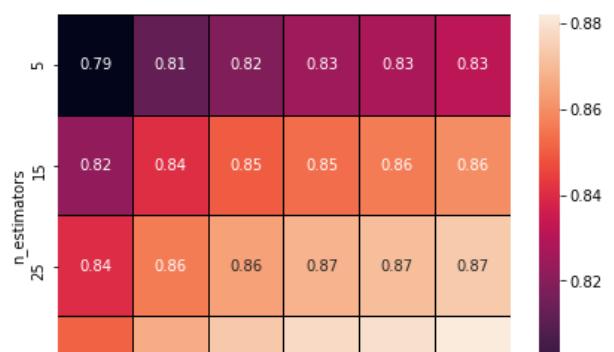
    plt.xlabel('num_leaves')
    plt.ylabel('n_estimators')
    plt.title(f"heatmap of train ROC AUC with learning rate: {lr} \n")
    plt.show()

```

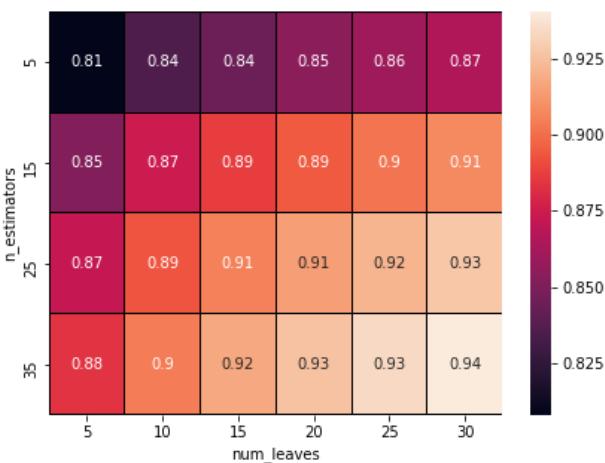
heatmap of train ROC AUC with learning rate: 0.1



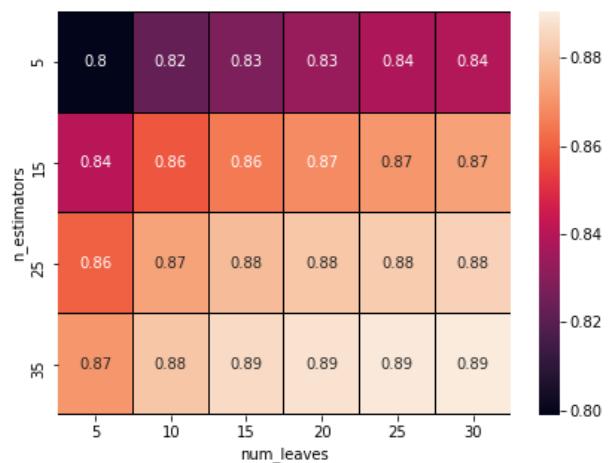
heatmap of train ROC AUC with learning rate: 0.1



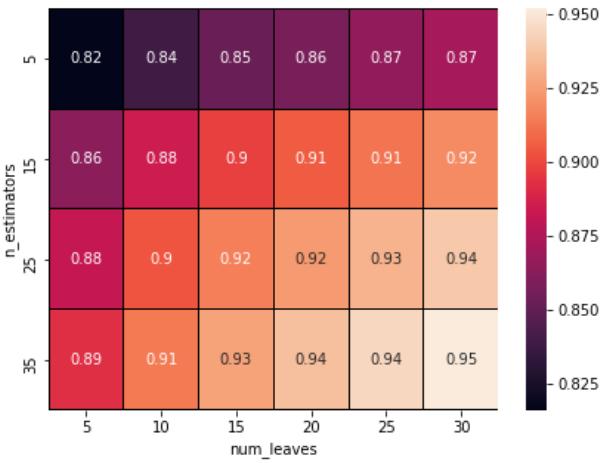
heatmap of train ROC AUC with learning rate: 0.2



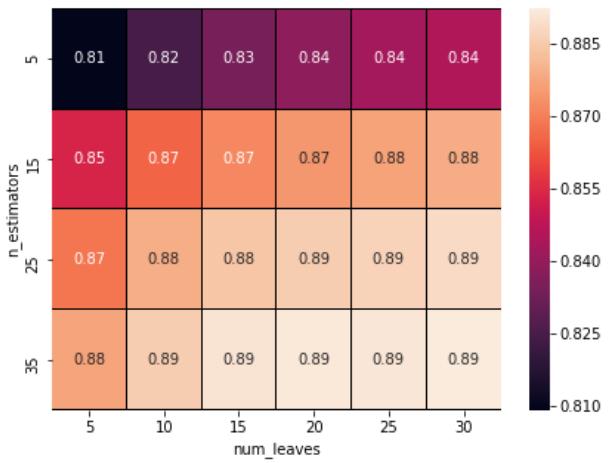
heatmap of train ROC AUC with learning rate: 0.2



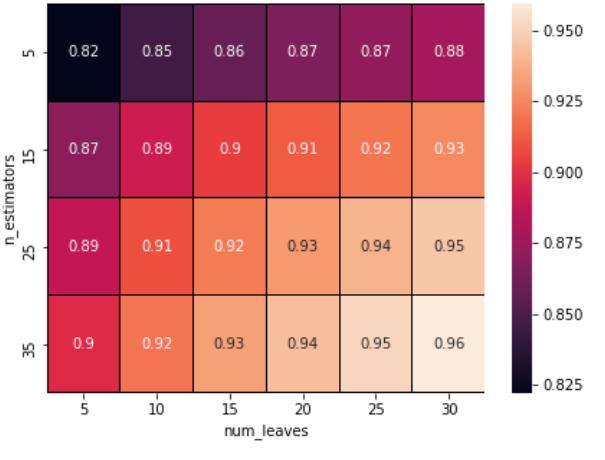
heatmap of train ROC AUC with learning rate: 0.3



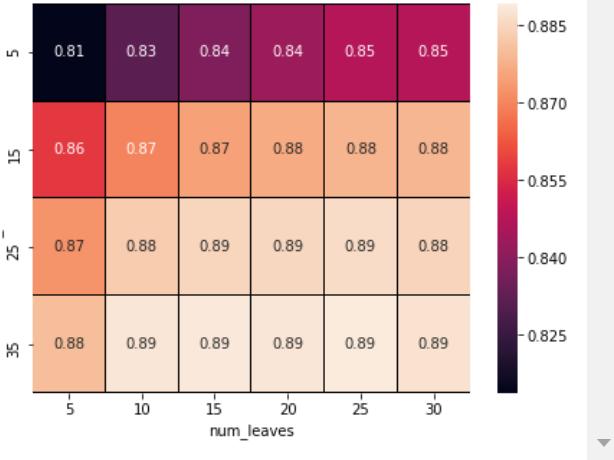
heatmap of train ROC AUC with learning rate: 0.3

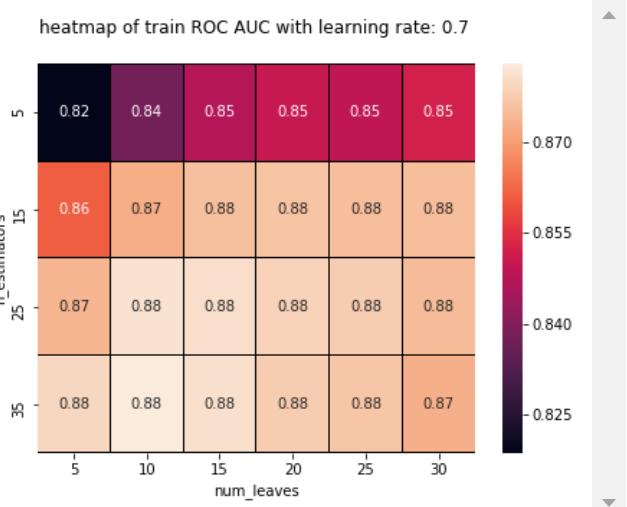
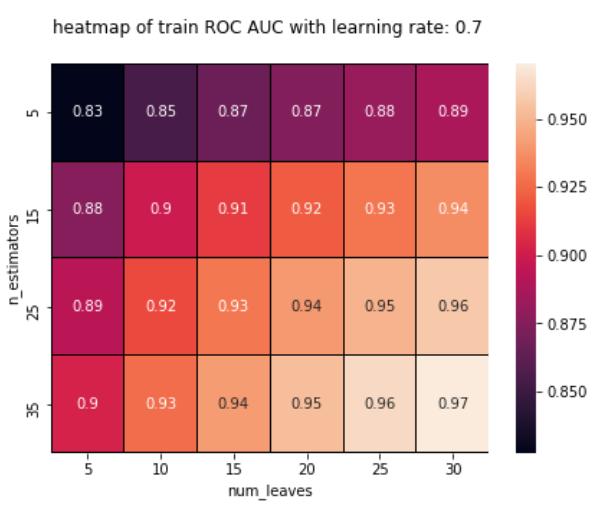
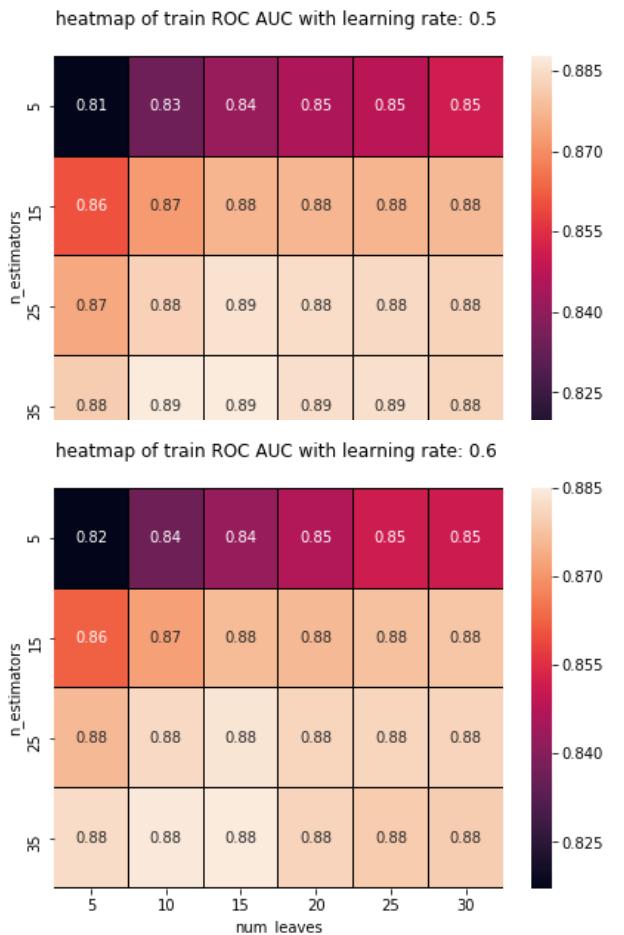
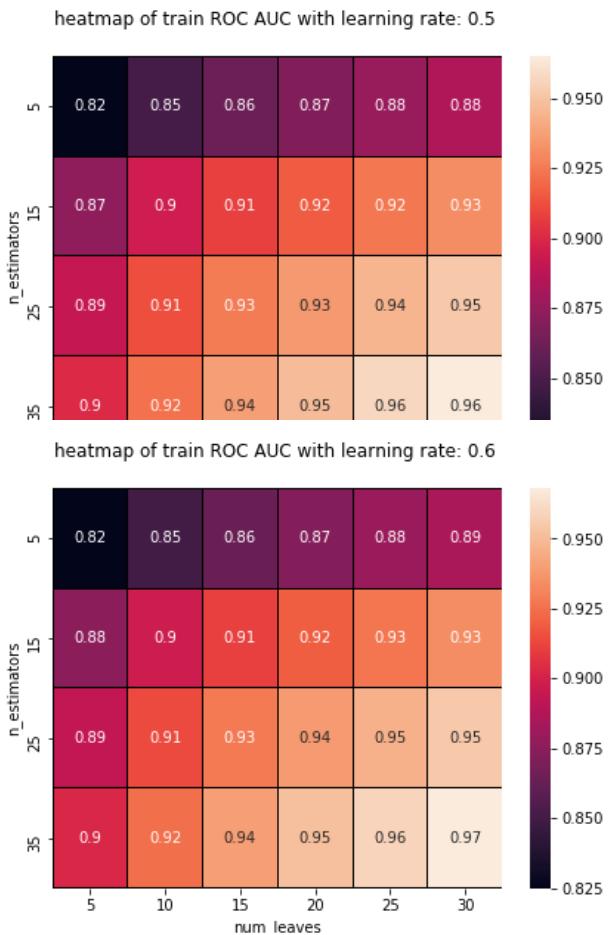


heatmap of train ROC AUC with learning rate: 0.4

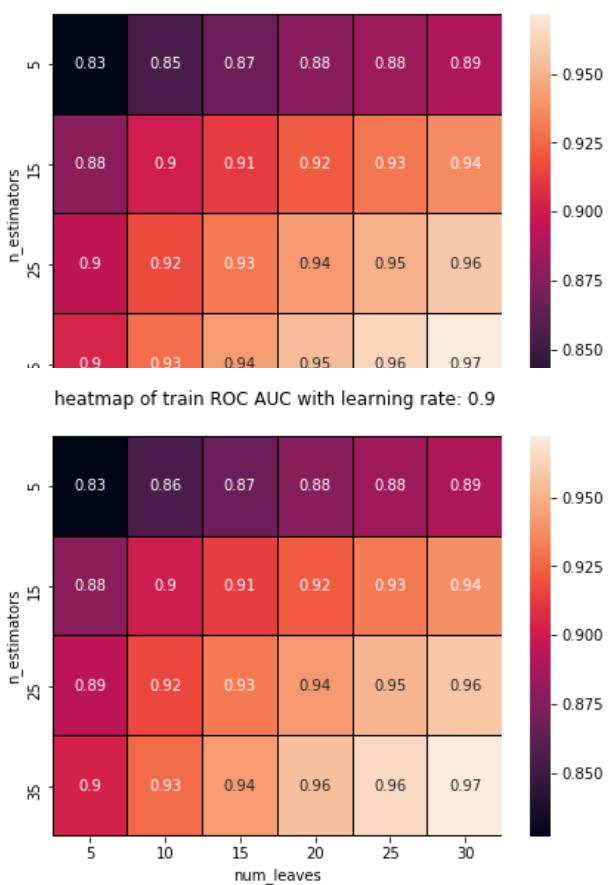


heatmap of train ROC AUC with learning rate: 0.4

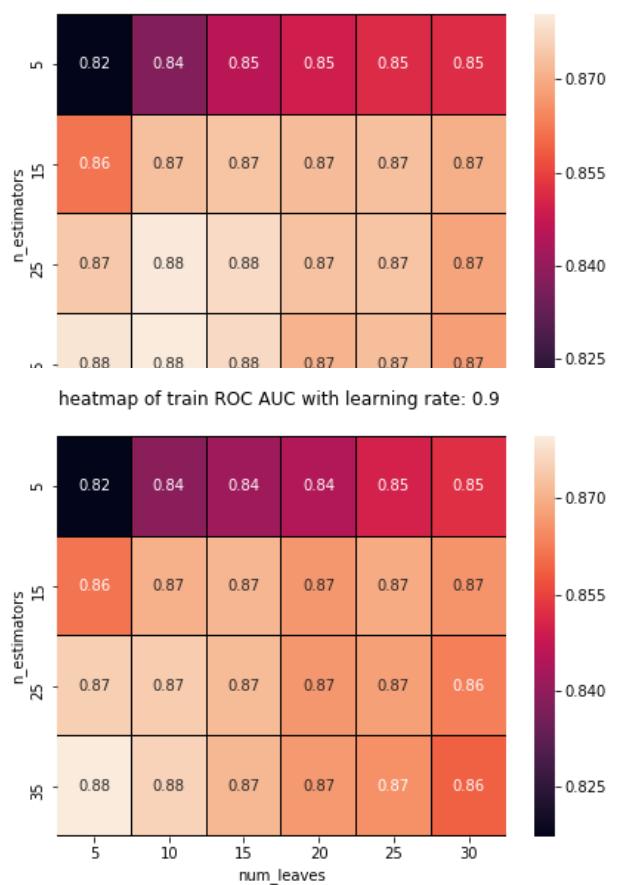




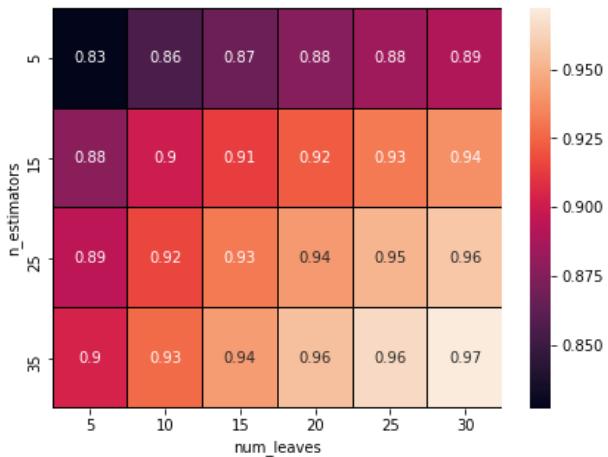
heatmap of train ROC AUC with learning rate: 0.8



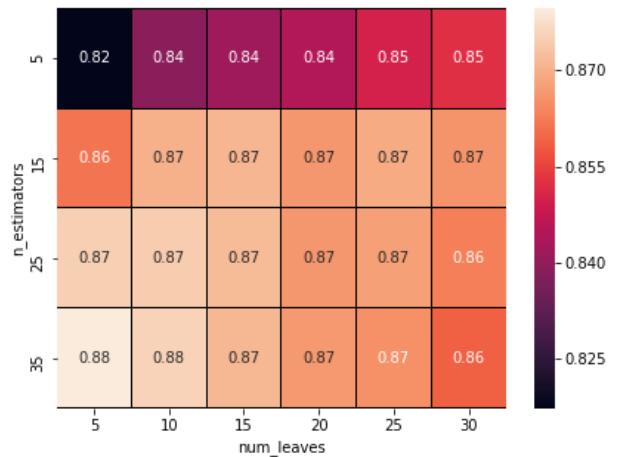
heatmap of train ROC AUC with learning rate: 0.8



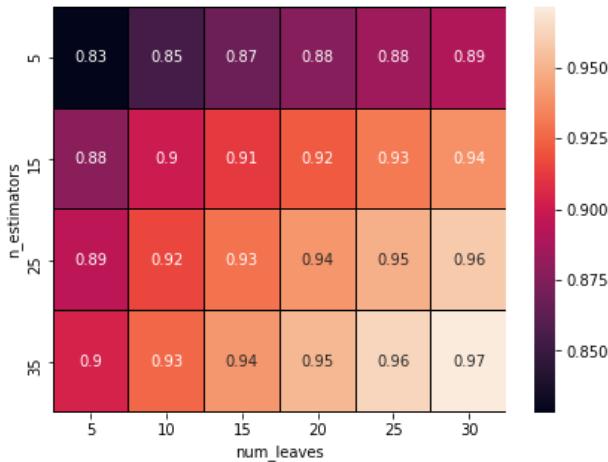
heatmap of train ROC AUC with learning rate: 0.9



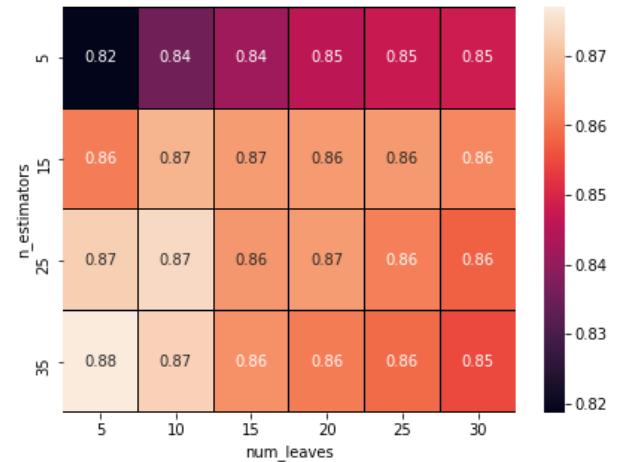
heatmap of train ROC AUC with learning rate: 0.9



heatmap of train ROC AUC with learning rate: 1



heatmap of train ROC AUC with learning rate: 1



best hyperparameter

```
In [67]: # best estimators
print(f"best hyperparameters are: {gsCV.best_params_}")
```

```
best hyperparameters are: {'learning_rate': 0.3, 'n_estimators': 35, 'num_leaves': 30}
```

```
In [68]: #testing
from sklearn.metrics import roc_auc_score
```

```

from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve

#probability score for ROC_AUC score
y_train_pred_proba=gsCV.predict_proba(train_tf_idf_w2v[:60000])[:,1]
y_test_pred_proba=gsCV.predict_proba(test_tf_idf_w2v[:20000])[:,1]

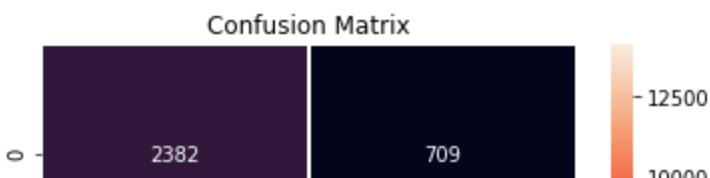
train_roc_score=roc_auc_score(y_train[:60000],y_train_pred_proba)
test_roc_score=roc_auc_score(y_test[:20000],y_test_pred_proba)

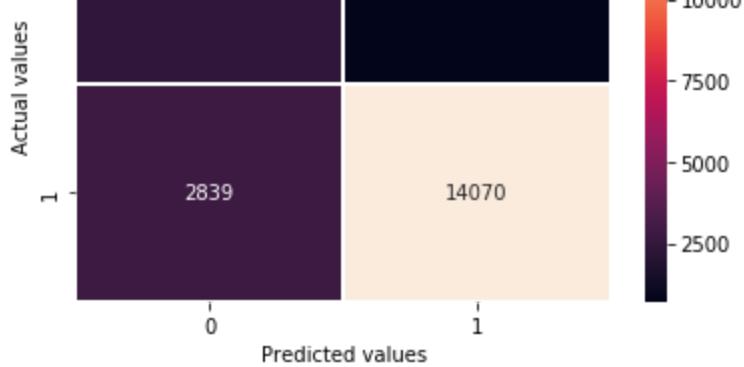
#ploting confusion matrix
y_pred=gsCV.predict(test_tf_idf_w2v[:20000])
sn.heatmap(confusion_matrix(y_test[:20000],y_pred), annot=True, fmt="d", linewidths=.5)
plt.title('Confusion Matrix')
plt.xlabel('Predicted values')
plt.ylabel('Actual values')
plt.show()
print("\n\nclassification report:\n",classification_report(y_test[:20000],y_pred))

# ROC Curve (reference:stack overflow with little modification)
train_fpr, train_tpr, train_threshold =roc_curve(y_train[:60000], y_train_pred_proba)
test_fpr, test_tpr, test_threshold =roc_curve(y_test[:20000], y_test_pred_proba)

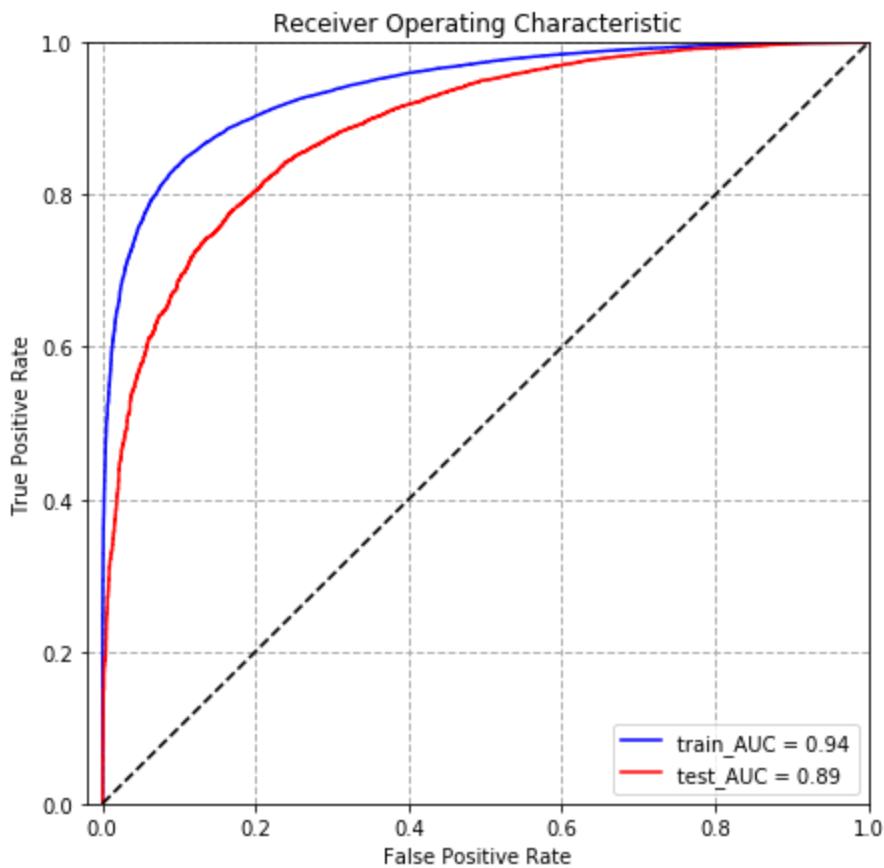
#ploting ROC curve
plt.figure(figsize=(7,7))
plt.title('Receiver Operating Characteristic')
plt.plot(train_fpr, train_tpr, 'b', label = 'train_AUC = %0.2f' % train_roc_score)
plt.plot(test_fpr, test_tpr, 'r', label = 'test_AUC = %0.2f' % test_roc_score)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([-0.02, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.grid(linestyle='--', linewidth=1)
plt.show()

```





```
classification report:
              precision    recall   f1-score   support
             0       0.46      0.77      0.57     3091
             1       0.95      0.83      0.89    16909
        micro avg       0.82      0.82      0.82     20000
        macro avg       0.70      0.80      0.73     20000
    weighted avg       0.88      0.82      0.84     20000
```



[6] Conclusions

```
In [38]: # Random Forest

from prettytable import PrettyTable
x=PrettyTable(field_names=["vectoriser", "max_depth", "n_estimator", "train_AUC",
"Generalization AUC"])

x.add_row(["BOW", 16, 1000, 0.96, 0.93])
```

```

x.add_row(["TF IDF", 20, 1000, 0.97, 0.93])
x.add_row(["W2V", 20, 1000, 1.0, 0.91])
x.add_row(["TFIDF W2V", 20, 1000, 1, 0.88])

print(35*"*", "Random Forest", 35*)
print(x, "\n")

***** Random Forest *****

# LightGBM

from prettytable import PrettyTable
y=PrettyTable(field_names=["vectoriser","learning_rate","max_leaves","n_estimator","train_AUC","Generalization AUC"])

y.add_row(["BOW", 0.1, 30, 500, 0.99, 0.95])
y.add_row(["TF IDF", 0.4, 30, 35, 0.97, 0.93])
y.add_row(["W2V", 0.3, 25, 35, 0.95, 0.91])
y.add_row(["TFIDF W2V", 0.3, 30, 35, 0.94, 0.89])

print(35*"*", "LightGBM", 35*)
print(y)

```

Random Forest *****						

vectoriser max_depth n_estimator train_AUC Generalization AUC						
+-----+-----+-----+-----+-----+						
BOW 16 1000 0.96 0.93						
TF IDF 20 1000 0.97 0.93						
W2V 20 1000 1.0 0.91						
TFIDF W2V 20 1000 1 0.88						
+-----+-----+-----+-----+-----+						

LightGBM *****						

vectoriser learning_rate max_leaves n_estimator train_AUC Generalization AUC						
+-----+-----+-----+-----+-----+						
BOW 0.1 30 500 0.99						
0.95						
TF IDF 0.4 30 35 0.97						
0.93						
W2V 0.3 25 35 0.95						
0.91						
TFIDF W2V 0.3 30 35 0.94						
0.89						
+-----+-----+-----+-----+-----+						
- - - - -						

Observation

- Best generalization ROC_AUC is achieved by LightGBM model using BOW vectoriser with hyperparameter: {"learning_rate": 0.1, "max_leaves": 30, "n_estimator": 500}
- Training using LightGBM is faster than simple XGBoost and RF.
- Ensemble model works best among all the previous models but takes more time to train.

