

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 [219]:

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

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

```

import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os
from sklearn import tree
import graphviz
from sklearn.model_selection import RandomizedSearchCV

```

In [197]:

```

# 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 500000""", 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 (500000, 10)

Out[197]:

	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

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	1	1219017600

In [198]:

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

In [199]:

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

(80668, 7)

Out[199]:

	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 [200]:

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

Out[200]:

	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 [201]:

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

Out[201]:

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 [202]:

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

Out[202]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
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 [203]:

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

```
sorted_data_filtered_duplicates_value = filtered_data, axis = 0, ascending = True, inplace = True, na_posi  
cksort', na_position='last')
```

In [204]:

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

Out[204]:

(348262, 10)

In [205]:

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

Out[205]:

69.6524

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 [206]:

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

Out[206]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5	12248926
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2	4	12128832

In [207]:

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

In [208]:

```
#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()
```

```
(348260, 10)
```

```
Out[208]:
```

```
1    293516
0     54744
Name: Score, dtype: int64
```

segragating datapoints w.r.t calss labels and sampling optimum number of data points

```
In [209]:
```

```
zero_class=final[final.Score==0]
print(zero_class['Score'].value_counts())
print(zero_class.shape)
one_class=final[final.Score==1]
print(one_class['Score'].value_counts())
print(one_class.shape)
```

```
0     54744
Name: Score, dtype: int64
(54744, 10)
1     293516
Name: Score, dtype: int64
(293516, 10)
```

```
In [210]:
```

```
one_class1=one_class.sample(n=50000)
zero_class1=zero_class.sample(n=50000)
print(zero_class1.shape)
print(one_class1.shape)
combined_frame=pd.concat([zero_class1,one_class1])
print(combined_frame.shape)
```

```
(50000, 10)
(50000, 10)
(100000, 10)
```

```
In [211]:
```

```
final_new_frame=combined_frame.sample(frac=1)
```

```
In [212]:
```

```
print(type(final_new_frame))
print(final_new_frame.shape)
print(final_new_frame['Score'].value_counts())
```

```
<class 'pandas.core.frame.DataFrame'>
(100000, 10)
1     50000
0     50000
Name: Score, dtype: int64
```

```
In [213]:
```

```
# 1.11 -this here cotinuation https://stackoverflow.com/a/47091490/4084039
import re
from bs4 import BeautifulSoup

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)
```

```

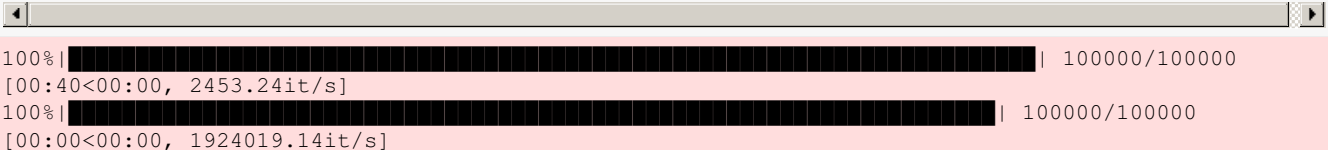
phrase = re.sub(r"\s", " is", phrase)
phrase = re.sub(r"\d", " would", phrase)
phrase = re.sub(r"\ll", " will", phrase)
phrase = re.sub(r"\t", " not", phrase)
phrase = re.sub(r"\ve", " have", phrase)
phrase = re.sub(r"\m", " am", phrase)
return phrase

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

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

j=0
for i in tqdm(preprocessed_reviews):
    j=j+1
print(j)

```



```

100%|████████████████████████████████████████████████████████████████████████████████| 100000/100000
[00:40<00:00, 2453.24it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 100000/100000
[00:00<00:00, 1924019.14it/s]

```

100000

[4] Featurization

[4.1] BAG OF WORDS

In [291]:

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(preprocessed_reviews, final_new_frame['Score'],
test_size=0.33)
```

In [292]:

```
#Bow
vectorizer = CountVectorizer()
vectorizer.fit(X_train)

X_train_bow = vectorizer.transform(X_train)
X_test_bow = vectorizer.transform(X_test)
```

[4.3] TF-IDF

In [304]:

```
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(X_train)

X_train_tfidf = tf_idf_vect.transform(X_train)
X_test_tfidf = tf_idf_vect.transform(X_test)
```

[4.4] Word2Vec

In [215]:

```
# 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 [0]:

```
# Using Google News Word2Vectors

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

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


◀ ▶

=====

```
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 3817
```

Тн [0] .

```
# 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_model.wv and word in tfidf_feat:
            vec = w2v_model.wv[word]
            #
            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors.append(sent_vec)
    row += 1
```

[illegible]

[5] Assignment 8: Decision Trees

1. Apply Decision Trees 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 (best `depth` in range [4,6, 8, 9,10,12,14,17] , and the best `min_samples_split` in range [2,10,20,30,40,50])

- 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. Graphviz

- Visualize your decision tree with Graphviz. It helps you to understand how a decision is being made, given a new vector.
- Since feature names are not obtained from word2vec related models, visualize only BOW & TFIDF decision trees using Graphviz
- Make sure to print the words in each node of the decision tree instead of printing its index.
- Just for visualization purpose, limit `max_depth` to 2 or 3 and either embed the generated images of graphviz in your notebook, or directly upload them as .png files.

4. Feature importance

- Find the top 20 important features from both feature sets **Set 1** and **Set 2** using `feature_importances_` method of **Decision Tree Classifier** and print their corresponding feature names

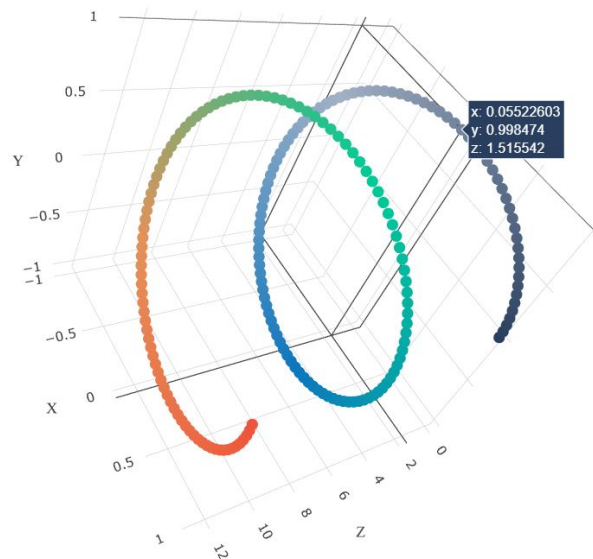
5. Feature engineering

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

6. Representation of results

You need to plot the performance of model both on train data and across validation data for each binary parameter.

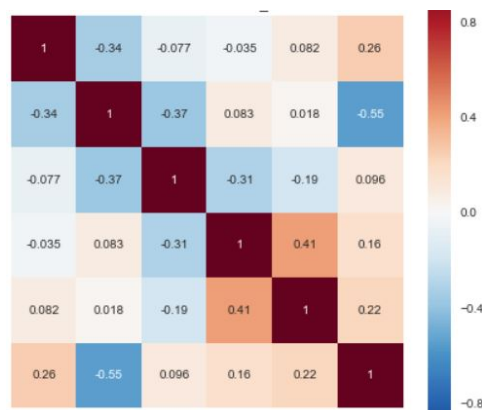
- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure



with X-axis as **min_sample_split**, 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](#)

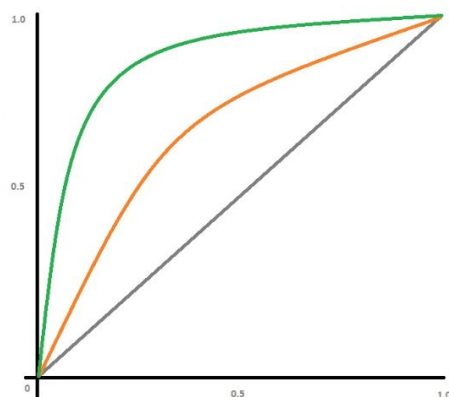
or

- 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 **min_sample_split**, 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](#).

7. Conclusion

- [You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this \[python table library\]\(#\) link](#)

Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on your train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](#).

Applying Decision Trees

[5.1] Applying Decision Trees on BOW, SET 1

In [293]:

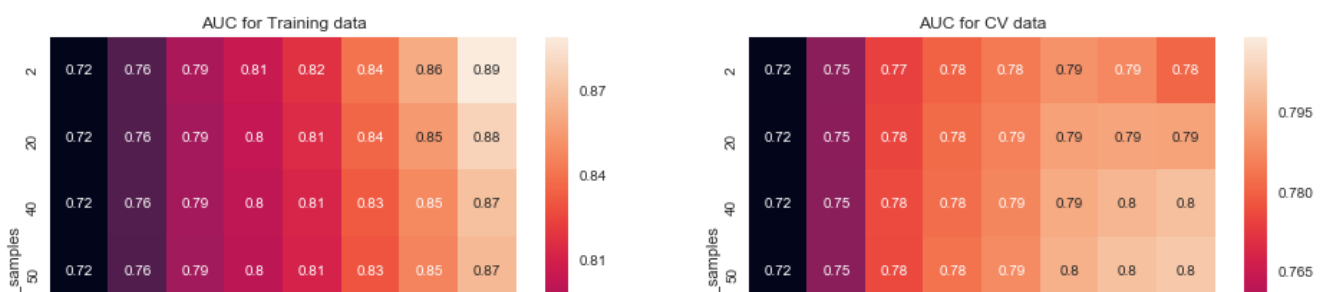
```
from sklearn.model_selection import GridSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import roc_auc_score

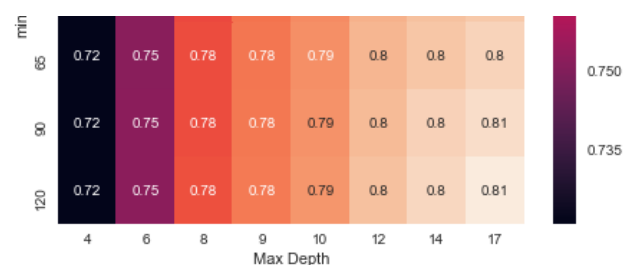
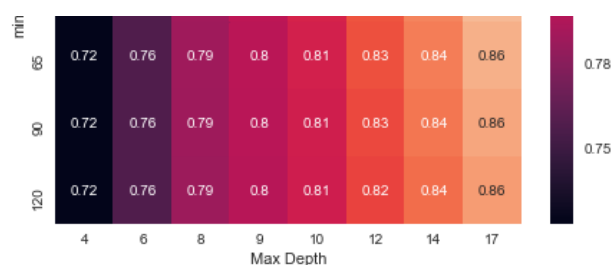
DT = DecisionTreeClassifier(class_weight='balanced')
parameters = {'max_depth':[4,6, 8, 9,10,12,14,17], 'min_samples_split':[2,20,40,50,65,90,120]}
clf = GridSearchCV(DT, parameters, cv=3, scoring='roc_auc', return_train_score=True)
clf.fit(X_train_bow, y_train)
print(clf.best_estimator_)
```

```
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=17,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=120,
                        min_weight_fraction_leaf=0.0, presort=False,
                        random_state=None, splitter='best')
```

In [294]:

```
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
min_samples_list = list(clf.cv_results_['param_min_samples_split'].data)
train_auc_score=clf.cv_results_['mean_train_score']
cv_auc_score=clf.cv_results_['mean_test_score']
train_data=pd.DataFrame(data={'min_samples_split':min_samples_list, 'Max Depth':max_depth_list,
                              'AUC':train_auc_score})
cv_data = pd.DataFrame(data={'Estimators':min_samples_list, 'Max Depth':max_depth_list,
                              'AUC':cv_auc_score})
sns.set_style("whitegrid")
plt.figure(figsize=(16,6))
plt.subplot(1,2,1)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
                              'AUC':train_auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')
plt.subplot(1,2,2)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
                              'AUC':cv_auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')
plt.show()
```





In [295]:

```
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
import numpy as np

x1 = min_samples_list
y1 = max_depth_list
z1 = train_Auc_score

x2 = min_samples_list
y2 = max_depth_list
z2 = cv_Auc_score

# https://plot.ly/python/3d-axes/
trace1 = go.Scatter3d(x=x1,y=y1,z=z1, name = 'train')
trace2 = go.Scatter3d(x=x2,y=y2,z=z2, name = 'Cross validation')
data = [trace1, trace2]

layout = go.Layout(scene = dict(
    xaxis = dict(title='min_samples_split'),
    yaxis = dict(title='max_depth'),
    zaxis = dict(title='AUC'),))

fig = go.Figure(data=data, layout=layout)
offline.iplot(fig, filename='3d-scatter-colorscale')
```

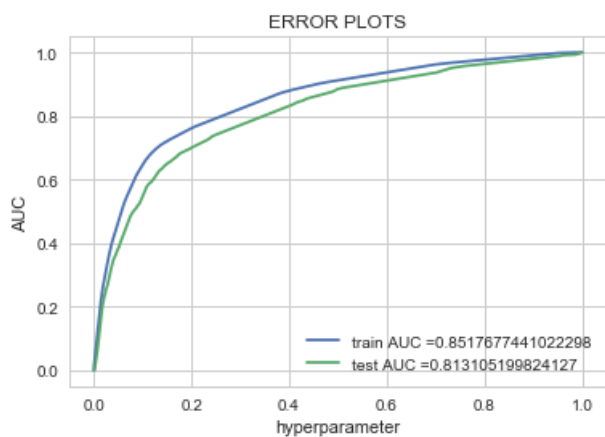
In [296]:

```
from sklearn.metrics import roc_curve, auc
```

```
DT = DecisionTreeClassifier(max_depth=17,min_samples_split=120,class_weight='balanced')
DT.fit(X_train_bow, y_train)

train_fpr, train_tpr, thresholds = roc_curve(y_train,DT.predict_proba(X_train_bow)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test,DT.predict_proba(X_test_bow)[: ,1])

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



In [297]:

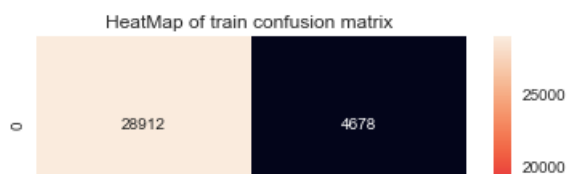
```
from sklearn.metrics import confusion_matrix
import seaborn as sn

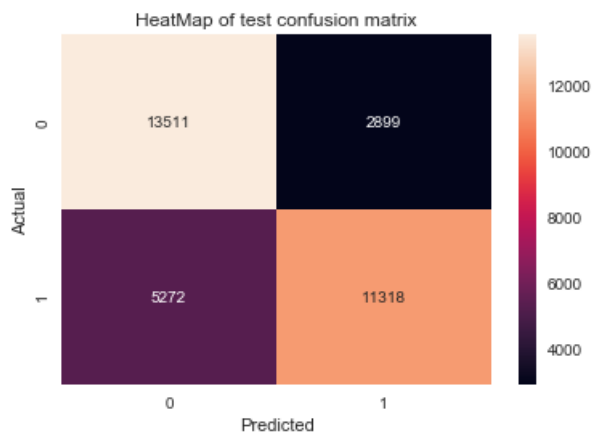
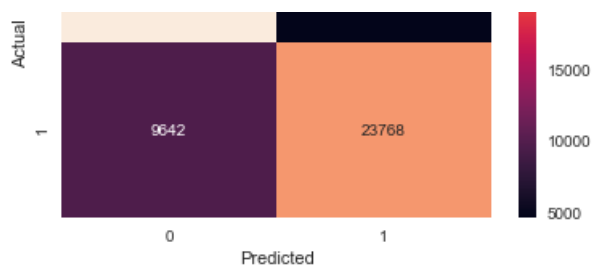
print("Train confusion matrix")
x=confusion_matrix(y_train, DT.predict(X_train_bow))
y=confusion_matrix(y_test, DT.predict(X_test_bow))
print(x)
print("Test confusion matrix")
print(y)

ax = plt.axes()
sns.heatmap(x, ax = ax,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
ax.set_title("HeatMap of train confusion matrix ")
plt.show()

bx = plt.axes()
sns.heatmap(y, ax = bx,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
bx.set_title("HeatMap of test confusion matrix")
plt.show()
```

```
Train confusion matrix
[[28912  4678]
 [ 9642 23768]]
Test confusion matrix
[[13511  2899]
 [ 5272 11318]]
```





[5.1.1] Top 20 important features from SET 1

In [298]:

```
FI=DT.feature_importances_

FeatInd=np.argsort(FI)
Fnames=vectorizer.get_feature_names()
topimp=FeatInd[-20:]
imp=[]
for i in topimp:
    imp.append(Fnames[i])

print("Top 20 important features")
print("=====")
print("")
print(imp)
```

Top 20 important features
=====

```
['however', 'horrible', 'money', 'nice', 'worst', 'easy', 'favorite', 'thought', 'bad',
'wonderful', 'excellent', 'loves', 'perfect', 'good', 'disappointed', 'love', 'delicious', 'best',
'great', 'not']
```

[5.1.2] Graphviz visualization of Decision Tree on BOW, SET 1

In [0]:

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

In [299]:

```
DT = DecisionTreeClassifier(max_depth=3,min_samples_split=2,class_weight='balanced')
DT.fit(X_train_bow, y_train)

dot_data = tree.export_graphviz(DT, out_file=None,feature_names=Fnames, filled=True, rounded=True,
special_characters=True)
graph = graphviz.Source(dot_data)
graph
```

Out[299]:

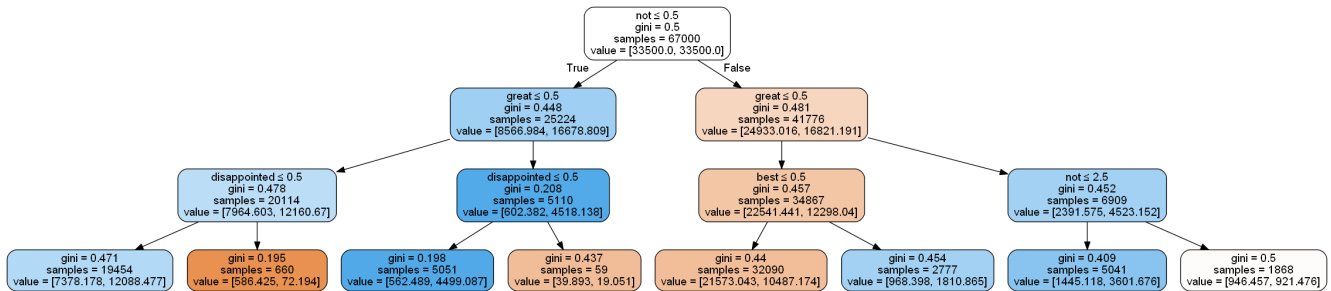
In [300]:

```
#Code refered from https://stackoverflow.com/questions/27817994/visualizing-decision-tree-in-scik
it-learn

png_bytes = graph.pipe(format='png')
with open('dtree_bow_graph.png', 'wb') as f:
    f.write(png_bytes)

from IPython.display import Image
Image(png_bytes)
```

Out[300]:



Review length feature added

In [315]:

```
from scipy import sparse
review_length_train=[]

for eachreview in X_train:
    x=len(eachreview)
    review_length_train.append(x)
```

In [316]:

```
review_length_test=[]

for eachreview in X_test:
    x=len(eachreview)
    review_length_test.append(x)

print(review_length_test[0])
```

104

In [317]:

```
print(type(X_train_bow))
print(X_train_bow.shape)
```

```
<class 'scipy.sparse.csr.csr_matrix'>
(67000, 50919)
```

In [318]:

```
X_train_bow_f1=sparse.hstack((X_train_bow,np.array(review_length_train)[: ,None]))
```

In [319]:

```
print(type(X_train_bow_f1))
print(X_train_bow_f1.shape)
```



```
<class 'scipy.sparse.coo.coo_matrix'>
(67000, 50920)
```

In [320]:

```
print(type(X_test_bow))
print(X_test_bow.shape)
```

```
<class 'scipy.sparse.csr.csr_matrix'>
(33000, 50919)
```

In [321]:

```
X_test_bow_f1=sparse.hstack((X_test_bow,np.array(review_length_test)[: ,None]))
print(type(X_test_bow_f1))
print(X_test_bow_f1.shape)
```

```
<class 'scipy.sparse.coo.coo_matrix'>
(33000, 50920)
```

In []:

```
X_train_bow_f1
X_test_bow_f1
```

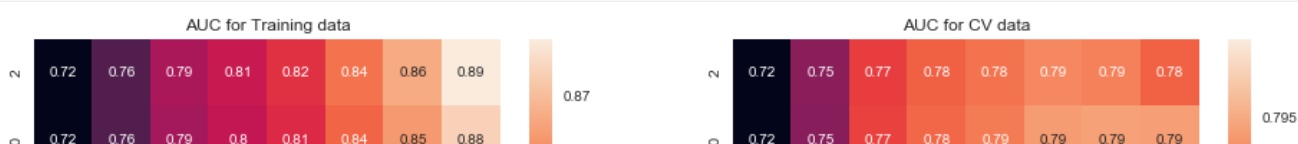
In [322]:

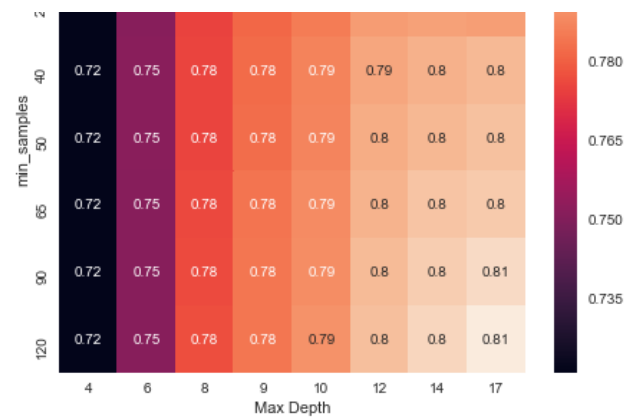
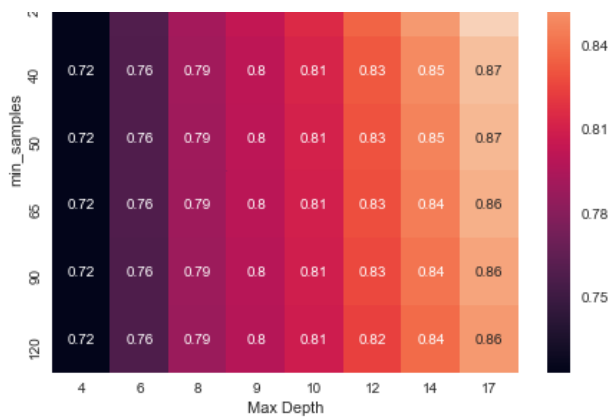
```
DT = DecisionTreeClassifier(class_weight='balanced')
parameters ={'max_depth':[4,6, 8, 9,10,12,14,17], 'min_samples_split':[2,20,40,50,65,90,120]}
clf = GridSearchCV(DT, parameters, cv=3, scoring='roc_auc',return_train_score=True)
clf.fit(X_train_bow_f1, y_train)
print(clf.best_estimator_)
```

```
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=17,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=120,
                        min_weight_fraction_leaf=0.0, presort=False,
                        random_state=None, splitter='best')
```

In [323]:

```
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
min_samples_list = list(clf.cv_results_['param_min_samples_split'].data)
train_Auc_score=clf.cv_results_['mean_train_score']
cv_Auc_score=clf.cv_results_['mean_test_score']
train_data=pd.DataFrame(data={'min_samples_split':min_samples_list, 'Max Depth':max_depth_list,
'AUC':train_Auc_score})
cv_data = pd.DataFrame(data={'Estimators':min_samples_list, 'Max Depth':max_depth_list,
'AUC':cv_Auc_score})
sns.set_style("whitegrid")
plt.figure(figsize=(16,6))
plt.subplot(1,2,1)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
'AUC':train_Auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')
plt.subplot(1,2,2)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
'AUC':cv_Auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')
plt.show()
```





In [324]:

```
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
import numpy as np

x1 = min_samples_list
y1 = max_depth_list
z1 = train_Auc_score

x2 = min_samples_list
y2 = max_depth_list
z2 = cv_Auc_score

# https://plot.ly/python/3d-axes/
trace1 = go.Scatter3d(x=x1,y=y1,z=z1, name = 'train')
trace2 = go.Scatter3d(x=x2,y=y2,z=z2, name = 'Cross validation')
data = [trace1, trace2]

layout = go.Layout(scene = dict(
    xaxis = dict(title='min_samples_split'),
    yaxis = dict(title='max_depth'),
    zaxis = dict(title='AUC'),))

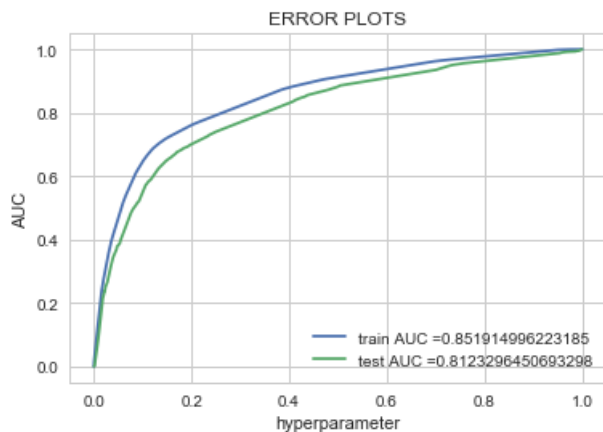
fig = go.Figure(data=data, layout=layout)
offline.iplot(fig, filename='3d-scatter-colorscale')
```

In [325]:

```
DT = DecisionTreeClassifier(max_depth=17,min_samples_split=120,class_weight='balanced')
DT.fit(X_train_bow_f1, y_train)

train_fpr, train_tpr, thresholds = roc_curve(y_train,DT.predict_proba(X_train_bow_f1)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test,DT.predict_proba(X_test_bow_f1)[: ,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("hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [326]:

```
from sklearn.metrics import confusion_matrix
import seaborn as sn

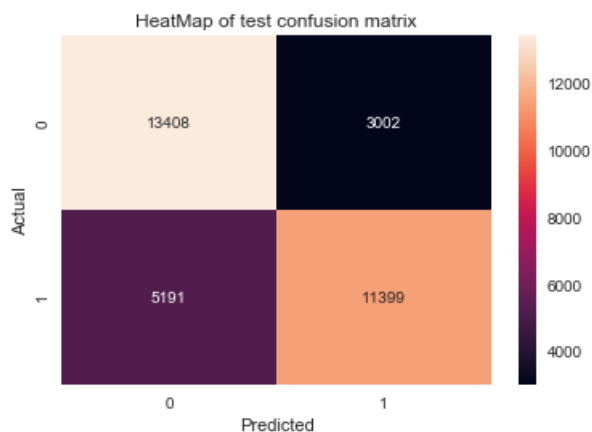
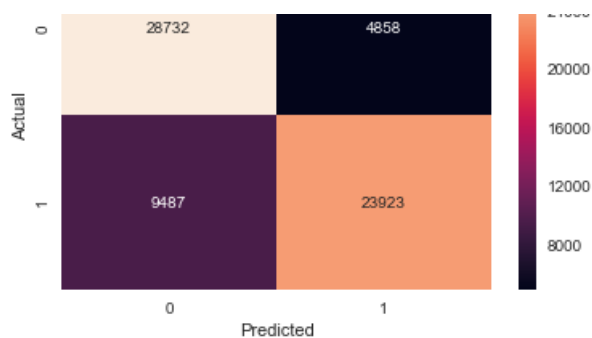
print("Train confusion matrix")
x=confusion_matrix(y_train, DT.predict(X_train_bow_f1))
y=confusion_matrix(y_test, DT.predict(X_test_bow_f1))
print(x)
print("Test confusion matrix")
print(y)

ax = plt.axes()
sns.heatmap(x, ax = ax,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
ax.set_title("HeatMap of train confusion matrix ")
plt.show()

bx = plt.axes()
sns.heatmap(y, ax = bx,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
bx.set_title("HeatMap of test confusion matrix")
plt.show()
```

```
Train confusion matrix
[[28732  4858]
 [ 9487 23923]]
Test confusion matrix
[[13408  3002]
 [ 5191 11399]]
```





[5.2] Applying Decision Trees on TFIDF, SET 2

In [305]:

```
DT = DecisionTreeClassifier(class_weight='balanced')
parameters = {'max_depth':[4,6,9,12,14,17], 'min_samples_split':[2,20,40,100,150]}
clf = GridSearchCV(DT, parameters, cv=3, scoring='roc_auc', return_train_score=True)
clf.fit(X_train_tfidf, y_train)
print(clf.best_estimator_)
```

```
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=17,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=150,
                        min_weight_fraction_leaf=0.0, presort=False,
                        random_state=None, splitter='best')
```

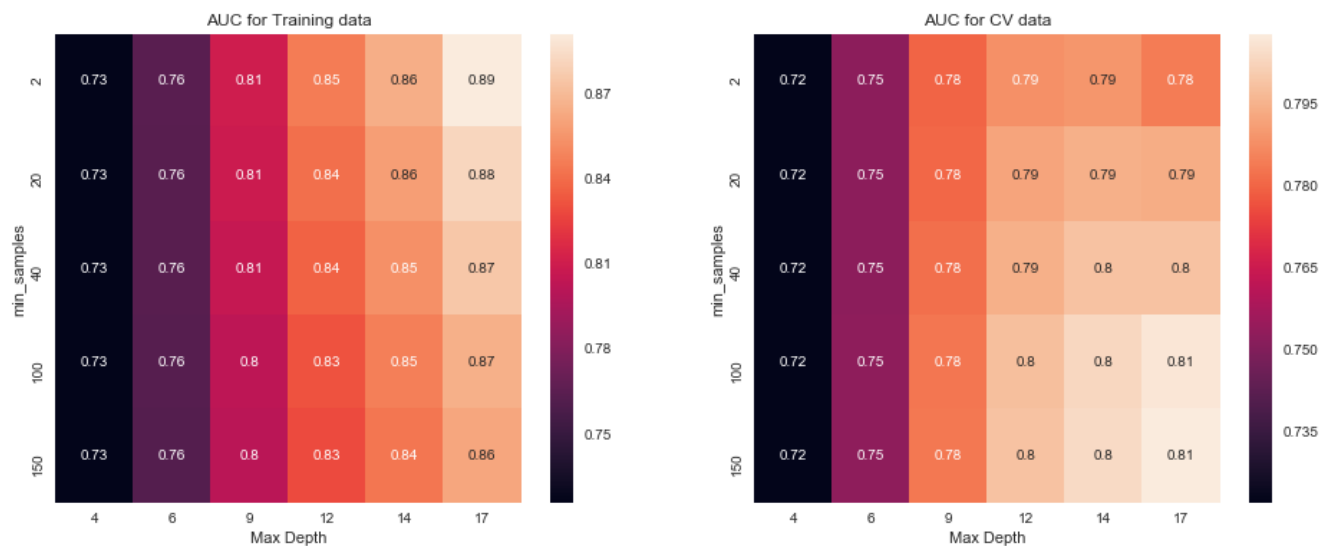
In [306]:

```
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
min_samples_list = list(clf.cv_results_['param_min_samples_split'].data)
train_auc_score = clf.cv_results_['mean_train_score']
cv_auc_score = clf.cv_results_['mean_test_score']
train_data = pd.DataFrame(data={'min_samples_split':min_samples_list, 'Max Depth':max_depth_list,
                              'AUC':train_auc_score})
cv_data = pd.DataFrame(data={'Estimators':min_samples_list, 'Max Depth':max_depth_list,
                              'AUC':cv_auc_score})
```

In [307]:

```
sns.set_style("whitegrid")
plt.figure(figsize=(16,6))
plt.subplot(1,2,1)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
                          'AUC':train_auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')
plt.subplot(1,2,2)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
                          'AUC':cv_auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
```

```
sns.heatmap(data, annot=True).set_title('AUC for CV data')
plt.show()
```



In [308]:

```
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
import numpy as np

x1 = min_samples_list
y1 = max_depth_list
z1 = train_Auc_score

x2 = min_samples_list
y2 = max_depth_list
z2 = cv_Auc_score

# https://plot.ly/python/3d-axes/
trace1 = go.Scatter3d(x=x1,y=y1,z=z1, name = 'train')
trace2 = go.Scatter3d(x=x2,y=y2,z=z2, name = 'Cross validation')
data = [trace1, trace2]

layout = go.Layout(scene = dict(
    xaxis = dict(title='min_samples_split'),
    yaxis = dict(title='max_depth'),
    zaxis = dict(title='AUC'),))

fig = go.Figure(data=data, layout=layout)
offline.iplot(fig, filename='3d-scatter-colorscale')
```

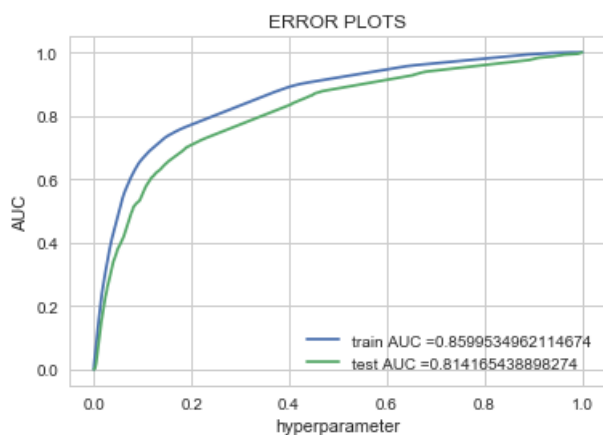
In [309]:

```
from sklearn.metrics import roc_curve, auc

DT = DecisionTreeClassifier(max_depth=17,min_samples_split=150,class_weight='balanced')
DT.fit(X_train_tfidf, y_train)

train_fpr, train_tpr, thresholds = roc_curve(y_train,DT.predict_proba(X_train_tfidf)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test,DT.predict_proba(X_test_tfidf)[:,1])

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



In [310]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
x=confusion_matrix(y_train, DT.predict(X_train_tfidf))
y=confusion_matrix(y_test, DT.predict(X_test_tfidf))
print(x)
print("Test confusion matrix")
print(y)
```

```
Train confusion matrix
[[28573  5017]
 [ 8852 24558]]
Test confusion matrix
[[13328  3082]
 [ 5017 11573]]
```

In [311]:

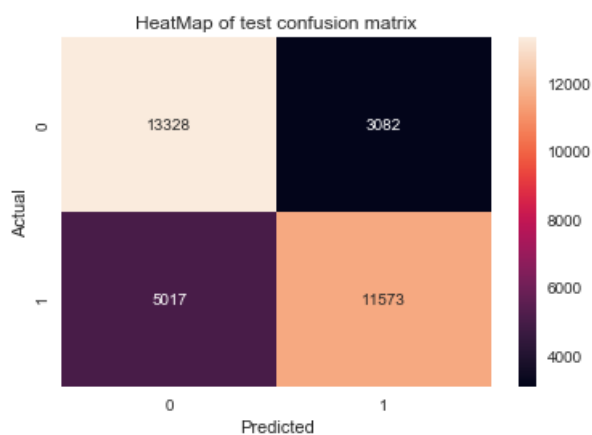
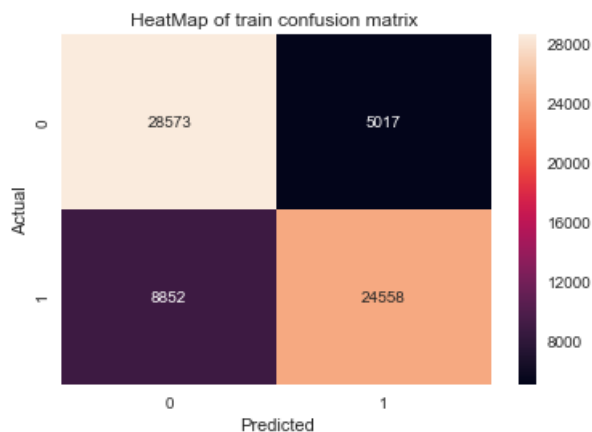
```
import seaborn as sn

ax = plt.axes()
sns.heatmap(x, ax = ax,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
ax.set_title("HeatMap of train confusion matrix ")
plt.show()
```

```

bx = plt.axes()
sns.heatmap(y, ax = bx,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
bx.set_title("HeatMap of test confusion matrix")
plt.show()

```



[5.2.1] Top 20 important features from SET 2

In [312]:

```
FI=DT.feature_importances_
```

```
FeatInd=np.argsort(FI)
```

```
Fnames=tf_idf_vect.get_feature_names()
```

```
topimp=FeatInd[-20:]
```

```
imp=[]
```

```
for i in topimp:
    imp.append(Fnames[i])
```

```
print("Top 20 important features")
```

```
print("=====")
```

```
print("")
```

```
print(imp)
```

Top 20 important features

=====

```
['horrible', 'money', 'easy', 'worst', 'bad', 'thought', 'nice', 'not good', 'favorite',
'wonderful', 'excellent', 'loves', 'perfect', 'disappointed', 'good', 'love', 'delicious', 'best',
'great', 'not']
```

[5.2.2] Graphviz visualization of Decision Tree on TFIDF, SET 2

In [01]:

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

In [313]:

```
DT = DecisionTreeClassifier(max_depth=3,min_samples_split=2,class_weight='balanced')
DT.fit(X_train_tfidf, y_train)

dot_data = tree.export_graphviz(DT, out_file=None,feature_names=tf_idf_vect.get_feature_names(),
filled=True, rounded=True, special_characters=True)
graph = graphviz.Source(dot_data)
graph
```

Out [313]:

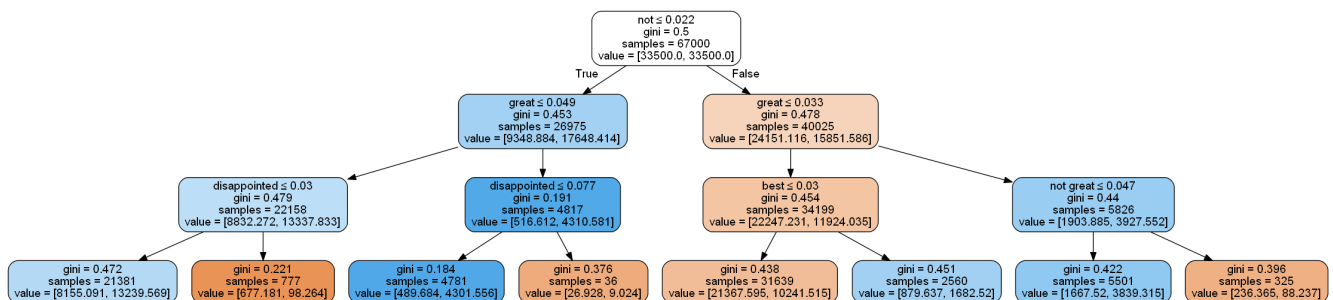
In [314]:

```
#Code refered from https://stackoverflow.com/questions/27817994/visualizing-decision-tree-in-scik
it-learn

png_bytes = graph.pipe(format='png')
with open('dtree_tfidf_graph.png','wb') as f:
    f.write(png_bytes)

from IPython.display import Image
Image(png_bytes)
```

Out [314]:



[5.3] Applying Decision Trees on AVG W2V, SET 3

In [270]:

```
list_of_sentence_train=[]
for sentence in X_train:
    list_of_sentence_train.append(sentence.split())

w2v_model=Word2Vec(list_of_sentence_train,min_count=5,size=50, workers=4)
```

In [271]:

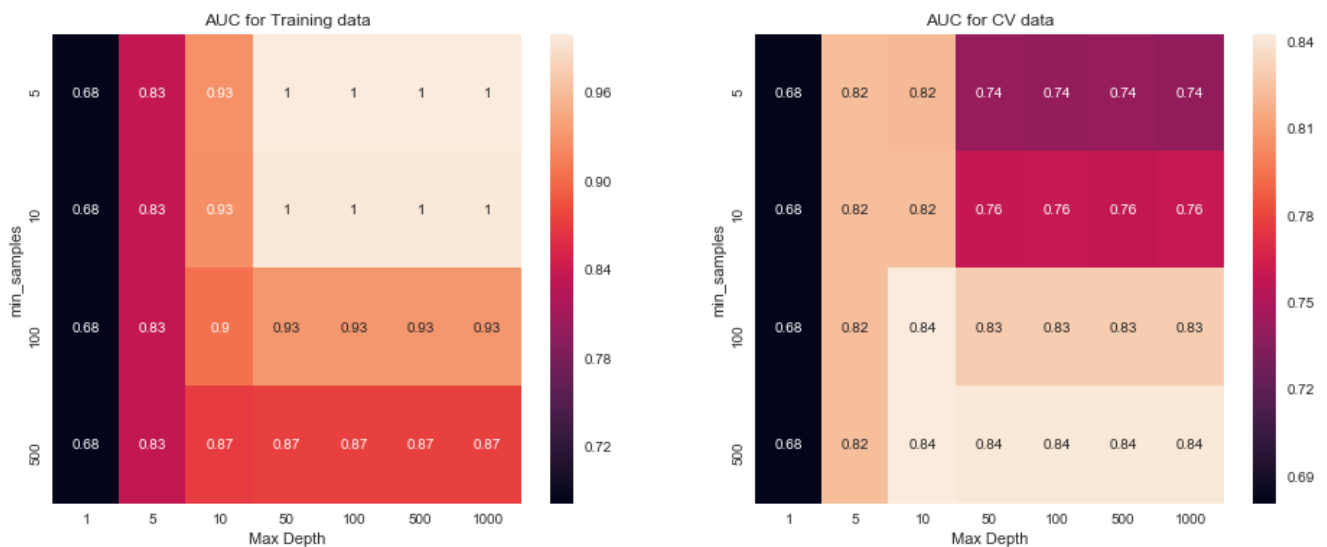
```
w2v_words = list(w2v_model.wv.vocab)
```

In [272]:

```
sent_vectors_train = []
for sent in tqdm(list_of_sentence_train):
    sent_vec = np.zeros(50)
    cnt_words = 0;
    for word in sent:
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
```


In [276]:

```
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
min_samples_list = list(clf.cv_results_['param_min_samples_split'].data)
train_Auc_score=clf.cv_results_['mean_train_score']
cv_Auc_score=clf.cv_results_['mean_test_score']
train_data=pd.DataFrame(data={'min_samples_split':min_samples_list, 'Max Depth':max_depth_list,
'AUC':train_Auc_score})
cv_data = pd.DataFrame(data={'Estimators':min_samples_list, 'Max Depth':max_depth_list,
'AUC':cv_Auc_score})
sns.set_style("whitegrid")
plt.figure(figsize=(16,6))
plt.subplot(1,2,1)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
'AUC':train_Auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')
plt.subplot(1,2,2)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
'AUC':cv_Auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')
plt.show()
```



In [277]:

```
x1 = min_samples_list
y1 = max_depth_list
z1 = train_Auc_score

x2 = min_samples_list
y2 = max_depth_list
z2 = cv_Auc_score

# https://plot.ly/python/3d-axes/
trace1 = go.Scatter3d(x=x1,y=y1,z=z1, name = 'train')
trace2 = go.Scatter3d(x=x2,y=y2,z=z2, name = 'Cross validation')
data = [trace1, trace2]

layout = go.Layout(scene = dict(
    xaxis = dict(title='min_samples_split'),
    yaxis = dict(title='max_depth'),
    zaxis = dict(title='AUC'),))

fig = go.Figure(data=data, layout=layout)
offline.iplot(fig, filename='3d-scatter-colorscale')
```

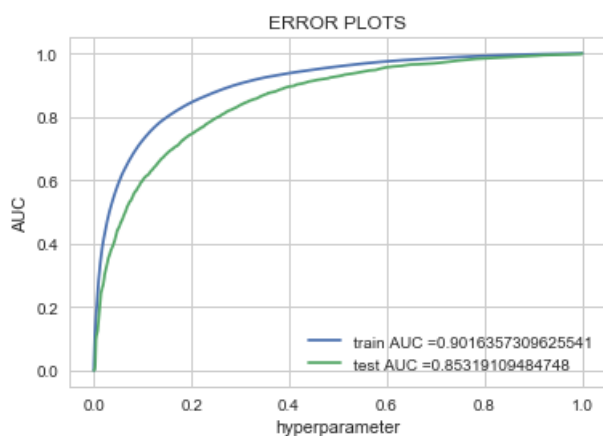
In [278]:

```
from sklearn.metrics import roc_curve, auc

DT = DecisionTreeClassifier(max_depth=10,min_samples_split=100,class_weight='balanced')
DT.fit(sent_vectors_train, y_train)

train_fpr, train_tpr, thresholds = roc_curve(y_train,DT.predict_proba(sent_vectors_train)[:,:1])
test_fpr, test_tpr, thresholds = roc_curve(y_test,DT.predict_proba(sent_vectors_test)[:,: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("hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [279]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
x=confusion_matrix(y_train, DT.predict(sent_vectors_train))
y=confusion_matrix(y_test, DT.predict(sent_vectors_test))
print(x)
print("Test confusion matrix")
print(y)

import seaborn as sn

ax = plt.axes()
```

```
sns.heatmap(x, ax = ax,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
ax.set_title("HeatMap of train confusion matrix ")
plt.show()

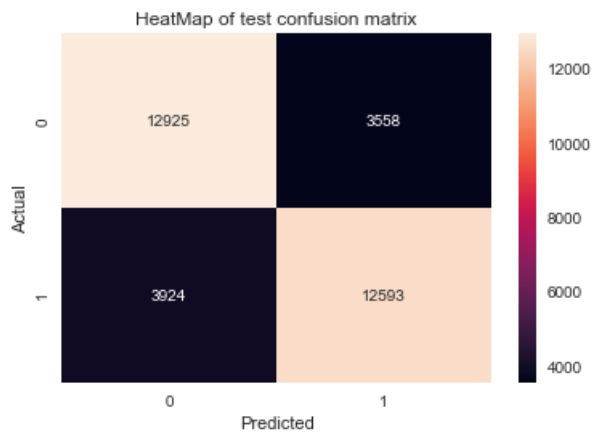
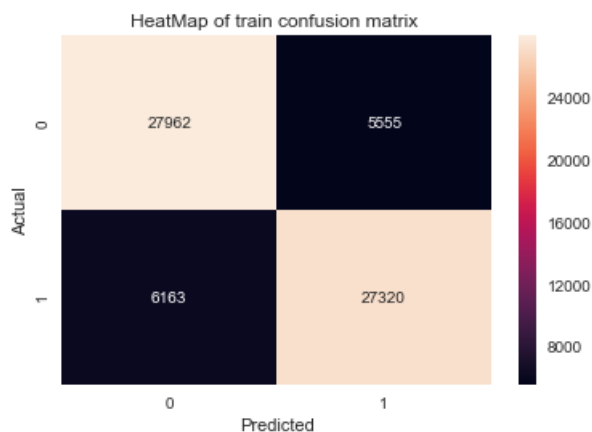
bx = plt.axes()
sns.heatmap(y, ax = bx,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
bx.set_title("HeatMap of test confusion matrix")
plt.show()
```

Train confusion matrix

```
[[27962  5555]
 [ 6163 27320]]
```

Test confusion matrix

```
[[12925  3558]
 [ 3924 12593]]
```



[5.4] Applying Decision Trees on TFIDF W2V, SET 4

In [328]:

```
from gensim.models import Word2Vec
from gensim.models import KeyedVectors

list_of_sentence_train=[]
for sentence in X_train:
    list_of_sentence_train.append(sentence.split())

w2v_model=Word2Vec(list_of_sentence_train,min_count=5,size=50, workers=4)
```

In [329]:

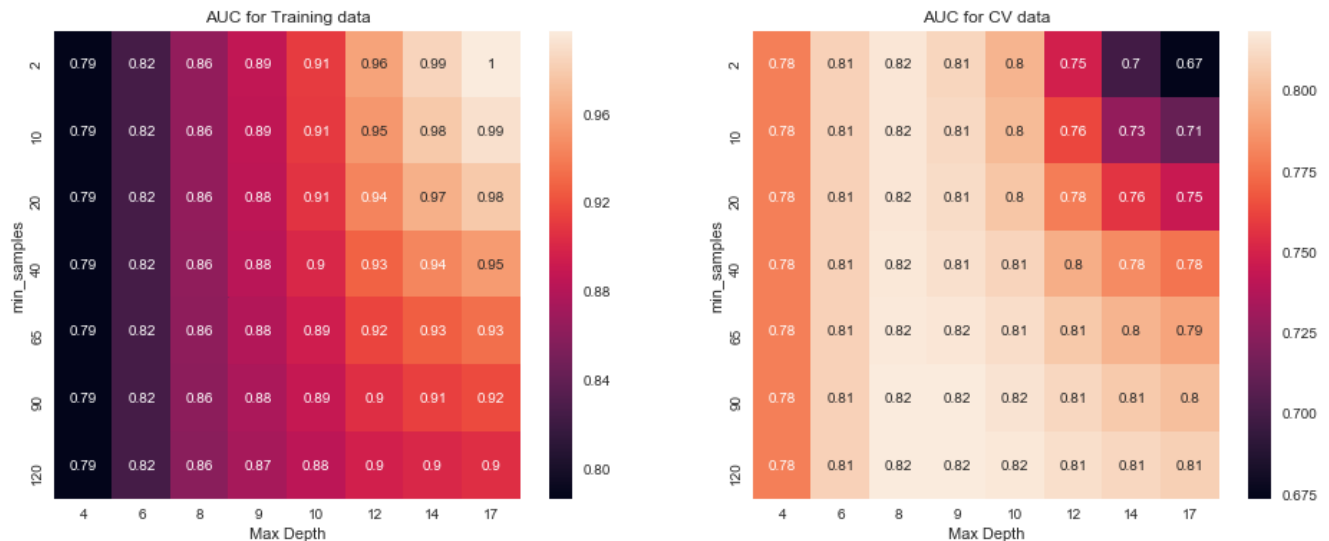
```
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(X_train)
```


In [284]:

```
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
min_samples_list = list(clf.cv_results_['param_min_samples_split'].data)
train_Auc_score=clf.cv_results_['mean_train_score']
cv_Auc_score=clf.cv_results_['mean_test_score']
train_data=pd.DataFrame(data={'min_samples_split':min_samples_list, 'Max Depth':max_depth_list,
'AUC':train_Auc_score})
cv_data = pd.DataFrame(data={'Estimators':min_samples_list, 'Max Depth':max_depth_list,
'AUC':cv_Auc_score})
```

In [285]:

```
sns.set_style("whitegrid")
plt.figure(figsize=(16,6))
plt.subplot(1,2,1)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
'AUC':train_Auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')
plt.subplot(1,2,2)
data = pd.DataFrame(data={'min_samples':min_samples_list, 'Max Depth':max_depth_list,
'AUC':cv_Auc_score})
data = data.pivot(index='min_samples', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')
plt.show()
```



In [286]:

```
x1 = min_samples_list
y1 = max_depth_list
z1 = train_Auc_score

x2 = min_samples_list
y2 = max_depth_list
z2 = cv_Auc_score

# https://plot.ly/python/3d-axes/
trace1 = go.Scatter3d(x=x1,y=y1,z=z1, name = 'train')
trace2 = go.Scatter3d(x=x2,y=y2,z=z2, name = 'Cross validation')
data = [trace1, trace2]

layout = go.Layout(scene = dict(
    xaxis = dict(title='min_samples_split'),
    yaxis = dict(title='max_depth'),
    zaxis = dict(title='AUC'),))

fig = go.Figure(data=data, layout=layout)
offline.iplot(fig, filename='3d-scatter-colorscale')
```

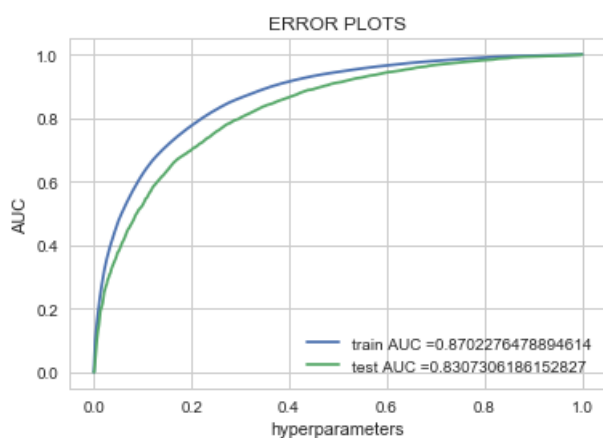
In [334]:

```
from sklearn.metrics import roc_curve, auc

DT = DecisionTreeClassifier(max_depth=9,min_samples_split=120,class_weight='balanced')
DT.fit(tfidf_sent_vectors, y_train)

train_fpr, train_tpr, thresholds = roc_curve(y_train,DT.predict_proba(tfidf_sent_vectors)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test,DT.predict_proba(tfidf_sent_vectors_test)[: ,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("hyperparameters")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [335]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
x=confusion_matrix(y_train, DT.predict(tfidf_sent_vectors))
y=confusion_matrix(y_test, DT.predict(tfidf_sent_vectors_test))
print(x)
print("Test confusion matrix")
```

```

print(y)

import seaborn as sn

ax = plt.axes()
sns.heatmap(x, ax = ax,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
ax.set_title("HeatMap of train confusion matrix ")
plt.show()

bx = plt.axes()
sns.heatmap(y, ax = bx,annot=True, fmt="d")
plt.xlabel("Predicted")
plt.ylabel("Actual")
bx.set_title("HeatMap of test confusion matrix")
plt.show()

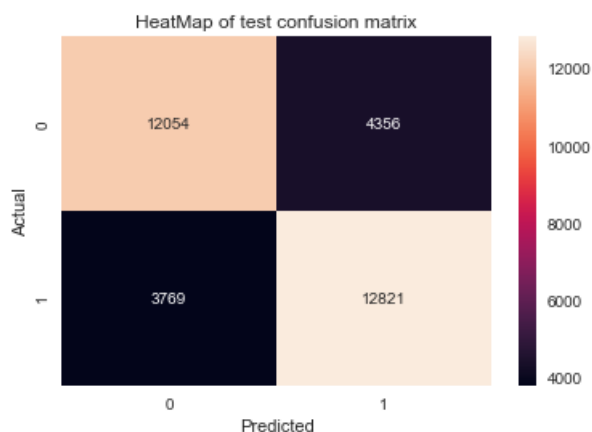
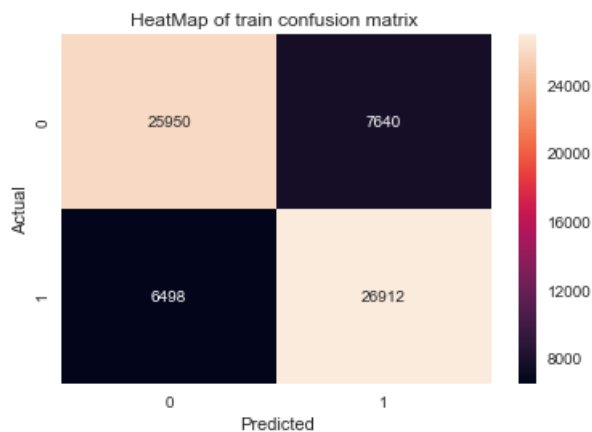
```

Train confusion matrix

```
[[25950  7640]
 [ 6498 26912]]
```

Test confusion matrix

```
[[12054  4356]
 [ 3769 12821]]
```



[6] Conclusions

In [337]:

```

from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Model","Vectorizer","Hyper Parameter(max_depth)","Hyper
Parameter(min_samples_split)","AUC"]
x.add_row(["DecisionTreeClassifier","BOW","17","120","0.8131"])
x.add_row(["DecisionTreeClassifier","TF-IDF","17","150","0.8141"])

```



```
x.add_row(["DecisionTreeClassifier", "AVG W2V", "10", "100", "0.8531"])
x.add_row(["DecisionTreeClassifier", "TFIDF W2V", "9", "120", "0.8307"])
```

```
print(x)
print("\n")
print("\n")
print("Feature engineered output after adding review length to BOW vectorized data:=")
print("\n")
y = PrettyTable()
y.field_names = ["Model", "Vectorizer", "Hyper Parameter(max_depth)", "Hyper
Parameter(min_samples_split)", "AUC"]
y.add_row(["DecisionTreeClassifier", "BOW", "17", "120", "0.8123"])
print(y)
```

```
+-----+-----+-----+-----+
+-----+
|           Model           | Vectorizer | Hyper Parameter(max_depth) | Hyper
Parameter(min_samples_split) |   AUC   |
+-----+-----+-----+-----+
+-----+
| DecisionTreeClassifier |    BOW    |           17              |           120
| 0.8131 |
| DecisionTreeClassifier |   TF-IDF   |           17              |           150
| 0.8141 |
| DecisionTreeClassifier |   AVG W2V   |           10              |           100
| 0.8531 |
| DecisionTreeClassifier | TFIDF W2V   |            9              |           120
| 0.8307 |
+-----+-----+-----+-----+
+-----+
```

Feature engineered output after adding review length to BOW vectorized data:=

```
+-----+-----+-----+-----+
+-----+
|           Model           | Vectorizer | Hyper Parameter(max_depth) | Hyper
Parameter(min_samples_split) |   AUC   |
+-----+-----+-----+-----+
+-----+
| DecisionTreeClassifier |    BOW    |           17              |           120
| 0.8123 |
+-----+-----+-----+-----+
+-----+
```

Observations:-

Adding additional feature review length does not seem to affect the classifier much the AUC score of the model was observed to remain the same upon adding the feature

Feature importance was observed using the DecisionTreeClassifier Attribute

The behaviour of the models corresponding to all the vectorizers was similar best performance was observed with AVG W2V vectorizers