

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 [69]:

```
%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 [70]:

```

# 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 """, 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[70]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	1	1303862400
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	0	1346976000

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

In [71]:

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

In [72]:

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

(80668, 7)

Out[72]:

	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 [73]:

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

Out[73]:

	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 [74]:

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

Out[74]:

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 [75]:

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

Out[75]:

	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 [76]:

```
#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 [77]:

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

Out[77]:

(364173, 10)

In [78]:

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

Out[78]:

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 [79]:

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

display.head()
```

Out[79]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5	12248928
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2	4	12128832

In [80]:

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

In [81]:

```
#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[81]:
```

```
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 [82]:
```

```
# 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 car as we're driving along a
nd he always can sing the refrain. he's learned about whales, India, drooping roses: i love all t
he 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, b
ut geez, 2 years expired!!! I'm hoping to find local San Diego area shoppe that carries pods so t
hat 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 someting a dog would ever fi
nd 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 v
irgin 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.<br /><br />Thick, delicious. Perfect. 3 ingredients: Water, Maltitol, Natural Maple
Flavor. PERIOD. No chemicals. No garbage.<br /><br />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.<br /><br />I use this as my SWEETENER in baking: cheesecakes, white brownies,
muffins, pumpkin pies, etc... Unbelievably delicious...<br /><br />Can you tell I like it? :)
=====
```

In [83]:

```
# 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 [84]:

```
# 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 [85]:

```
# 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 [86]:

```
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 70 is it was poisonous until they figured out a way to fix that. I still like it but it could be better.

=====

In [87]:

```
#remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
sent_0 = re.sub(r"\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 [88]:

```
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub(r'[^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 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 70 is it was poisonous until they figured out a way to fix that I still like it but it could be better

In [89]:

```
# 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', '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"]])

```

In [90]:

```

# Combining all the above stundents
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
[08:01<00:00, 756.50it/s]

```

In [91]:

```
preprocessed_reviews[1500]
```

Out[91]:

```
'great ingredients although chicken rather chicken broth thing not think belongs canola oil canola
rapeseed not someting dog would ever find nature find rapeseed nature eat would poison today food
industries convinced masses canola oil safe even better oil olive virgin coconut facts though say
otherwise late poisonous figured way fix still like could better'
```

In [92]:

```
data = preprocessed_reviews[:100000]
scores = final["Score"][:100000]
```

In [93]:

```

from sklearn.model_selection import train_test_split

data_train,data_test,scores_train,scores_test = train_test_split(data,scores,shuffle = False,random
_state = 42,test_size = 0.2)
data_train,data_cv,scores_train,scores_cv = train_test_split(data_train,scores_train,shuffle =
False,random_state = 42,test_size = 0.25)

```

[3.2] Preprocessing Review Summary

In [94]:

```
## Similarly you can do preprocessing for review summary also.
```

[4] Featurization

[4.1] BAG OF WORDS

In [95]:

```
# BoW
bow_vect = CountVectorizer() #in scikit-learn
bow_vect.fit(data_train)
bow_data_train = bow_vect.fit_transform(data_train)
bow_data_cv = bow_vect.transform(data_cv)
bow_data_test = bow_vect.transform(data_test)
```

[4.2] TF-IDF

In [96]:

```
# tf-idf
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(data_train)
tf_idf_data_train = tf_idf_vect.fit_transform(data_train)
tf_idf_data_cv = tf_idf_vect.transform(data_cv)
tf_idf_data_test = tf_idf_vect.transform(data_test)
```

[4.3] Word2Vec

In [97]:

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

w2v_model=Word2Vec(X_train,min_count=5,size=50, workers=4)
w2v_words = list(w2v_model.wv.vocab)
```

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

[4.4.1.1] Avg W2v

In [98]:

```
def avg_W2V(list_of_sentence, w2v_model, w2v_words):
    # average Word2Vec
    # compute average word2vec for each review.
    sent_vectors = []; # the avg-w2v for each sentence/review is stored in this list
    for sent in tqdm(list_of_sentence): # for each review/sentence
        sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change t
his to 300 if you use google's w2v
        cnt_words = 0; # num of words with a valid vector in the sentence/review
        for word in sent.split(): # for each word in a review/sentence
            if word in w2v_words:
                vec = w2v_model.wv[word]
                sent_vec += vec
                cnt_words += 1
        if cnt_words != 0:
            sent_vec /= cnt_words
        sent_vectors.append(sent_vec)
    return sent_vectors

avgw2v_data_train = avg_W2V(data_train, w2v_model, w2v_words)
avgw2v_data_cv = avg_W2V(data_cv, w2v_model, w2v_words)
avgw2v_data_test = avg_W2V(data_test, w2v_model, w2v_words)
```

```
100%|███████████| 60000/60000 [06:  
55<00:00, 144.32it/s]
```

```
100%|███████████| 20000/20000 [02:
```

[illegible]

[4.4.1.2] TFIDF weighted W2v

In [99]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(data_train)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
tf_idf_features = model.get_feature_names()
```

In [100]:

```
# TF-IDF weighted Word2Vec

def tf_idf_w2v(list_of_sentence, w2v_model, w2v_words, tfidf_feat, dictionary):
    tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list

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

    return tfidf_sent_vectors

tf_idf_w2v_data_train = tf_idf_w2v(data_train, w2v_model, w2v_words, tf_idf_features, dictionary)
tf_idf_w2v_data_cv = tf_idf_w2v(data_cv, w2v_model, w2v_words, tf_idf_features, dictionary)
tf_idf_w2v_data_test = tf_idf_w2v(data_test, w2v_model, w2v_words, tf_idf_features, dictionary)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 60000/60000  
[1:58:12<00:00, 8.46it/s]  
100%|██████████████████████████████████████████████████████████████████████████| 20000/20000 [42  
:32<00:00, 7.83it/s]  
100%|██████████████████████████████████████████████████████████████████████████| 20000/20000 [38  
:52<00:00, 8.57it/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 paramter 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 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 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-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method fit_transform() on you train data, and apply the method transform() on cv/test data.
4. For more details please go through this [link](#).

Applying Decision Trees

In [101]:

```
from sklearn import tree

def get_AUC(X_train,y_train,X_cv,y_cv,list_depth,sample_split):
    """This function apply decision tree classifier
    on train and cv data and return AUC values for train and cross validation"""
    auc_train = []
    auc_cv = []
    # applying Decision Tree on list of hyper parameters to find best alpha using simple loop
    for depth in list_depth:
        clf = tree.DecisionTreeClassifier(max_depth = depth,min_samples_split = sample,class_weight
= "balanced",random_state = 42)
        clf.fit(X_train, y_train)
        prob_train = clf.predict_proba(X_train)
        fpr, tpr, threshold = roc_curve(y_train, prob_train[:, 1])
        auc_train.append(auc(fpr,tpr))
        prob_cv = clf.predict_proba(X_cv)
        fpr, tpr, threshold = roc_curve(y_cv, prob_cv[:, 1])
        auc_cv.append(auc(fpr,tpr))
```

```

return auc_train, auc_cv

def plot_AUC_Curves(auc_train, auc_cv, depth, title):
    """This function plots the auc curves for the given auc values and alpha"""
    sns.set_style("whitegrid", {'axes.grid' : False})
    plt.plot(depth, auc_train, "b-", label = "AUC_Train")
    plt.plot(depth, auc_cv, "r-", label = "AUC_Validation")
    plt.scatter(depth, auc_train)
    plt.scatter(depth, auc_cv)
    plt.legend()
    plt.xlabel("Hyper Parameter (Depth)")
    plt.ylabel("AUC")
    plt.title(title)
    plt.show()

def apply_roc_curve(X_train, y_train, X_test, y_test, max_depth, min_sample_split):
    """This function apply DecisionTree model on train and predict labels for test data
    and also find FPR and TPR for train and test data.
    Returns the predicted labels, FPR and TPR values"""
    clf = tree.DecisionTreeClassifier(max_depth = max_depth, min_samples_split = min_sample_split, random_state = 42)
    clf.fit(X_train, y_train)
    prob_train = clf.predict_proba(X_train)
    fpr_train, tpr_train, threshold = roc_curve(y_train, prob_train[:, 1])
    prob_test = clf.predict_proba(X_test)
    fpr_test, tpr_test, threshold = roc_curve(y_test, prob_test[:, 1])

    # predict the class labels
    pred_train = clf.predict(X_train)
    pred_test = clf.predict(X_test)
    return fpr_train, tpr_train, fpr_test, tpr_test, pred_train, pred_test, clf

def plot_roc_curve(fpr_train, tpr_train, fpr_test, tpr_test):
    """This function plot the roc curves for train and test data"""
    # plot ROC curves for train and test data
    plt.plot(fpr_train, tpr_train, "g-", label = "AUC_Train : "+str(auc(fpr_train, tpr_train)))
    plt.plot(fpr_test, tpr_test, "r-", label = "AUC_Test : "+str(auc(fpr_test, tpr_test)))
    plt.plot([0, 1], [0, 1], "b-")
    plt.legend(loc="lower right")
    plt.xlabel("False Positive Rate")
    plt.ylabel("True Positive Rate")
    plt.title("ROC Curve")
    plt.show()

def plot_Confusion_Matrix(actual_labels, predict_labels, title):
    """This function plot the confusion matrix"""
    # Reference : https://seaborn.pydata.org/generated/seaborn.heatmap.html
    cm = confusion_matrix(actual_labels, predict_labels)
    classNames = ['NO', 'YES']
    cm_data = pd.DataFrame(cm, index = classNames,
                           columns = classNames)
    plt.figure(figsize = (5, 4))
    sns.heatmap(cm_data, annot=True, fmt="d")
    plt.title(title)
    plt.ylabel('Actual label')
    plt.xlabel('Predicted label')
    plt.show()

```

[5.1] Applying Decision Trees on BOW, SET 1

In [102]:

```

# max_depth
max_depths = [1, 5, 10, 50, 100, 500, 1000]
# min_sample_split
min_samples = [5, 10, 100, 500]
bow_max_depth = 0
bow_min_samples_split = 0

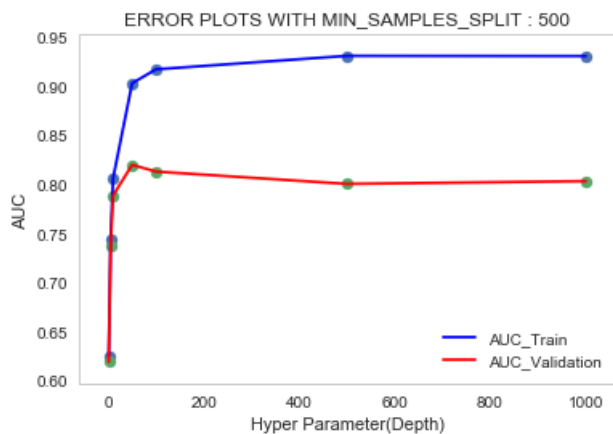
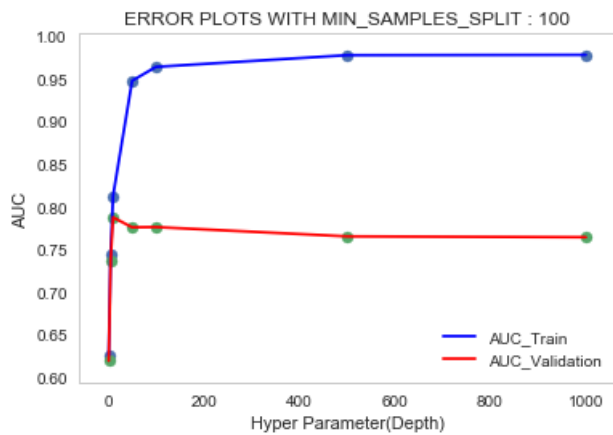
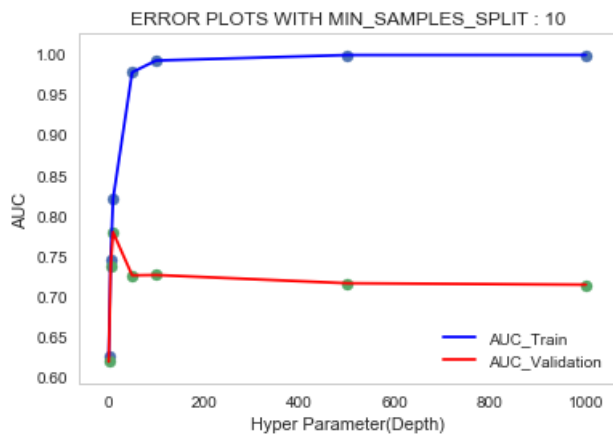
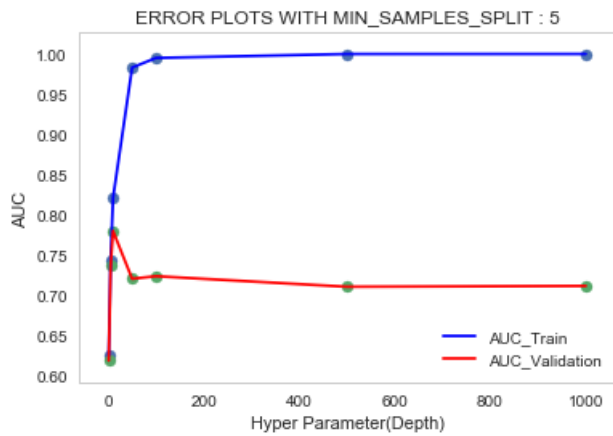
for sample in min_samples:
    # AUC
    auc_train, auc_cv = get_AUC(bow_data_train, scores_train, bow_data_cv, scores_cv, max_depths, sample)
    # Plot AUC curves
    plot_AUC_Curves(auc_train, auc_cv, max_depths, "ERROR PLOTS WITH MIN_SAMPLES_SPLIT : "+str(sample))

```

```

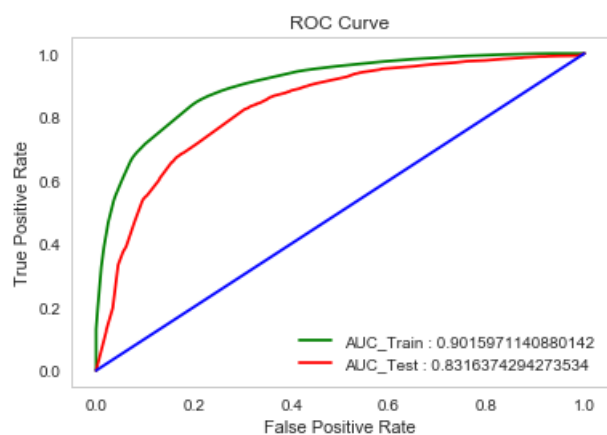
)
max_depth = max_depths[auc_cv.index(max(auc_cv))]
min_samples_split = sample
if max_depth >= bow_max_depth:
    bow_max_depth = max_depth
    bow_min_samples_split = min_samples_split

```



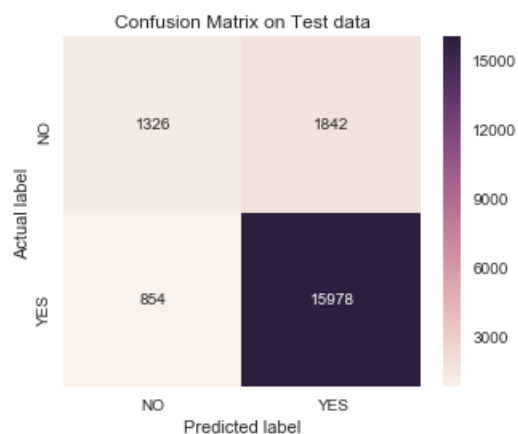
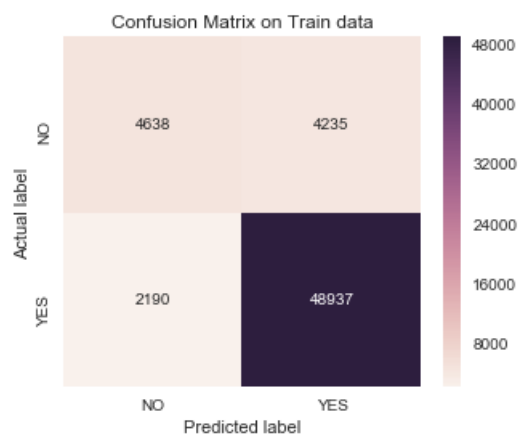
In [103]:

```
# roc
fpr_train,tpr_train,fpr_test,tpr_test,pred_train,pred_test,clf =
apply_roc_curve(bow_data_train,scores_train,bow_data_test,scores_test,bow_max_depth,bow_min_samples
_split)
# plot roc
plot_roc_curve(fpr_train,tpr_train,fpr_test,tpr_test)
```



In [104]:

```
# Confusion matrix
plot_Confusion_Matrix(scores_train,pred_train,"Confusion Matrix on Train data")
plot_Confusion_Matrix(scores_test,pred_test,"Confusion Matrix on Test data")
```



In [105]:

```
# AUC
bow_auc = auc(fpr_test,tpr_test)
print("AUC = ",round(bow_auc,2))
```

```
print("Max_depth = ",bow_max_depth)
print("Min_samples_split = ",bow_min_samples_split)
```

```
AUC = 0.83
Max_depth = 50
Min_samples_split = 500
```

[5.1.1] Top 20 important features from SET 1

In [106]:

```
features = clf.feature_importances_
top20_important_indices = list(features.argsort()[-20:])
top20_important_indices.reverse()
rank = np.array(range(1,21))
top20_important_features = np.take(bow_vect.get_feature_names(),top20_important_indices)
gini_importance = np.take(features,top20_important_indices)
top20_important_features_details = pd.DataFrame(data = {'Rank' : rank,'Feature' :
top20_important_features,'Gini_Importance' : gini_importance})
print(top20_important_features_details)
```

	Feature	Gini_Importance	Rank
0	not	0.099506	1
1	great	0.054575	2
2	disappointed	0.053912	3
3	money	0.042002	4
4	best	0.036840	5
5	worst	0.033615	6
6	love	0.026004	7
7	delicious	0.022174	8
8	terrible	0.021446	9
9	threw	0.020735	10
10	good	0.020539	11
11	waste	0.019217	12
12	refund	0.016199	13
13	disappointing	0.015862	14
14	loves	0.014145	15
15	awful	0.013687	16
16	perfect	0.011646	17
17	favorite	0.011282	18
18	find	0.011161	19
19	wonderful	0.009891	20

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

In [107]:

```
# Please write all the code with proper documentation
import graphviz
dot_data = tree.export_graphviz(clf,out_file= "tree_bow.pdf",feature_names=bow_vect.get_feature_names(),max_depth=3)
graph = graphviz.Source(dot_data)

with open("tree_bow.pdf") as f:
    dot_graph=f.read()
graphviz.Source(dot_graph)
```

Out[107]:



[5.2] Applying Decision Trees on TFIDF, SET 2

In [108]:

```
# max_depth
max_depths = [1,5,10,50,100,500,1000]
# min_sample_split
```



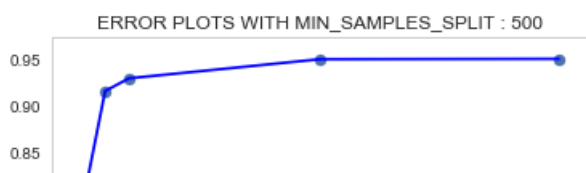
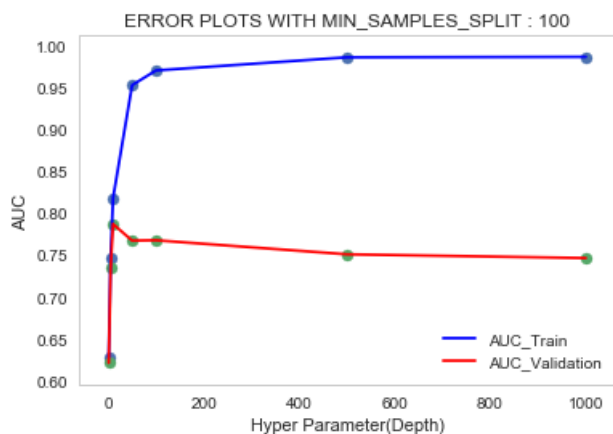
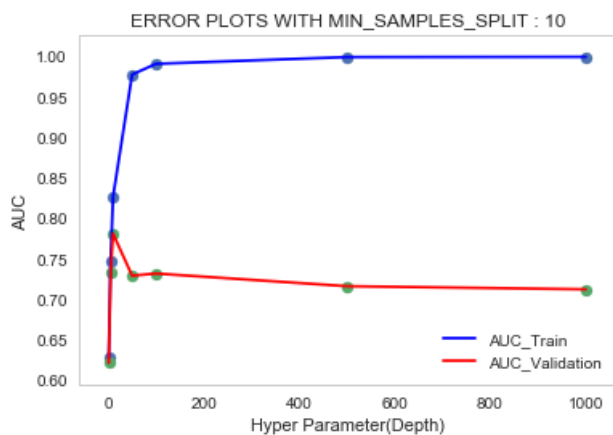
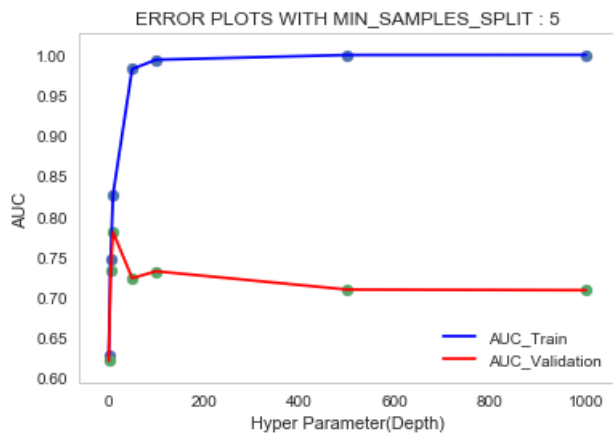
```

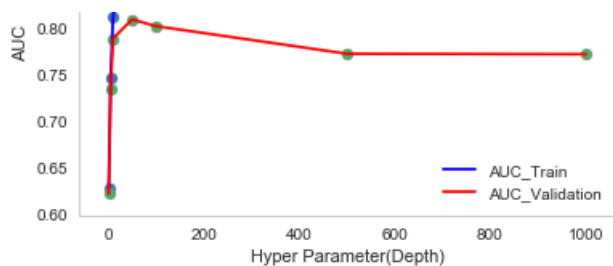
min_samples = [5,10,100,500]
tf_idf_max_depth = 0
tf_idf_min_samples_split = 0

for sample in min_samples:
    # AUC
    auc_train,auc_cv = get_AUC(tf_idf_data_train,scores_train,tf_idf_data_cv,scores_cv,max_depths,sample)
    # Plot AUC curves
    plot_AUC_Curves(auc_train,auc_cv,max_depths,"ERROR PLOTS WITH MIN_SAMPLES_SPLIT : "+str(sample)
)

    max_depth = max_depths[auc_cv.index(max(auc_cv))]
    min_samples_split = sample
    if max_depth >= tf_idf_max_depth:
        tf_idf_max_depth = max_depth
        tf_idf_min_samples_split = min_samples_split

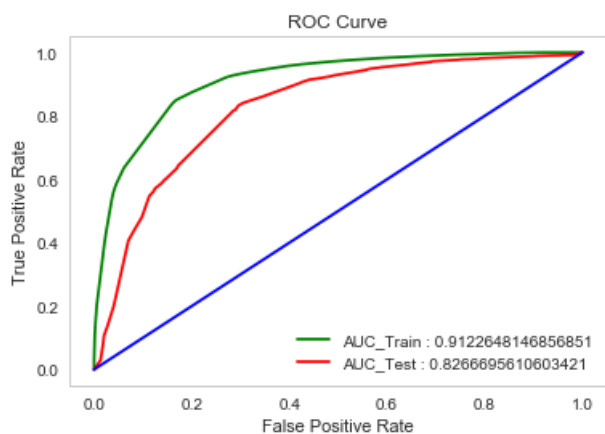
```





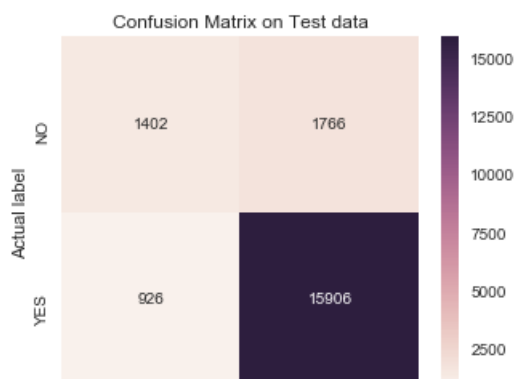
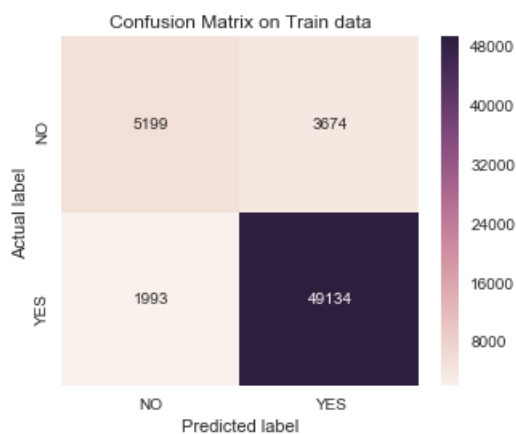
In [109]:

```
# roc
fpr_train, tpr_train, fpr_test, tpr_test, pred_train, pred_test, clf = apply_roc_curve(tf_idf_data_train,
,scores_train, tf_idf_data_test, scores_test, tf_idf_max_depth, tf_idf_min_samples_split)
# plot roc
plot_roc_curve(fpr_train, tpr_train, fpr_test, tpr_test)
```



In [110]:

```
# Confusion matrix
plot_confusion_matrix(scores_train, pred_train, "Confusion Matrix on Train data")
plot_confusion_matrix(scores_test, pred_test, "Confusion Matrix on Test data")
```



NO

YES

Predicted label

In [111]:

```
# AUC
tf_idf_auc = auc(fpr_test, tpr_test)
print("AUC = ", round(tf_idf_auc, 2))
print("Max_depth = ", tf_idf_max_depth)
print("Min_samples_split = ", tf_idf_min_samples_split)
```

```
AUC = 0.83
Max_depth = 50
Min_samples_split = 500
```

[5.2.1] Top 20 important features from SET 2

In [112]:

```
features = clf.feature_importances_  
top20_important_indices = list(features.argsort() [-20:])  
top20_important_indices.reverse()  
rank = np.array(range(1,21))  
top20_important_features = np.take(tf_idf_vect.get_feature_names(),top20_important_indices)  
gini_importance = np.take(features,top20_important_indices)  
top20_important_features_details = pd.DataFrame(data = {'Rank' : rank,'Feature' :  
top20_important_features,'Gini_Importance' : gini_importance})  
print(top20_important_features_details)
```

	Feature	Gini_Importance	Rank
0	not	0.101129	1
1	disappointed	0.046651	2
2	great	0.045408	3
3	worst	0.030670	4
4	not buy	0.027788	5
5	money	0.025470	6
6	threw	0.025416	7
7	waste money	0.025350	8
8	best	0.024318	9
9	not recommend	0.022227	10
10	refund	0.020327	11
11	terrible	0.020058	12
12	awful	0.016179	13
13	not worth	0.015372	14
14	love	0.015173	15
15	not disappointed	0.014343	16
16	horrible	0.014313	17
17	disappointing	0.012528	18
18	not even	0.012501	19
19	delicious	0.011968	20

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

In [113]:

```
import graphviz
dot_data = tree.export_graphviz(clf,out_file= "tree_tfidf.pdf",feature_names=tf_idf_vect.get_feature_names(),max_depth=3)
graph = graphviz.Source(dot_data)

with open("tree_bow.pdf") as f:
    dot_graph=f.read()
graphviz.Source(dot_graph)
```

Out[113]:

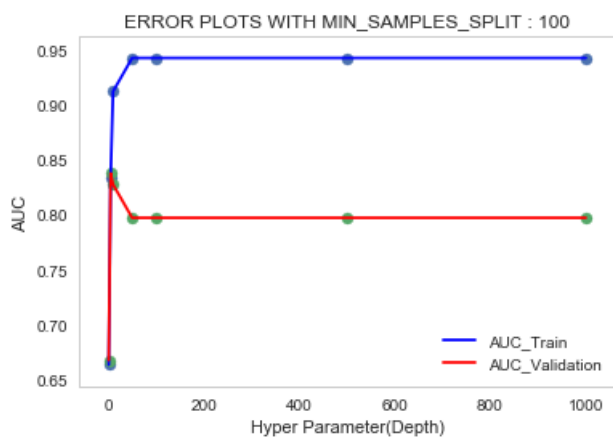
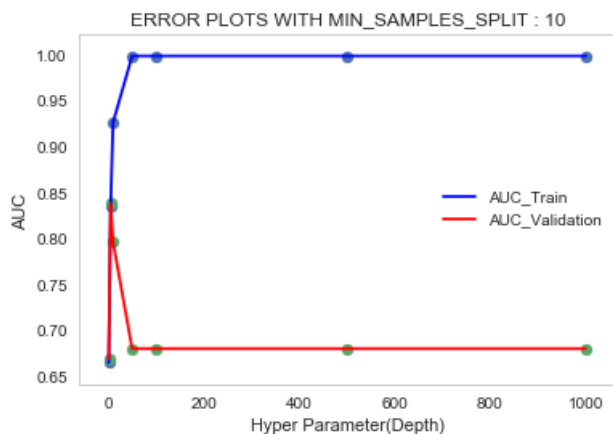
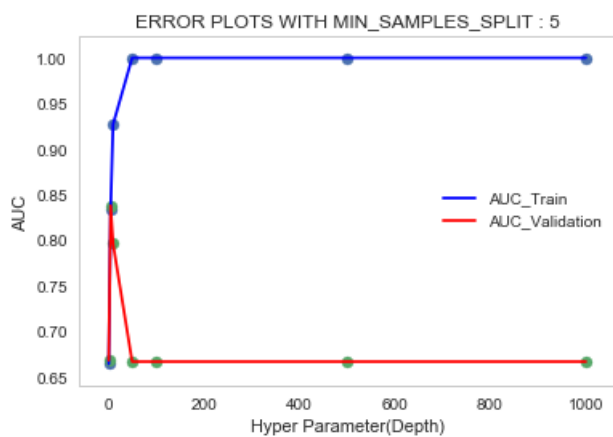
[5.3] Applying Decision Trees on AVG W2V, SET 3

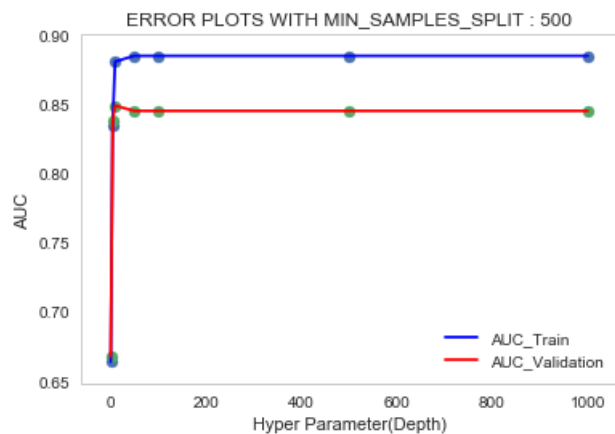
In [114]:

```
# max_depth
max_depths = [1,5,10,50,100,500,1000]
# min_sample_split
min_samples = [5,10,100,500]
avgw2v_max_depth = 0
avgw2v_min_samples_split = 0

for sample in min_samples:
    # AUC
    auc_train,auc_cv = get_AUC(avgw2v_data_train,scores_train,avgw2v_data_cv,scores_cv,max_depths,sample)
    # Plot AUC curves
    plot_AUC_Curves(auc_train,auc_cv,max_depths,"ERROR PLOTS WITH MIN_SAMPLES_SPLIT : "+str(sample))
)

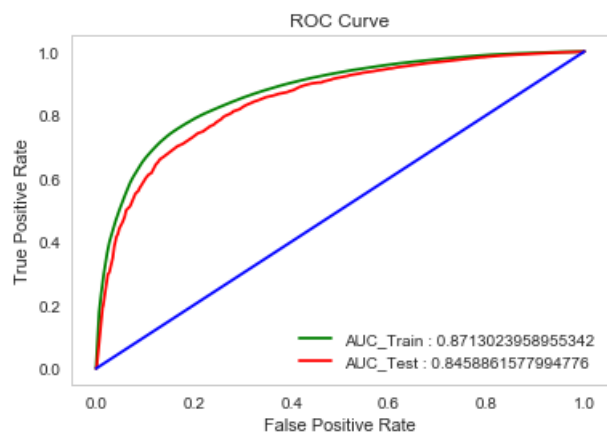
max_depth = max_depths[auc_cv.index(max(auc_cv))]
min_samples_split = sample
if max_depth >= avgw2v_max_depth:
    avgw2v_max_depth = max_depth
    avgw2v_min_samples_split = min_samples_split
```





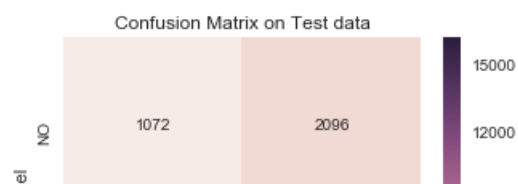
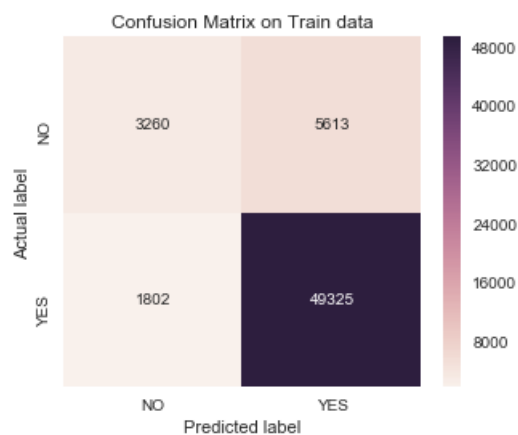
In [115]:

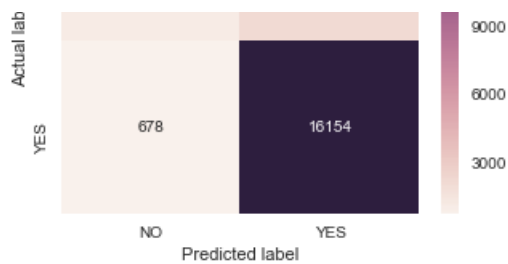
```
# roc
fpr_train,tpr_train,fpr_test,tpr_test,pred_train,pred_test,clf = apply_roc_curve(avgw2v_data_train
,scores_train,avgw2v_data_test,scores_test,avgw2v_max_depth,avgw2v_min_samples_split)
# plot roc
plot_roc_curve(fpr_train,tpr_train,fpr_test,tpr_test)
```



In [116]:

```
# Confusion matrix
plot_Confusion_Matrix(scores_train,pred_train,"Confusion Matrix on Train data")
plot_Confusion_Matrix(scores_test,pred_test,"Confusion Matrix on Test data")
```





In [117]:

```
# AUC
avgw2v_auc = auc(fpr_test,tpr_test)
print("AUC = ",round(avgw2v_auc,2))
print("Max_depth = ",avgw2v_max_depth)
print("Min_samples_split = ",avgw2v_min_samples_split)
```

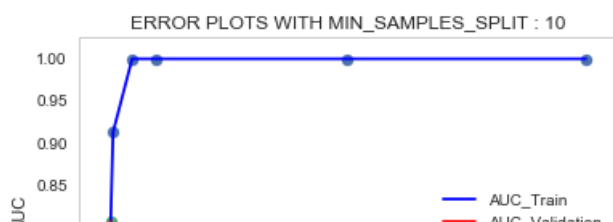
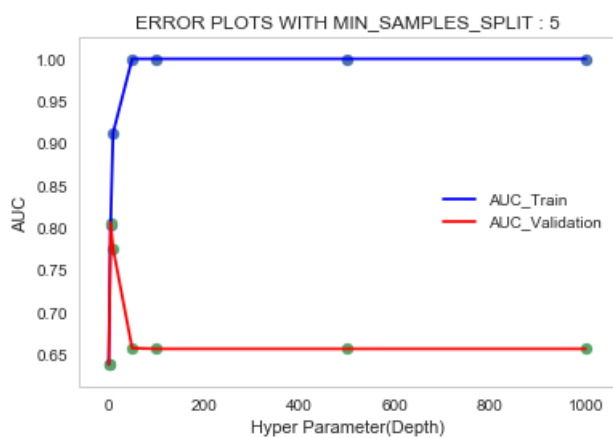
```
AUC = 0.85
Max_depth = 10
Min_samples_split = 500
```

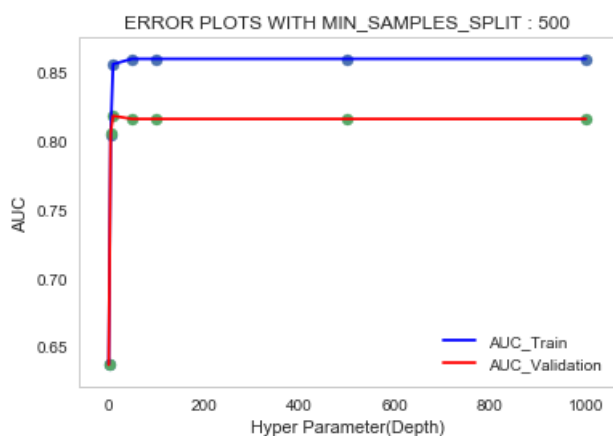
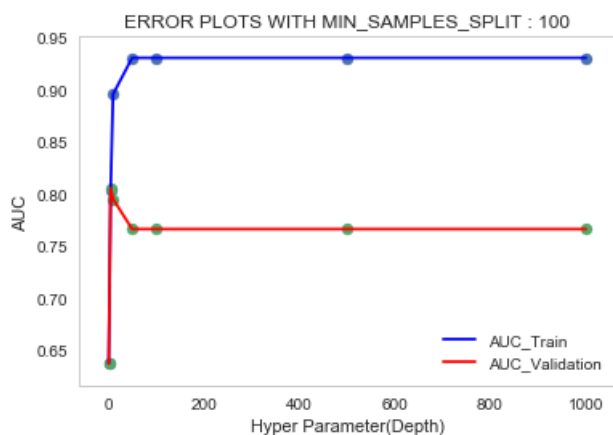
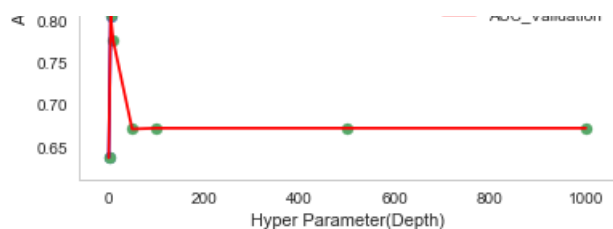
[5.4] Applying Decision Trees on TFIDF W2V, SET 4

In [118]:

```
# max_depth
max_depths = [1,5,10,50,100,500,1000]
# min_sample_split
min_samples = [5,10,100,500]
tf_idf_w2v_max_depth = 0
tf_idf_w2v_min_samples_split = 0

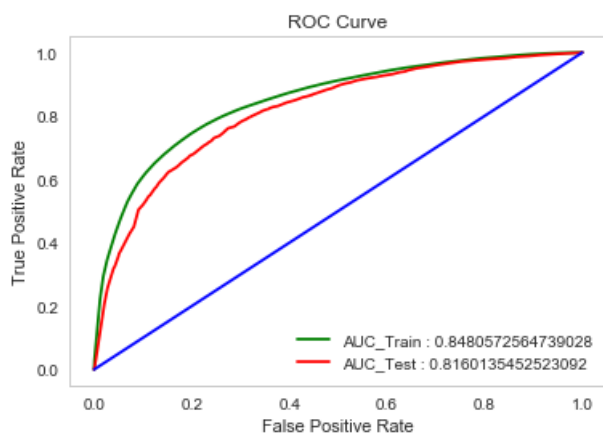
for sample in min_samples:
    # AUC
    auc_train,auc_cv =
get_AUC(tf_idf_w2v_data_train,scores_train,tf_idf_w2v_data_cv,scores_cv,max_depths,sample)
    # Plot AUC curves
    plot_AUC_Curves(auc_train,auc_cv,max_depths,"ERROR PLOTS WITH MIN_SAMPLES_SPLIT : "+str(sample)
)
    max_depth = max_depths[auc_cv.index(max(auc_cv))]
    min_samples_split = sample
    if max_depth >= tf_idf_w2v_max_depth:
        tf_idf_w2v_max_depth = max_depth
        tf_idf_w2v_min_samples_split = min_samples_split
```





In [119]:

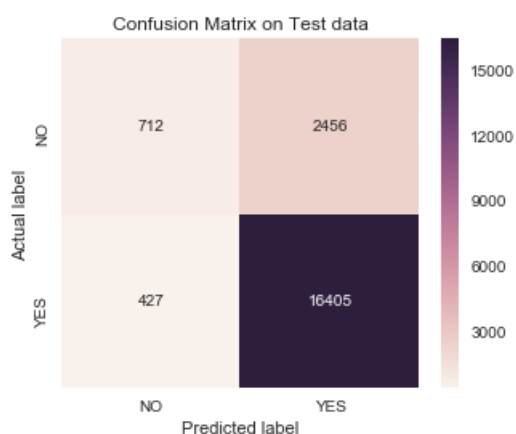
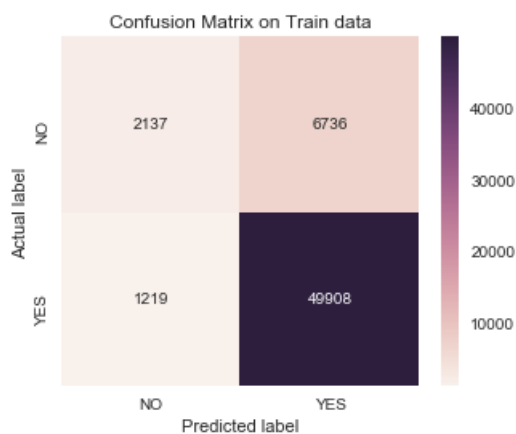
```
# roc
fpr_train,tpr_train,fpr_test,tpr_test,pred_train,pred_test,clf =
apply_roc_curve(tf_idf_w2v_data_train,scores_train,tf_idf_w2v_data_test,scores_test,tf_idf_w2v_max_
depth,tf_idf_w2v_min_samples_split)
# plot roc
plot_roc_curve(fpr_train,tpr_train,fpr_test,tpr_test)
```



In [120]:

```
# Confusion matrix
plot Confusion Matrix(scores_train,pred_train,"Confusion Matrix on Train data")
```

```
plot_Confusion_Matrix(scores_test,pred_test,"Confusion Matrix on Test data")
```



In [121]:

```
# AUC
tf_idf_w2v_auc = auc(fpr_test, tpr_test)
print("AUC = ", round(tf_idf_w2v_auc, 2))
print("Max_depth = ", tf_idf_w2v_max_depth)
print("Min_samples_split = ", tf_idf_w2v_min_samples_split)
```

```
AUC = 0.82
Max_depth = 10
Min_samples_split = 500
```

[6] Conclusions

In [124]:

```
from prettytable import PrettyTable

table = PrettyTable()
table.field_names = ["Vectorization", "Max_Depth", "Min_Samples_Split", "AUC"]
table.add_row(["BOW", bow_max_depth, bow_min_samples_split, round(bow_auc, 2)])
table.add_row(["TF-IDF", tf_idf_max_depth, tf_idf_min_samples_split, round(tf_idf_auc, 2)])
table.add_row(["Avg-W2V", avgw2v_max_depth, avgw2v_min_samples_split, round(avgw2v_auc, 2)])
table.add_row(["TF-IDF W2V", tf_idf_w2v_max_depth, tf_idf_w2v_min_samples_split, round(tf_idf_w2v_auc, 2)])
print(table.get_string(title="Results"))
```

```
+-----+
|               Results               |
+-----+-----+-----+-----+
| Vectorization | Max_Depth | Min_Samples_Split | AUC |
+-----+-----+-----+-----+
| BOW           | 50        | 500               | 0.83 |
| TF-IDF        | 50        | 500               | 0.83 |
| Avg-W2V       | 10        | 500               | 0.85 |
| TF-IDF W2V    | 10        | 500               | 0.85 |
+-----+-----+-----+-----+
```


	Avg-W2V		10		500		0.85	
	TF-IDF W2V		10		500		0.82	
+-----+-----+-----+-----+								

In []: