

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
```

```
# using SQLite Table to read data.

db_path = "/home/monodeepdas112/Datasets/amazon-fine-food-reviews/database.sqlite"
con = sqlite3.connect(db_path)

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

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

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

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

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

Out[2]:

Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
2	3	B000LQOCH0	ABXLMWJXXAIN	1	1	1	1219017600	Natania says it all

In [3]:

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

In [4]:

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

(80668, 7)

Out[4]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
0	#oc-R115TNMSPFT9I7	B007Y59HVM	Breyton	1331510400	2	Overall its just OK when considering the price...	2
1	#oc-R11D9D7SHXIB9	B005HG9ET0	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms, u...	3
2	#oc-R11DNU2NBKQ23Z	B007Y59HVM	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately not ...	2
3	#oc-R11O5J5ZVQE25C	B005HG9ET0	Penguin Chick	1346889600	5	This will be the bottle that you grab from the...	3
4	#oc-R12KPBODL2B5ZD	B007OSBE1U	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2

In [5]:

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

Out[5]:

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

In [6]:

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

Out[6]:

393063

[2] Exploratory Data Analysis

[2.1] Data Cleaning: Deduplication

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

In [7]:

```
display = pd.read_sql_query("""
```

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

Out [7]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summ
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE

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

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

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

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

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

In [8]:

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

In [9]:

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

Out [9]:

(87775, 10)

In [10]:

```
#Checking to see how much % of data still remains
```

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

Out[10]:

87.775

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

In [11]:

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

display.head()
```

Out[11]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5	1224892800	Bought This for My Son at College
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2	4	1212883200	Pure cocoa taste with crunchy almonds inside

In [12]:

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

In [13]:

```
#Before starting the next phase of preprocessing lets see the number of entries left
print(final.shape)

#How many positive and negative reviews are present in our dataset?
final['Score'].value_counts()
```

(87773, 10)

Out[13]:

```
1    73592
0    14181
Name: Score, dtype: int64
```

[3] Preprocessing

[3.1]. Preprocessing Review Text

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

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

1. Begin by removing the html tags

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

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

In [14]:

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

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

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

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

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its ver y hard to find any chicken products made in the USA but they are out there, but this one isnt. It s too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

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

=====

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

=====

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

=====

In [15]:

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

print(sent_0)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its ver y hard to find any chicken products made in the USA but they are out there, but this one isnt. It s too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

In [16]:

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

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

```

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

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

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

```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its ver y hard to find any chicken products made in the USA but they are out there, but this one isnt. It s too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

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

=====

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

=====

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

In [17]:

```

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

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

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

```

In [18]:

```

sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)

```

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

=====

In [19]:

```

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

```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its ver y hard to find any chicken products made in the USA but they are out there, but this one isnt. It s too bad too because its a good product but I wont take any chances till they know what is going

on with the china imports.

In [20]:

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

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

In [21]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have revmoved in the 1st step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "y
ou're", "you've", \
                "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
'himself', \
                'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them',
'their', \
                'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll",
'these', 'those', \
                'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
'do', 'does', \
                'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', '
while', 'of', \
                'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during',
'before', 'after', \
                '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', "d
oesn'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 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%|██████████| 87773/87773 [00:35<00:00, 2481.07it/s]

In [23]:

```
preprocessed_reviews[1500]
```

Out[23]:


```
'way hot blood took bite jig lol'
```

[3.2] Preprocessing Review Summary

In [24]:

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

[4] Featurization

[4.1] BAG OF WORDS

In [25]:

```
# #BoW
# count_vect = CountVectorizer() #in scikit-learn
# count_vect.fit(preprocessed_reviews)
# print("some feature names ", count_vect.get_feature_names()[0:10])
# print('='*50)

# final_counts = count_vect.transform(preprocessed_reviews)
# print("the type of count vectorizer ",type(final_counts))
# print("the shape of out text BOW vectorizer ",final_counts.get_shape())
# print("the number of unique words ", final_counts.get_shape()[1])
```

[4.2] Bi-Grams and n-Grams.

In [26]:

```
# #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)
# tf_idf_vect.fit(preprocessed_reviews)
# print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names()[0:10])
# print('='*50)

# final_tf_idf = tf_idf_vect.transform(preprocessed_reviews)
# print("the type of count vectorizer ",type(final_tf_idf))
# print("the shape of out text TFIDF vectorizer ",final_tf_idf.get_shape())
# print("the number of unique words including both unigrams and bigrams ",
final_tf_idf.get_shape()[1])
```

[4.4] Word2Vec

In [28]:

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

In [29]:

```
# # Using Google News Word2Vectors

# # in this project we are using a pretrained model by google
# # its 3.3G file, once you load this into your memory
# # it occupies ~9Gb, so please do this step only if you have >12G of ram
# # we will provide a pickle file wich contains a dict ,
# # and it contains all our courpus words as keys and model[word] as values
# # To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
# # from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS2lpQmM/edit
# # it's 1.9GB in size.

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

# is_your_ram_gt_16g=False
# want_to_use_google_w2v = False
# want_to_train_w2v = True

# if want_to_train_w2v:
#     # min_count = 5 considers only words that occurred atleast 5 times
#     w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
#     print(w2v_model.wv.most_similar('great'))
#     print('='*50)
#     print(w2v_model.wv.most_similar('worst'))

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

In [30]:

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

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

[4.4.1.1] Avg W2v

In [31]:

```
# # 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 th
is to 300 if you use google's w2v
#     cnt_words = 0; # num of words with a valid vector in the sentence/review
#     for word in sent: # for each word in a review/sentence
#         if word in w2v_words:
#             vec = w2v_model.wv[word]
#             sent_vec += vec
#             cnt_words += 1
#     if cnt_words != 0:
#         sent_vec /= cnt_words
```

```

# sent_vec /= cnt_words
# sent_vectors.append(sent_vec)
# print(len(sent_vectors))
# print(len(sent_vectors[0]))

```

[4.4.1.2] TFIDF weighted W2v

In [32]:

```

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

```

In [33]:

```

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

# tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
# row=0;
# for sent in tqdm(list_of_sentence): # for each review/sentence
#     sent_vec = np.zeros(50) # as word vectors are of zero length
#     weight_sum = 0; # num of words with a valid vector in the sentence/review
#     for word in sent: # for each word in a review/sentence
#         if word in w2v_model.wv:
#             vec = w2v_model.wv[word]
#             tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
#             # to reduce the computation we are
#             # dictionary[word] = idf value of word in whole corpus
#             # sent.count(word) = tf value of word in this review
#             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
#             sent_vec += (vec * tf_idf)
#             weight_sum += tf_idf
#     if weight_sum != 0:
#         sent_vec /= weight_sum
#     tfidf_sent_vectors.append(sent_vec)
#     row += 1

```

[5] Assignment 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

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

Applying Decision Trees

In [34]:

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve, auc
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import train_test_split
from sklearn.model_selection import StratifiedKFold
from sklearn.tree import DecisionTreeClassifier
import pprint
import os.path
import pickle
import math

import warnings
warnings.filterwarnings('ignore')
```

[5.0.0] Splitting up the Dataset into D_train and D_test

In [35]:

```
num_data_points = 100000
```

In [36]:

```
Dx_train, Dx_test, Dy_train, Dy_test = train_test_split(preprocessed_reviews[:num_data_points],
final['Score'].tolist()[:num_data_points], test_size=0.30, random_state=42)
```

In [37]:

```
prettytable_data = []
```

[5.0.1] Defining some functions to increase code reusability and readability

In [38]:

```
'''Creating Custom Vectorizers for TFIDF - W2Vec and Avg - W2Vec'''
class Tfidf_W2Vec_Vectorizer(object):
    def __init__(self, w2vec_model):
        if(w2v_model is None):
            raise Exception('Word 2 Vector model passed to Tfidf_W2Vec Vectorizer is None !')
        self.tfidf = TfidfVectorizer(max_features=300)
        self.dictionary = None
        self.tfidf_feat = None

        self.word2vec = w2vec_model

    def fit(self, X):
        '''X : list'''
        #Initializing the TFIDF Vectorizer
        self.tfidf.fit_transform(X)
        # we are converting a dictionary with word as a key, and the idf as a value
        self.dictionary = dict(zip(self.tfidf.get_feature_names(), list(self.tfidf.idf_)))
        self.tfidf_feat = self.tfidf.get_feature_names()

        return self

    def transform(self, X):
        '''X : list'''
        return np.array([
            np.mean([self.word2vec[w] * self.dictionary[word]*(X.count(word)/len(X))
                     for w in words if w in self.word2vec and w in self.tfidf_feat] or
                     [np.zeros(300)], axis=0)
            for words in X
        ])

class Avg_W2Vec_Vectorizer(object):
    def __init__(self, w2vec_model):
        if(w2v_model is None):
            raise Exception('Word 2 Vector model passed to Avg_W2Vec Vectorizer is None !')
        self.word2vec = w2vec_model

    def fit(self, X):
        return self

    def transform(self, X):
        '''X : list'''
        return np.array([
            np.mean([self.word2vec[w] for w in words if w in self.word2vec]
                    or [np.zeros(300)], axis=0)
            for words in X
        ])
1)
```

In [39]:

```
def get_vectorizer(vectorizer, train, W2V_model=None):
    if(vectorizer=='BOW'):
        vectorizer = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
    if(vectorizer=='TFIDF'):
        vectorizer = TfidfVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
    if(vectorizer=='TFIDF-W2Vec'):
        vectorizer = Tfidf_W2Vec_Vectorizer(W2V_model)
    if(vectorizer=='Avg-W2Vec'):
        vectorizer = Avg_W2Vec_Vectorizer(W2V_model)

    vectorizer.fit(train)
    return vectorizer
```

In [40]:

```
'''Perform Simple Cross Validation'''
def perform_hyperparameter_tuning(X, Y, vectorizer, results_path, retrain=False, W2V_model=None):
    #If the pandas dataframe with the hyperparameter info exists then return it

    if(retrain==False):
```

```

# If Cross Validation results exists then return them
if(os.path.exists(results_path)):
    return pd.read_csv(results_path)
else:
    # If no data exists but retrain=False then mention accordingly
    print('Retrain is set to be False but no Cross Validation Results DataFrame was found !
\nPlease set retrain to True.')
else:
    # else perform hyperparameter tuning
    print('Performing Hyperparameter Tuning...\n')
    # regularization parameter
    hyperparameters = {
        'dt_max_depth' : [1, 5, 10, 50, 100, 500, 1000],
        'dt_min_samples_split' : [5, 10, 100, 500]
    }

    max_depth = []
    min_samples_split = []

    train_scores = []
    test_scores = []

    train_mean_score = []
    test_mean_score = []

    # Initializing KFold
    skf = StratifiedKFold(n_splits=3)
    X = np.array(X)
    Y = np.array(Y)

    for depth in hyperparameters['dt_max_depth']:
        for min_samples in hyperparameters['dt_min_samples_split']:

            #Performing Cross Validation
            for train_index, test_index in skf.split(X, Y):
                Dx_train, Dx_cv = X[train_index], X[test_index]
                Dy_train, Dy_cv = Y[train_index], Y[test_index]

                #Initializing the Vectorizer
                vectorizer = get_vectorizer(vectorizer, Dx_train.tolist(), W2V_model)

                #Transforming the data to features
                x_train = vectorizer.transform(Dx_train.tolist())
                x_cv = vectorizer.transform(Dx_cv.tolist())

                #Initializing the LR model
                dt_clf = DecisionTreeClassifier(max_depth=depth, min_samples_split=min_samples)

                # Fit the model
                dt_clf.fit(x_train, Dy_train)

                #Prediction
                train_results = dt_clf.predict_proba(x_train)
                cv_results = dt_clf.predict_proba(x_cv)

                try:
                    train_score = roc_auc_score(Dy_train, train_results[:, 1])
                    test_score = roc_auc_score(Dy_cv, cv_results[:, 1])

                    #storing the results to form a dataframe
                    train_scores.append(train_score)
                    test_scores.append(test_score)

                except Exception as e:
                    print('Error Case : ', e)
                    print(('Actual, Predicted'))
                    [print((Dy_cv[i], cv_results[i, 1])) for i in range(len(Dy_cv))]

                print('CV iteration : depth={0}, min_samples={1}, train_score={2}, test_score={
3}'

                    .format(depth, min_samples, train_score, test_score))

            train_mean_score.append(sum(train_scores)/len(train_scores))
            test_mean_score.append(sum(test_scores)/len(test_scores))

        max_depth.append(depth)
        min_samples_split.append(min_samples)

```

```

        print('CV : depth={0}, min_samples_split={1}, train_score={2}, test_score={3}'
              .format(depth, min_samples, sum(train_scores)/len(train_scores), sum(test_scores)/len(test_scores)))

    train_scores = []
    test_scores = []

    # Creating a DataFrame from the saved data for visualization
    results_df = pd.DataFrame({
        'max_depth' : max_depth,
        'min_samples_split' : min_samples_split,
        'train_score' : train_mean_score,
        'test_score' : test_mean_score
    })

    #writing the results to csv after performing hyperparameter tuning
    try:
        results_df.to_csv(results_path)
    except Exception as ex:
        print(str(ex), "\nError occured while converting DataFrame to CSV after cross validation.")

    return results_df

```

In [41]:

```

def analyse_results(df):
    # plotting error curves
    fig = plt.figure(figsize=(15, 15))
    ax = fig.gca()

    unique_min_samples = np.unique(df['min_samples_split'].values)
    c = 1
    for i in unique_min_samples:
        mini = df.loc[df['min_samples_split'] == i]
        plt.subplot(len(unique_min_samples)//2, len(unique_min_samples)//2, c)
        plt.plot([math.log10(i) for i in mini.max_depth.tolist()], mini.train_score.tolist(), '-o',
c='r', label='Train AUC')
        plt.plot([math.log10(i) for i in mini.max_depth.tolist()], mini.test_score.tolist(), '-o',
c='b', label='Validation AUC')
        plt.grid(True)
        plt.xlabel('log10 of Hyperparameter : max_depth')
        plt.ylabel('Area Under ROC Curve')
        plt.title('AUC ROC : min_samples_split = {0}'.format(i))
        plt.legend(loc='best')
        c = c + 1

    plt.show()

    # return the best parameters
    mmax = 0
    ind_max = 0
    for index, row in df.iterrows():
        if(row['test_score']>mmax):
            mmax=row['test_score']
            ind_max = index

    best_params = {
        'max_depth': df.loc[ind_max, 'max_depth'],
        'min_samples_split':df.loc[ind_max, 'min_samples_split']
    }

    return best_params

```

In [42]:

```

def retrain_with_best_params(data, labels, best_params, vec_name, model_path, word2vec):
    if(os.path.exists(model_path)):
        print('Loading Model....')
        with open(model_path, 'rb') as input_file:
            dt_clf = pickle.load(input_file)
    else:
        dt_clf = DecisionTreeClassifier(max_depth=best_params['max_depth'], min_samples_split=best_params['min samples split'])

```

```

print('Initializing Vectorizer')
vectorizer = get_vectorizer(vectorizer=vec_name, train=data, W2V_model=word2vec)
print('Training Model....')
dt_clf.fit(vectorizer.transform(data), np.array(labels))

print('Saving Trained Model....')
with open(model_path, 'wb') as file:
    pickle.dump(dt_clf, file)
return dt_clf

```

In [43]:

```

def plot_confusion_matrix(model, data, labels, dataset_label):
    pred = model.predict(data)
    conf_mat = confusion_matrix(labels, pred)

    strings = strings = np.asarray([[ 'TN = ', 'FP = '],
                                     [ 'FN = ', 'TP = ']])

    labels = (np.asarray(["{0}{1}".format(string, value)
                          for string, value in zip(strings.flatten(),
                                                    conf_mat.flatten())])
              ).reshape(2, 2)

    fig, ax = plt.subplots()
    ax.set(xlabel='Predicted', ylabel='Actual', title='Confusion Matrix : {0}'.format(dataset_label))
    sns.heatmap(conf_mat, annot=labels, fmt="", cmap='YlGnBu', ax=ax)
    ax.set_xlabel('Predicted')
    ax.set_ylabel('Actual')
    ax.set_xticklabels(['False', 'True'])
    ax.set_yticklabels(['False', 'True'])
    plt.show()

```

In [44]:

```

def plot_AUC_ROC(model, vectorizer, Dx_train, Dx_test, Dy_train, Dy_test):

    #predicting probability of Dx_test, Dx_train
    test_score = model.predict_proba(vectorizer.transform(Dx_test))
    train_score = model.predict_proba(vectorizer.transform(Dx_train))

    #Finding out the ROC_AUC_SCORE
    train_roc_auc_score = roc_auc_score(np.array(Dy_train), train_score[:, 1])
    print('Area Under the Curve for Train : ', train_roc_auc_score)
    test_roc_auc_score = roc_auc_score(np.array(Dy_test), test_score[:, 1])
    print('Area Under the Curve for Test : ', test_roc_auc_score)

    #Plotting with matplotlib.pyplot
    #ROC Curve for D-train
    train_fpr, train_tpr, thresholds = roc_curve(np.array(Dy_train), train_score[:, 1])
    plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))

    # ROC Curve for D-test
    test_fpr, test_tpr, thresholds = roc_curve(np.array(Dy_test), test_score[:, 1])
    plt.plot(test_fpr, test_tpr, label="train AUC =" + str(auc(test_fpr, test_tpr)))

    plt.legend()
    plt.xlabel("FPR : False Positive Ratio")
    plt.ylabel("TPF : True Positive Ratio")
    plt.title("Area Under ROC Curve")
    plt.show()

    plot_confusion_matrix(model, vectorizer.transform(Dx_train), np.array(Dy_train), 'Training')
    plot_confusion_matrix(model, vectorizer.transform(Dx_test), np.array(Dy_test), 'Testing')
    return train_roc_auc_score, test_roc_auc_score

```

[5.1] Applying Decision Trees on BOW, SET 1

In [45]:


```
# Please write all the code with proper documentation
csv_path = 'saved_models/Assignment8/BOW_dtree_results.csv'
cv_results = perform_hyperparameter_tuning(X=Dx_train, Y=Dy_train, vectorizer='BOW',
                                          results_path=csv_path, retrain=False, W2V_model=None)

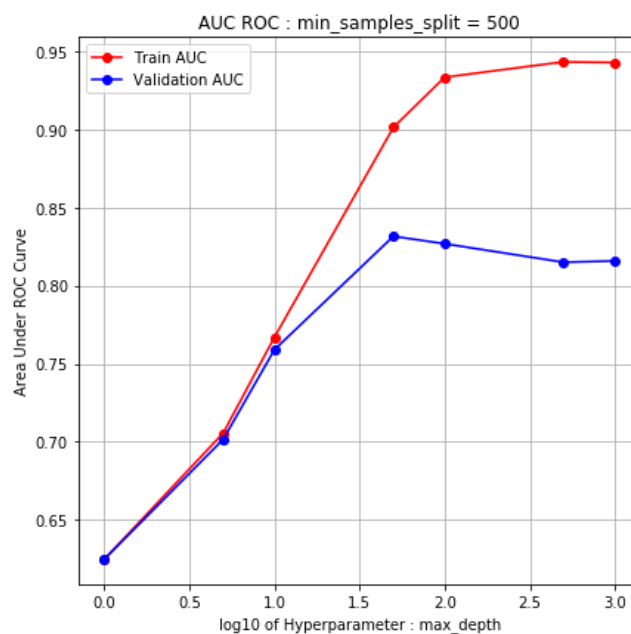
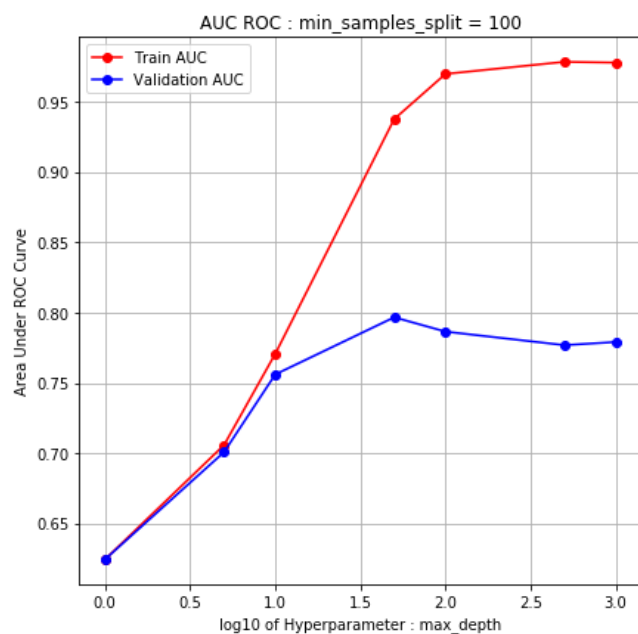
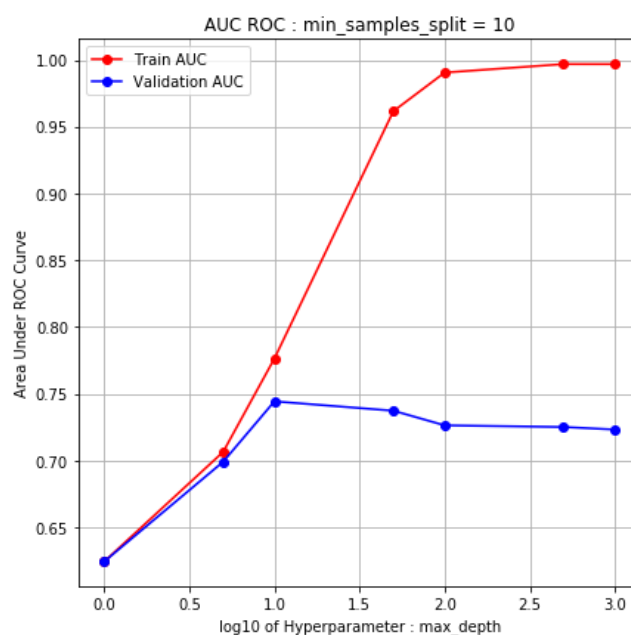
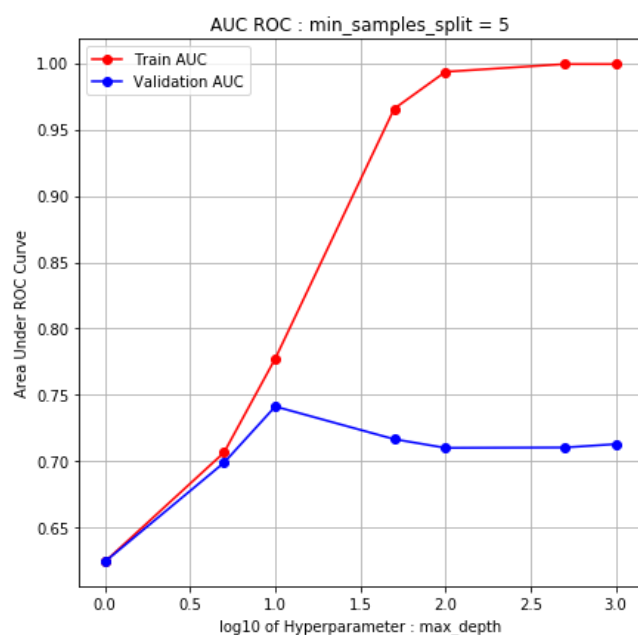
# Analysing best parameters
best_parameters = analyse_results(cv_results)
pprint.pprint(best_parameters)

# retraining the model with best parameters
model_path = 'saved_models/Assignment8/{0}_dtree.pkl'.format('BOW')
clf = retrain_with_best_params(Dx_train, Dy_train, best_parameters, 'BOW', model_path, None)

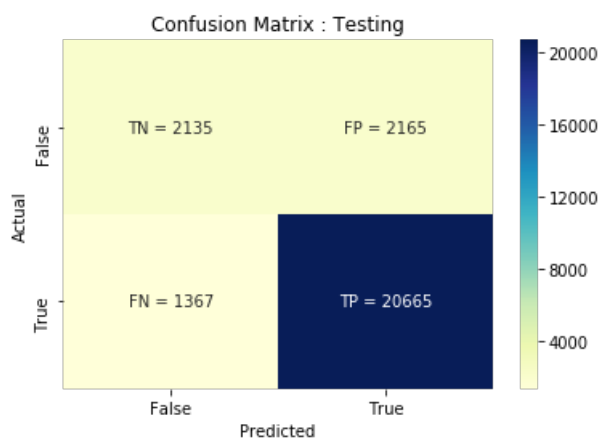
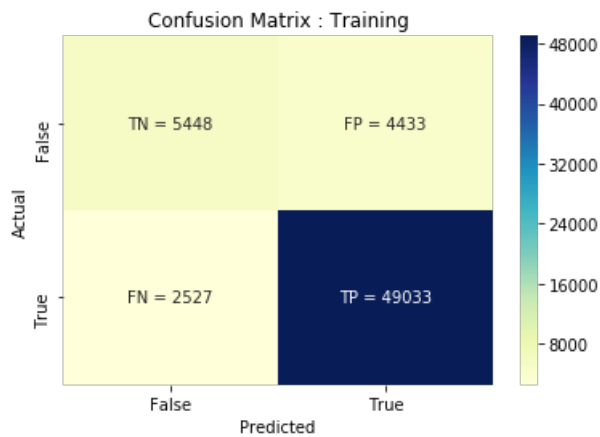
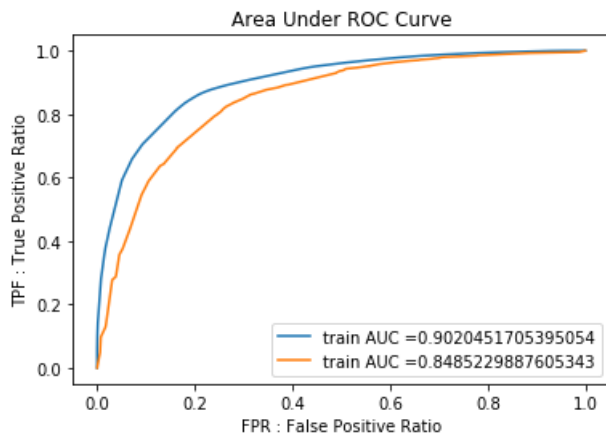
print('Retraining Vectorizer with Dx_train')
vectorizer_obj = get_vectorizer(W2V_model = None, train=Dx_train, vectorizer='BOW')

# plotting AUC ROC
train_score, test_score = plot_AUC_ROC(clf, vectorizer_obj, Dx_train, Dx_test, Dy_train, Dy_test)

# appending the data results
prettytable_data.append(['BOW', 'Decision Trees', best_parameters['max_depth'], best_parameters['min_samples_split'], train_score, test_score])
```



```
{'max_depth': 50, 'min_samples_split': 500}
Loading Model....
Retraining Vectorizer with Dx_train
Area Under the Curve for Train : 0.9020451705395054
Area Under the Curve for Test : 0.8485229887605343
```



[5.1.1] Top 20 important features from SET 1

In [46]:

```
clf1 = DecisionTreeClassifier(random_state=0)
vectorizer_obj1 = get_vectorizer(W2V_model = None, train=Dx_train, vectorizer='BOW')
clf1.fit(vectorizer_obj1.transform(Dx_train), np.array(Dy_train))
```

Out[46]:

```
DecisionTreeClassifier(class_weight=None, 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=0,
splitter='best')
```

In [47]:

```
feature_importance = clf1.feature_importances_
features = vectorizer_obj1.get_feature_names()
features_with_names = [(features[i], feature_importance[i]) for i in range(feature_importance.shape
[0]) if feature_importance[i]>0]
features_with_names.sort(key=lambda x: x[1], reverse=True)
features_with_names[:20]
```

Out[47]:

```
[('not', 0.03789421537218232),
 ('great', 0.024372247031116144),
 ('disappointed', 0.01590932693459837),
 ('not buy', 0.015736499956928187),
 ('money', 0.013322794154553919),
 ('worst', 0.0119890997033828),
 ('horrible', 0.011930031773261929),
 ('good', 0.01027286132253695),
 ('return', 0.009826070015016357),
 ('best', 0.009476224457512164),
 ('delicious', 0.009163800137051908),
 ('not disappointed', 0.00906987831120911),
 ('not recommend', 0.008992242315023737),
 ('awful', 0.008884596965140311),
 ('love', 0.00831565566432052),
 ('not worth', 0.008000935963592432),
 ('not good', 0.006321453599875989),
 ('loves', 0.005971119116773436),
 ('bad', 0.00585598714676059),
 ('terrible', 0.005721904233097538)]
```

In [48]:

```
del clf1
del vectorizer_obj1
del feature_importance
del features
del features_with_names
```

[5.1.2] Graphviz visualization of Decision Tree on BOW, SET 2

In [50]:

```
from sklearn import tree
from graphviz import Source
```

In [51]:

```
dot_data = tree.export_graphviz(clf, filled=True)
graph = Source(dot_data)
graph
```

Out[51]:



[5.2] Applying Decision Trees on TFIDF, SET 2

In [52]:

```
# Please write all the code with proper documentation
csv_path = 'saved_models/Assignment8/TFIDF_dtree_results.csv'
cv_results = perform_hyperparameter_tuning(X=Dx_train, Y=Dy_train, vectorizer='TFIDF',
                                           results_path=csv_path, retrain=False, W2V_model=None)

# Analysing best parameters
best_parameters = analyse_results(cv_results)
pprint.pprint(best_parameters)

# retraining the model with best parameters
model_path = 'saved_models/Assignment8/{0}_dtree.pkl'.format('TFIDF')
clf = retrain_with_best_params(Dx_train, Dy_train, best_parameters, 'TFIDF', model_path, None)
```

```

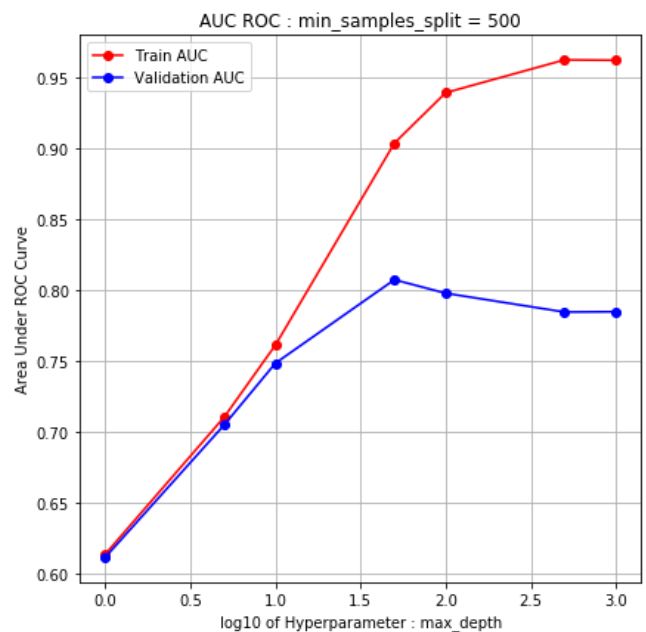
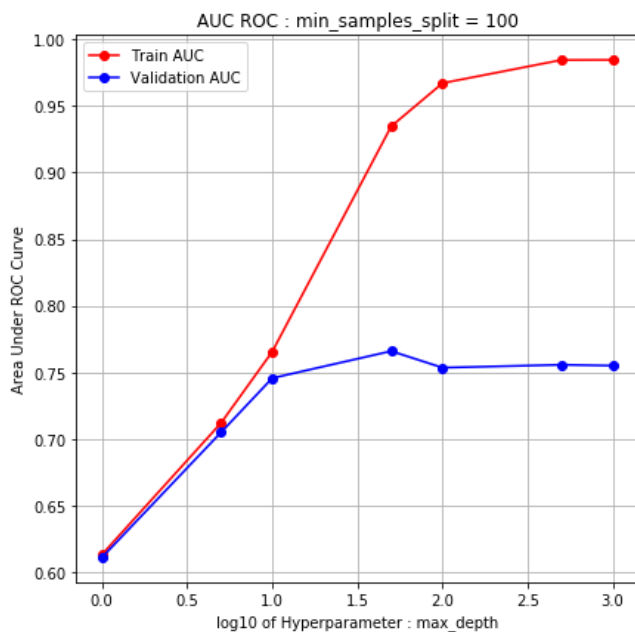
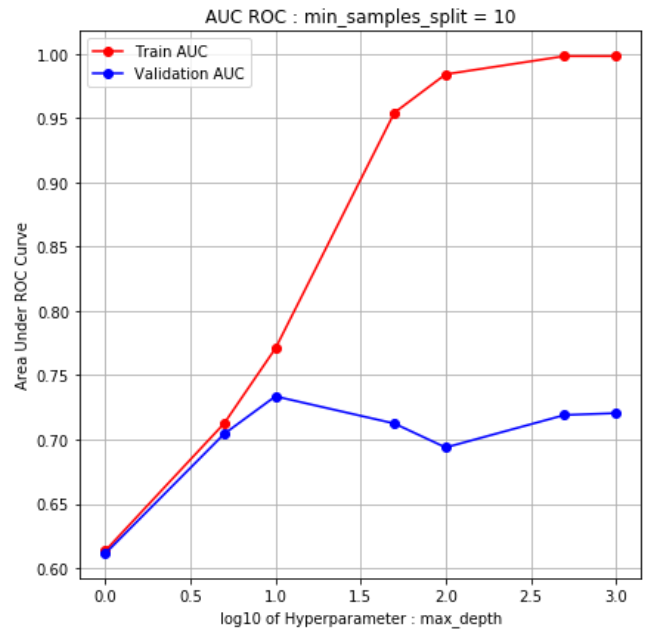
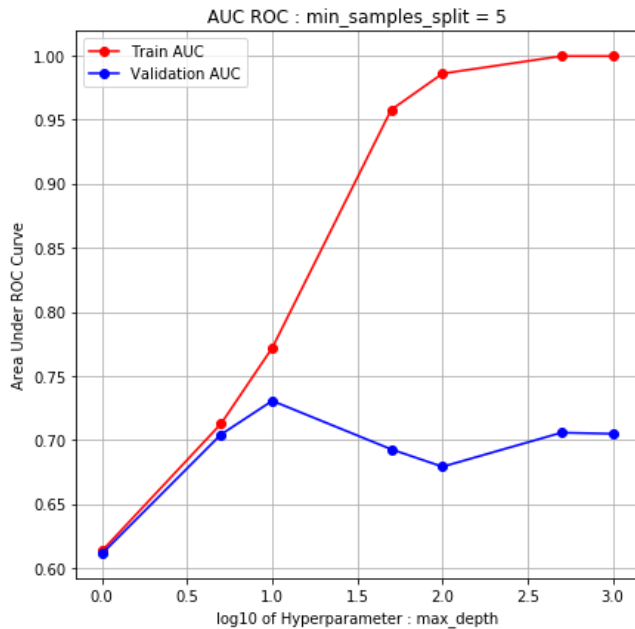
train_with_best_params(Dx_train, Dy_train, best_parameters, 'TFIDF', model_path, None,
                        'TFIDF')

print('Retraining Vectorizer with Dx_train')
vectorizer_obj = get_vectorizer(W2V_model = None, train=Dx_train, vectorizer='TFIDF')

# plotting AUC ROC
train_score, test_score = plot_AUC_ROC(clf, vectorizer_obj, Dx_train, Dx_test, Dy_train, Dy_test)

# appending the data results
prettytable_data.append(['TFIDF', 'Decision Trees', best_parameters['max_depth'], best_parameters[
'min_samples_split'], train_score, test_score])

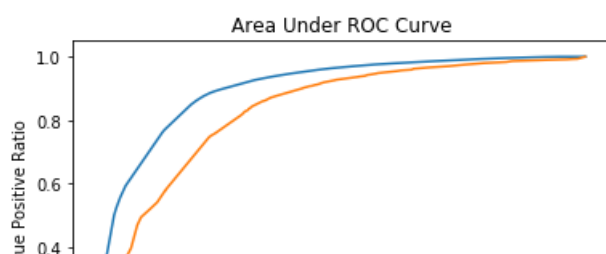
```

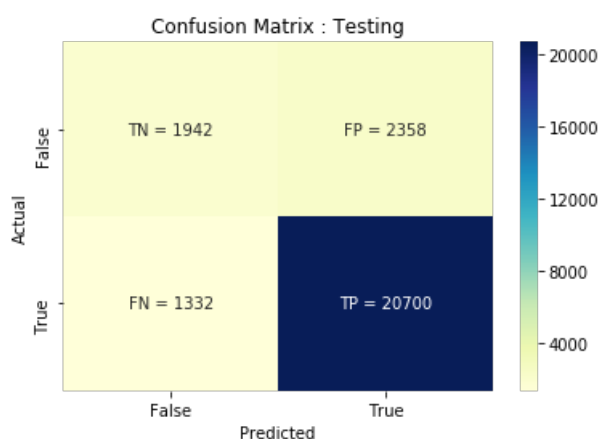
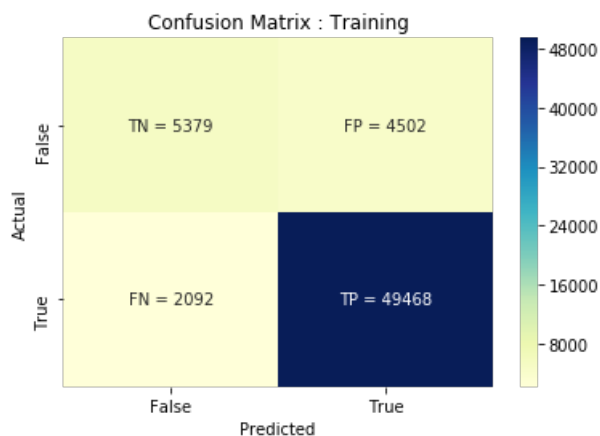
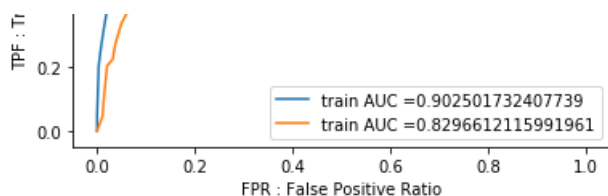


```

{'max_depth': 50, 'min_samples_split': 500}
Loading Model....
Retraining Vectorizer with Dx_train
Area Under the Curve for Train : 0.902501732407739
Area Under the Curve for Test : 0.8296612115991961

```





[5.2.1] Top 20 important features from SET 2

In [53]:

```
clf1 = DecisionTreeClassifier(random_state=0)
vectorizer_obj1 = get_vectorizer(W2V_model = None, train=Dx_train, vectorizer='TFIDF')
clf1.fit(vectorizer_obj1.transform(Dx_train), np.array(Dy_train))
```

Out[53]:

```
DecisionTreeClassifier(class_weight=None, 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=0,
splitter='best')
```

In [54]:

```
feature_importance = clf1.feature_importances_
features = vectorizer_obj1.get_feature_names()
features_with_names = [(features[i], feature_importance[i]) for i in range(feature_importance.shape
[0]) if feature_importance[i]>0]
features_with_names.sort(key=lambda x: x[1], reverse=True)
features_with_names[:20]
```

Out[54]:

```
[('not', 0.052857740177687275),
 ('', 0.004642101510275025)]
```

```
(
    'great', 0.02464312151077595),
    ('disappointed', 0.018193155360967064),
    ('worst', 0.014167732198852845),
    ('money', 0.013095772156989174),
    ('not buy', 0.012989791164653671),
    ('awful', 0.01196214528199168),
    ('return', 0.011809157283450666),
    ('good', 0.011741300079850681),
    ('horrible', 0.011704491465894641),
    ('love', 0.010272241989765755),
    ('best', 0.0100960915665325),
    ('delicious', 0.00884606501470341),
    ('bad', 0.008365333667113348),
    ('not worth', 0.008049574196986922),
    ('waste money', 0.008007015390045853),
    ('like', 0.007860762110616981),
    ('not recommend', 0.007835141539032498),
    ('not disappointed', 0.007081061149136909),
    ('taste', 0.006445505143929829)]
```

In [55]:

```
del clf1
del vectorizer_obj1
del feature_importance
del features
del features_with_names
```

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

In [56]:

```
from sklearn import tree
from graphviz import Source
```

In [57]:

```
dot_data = tree.export_graphviz(clf, filled=True)
graph = Source(dot_data)
graph
```

Out[57]:

Preparing/Training Google Word2Vec

In [58]:

```
is_your_ram_gt_16g=True
want_to_use_google_w2v = False
want_to_train_w2v = True

path_to_word2vec = '/home/monodeepdas112/Datasets/GoogleNews-vectors-negative300.bin.gz'

if want_to_train_w2v:

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

    # min_count = 5 considers only words that occurred atleast 5 times
    w2v_model=Word2Vec(list_of_sentences,min_count=5,size=300,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(path_to_word2vec):
    print('Preparing to load pre-trained Word2Vec model !')
    w2v_model=KeyedVectors.load_word2vec_format(path_to_word2vec, binary=True, class_weight = '
balanced')
    print('Successfully loaded model into memory !!!')
    print('Words similar to "similar" : ', w2v_model.wv.most_similar('great'))
    print('Words similar to "worst" : ',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 ")

```

```

[('fantastic', 0.7636513710021973), ('excellent', 0.7470881938934326), ('awesome',
0.7396433353424072), ('terrific', 0.7010406255722046), ('wonderful', 0.6899394989013672), ('good',
0.665637731552124), ('amazing', 0.6234958171844482), ('fabulous', 0.6180654764175415), ('perfect',
0.5894500613212585), ('nice', 0.5754525661468506)]
=====
[('greatest', 0.7833311557769775), ('tastiest', 0.6857161521911621), ('best', 0.6853801012039185),
('nastiest', 0.6730923652648926), ('disgusting', 0.6282191276550293), ('awful',
0.6015647649765015), ('smoothest', 0.592754602432251), ('weakest', 0.5881092548370361),
('horrible', 0.5868627429008484), ('coolest', 0.5763354897499084)]

```

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

In [59]:

```

# Please write all the code with proper documentation
csv_path = 'saved_models/Assignment8/Avg-W2Vec_dtree_results.csv'
cv_results = perform_hyperparameter_tuning(X=Dx_train, Y=Dy_train, vectorizer='Avg-W2Vec',
                                           results_path=csv_path, retrain=False, W2V_model=w2v_model)

# Analysing best parameters
best_parameters = analyse_results(cv_results)
pprint.pprint(best_parameters)

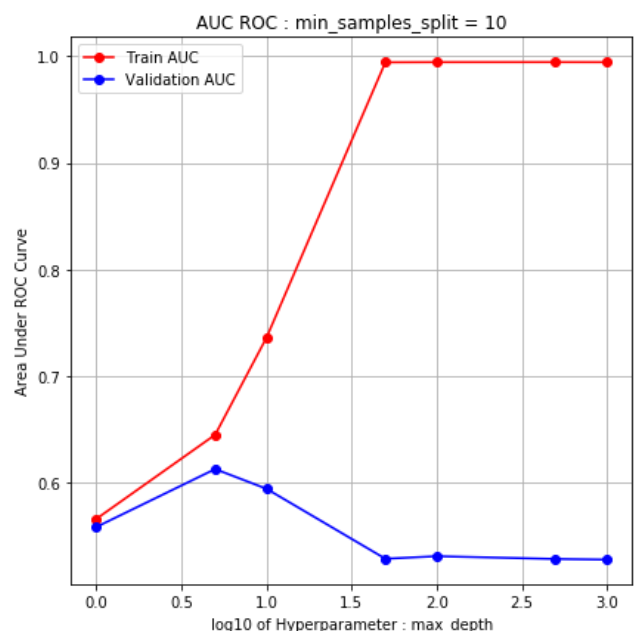
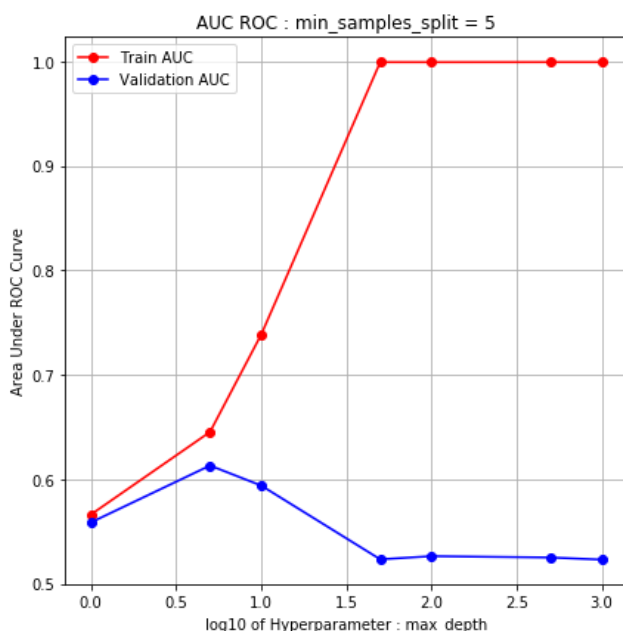
# retraining the model with best parameters
model_path = 'saved_models/Assignment8/{0}_dtree.pkl'.format('Avg-W2Vec')
clf = retrain_with_best_params(Dx_train, Dy_train, best_parameters, 'Avg-W2Vec', model_path, w2v_model)

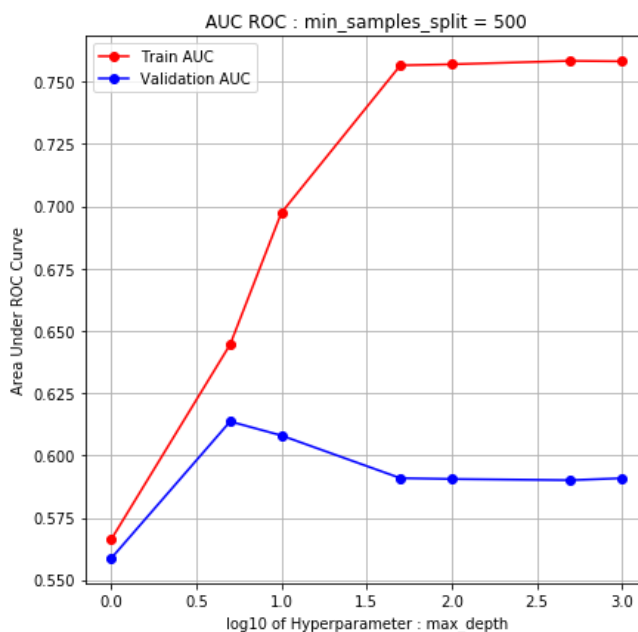
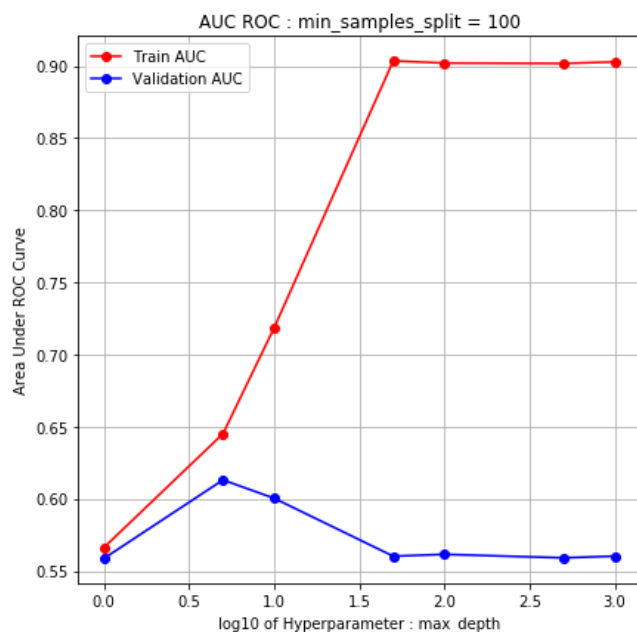
print('Retraining Vectorizer with Dx_train')
vectorizer_obj = get_vectorizer(W2V_model = w2v_model, train=Dx_train, vectorizer='Avg-W2Vec')

# plotting AUC ROC
train_score, test_score = plot_AUC_ROC(clf, vectorizer_obj, Dx_train, Dx_test, Dy_train, Dy_test)

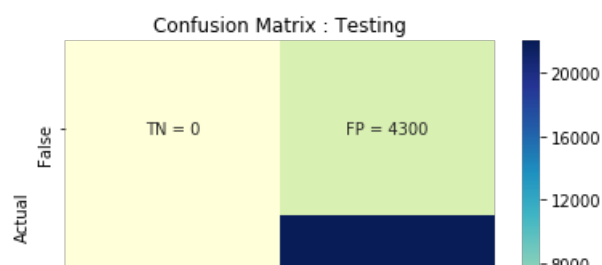
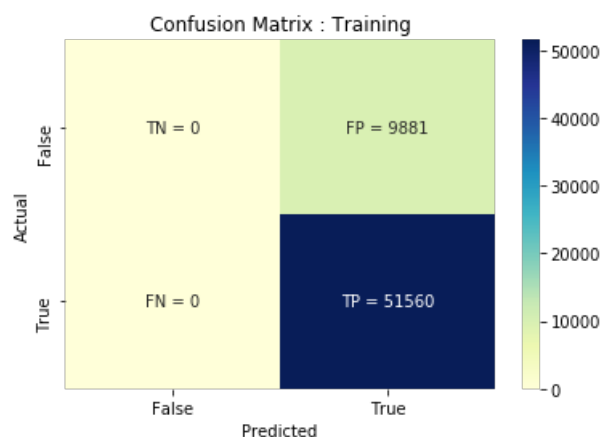
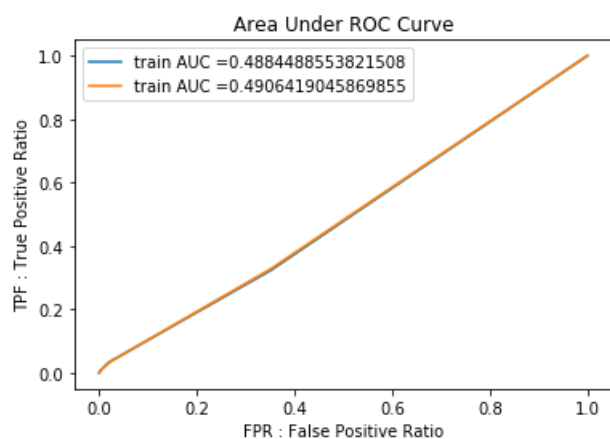
# appending the data results
prettytable_data.append(['Avg-W2Vec', 'Decision Trees', best_parameters['max_depth'],
best_parameters['min_samples_split'], train_score, test_score])

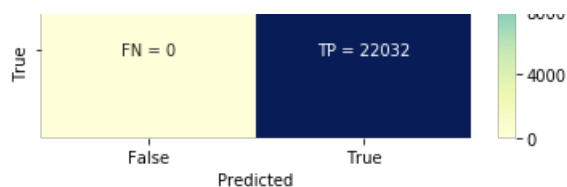
```





```
{'max_depth': 5, 'min_samples_split': 500}
Loading Model....
Retraining Vectorizer with Dx_train
Area Under the Curve for Train : 0.4884488553821508
Area Under the Curve for Test : 0.4906419045869855
```





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

In [60]:

```
# Please write all the code with proper documentation
csv_path = 'saved_models/Assignment8/TFIDF-W2Vec_dtree_results.csv'
cv_results = perform_hyperparameter_tuning(X=Dx_train, Y=Dy_train, vectorizer='TFIDF-W2Vec',
                                           results_path=csv_path, retrain=False, W2V_model=w2v_model)

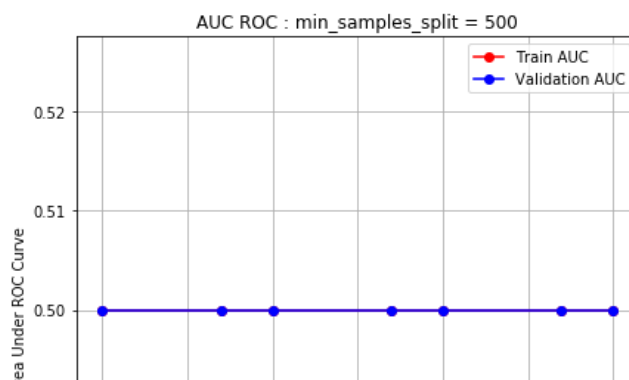
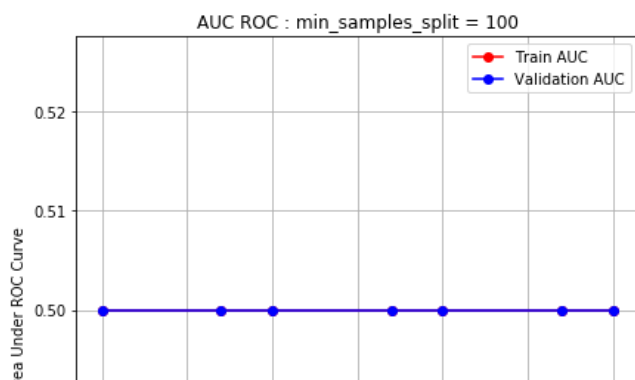
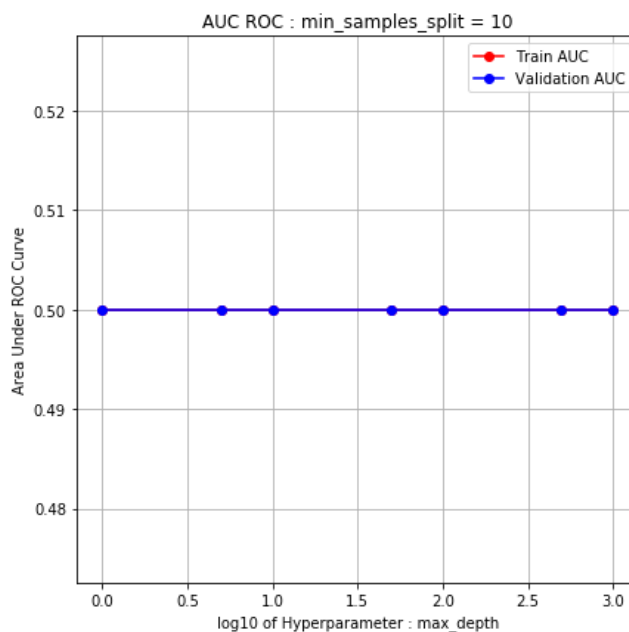
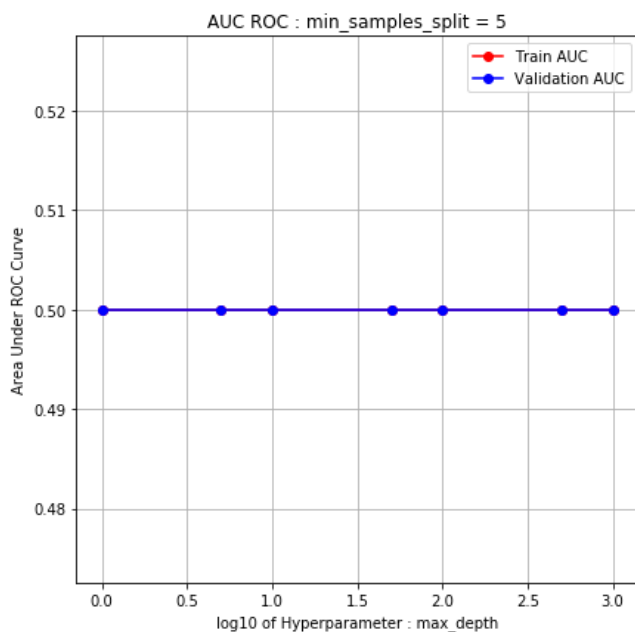
# Analysing best parameters
best_parameters = analyse_results(cv_results)
pprint.pprint(best_parameters)

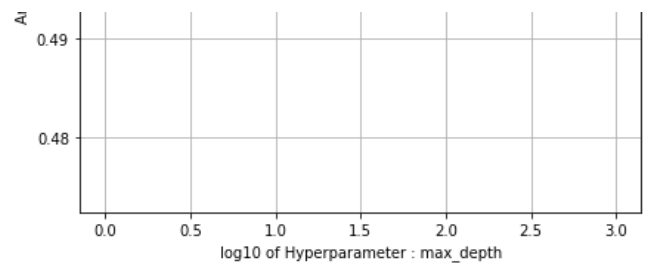
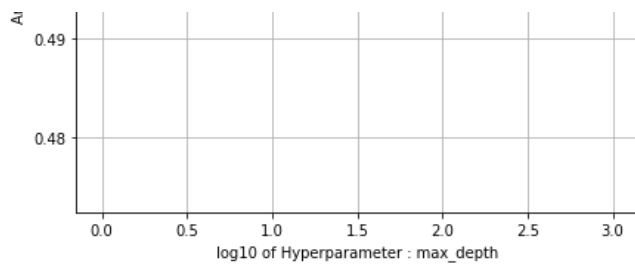
# retraining the model with best parameters
model_path = 'saved_models/Assignment8/{0}_dtree.pkl'.format('TFIDF-W2Vec')
clf = retrain_with_best_params(Dx_train, Dy_train, best_parameters, 'TFIDF-W2Vec', model_path, w2v_model)

print('Retraining Vectorizer with Dx_train')
vectorizer_obj = get_vectorizer(W2V_model = w2v_model, train=Dx_train, vectorizer='TFIDF-W2Vec')

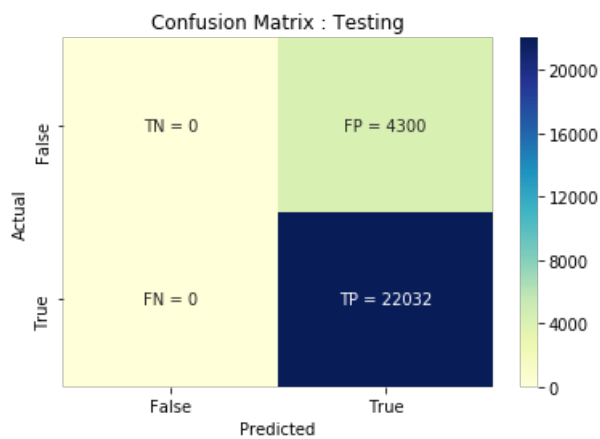
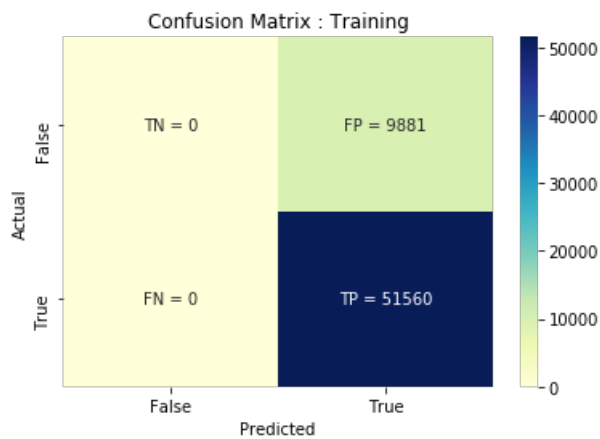
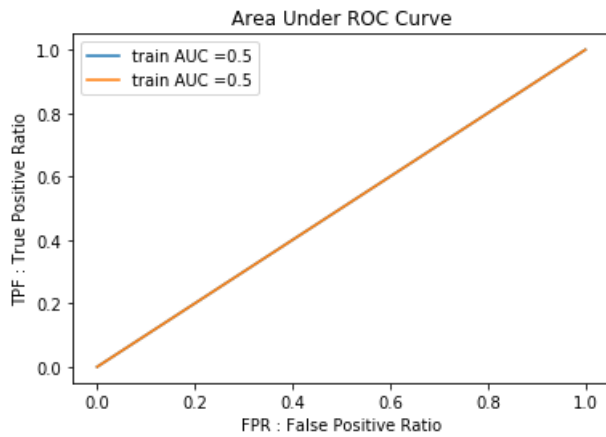
# plotting AUC ROC
train_score, test_score = plot_AUC_ROC(clf, vectorizer_obj, Dx_train, Dx_test, Dy_train, Dy_test)

# appending the data results
prettytable_data.append(['TFIDF-W2Vec', 'Decision Trees', best_parameters['max_depth'],
                        best_parameters['min_samples_split'], train_score, test_score])
```





```
{'max_depth': 1, 'min_samples_split': 5}
Loading Model....
Retraining Vectorizer with Dx_train
Area Under the Curve for Train : 0.5
Area Under the Curve for Test : 0.5
```



[6] Conclusions

In [61]:

```
from prettytable import PrettyTable
```

In [63]:

```
# Please compare all your models using Prettytable library
```

```
x = PrettyTable()
```

```
x.field_names = ["Vectorizer", "Model", "max_depth", "min_samples_split", "Train AUC", "Test AUC"]  
[x.add_row(i) for i in prettytable_data]  
print(x)
```

```
+-----+-----+-----+-----+-----+-----+  
---+  
| Vectorizer | Model | max_depth | min_samples_split | Train AUC | Test AUC |  
|-----+-----+-----+-----+-----+-----+  
---+  
| BOW | Decision Trees | 50 | 500 | 0.9020451705395054 |  
0.8485229887605343 |  
| TFIDF | Decision Trees | 50 | 500 | 0.902501732407739 |  
0.8296612115991961 |  
| Avg-W2Vec | Decision Trees | 5 | 500 | 0.4884488553821508 |  
0.4906419045869855 |  
| TFIDF-W2Vec | Decision Trees | 1 | 5 | 0.5 | 0.5 |  
|-----+-----+-----+-----+-----+-----+  
---+  
◀────────────────────────────────────────────────────────────────────────────────▶▶
```

[Save to PDF](#)