

Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

Objective:

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

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

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

[1]. Reading Data

[1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

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

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

```

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

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

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

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

from tqdm import tqdm
import os

```

```

In [2]: # using SQLite Table to read data.
con = sqlite3.connect('database.sqlite')

# filtering only positive and negative reviews i.e.
# not taking into consideration those reviews with Score=3
# SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 50
0000 data points
# you can change the number to any other number based on your computing
power

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

```

```

filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score
!= 3 """, 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 (525814, 10)

Out[2]:

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

```
In [3]: display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

```
In [4]: print(display.shape)
display.head()
```

```
(80668, 7)
```

Out[4]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
0	#oc-R115TNMSPFT9I7	B005ZBZLT4	Breyton	1331510400	2	Overall its just OK when considering the price...	2
1	#oc-R11D9D7SHXIJB9	B005HG9ESG	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms, u...	3
2	#oc-R11DNU2NBKQ23Z	B005ZBZLT4	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately not ...	2
3	#oc-R11O5J5ZVQE25C	B005HG9ESG	Penguin Chick	1346889600	5	This will be the bottle that you grab from the...	3
4	#oc-R12KPBODL2B5ZD	B007OSBEV0	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2

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

Out[5]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
80638	AZY10LLTJ71NX	B001ATMQK2	undertheshrine "undertheshrine"	1296691200	5	I bought this 6 pack because for the price tha...	5

In [6]: `display['COUNT(*)'].sum()`

Out[6]: 393063

[2] Exploratory Data Analysis

[2.1] Data Cleaning: Deduplication

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

In [7]: `display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()`

Out[7]:

Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenon
----	-----------	--------	-------------	----------------------	------------------

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenon
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	

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

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

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

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

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

```
In [8]: #Sorting data according to ProductId in ascending order
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True,
inplace=False, kind='quicksort', na_position='last')
```

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

```
Out[9]: (364173, 10)
```

```
In [10]: #Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[10]: 69.25890143662969
```

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

```
In [11]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)
```



```
display.head()
```

Out[11]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenom
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	

```
In [12]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [13]: #Before starting the next phase of preprocessing lets see the number of  
entries left  
print(final.shape)  
  
#How many positive and negative reviews are present in our dataset?  
final['Score'].value_counts()  
  
(364171, 10)
```

Out[13]:

1	307061
0	57110

Name: Score, dtype: int64

[3] Preprocessing

[3.1]. Preprocessing Review Text

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

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

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

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

```
In [14]: # printing some random reviews
sent_0 = final['Text'].values[0]
print(sent_0)
print("="*50)

sent_1000 = final['Text'].values[1000]
print(sent_1000)
print("="*50)

sent_1500 = final['Text'].values[1500]
print(sent_1500)
print("="*50)

sent_4900 = final['Text'].values[4900]
print(sent_4900)
print("="*50)
```

```
this witty little book makes my son laugh at loud. i recite it in the c
ar as we're driving along and he always can sing the refrain. he's lear
ned about whales, India, drooping roses: i love all the new words this
```

book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

=====

I was really looking forward to these pods based on the reviews. Starbucks is good, but I prefer bolder taste.... imagine my surprise when I ordered 2 boxes - both were expired! One expired back in 2005 for gosh sakes. I admit that Amazon agreed to credit me for cost plus part of shipping, but geez, 2 years expired!!! I'm hoping to find local San Diego area shoppe that carries pods so that I can try something different than starbucks.

=====

Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70's it was poisonous until they figured out a way to fix that. I still like it but it could be better.

=====

Can't do sugar. Have tried scores of SF Syrups. NONE of them can touch the excellence of this product.

Thick, delicious. Perfect. 3 ingredients: Water, Maltitol, Natural Maple Flavor. PERIOD. No chemicals. No garbage.

Have numerous friends & family members hooked on this stuff. My husband & son, who do NOT like "sugar free" prefer this over major label regular syrup.

I use this as my SWEETENER in baking: cheesecakes, white brownies, muffins, pumpkin pies, etc... Unbelievably delicious...

Can you tell I like it? :)

=====

```
In [15]: # remove urls from text python: https://stackoverflow.com/a/40823105/4084039
sent_0 = re.sub(r"http\S+", "", sent_0)
sent_1000 = re.sub(r"http\S+", "", sent_1000)
sent_150 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)
```

```
print(sent_0)
```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

```
In [16]: # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all-tags-from-an-element  
from bs4 import BeautifulSoup
```

```
soup = BeautifulSoup(sent_0, 'lxml')  
text = soup.get_text()  
print(text)  
print("="*50)
```

```
soup = BeautifulSoup(sent_1000, 'lxml')  
text = soup.get_text()  
print(text)  
print("="*50)
```

```
soup = BeautifulSoup(sent_1500, 'lxml')  
text = soup.get_text()  
print(text)  
print("="*50)
```

```
soup = BeautifulSoup(sent_4900, 'lxml')  
text = soup.get_text()  
print(text)
```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

=====

I was really looking forward to these pods based on the reviews. Starbucks is good, but I prefer bolder taste.... imagine my surprise when I ordered 2 boxes - both were expired! One expired back in 2005 for gosh sakes. I admit that Amazon agreed to credit me for cost plus part of shipping, but geez, 2 years expired!!! I'm hoping to find local San Diego area shoppe that carries pods so that I can try something different than Starbucks.

=====

Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70's it was poisonous until they figured out a way to fix that. I still like it but it could be better.

=====

Can't do sugar. Have tried scores of SF Syrups. NONE of them can touch the excellence of this product. Thick, delicious. Perfect. 3 ingredients: Water, Maltitol, Natural Maple Flavor. PERIOD. No chemicals. No garbage. Have numerous friends & family members hooked on this stuff. My husband & son, who do NOT like "sugar free" prefer this over major label regular syrup. I use this as my SWEETENER in baking: cheesecakes, white brownies, muffins, pumpkin pies, etc... Unbelievably delicious... Can you tell I like it? :)

In [17]: `# https://stackoverflow.com/a/47091490/4084039`

```
import re

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

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"\ 're", " are", phrase)
    phrase = re.sub(r"\ 's", " is", phrase)
```

```

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

```

```

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

```

Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today is Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70s it was poisonous until they figured out a way to fix that. I still like it but it could be better.

=====

```

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

```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

```

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

```

Great ingredients although chicken should have been 1st rather than chicken broth the only thing I do not think belongs in it is Canola oil Ca

nola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it it would poison them Today is Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut facts though say otherwise Until the late 70s it was poisonous until they figured out a way to fix that I still like it but it could be better

```
In [21]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have been 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 [22]: # Combining all the above students
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentence in tqdm(final['Text'].values):
    sentence = re.sub(r"http\S+", "", sentence)
    sentence = BeautifulSoup(sentence, 'lxml').get_text()
    sentence = decontracted(sentence)
    sentence = re.sub("\S*\d\S*", "", sentence).strip()
    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower()
    not in stopwords)
    preprocessed_reviews.append(sentence.strip())
```

```
100%|██████████| 364171/364171 [01:43<00:00, 3530.30it/s]
```

```
In [23]: preprocessed_reviews[1500]
```

```
Out[23]: 'great ingredients although chicken rather chicken broth thing not thin
k belongs canola oil canola rapeseed not someting dog would ever find n
ature find rapeseed nature eat would poison today food industries convi
nced masses canola oil safe even better oil olive virgin coconut facts
though say otherwise late poisonous figured way fix still like could be
tter'
```

[3.2] Preprocessing Review Summary

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


[4] Featurization

```
In [40]: # split data points to avoid data leakage
from sklearn.model_selection import train_test_split

n_samples= 100000
x_sampled_data= preprocessed_reviews[:n_samples]
y_sampled_data= final.Score[:n_samples]

X_train, x_test, Y_train, y_test= train_test_split(x_sampled_data, y_sampled_data, test_size= .33)
```

[4.1] BAG OF WORDS

```
In [36]: #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("the type of count vectorizer ",type(X_train_bow))
print("the shape of out text BOW vectorizer ",X_train_bow.get_shape())
print('='*50)

X_test_bow = count_vect.transform(x_test)
print("the type of count vectorizer ",type(X_test_bow))
print("the shape of out text BOW vectorizer ",X_test_bow.get_shape())

some feature names  ['aa', 'aaa', 'aaaa', 'aaaaa', 'aaaaaa', 'aaaaaaaa', 'aaaaaaaaa', 'aaaaaaa', 'aaaaaa', 'aaaaa']
gghh', 'aaaaawesome', 'aaaah', 'aaaahhhhhhhhhhh', 'aaaallll']
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (67000, 49319)
=====
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (33000, 49319)
```

[4.2] Bi-Grams and n-Grams.

```
In [ ]: #bi-gram, tri-gram and n-gram

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

# you can choose these numebrs min_df=10, max_features=5000, of your choice
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
final_bigram_counts = count_vect.fit_transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_bigram_counts))
print("the shape of out text BOW vectorizer ",final_bigram_counts.get_shape())
print("the number of unique words including both unigrams and bigrams ", final_bigram_counts.get_shape()[1])
```

[4.3] TF-IDF

```
In [27]: tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
X_train_tfidf= tf_idf_vect.fit_transform(X_train)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names()[0:10])
print("the type of count vectorizer ",type(X_train_tfidf))
print("the shape of out text TFIDF vectorizer ",X_train_tfidf.get_shape())
print('='*50)

X_test_tfidf= tf_idf_vect.transform(x_test)
print("the type of count vectorizer ",type(X_test_tfidf))
print("the shape of out text TFIDF vectorizer ",X_test_tfidf.get_shape())
```

some sample features(unique words in the corpus) ['ab', 'abandon', 'ab

```

c', 'abdominal', 'ability', 'able', 'able break', 'able buy', 'able che
w', 'able drink']
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (67000, 37841)
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (33000, 37841)

```

[4.4] Word2Vec

```

In [41]: # Train your own Word2Vec model using your own text corpus
i=0
list_of_sentence=[]
for sentence in X_train:
    list_of_sentence.append(sentence.split())

```

```

In [42]: # Using Google News Word2Vectors

# in this project we are using a pretrained model by google
# its 3.3G file, once you load this into your memory
# it occupies ~9Gb, so please do this step only if you have >12G of ram
# we will provide a pickle file wich contains a dict ,
# and it contains all our courpus words as keys and model[word] as val
ues
# To use this code-snippet, download "GoogleNews-vectors-negative300.bi
n"
# from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edi
t
# it's 1.9GB in size.

# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17
SRFAzZPY
# 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 at least 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 google's word2vec file, keep want_to_train_w2v = True, to train your own w2v ")

```

```

[('terrific', 0.8466930389404297), ('excellent', 0.8420851230621338),
 ('awesome', 0.8315905928611755), ('fantastic', 0.8231490850448608), ('good', 0.8185310959815979), ('wonderful', 0.8021679520606995), ('perfect', 0.7707405090332031), ('amazing', 0.7036599516868591), ('fabulous', 0.6987220048904419), ('nice', 0.6965835690498352)]
=====

```

```

[('greatest', 0.7618408799171448), ('best', 0.7131766080856323), ('naughtiest', 0.7020399570465088), ('disgusting', 0.6772739291191101), ('tastiest', 0.6649549007415771), ('worse', 0.6477916836738586), ('displeasure', 0.6312363743782043), ('spoiled', 0.6199416518211365), ('terrible', 0.6138170957565308), ('horrible', 0.6043893098831177)]

```

```

In [43]: 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 15730
sample words ['thanks', 'order', 'great', 'job', 'quick', 'service',
'good', 'going', 'food', 'cat', 'future', 'not', 'usually', 'drinker',

```

```
'flavored', 'tea', 'fact', 'one', 'really', 'drink', 'bought', 'mango',  
'ceylon', 'several', 'years', 'ago', 'local', 'nature', 'store', 'whi  
m', 'fell', 'love', 'kids', 'since', 'told', 'flavor', 'discontinued',  
'believe', 'supplier', 'manufacturer', 'thank', 'goodness', 'online',  
'retailers', 'box', 'package', 'cheapest', 'could', 'find', 'time']
```

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

[4.4.1.1] Avg W2v

```
In [44]: # average Word2Vec  
# compute average word2vec for each review.  
X_tr_AvgW2V = []; # the avg-w2v for each sentence/review is stored in t  
his list  
for sent in tqdm(list_of_sentence): # for each review/sentence  
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo  
u might need to change this to 300 if you use google's w2v  
    cnt_words = 0; # num of words with a valid vector in the sentence/re  
view  
    for word in sent: # for each word in a review/sentence  
        if word in w2v_words:  
            vec = w2v_model.wv[word]  
            sent_vec += vec  
            cnt_words += 1  
    if cnt_words != 0:  
        sent_vec /= cnt_words  
    X_tr_AvgW2V.append(sent_vec)  
print(len(X_tr_AvgW2V))  
print(len(X_tr_AvgW2V[0]))
```

```
100%|██████████| 67000/67000 [01:35<00:00, 703.39it/s]
```

```
67000
```

```
50
```

```
In [45]: list_of_sentence_test=[]
        for sentence in x_test:
            list_of_sentence_test.append(sentence.split())
```

```
In [46]: # average Word2Vec
        # compute average word2vec for each review.
        X_test_AvgW2V = []; # the avg-w2v for each sentence/review is stored in
                             # this list
        for sent in tqdm(list_of_sentence_test): # for each review/sentence
            sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
            u might need to change this to 300 if you use google's w2v
            cnt_words = 0; # num of words with a valid vector in the sentence/re
            view
            for word in sent: # for each word in a review/sentence
                if word in w2v_words:
                    vec = w2v_model.wv[word]
                    sent_vec += vec
                    cnt_words += 1
            if cnt_words != 0:
                sent_vec /= cnt_words
            X_test_AvgW2V.append(sent_vec)
        print(len(X_test_AvgW2V))
        print(len(X_test_AvgW2V[0]))
```

```
100%|██████████| 33000/33000 [00:50<00:00, 652.71it/s]
```

```
33000
50
```

[4.4.1.2] TFIDF weighted W2v

```
In [55]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
        model = TfidfVectorizer()
        tf_idf_matrix = model.fit_transform(x_sampled_data)
        # we are converting a dictionary with word as a key, and the idf as a v
        alue
        dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

```
In [56]: # 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 ce
ll_val = tfidf

X_tr_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/r
review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
            #
            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole corpus
            # sent.count(word) = tf 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
    X_tr_tfidf_sent_vectors.append(sent_vec)
    row += 1
```

100%|██████████| 67000/67000 [26:35<00:00, 39.05it/s]

```
In [57]: # 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 ce
ll_val = tfidf

X_test_tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/revie
w is stored in this list
row=0;
for sent in tqdm(list_of_sentence_test): # for each review/sentence
```

```

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

```

100%|██████████| 33000/33000 [12:31<00:00, 43.91it/s]

[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 parameter tuning (best `depth` in range [1, 5, 10, 50, 100, 500, 100], and the best `min_samples_split` in range [5, 10, 100, 500])

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Find the best hyper parameter using k-fold cross validation or simple cross validation data

- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. 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 cross validation data for each hyper parameter, like shown in the figure.
 Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.



- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).



7. [Conclusion](#)

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



Note: Data Leakage

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

Applying Decision Trees

```
In [63]: # importing libraries
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_curve, auc
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import GridSearchCV
from sklearn.tree import DecisionTreeClassifier

import math
import warnings
warnings.filterwarnings("ignore")
```

[5.1] Applying Decision Trees on BOW, SET 1

```
In [38]: # grid search CV
depth=[1, 5, 10, 50, 100, 500, 1000]
min_samples=[5, 10, 100, 500]
tuned_parameters = [{'max_depth':depth, 'min_samples_split': min_samples}]

#Using GridSearchCV
clf = GridSearchCV(DecisionTreeClassifier(), tuned_parameters, scoring
= 'roc_auc', cv=5)
clf.fit(X_train_bow, Y_train)
```

```
Out[38]: GridSearchCV(cv=5, error_score='raise-deprecating',
      estimator=DecisionTreeClassifier(class_weight=None, criterion='g
      ini', max_depth=None,
      max_features=None, max_leaf_nodes=None,
      min_impurity_decrease=0.0, min_impurity_split=None,
      min_samples_leaf=1, min_samples_split=2,
      min_weight_fraction_leaf=0.0, presort=False, random_state=N
      one,
      splitter='best'),
      fit_params=None, iid='warn', n_jobs=None,
      param_grid=[{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_s
      amples_split': [5, 10, 100, 500]}],
      pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
      scoring='roc_auc', verbose=0)
```

```
In [41]: print('Best hyper parameter: ', clf.best_params_)
print('Model Score: ', clf.best_score_)
print('Model estimator: ', clf.best_estimator_)
```

```
Best hyper parameter: {'max_depth': 50, 'min_samples_split': 500}
Model Score:  0.8259173751582736
Model estimator: DecisionTreeClassifier(class_weight=None, criterion
='gini', max_depth=50,
      max_features=None, max_leaf_nodes=None,
```

```

min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=500,
min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
splitter='best')

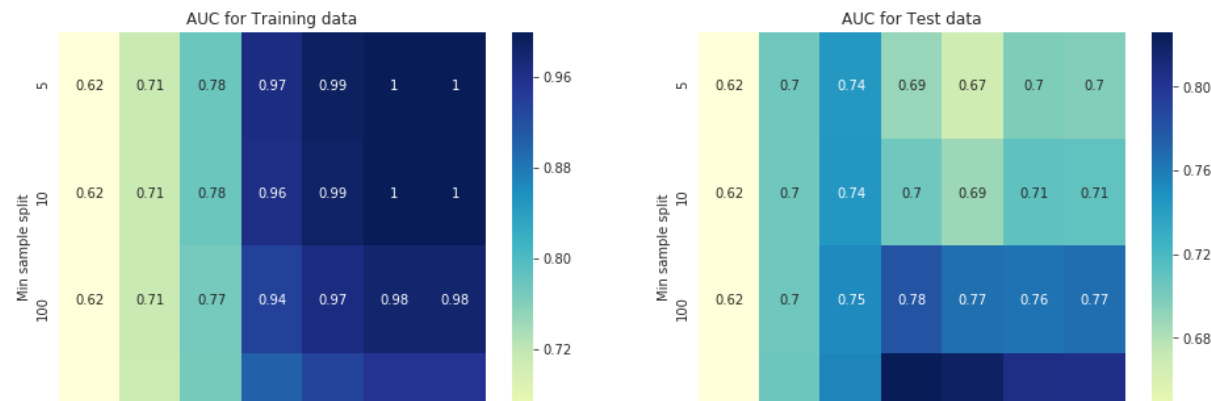
```

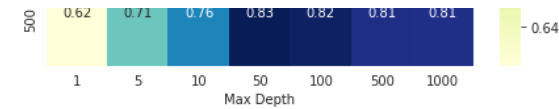
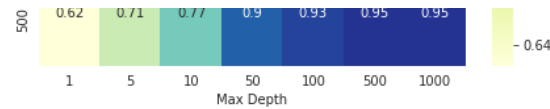
```

In [45]: max_depth_list = list(clf.cv_results_['param_max_depth'].data)
min_sample_split_list = list(clf.cv_results_['param_min_samples_split'].data)

sns.set_style("whitegrid")
plt.figure(figsize=(16,6))
plt.subplot(1,2,1)
data = pd.DataFrame(data={'Min sample split':min_sample_split_list, 'Max Depth':max_depth_list, 'AUC':clf.cv_results_['mean_train_score']})
data = data.pivot(index='Min sample split', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True, cmap="YlGnBu").set_title('AUC for Training data')
plt.subplot(1,2,2)
data = pd.DataFrame(data={'Min sample split':min_sample_split_list, 'Max Depth':max_depth_list, 'AUC':clf.cv_results_['mean_test_score']})
data = data.pivot(index='Min sample split', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True, cmap="YlGnBu").set_title('AUC for Test data')
plt.show()

```



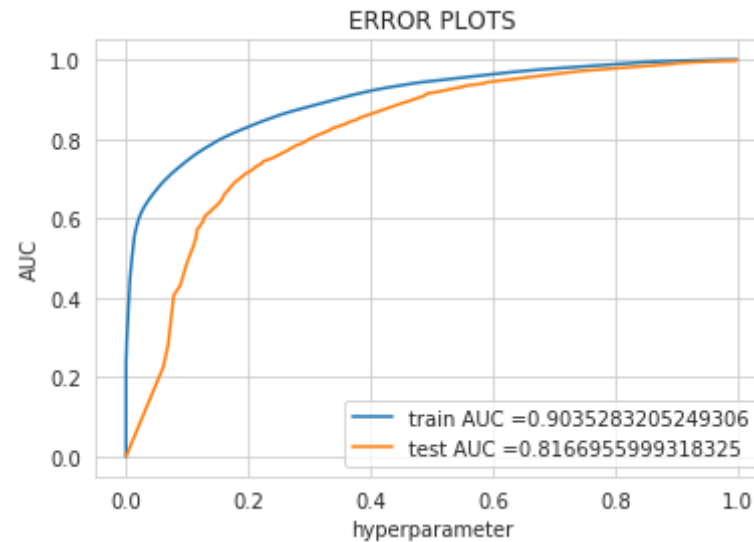


```
In [46]: model = DecisionTreeClassifier(max_depth= 50, min_samples_split=500, class_weight='balanced')
model.fit(X_train_bow, Y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(Y_train, model.predict_proba(X_train_bow)[:,-1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, model.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()

print('AUC: ', roc_auc_score(Y_train, model.predict(X_train_bow)))
```



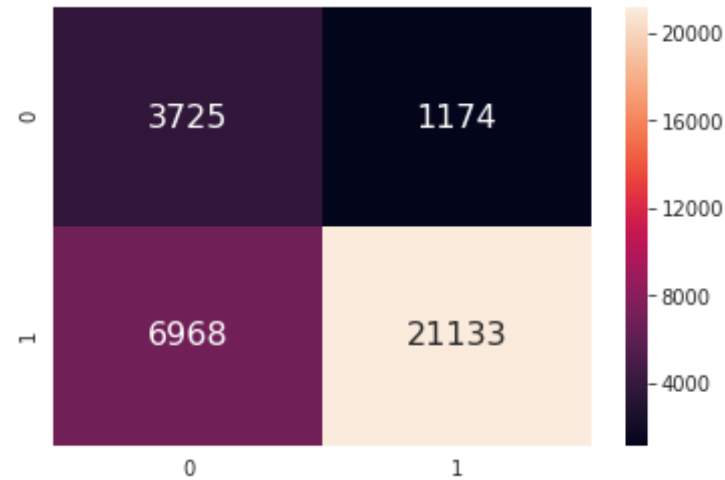
AUC: 0.8239492995726609

```
In [47]: from sklearn.metrics import confusion_matrix
print("Test confusion matrix")
print(confusion_matrix(y_test, model.predict(X_test_bow)))

# conprint('AUC: ',roc_auc_score(Y_train, m_nb.predict(X_tr_bow)))fusio
n matrix visualization using seaborn heatmap
df_test= pd.DataFrame(confusion_matrix(y_test, model.predict(X_test_bow
)))
sns.heatmap(df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

```
Test confusion matrix
[[ 3725  1174]
 [ 6968 21133]]
```

Out[47]: <matplotlib.axes._subplots.AxesSubplot at 0x7f8afde78ef0>



[5.1.1] Top 20 important features from SET 1

```
In [48]: w = count_vect.get_feature_names()
coef = model.feature_importances_
coeff_df = pd.DataFrame({'Word' : w, 'Coefficient' : coef})
coeff_df = coeff_df.sort_values(['Coefficient', 'Word'], ascending=[0,
1])
print(coeff_df.head(20).to_string(index=False))
```

Word	Coefficient
not	0.145390
great	0.089775
best	0.059412
delicious	0.040781
love	0.038556
disappointed	0.030840
good	0.025129
perfect	0.023498
loves	0.022366
wonderful	0.015745
excellent	0.014272
favorite	0.013016

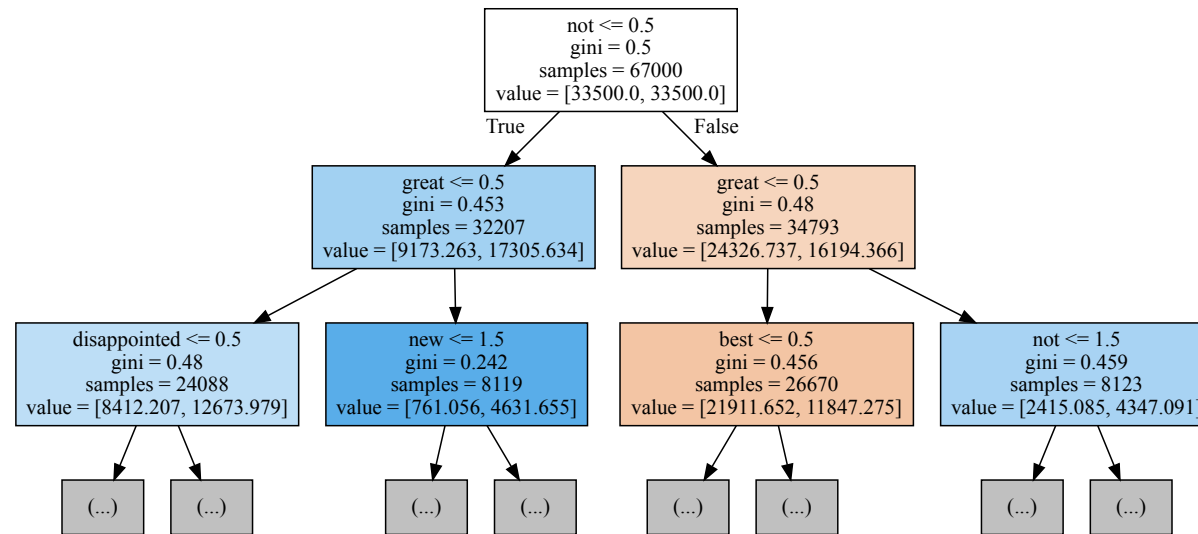
favorite	0.013040
money	0.011729
bad	0.011048
thought	0.010654
find	0.010015
unfortunately	0.009878
easy	0.009408
worst	0.009010
awful	0.008952

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

```
In [49]: from io import StringIO
from sklearn.tree.export import export_graphviz
from sklearn import tree
import graphviz
from IPython.display import SVG
```

```
In [50]: export_graphviz(model, out_file="tree.dot", feature_names = count_vect.
get_feature_names(), max_depth = 2, filled = True)
with open("tree.dot") as f:
    dot_graph = f.read()
graphviz.Source(dot_graph)
```

Out[50]:



[5.2] Applying Decision Trees on TFIDF, SET 2

```
In [51]: # grid search CV
depth=[1, 5, 10, 50, 100, 500, 1000]
min_samples=[5, 10, 50, 100, 500]
tuned_parameters = [{'max_depth':depth, 'min_samples_split': min_sample
s}]

#Using GridSearchCV
model = GridSearchCV(DecisionTreeClassifier(class_weight='balanced'), t
uned_parameters, scoring = 'roc_auc', cv=5)
model.fit(X_train_tfidf, Y_train)
```

```
Out[51]: GridSearchCV(cv=5, error_score='raise-deprecating',
                    estimator=DecisionTreeClassifier(class_weight='balanced', criterion='gini',
                    max_depth=None, max_features=None, max_leaf_nodes=None,
                    min_impurity_decrease=0.0, min_impurity_split=None,
                    min_samples_leaf=1, min_samples_split=2,
                    min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
```

```

        splitter='best'),
        fit_params=None, iid='warn', n_jobs=None,
        param_grid=[{'max_depth': [1, 5, 10, 50, 100, 500], 'min_samples
_split': [5, 10, 50, 100, 500]}],
        pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
        scoring='roc_auc', verbose=0)

```

```

In [52]: print('Best hyper parameter: ', model.best_params_)
         print('Model Score: ', model.best_score_)
         print('Model estimator: ', model.best_estimator_)

```

```

Best hyper parameter: {'max_depth': 50, 'min_samples_split': 500}
Model Score: 0.8027730607827726
Model estimator: DecisionTreeClassifier(class_weight='balanced', crite
rion='gini',
        max_depth=50, max_features=None, max_leaf_nodes=None,
        min_impurity_decrease=0.0, min_impurity_split=None,
        min_samples_leaf=1, min_samples_split=500,
        min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
        splitter='best')

```

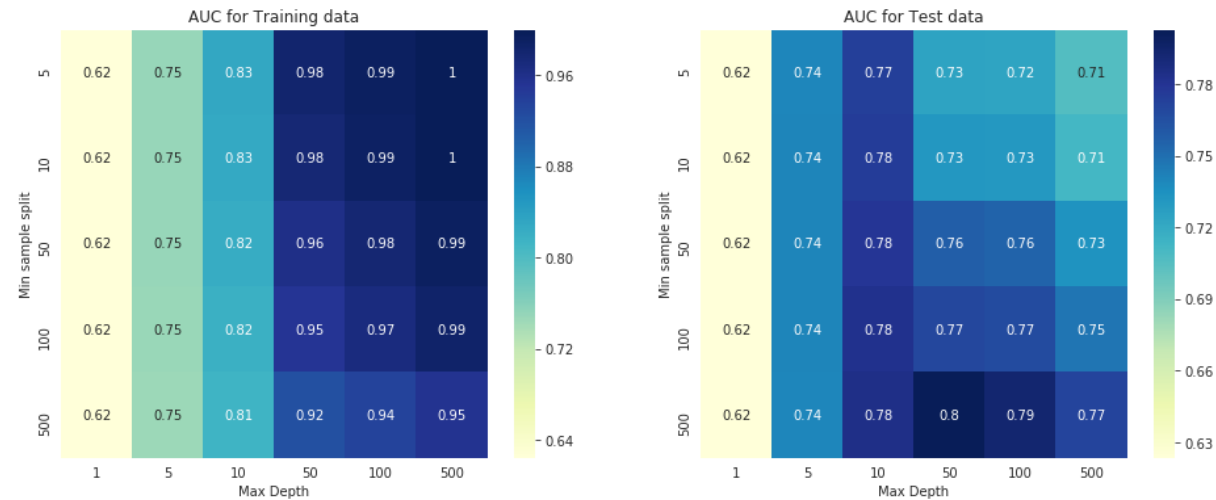
```

In [53]: max_depth_list = list(model.cv_results_['param_max_depth'].data)
         min_sample_split_list = list(model.cv_results_['param_min_samples_spli
t'].data)

         sns.set_style("whitegrid")
         plt.figure(figsize=(16,6))
         plt.subplot(1,2,1)
         data = pd.DataFrame(data={'Min sample split':min_sample_split_list, 'Ma
x Depth':max_depth_list, 'AUC':model.cv_results_['mean_train_score']})
         data = data.pivot(index='Min sample split', columns='Max Depth', values
='AUC')
         sns.heatmap(data, annot=True, cmap="YlGnBu").set_title('AUC for Trainin
g data')
         plt.subplot(1,2,2)
         data = pd.DataFrame(data={'Min sample split':min_sample_split_list, 'Ma
x Depth':max_depth_list, 'AUC':model.cv_results_['mean_test_score']})
         data = data.pivot(index='Min sample split', columns='Max Depth', values

```

```
='AUC')
sns.heatmap(data, annot=True, cmap="YlGnBu").set_title('AUC for Test da
ta')
plt.show()
```



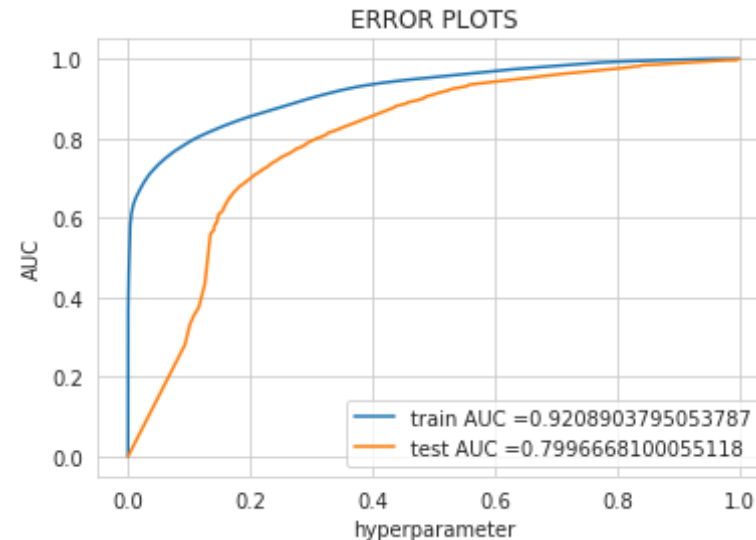
```
In [54]: model = DecisionTreeClassifier(max_depth= 50, min_samples_split=500, cl
ass_weight='balanced')
model.fit(X_train_tfidf, Y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probabilit
y estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(Y_train, model.predict_pro
ba(X_train_tfidf)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, model.predict_proba(
X_test_tfidf)[: ,1])

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, t
rain_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()

print('AUC: ',roc_auc_score(Y_train, model.predict(X_train_tfidf)))
```



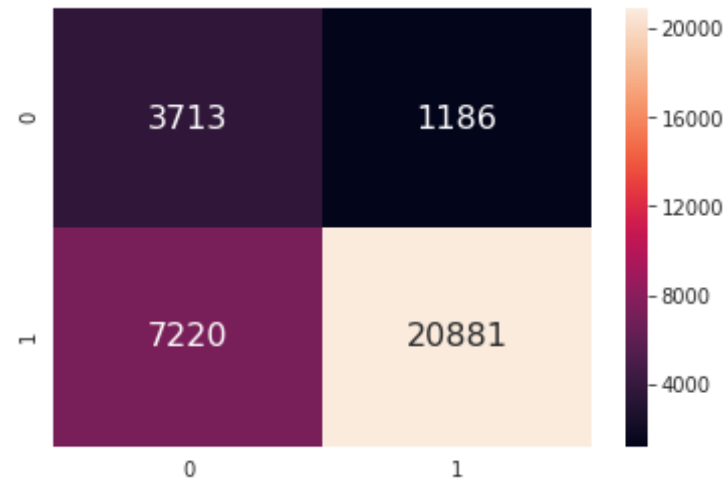
AUC: 0.8457787304006086

```
In [55]: from sklearn.metrics import confusion_matrix
print("Test confusion matrix")
print(confusion_matrix(y_test, model.predict(X_test_tfidf)))

# conprint('AUC: ',roc_auc_score(Y_train, m_nb.predict(X_tr_bow)))fusio
n matrix visualization using seaborn heatmap
df_test= pd.DataFrame(confusion_matrix(y_test, model.predict(X_test_tfi
df)))
sns.heatmap(df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

```
Test confusion matrix
[[ 3713  1186]
 [ 7220 20881]]
```

Out[55]: <matplotlib.axes._subplots.AxesSubplot at 0x7f8afb68beb8>



[5.2.1] Top 20 important features from SET 2

In [56]: `# Ref.: https://www.kaggle.com/laowingkin/amazon-fine-food-review-sentiment-analysis`

```
# Ref.:https://www.kaggle.com/laowingkin/amazon-fine-food-review-sentiment-analysis
w = tf_idf_vect.get_feature_names()
coef = model.feature_importances_.tolist()[0]
coeff_df = pd.DataFrame({'Word' : w, 'Coefficient' : coef})
coeff_df = coeff_df.sort_values(['Coefficient', 'Word'], ascending=[0, 1])
print(coeff_df.head(20).to_string(index=False))
```

Word	Coefficient
ab	0.0
abandon	0.0
abc	0.0
abdominal	0.0
ability	0.0
able	0.0
able break	0.0
able buy	0.0
able chew	0.0
able drink	0.0

able drink	0.0
able eat	0.0
able enjoy	0.0
able find	0.0
able get	0.0
able give	0.0
able go	0.0
able help	0.0
able keep	0.0
able locate	0.0
able make	0.0

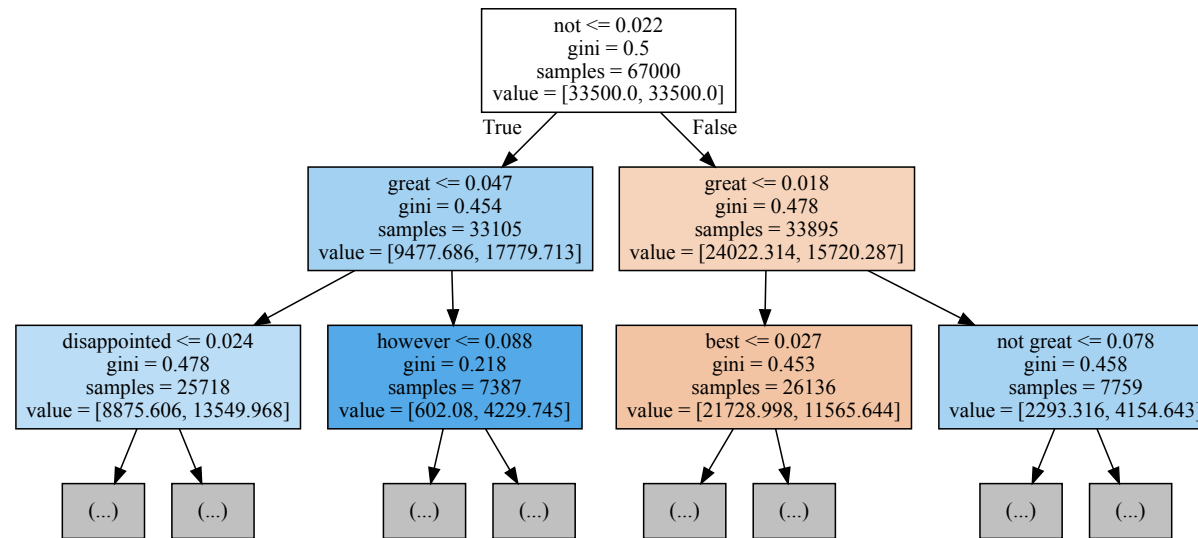
[5.2.2] Graphviz visualization of Decision Tree on TFIDF, SET 2

```
In [57]: model = DecisionTreeClassifier(max_depth= 50, min_samples_split=500, class_weight='balanced')
model.fit(X_train_tfidf, Y_train)
```

```
Out[57]: DecisionTreeClassifier(class_weight='balanced', criterion='gini',
                                max_depth=50, max_features=None, max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=1, min_samples_split=500,
                                min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                                splitter='best')
```

```
In [58]: export_graphviz(model, out_file="tree.dot", feature_names = tf_idf_vector.get_feature_names(), max_depth = 2, filled = True)
with open("tree.dot") as f:
    dot_graph = f.read()
graphviz.Source(dot_graph)
```

```
Out[58]:
```



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

```
In [48]: X_train_AvgW2v= np.array(X_tr_AvgW2V)
X_test_AvgW2v= np.array(X_test_AvgW2V)
```

```
In [49]: depth=[1, 5, 10, 50, 100, 500, 1000]
min_samples=[5, 10, 100, 500]
tuned_parameters = [{'max_depth':depth, 'min_samples_split': min_sample
s}]
```

#Using GridSearchCV

```
clf = GridSearchCV(DecisionTreeClassifier(class_weight='balanced'), tun
ed_parameters, scoring = 'roc_auc', cv=5)
clf.fit(X_train_AvgW2v, Y_train)
```

```
Out[49]: GridSearchCV(cv=5, error_score='raise-deprecating',
    estimator=DecisionTreeClassifier(class_weight='balanced', criter
ion='gini',
    max_depth=None, max_features=None, max_leaf_nodes=None,
    min_impurity_decrease=0.0, min_impurity_split=None,
    min_samples_leaf=1, min_samples_split=2,
```

```

        min_samples_leaf=1, min_samples_split=2,
        min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
        splitter='best'),
        fit_params=None, iid='warn', n_jobs=None,
        param_grid=[{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_s
amples_split': [5, 10, 100, 500]}],
        pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
        scoring='roc_auc', verbose=0)

```

```

In [50]: print('Best hyper parameter: ', clf.best_params_)
print('Model Score: ', clf.best_score_)
print('Model estimator: ', clf.best_estimator_)

```

```

Best hyper parameter: {'max_depth': 10, 'min_samples_split': 500}
Model Score: 0.8289073530474942
Model estimator: DecisionTreeClassifier(class_weight='balanced', crite
rion='gini',
        max_depth=10, max_features=None, max_leaf_nodes=None,
        min_impurity_decrease=0.0, min_impurity_split=None,
        min_samples_leaf=1, min_samples_split=500,
        min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
        splitter='best')

```

```

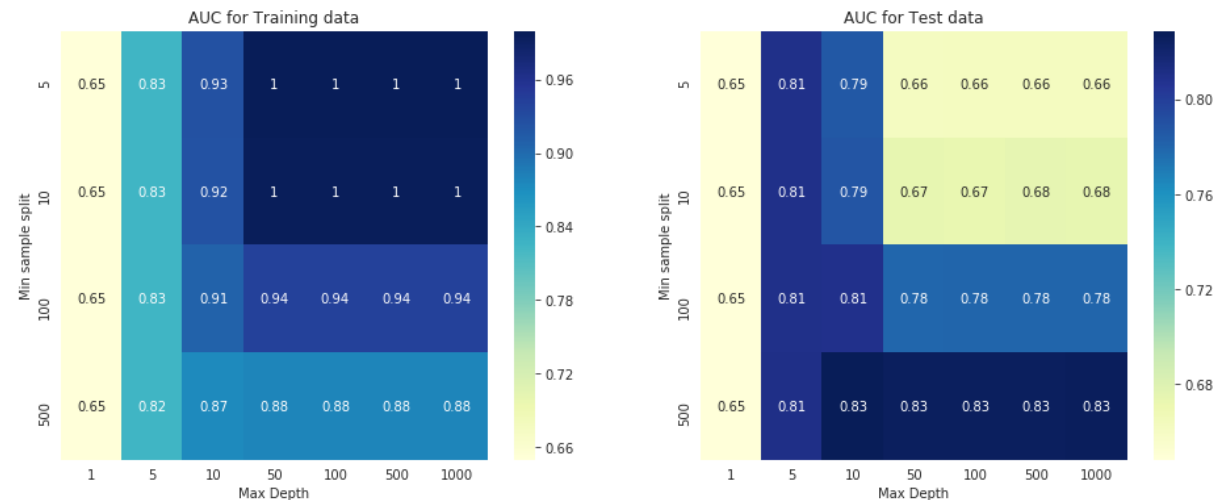
In [51]: max_depth_list = list(clf.cv_results_['param_max_depth'].data)
min_sample_split_list = list(clf.cv_results_['param_min_samples_split']
.data)

sns.set_style("whitegrid")
plt.figure(figsize=(16,6))
plt.subplot(1,2,1)
data = pd.DataFrame(data={'Min sample split':min_sample_split_list, 'Ma
x Depth':max_depth_list, 'AUC':clf.cv_results_['mean_train_score']})
data = data.pivot(index='Min sample split', columns='Max Depth', values
='AUC')
sns.heatmap(data, annot=True, cmap="YlGnBu").set_title('AUC for Trainin
g data')
plt.subplot(1,2,2)

```



```
data = pd.DataFrame(data={'Min sample split':min_sample_split_list, 'Max Depth':max_depth_list, 'AUC':clf.cv_results_['mean_test_score']})
data = data.pivot(index='Min sample split', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True, cmap="YlGnBu").set_title('AUC for Test data')
plt.show()
```



```
In [53]: model = DecisionTreeClassifier(max_depth= 10, min_samples_split=500, class_weight='balanced')
model.fit(X_train_AvgW2v, Y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(Y_train, model.predict_proba(X_train_AvgW2v)[:,-1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, model.predict_proba(X_test_AvgW2V)[:,-1])

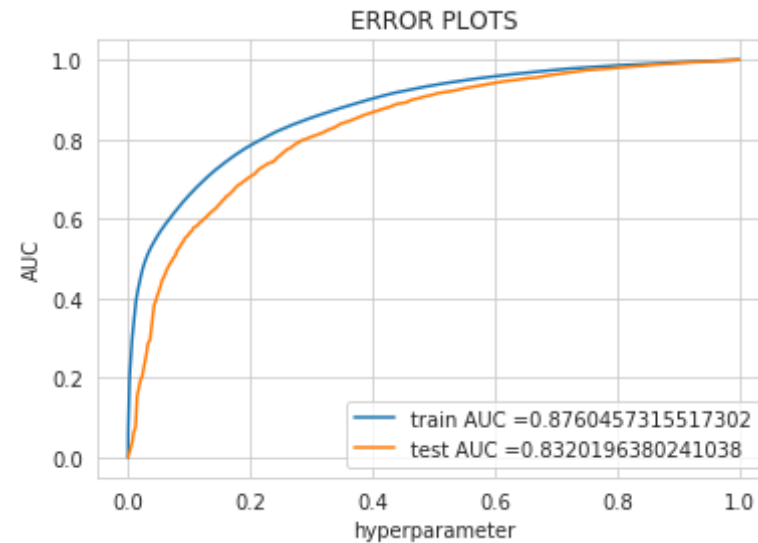
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
```

```

tpr))
plt.legend()
plt.xlabel("hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()

print('AUC: ',roc_auc_score(Y_train, model.predict(X_train_AvgW2v)))

```



AUC: 0.7927637249231676

```

In [54]: from sklearn.metrics import confusion_matrix
print("Test confusion matrix")
print(confusion_matrix(y_test, model.predict(X_test_AvgW2v)))

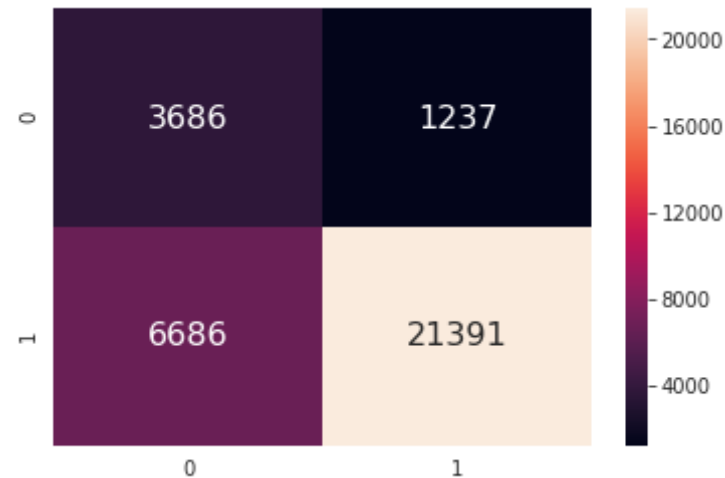
# conprint('AUC: ',roc_auc_score(Y_train, m_nb.predict(X_tr_bow)))fusio
n matrix visualization using seaborn heatmap
df_test= pd.DataFrame(confusion_matrix(y_test, model.predict(X_test_Avg
W2v)))
sns.heatmap(df_test, annot=True,annot_kws={"size": 16}, fmt='g')

```

Test confusion matrix
[[3686 1237]

```
[ 6686 21391]]
```

```
Out[54]: <matplotlib.axes._subplots.AxesSubplot at 0x7f81c94d57f0>
```



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

```
In [58]: X_train_tfidfW2V= np.array(X_tr_tfidf_sent_vectors)
X_test_tfidfW2V= np.array(X_test_tfidf_sent_vectors)
```

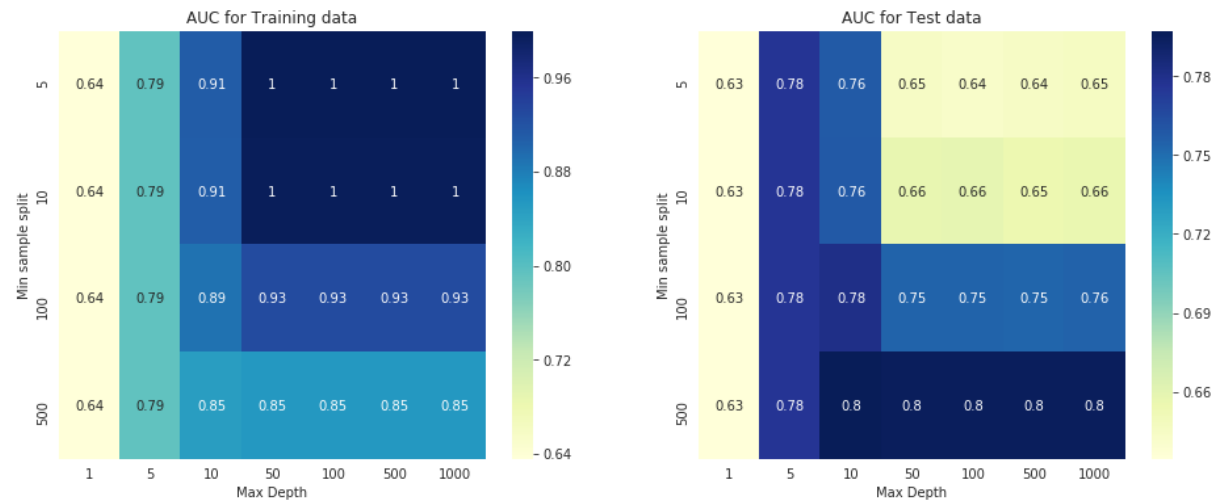
```
In [59]: depth=[1, 5, 10, 50, 100, 500, 1000]
min_samples=[5, 10, 100, 500]
tuned_parameters = [{'max_depth':depth, 'min_samples_split': min_samples}]

#Using GridSearchCV
model = GridSearchCV(DecisionTreeClassifier(class_weight='balanced'), tuned_parameters, scoring = 'roc_auc', cv=5)
model.fit(X_train_tfidfW2V, Y_train)
print('Best hyper parameter: ', model.best_params_)
print('Model Score: ', model.best_score_)
print('Model estimator: ', model.best_estimator_)
```

```
Best hyper parameter: {'max_depth': 10, 'min_samples_split': 500}
Model Score: 0.7972494691196318
Model estimator: DecisionTreeClassifier(class_weight='balanced', criterion='gini',
                                         max_depth=10, max_features=None, max_leaf_nodes=None,
                                         min_impurity_decrease=0.0, min_impurity_split=None,
                                         min_samples_leaf=1, min_samples_split=500,
                                         min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                                         splitter='best')
```

```
In [60]: max_depth_list = list(model.cv_results_['param_max_depth'].data)
min_sample_split_list = list(model.cv_results_['param_min_samples_split'].data)

sns.set_style("whitegrid")
plt.figure(figsize=(16,6))
plt.subplot(1,2,1)
data = pd.DataFrame(data={'Min sample split':min_sample_split_list, 'Max Depth':max_depth_list, 'AUC':model.cv_results_['mean_train_score']})
data = data.pivot(index='Min sample split', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True, cmap="YlGnBu").set_title('AUC for Training data')
plt.subplot(1,2,2)
data = pd.DataFrame(data={'Min sample split':min_sample_split_list, 'Max Depth':max_depth_list, 'AUC':model.cv_results_['mean_test_score']})
data = data.pivot(index='Min sample split', columns='Max Depth', values='AUC')
sns.heatmap(data, annot=True, cmap="YlGnBu").set_title('AUC for Test data')
plt.show()
```



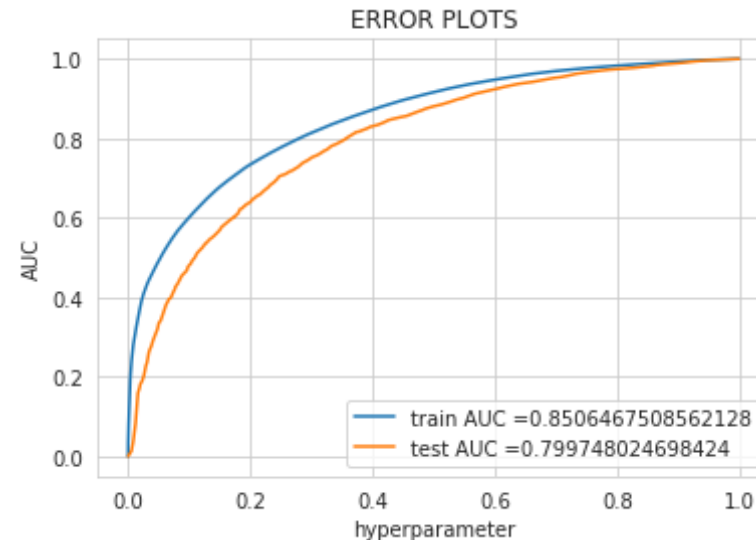
```
In [61]: model = DecisionTreeClassifier(max_depth= 10, min_samples_split=500, cl
ass_weight='balanced')
model.fit(X_train_tfidfW2V, Y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probabilit
y estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(Y_train, model.predict_pro
ba(X_train_tfidfW2V)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, model.predict_proba(
X_test_tfidfW2V)[: ,1])

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, t
rain_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()

print('AUC: ',roc_auc_score(Y_train, model.predict(X_train_tfidfW2V)))
```



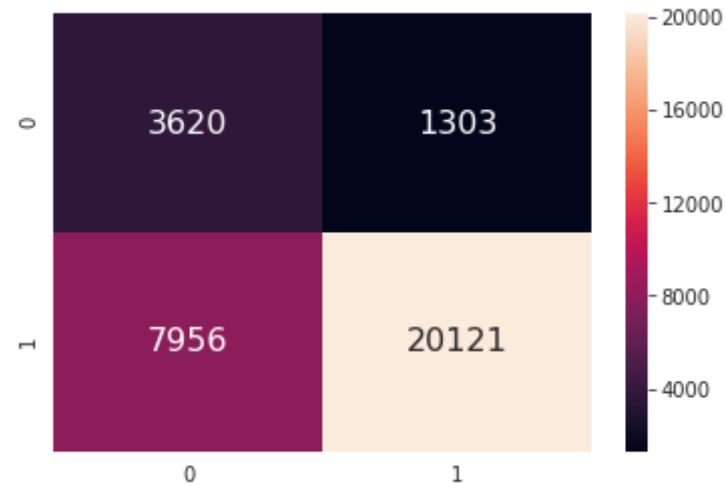
AUC: 0.7672198383401944

```
In [62]: from sklearn.metrics import confusion_matrix
print("Test confusion matrix")
print(confusion_matrix(y_test, model.predict(X_test_tfidfW2V)))

# conprint('AUC: ',roc_auc_score(Y_train, m_nb.predict(X_tr_bow)))fusio
n matrix visualization using seaborn heatmap
df_test= pd.DataFrame(confusion_matrix(y_test, model.predict(X_test_tfi
dfW2V)))
sns.heatmap(df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

```
Test confusion matrix
[[ 3620  1303]
 [ 7956 20121]]
```

Out[62]: <matplotlib.axes._subplots.AxesSubplot at 0x7f81c8159908>



[6] Conclusions

```
In [67]: # http://zetcode.com/python/prettytable/
from prettytable import PrettyTable

x = PrettyTable()
x.field_names= ["Vectorizer", "Hyper parameter(max_depth)", "Hyper para
meter(min_sample_split)", "AUC", "Score"]

x.add_row(["BOW", 50, 500, .82, .82])
x.add_row(["TFIDF", 50, 500, .84, .80])
x.add_row(["Avg W2V", 10, 500, .83, .79])
x.add_row(["TFIDF Avg W2V", 10, 500, .79, .76])
print(x)
```

```
+-----+-----+-----+-----+
| Vectorizer | Hyper parameter(max_depth) | Hyper parameter(min_samp
le_split) | AUC | Score |
+-----+-----+-----+-----+
| BOW       | 50 | 500 | 0.82 | 0.82 |
| TFIDF     | 50 | 500 | 0.84 | 0.80 |
| Avg W2V   | 10 | 500 | 0.83 | 0.79 |
| TFIDF Avg W2V | 10 | 500 | 0.79 | 0.76 |
```

	BOW			50		500
	0.82		0.82			
	TFIDF			50		500
	0.84		0.8			
	Avg W2V			10		500
	0.83		0.79			
	TFIDF Avg W2V			10		500
	0.79		0.76			
+	-----	+		+	-----	
-----	+	-----	+			

In []: