

Amazon Fine Food Reviews Analysis

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

EDA: <https://nycdatasience.com/blog/student-works/amazon-fine-foods-visualization/>
(<https://nycdatasience.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]: 1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5
6 import sqlite3
7 import pandas as pd
8 import numpy as np
9 import nltk
10 import string
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 from sklearn.feature_extraction.text import TfidfTransformer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.metrics import confusion_matrix
18 from sklearn import metrics
19 from sklearn.metrics import roc_curve, auc
20 from nltk.stem.porter import PorterStemmer
21
22 import re
23 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
24 import string
25 from nltk.corpus import stopwords
26 from nltk.stem import PorterStemmer
27 from nltk.stem.wordnet import WordNetLemmatizer
28
29 from gensim.models import Word2Vec
30 from gensim.models import KeyedVectors
31 import pickle
32
33 from tqdm import tqdm
34 import os
```

```
C:\Users\sujpanda\Anaconda3\lib\site-packages\gensim\utils.py:1212: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

```

In [2]: 1 # using SQLite Table to read data.
2 con = sqlite3.connect('C:\\Users\\sujpanda\\Desktop\\applied\\database.sqlite')
3
4 # filtering only positive and negative reviews i.e.
5 # not taking into consideration those reviews with Score=3
6 # SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 500000 d
7 # you can change the number to any other number based on your computing power
8
9 # filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score !=
10 # for tsne assignment you can take 5k data points
11
12 filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3
13
14 # Give reviews with Score>3 a positive rating(1), and reviews with a score<3
15 def partition(x):
16     if x < 3:
17         return 0
18     return 1
19
20 #changing reviews with score less than 3 to be positive and vice-versa
21 actualScore = filtered_data['Score']
22 positiveNegative = actualScore.map(partition)
23 filtered_data['Score'] = positiveNegative
24 print("Number of data points in our data", filtered_data.shape)
25 filtered_data.head(3)

```

Number of data points in our data (100000, 10)

Out[2]:

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

```

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

```

```
In [4]: 1 print(display.shape)
        2 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-R11D9D7SHXIJ9	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]: 1 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]: 1 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]: 1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5 ORDER BY ProductID
6 """, con)
7 display.head()
```

Out[7]:

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

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

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

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

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

The method used for the same was that we first sort the data according to ProductId and then just keep the first similar product review and delete the others. for eg. in the above just the review for ProductId=B000HDL1RQ remains. This method ensures that there is only one representative for

each product and deduplication without sorting would lead to possibility of different representatives still existing for the same product.

```
In [8]: 1 #Sorting data according to ProductId in ascending order
        2 sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, in
```

```
In [9]: 1 #Deduplication of entries
        2 final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"
        3 final.shape
```

Out[9]: (87775, 10)

```
In [10]: 1 #Checking to see how much % of data still remains
         2 (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]: 1 display= pd.read_sql_query("""
        2 SELECT *
        3 FROM Reviews
        4 WHERE Score != 3 AND Id=44737 OR Id=64422
        5 ORDER BY ProductID
        6 """, con)
        7
        8 display.head()
```

Out[11]:

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

```
In [12]: 1 final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [13]: 1 #Before starting the next phase of preprocessing Lets see the number of entri
2 print(final.shape)
3
4 #How many positive and negative reviews are present in our dataset?
5 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
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]:

```

1 # printing some random reviews
2 sent_0 = final['Text'].values[0]
3 print(sent_0)
4 print("="*50)
5
6 sent_1000 = final['Text'].values[1000]
7 print(sent_1000)
8 print("="*50)
9
10 sent_1500 = final['Text'].values[1500]
11 print(sent_1500)
12 print("="*50)
13
14 sent_4900 = final['Text'].values[4900]
15 print(sent_4900)
16 print("="*50)

```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its 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. You 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]:

```

1 # remove urls from text python: https://stackoverflow.com/a/40823105/4084039
2 sent_0 = re.sub(r"http\S+", "", sent_0)
3 sent_1000 = re.sub(r"http\S+", "", sent_1000)
4 sent_150 = re.sub(r"http\S+", "", sent_1500)
5 sent_4900 = re.sub(r"http\S+", "", sent_4900)
6
7 print(sent_0)

```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its 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]: 1 # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-re
2 from bs4 import BeautifulSoup
3
4 soup = BeautifulSoup(sent_0, 'lxml')
5 text = soup.get_text()
6 print(text)
7 print("="*50)
8
9 soup = BeautifulSoup(sent_1000, 'lxml')
10 text = soup.get_text()
11 print(text)
12 print("="*50)
13
14 soup = BeautifulSoup(sent_1500, 'lxml')
15 text = soup.get_text()
16 print(text)
17 print("="*50)
18
19 soup = BeautifulSoup(sent_4900, 'lxml')
20 text = soup.get_text()
21 print(text)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its 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. You 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]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\ 're", " are", phrase)
12    phrase = re.sub(r"\ 's", " is", phrase)
13    phrase = re.sub(r"\ 'd", " would", phrase)
14    phrase = re.sub(r"\ 'll", " will", phrase)
15    phrase = re.sub(r"\ 't", " not", phrase)
16    phrase = re.sub(r"\ 've", " have", phrase)
17    phrase = re.sub(r"\ 'm", " am", phrase)
18    return phrase
```

```
In [18]: 1 sent_1500 = decontracted(sent_1500)
2 print(sent_1500)
3 print("="*50)
```

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

```
In [19]: 1 #remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
2 sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
3 print(sent_0)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its 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]: 1 #remove spacial character: https://stackoverflow.com/a/5843547/4084039
2 sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
3 print(sent_1500)
```

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

```
In [21]: 1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'
3 # <br /><br /> ==> after the above steps, we are getting "br br"
4 # we are including them into stop words list
5 # instead of <br /> if we have <br/> these tags would have revmoved in the 1s
6
7 stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours',
8               "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
9               'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'i',
10              'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'is',
11              'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have',
12              'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'beca',
13              'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
14              'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on',
15              'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 't',
16              'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'th',
17              's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shoul",
18              've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
19              "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'm',
20              "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shou",
21              'won', "won't", 'wouldn', "wouldn't"])
```

```
In [22]: 1 # Combining all the above stundents
2 from tqdm import tqdm
3 preprocessed_reviews = []
4 # tqdm is for printing the status bar
5 for sentence in tqdm(final['Text'].values):
6     sentence = re.sub(r"http\S+", "", sentence)
7     sentence = BeautifulSoup(sentence, 'lxml').get_text()
8     sentence = decontracted(sentence)
9     sentence = re.sub("\S*\d\S*", "", sentence).strip()
10    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
11    # https://gist.github.com/sebleier/554280
12    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not
13    preprocessed_reviews.append(sentence.strip())
```

100%|██████████| 87773/87773 [01:02<00:00, 1396.85it/s]

```
In [23]: 1 preprocessed_reviews[1500]
```

Out[23]: 'way hot blood took bite jig lol'

[3.2] Preprocessing Review Summary

In [24]:

```

1 from tqdm import tqdm
2 preprocessed_summary = []
3 # tqdm is for printing the status bar
4 for sentence in tqdm(final['Summary'].values):
5     sentence = re.sub(r"http\S+", "", sentence)
6     sentence = BeautifulSoup(sentence, 'lxml').get_text()
7     sentence = decontracted(sentence)
8     sentence = re.sub("\S*\d\S*", "", sentence).strip()
9     sentence = re.sub('[^A-Za-z]+', ' ', sentence)
10    # https://gist.github.com/sebleier/554280
11    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not
12    preprocessed_summary.append(sentence.strip())

```

```

37%|██████████| 32715/87773 [00:11<00:22, 2405.44it/s]C:\Users\sujpanda\Anaconda3\lib\site-packages\bs4\__init__.py:219: UserWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.

```

```

' BeautifulSoup Soup.' % markup)

```

```

70%|██████████| 61326/87773 [00:22<00:09, 2647.88it/s]C:\Users\sujpanda\Anaconda3\lib\site-packages\bs4\__init__.py:219: UserWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.

```

```

' BeautifulSoup Soup.' % markup)

```

```

75%|██████████| 65472/87773 [00:23<00:07, 3079.02it/s]C:\Users\sujpanda\Anaconda3\lib\site-packages\bs4\__init__.py:219: UserWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.

```

```

' BeautifulSoup Soup.' % markup)

```

```

96%|██████████| 84147/87773 [00:30<00:01, 2489.15it/s]C:\Users\sujpanda\Anaconda3\lib\site-packages\bs4\__init__.py:219: UserWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.

```

```

' BeautifulSoup Soup.' % markup)

```

```

100%|██████████| 87773/87773 [00:32<00:00, 2730.40it/s]

```

[5] Assignment 8: Decision Trees

1. Apply Decision Trees on these feature sets

- **SET 1:** Review text, preprocessed one converted into vectors using (BOW)
- **SET 2:** Review text, preprocessed one converted into vectors using (TFIDF)
- **SET 3:** Review text, preprocessed one converted into vectors using (AVG W2v)
- **SET 4:** Review text, preprocessed one converted into vectors using (TFIDF W2v)

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

- Find the best hyper parameter which will give the maximum [AUC](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/>) 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 \(https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html) 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 \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](https://seaborn.pydata.org/generated/seaborn.heatmap.html).

 (<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

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 \(http://zetcode.com/python/prettytable/\)](http://zetcode.com/python/prettytable/)



Note: Data Leakage

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

Some utility functions

```

In [25]: 1  ## Some utility functions
2
3  def check_trade_off(X_train,X_test,y_train,y_test):
4
5      from sklearn.metrics import roc_curve
6      from sklearn.metrics import roc_auc_score
7
8      [{'max_depth': [1, 5, 10, 50, 100, 500, 1000]}]
9
10     C_range1 = ['1','5','10','50','100','500','1000']
11     C_range = [1, 5, 10, 50, 100, 500, 1000]
12     dummy_range = [1,2,3,4,5,6,7]
13
14     auc_scores = []
15     auc_train_scores = []
16
17     i = 0
18     for i in C_range:
19         clf =DecisionTreeClassifier(max_depth=i)
20
21         # fitting the model on crossvalidation train
22         clf.fit(X_train, y_train)
23
24
25         #evaluate AUC score.
26         probs = clf.predict_proba(X_test)
27         probs = probs[:, 1]
28         # calculate AUC
29         auc = roc_auc_score(y_test, probs)
30         print('AUC: %.3f' % auc)
31         auc_scores.append(auc)
32
33     print('#####')
34     print('AUC from train data #####')
35     i = 0
36     for i in C_range:
37         clf =DecisionTreeClassifier(max_depth=i)
38
39         # fitting the model on crossvalidation train
40         clf.fit(X_train, y_train)
41
42         #evaluate AUC score.
43         probs = clf.predict_proba(X_train)
44         probs = probs[:, 1]
45         # calculate AUC
46         auc = roc_auc_score(y_train, probs)
47         print('AUC: %.3f' % auc)
48         auc_train_scores.append(auc)
49
50     plt.plot(dummy_range, auc_scores,'r')
51     plt.plot(dummy_range, auc_train_scores,'b')
52     plt.xticks(dummy_range, C_range1, rotation='vertical')
53     for xy in zip(dummy_range, auc_scores):
54         plt.annotate('(%.3f, %.3f)' % xy, xy=xy, textcoords='data')
55     for xy in zip(dummy_range, auc_train_scores):
56         plt.annotate('(%.3f, %.3f)' % xy, xy=xy, textcoords='data')

```

```
57  
58  
59     plt.xlabel('max_depth')  
60     plt.ylabel('auc_scores')  
61     plt.show()  
62
```



```

In [26]: 1 def dt_results(maxDepth,minimumSampleSplit,X_train,X_test,y_train,y_test):
2         # roc curve and auc
3         from sklearn.metrics import roc_curve
4         from sklearn.metrics import roc_auc_score
5         from matplotlib import pyplot
6         # ===== Decision Tree=====
7         clf = DecisionTreeClassifier(max_depth=maxDepth,min_samples_split=minimum
8
9         # fitting the model
10        clf.fit(X_train, y_train)
11
12        # predict the response
13        pred = clf.predict(X_test)
14
15        # evaluate accuracy
16        acc = accuracy_score(y_test, pred) * 100
17        print('\nThe accuracy of the DT classifier for maxDepth = %d and min spli
18
19        probs = clf.predict_proba(X_test)
20        probs = probs[:, 1]
21        # calculate AUC
22        auc = roc_auc_score(y_test, probs)
23        print('AUC: %.3f' % auc)
24        # calculate roc curve
25        fpr, tpr, thresholds = roc_curve(y_test, probs)
26        #####Train data#####
27
28        clf.fit(X_train, y_train)
29        pred_train = clf.predict(X_train)
30        probs = clf.predict_proba(X_train)
31        probs = probs[:, 1]
32        # calculate AUC
33        auc = roc_auc_score(y_train, probs)
34        print('AUC: %.3f' % auc)
35        fpr1, tpr1, thresholds1 = roc_curve(y_train, probs)
36
37
38        # plot no skill
39        pyplot.plot([0, 1], [0, 1], linestyle='--')
40        # plot the roc curve for the model
41        pyplot.plot(fpr, tpr, marker='.',label='test')
42        pyplot.plot(fpr1, tpr1, marker='*',label='train')
43        pyplot.legend()
44        # show the plot
45        pyplot.show()
46        from sklearn.metrics import confusion_matrix
47        con_mat = confusion_matrix(y_test, pred, [0, 1])
48        con_mat_train = confusion_matrix(y_train,pred_train,[0,1])
49        return con_mat,con_mat_train,clf

```

```
In [27]: 1 def showHeatMap(con_mat):
2         class_label = ["negative", "positive"]
3         df_cm = pd.DataFrame(con_mat, index = class_label, columns = class_label)
4         sns.heatmap(df_cm, annot = True, fmt = "d")
5         plt.title("Confusion Matrix")
6         plt.xlabel("Predicted Label")
7         plt.ylabel("True Label")
8         plt.show()
```

Applying Decision Trees

[5.1] Applying Decision Trees on BOW, SET 1

```
In [28]: 1 from sklearn.cross_validation import train_test_split
2         from sklearn.tree import DecisionTreeClassifier
3         from sklearn.metrics import accuracy_score
4         from sklearn.cross_validation import cross_val_score
5         from collections import Counter
6         from sklearn.metrics import accuracy_score
7         from sklearn import cross_validation
8         from sklearn.grid_search import GridSearchCV
9         import warnings
10        warnings.filterwarnings("ignore")
```

C:\Users\sujpanda\Anaconda3\lib\site-packages\sklearn\cross_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.

"This module will be removed in 0.20.", DeprecationWarning)

C:\Users\sujpanda\Anaconda3\lib\site-packages\sklearn\grid_search.py:42: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. This module will be removed in 0.20.

DeprecationWarning)

```
In [143]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev
```

```

In [144]: 1 count_vect = CountVectorizer()
          2 final_counts = count_vect.fit_transform(X_1)
          3 final_test_count = count_vect.transform(X_test)
          4
          5 # split the train data set into cross validation train and cross validation t
          6 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_siz
          7
          8 final_counts_tr_cv = count_vect.transform(X_tr)
          9 final_test_count_cv = count_vect.transform(X_cv)
         10
         11 tuned_parameters = [{'max_depth': [1, 5, 10, 50, 100, 500, 1000]}, 'min_samples
         12
         13 #Using GridSearchCV
         14 model = GridSearchCV(DecisionTreeClassifier(), tuned_parameters, scoring = 'r
         15 model.fit(final_counts_tr_cv, y_tr)
         16
         17 print(model.best_estimator_)
         18 print(model.score(final_test_count_cv, y_cv))
         19
         20 check_trade_off(final_counts_tr_cv, final_test_count_cv, y_tr, y_cv)

```

```

DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=50,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=500,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')

```

```
0.8217225902108642
```

```
AUC: 0.617
```

```
AUC: 0.688
```

```
AUC: 0.751
```

```
AUC: 0.675
```

```
AUC: 0.673
```

```
AUC: 0.679
```

```
AUC: 0.680
```

```
#####
```

```
AUC from train data #####
```

```
AUC: 0.627
```

```
AUC: 0.706
```

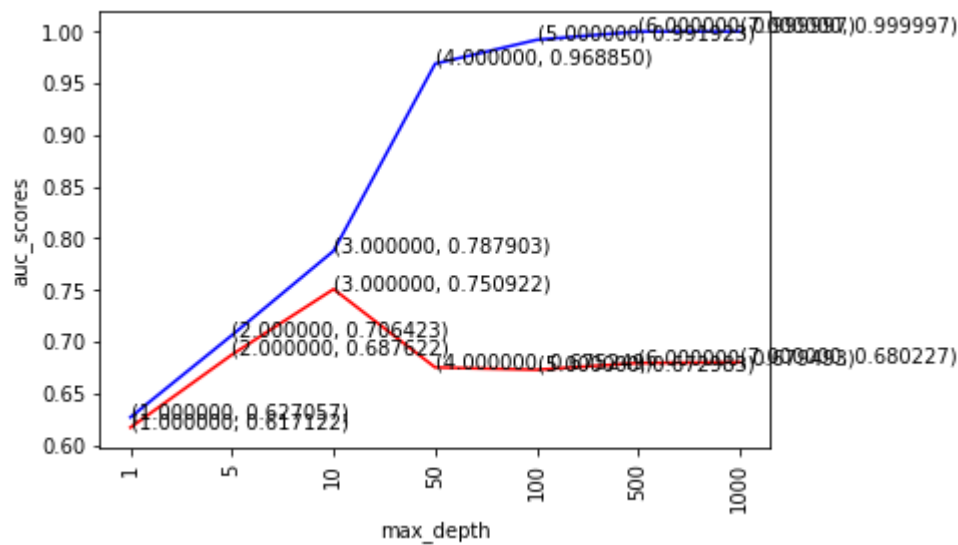
```
AUC: 0.788
```

```
AUC: 0.969
```

```
AUC: 0.992
```

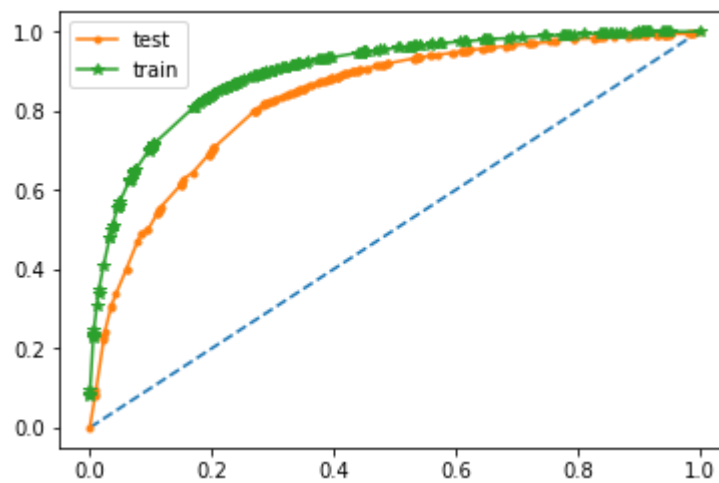
```
AUC: 1.000
```

```
AUC: 1.000
```



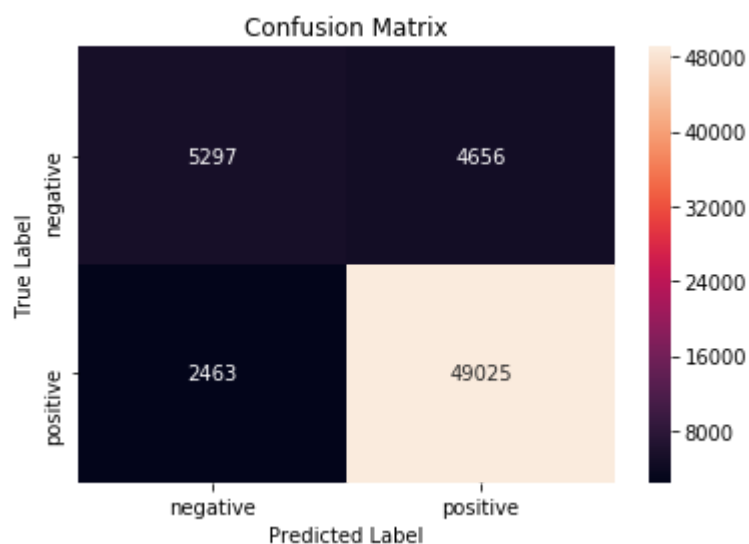
```
In [145]: 1 con_mat,con_mat_train,clf = dt_results(50,500,final_counts,final_test_count,y
```

The accuracy of the DT classifier for maxDepth = 50 and min split = 500 is 85.872702%
 AUC: 0.832
 AUC: 0.898



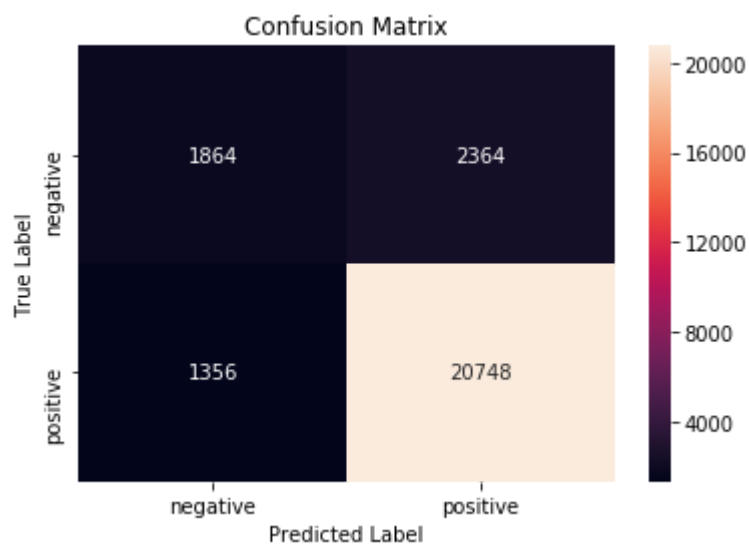
Observation: My model predicted with accuracy 85% with AUC: 0.833. There is no much difference between train and test ROC curve

```
In [146]: 1 showHeatMap(con_mat_train)
```



Observation: My model would have predicted 2463 + 4656 points wrongly even with train data

```
In [147]: 1 showHeatMap(con_mat)
```



Observation: My model misclassified 1356 + 2364 points

[5.1.1] Top 20 important features from SET 1

```
In [148]: 1 feature_names = np.array(count_vect.get_feature_names())
2 featureDict = dict(zip(feature_names, clf.feature_importances_))
3 sortedFeatures = sorted(featureDict.items(), key=lambda x: x[1], reverse=True)
4 print(type(sortedFeatures))
5 for i in range(0,20):
6     print(sortedFeatures[i])
```

```
<class 'list'>
('not', 0.10556541192374382)
('great', 0.06328011594676017)
('worst', 0.04198736350781279)
('disappointed', 0.041460848322431014)
('money', 0.03469326192544119)
('horrible', 0.02755088753016225)
('return', 0.02734130622754287)
('best', 0.02631809150225828)
('delicious', 0.02371161236410449)
('love', 0.02230452293131868)
('terrible', 0.021929296168877045)
('good', 0.019721769822422328)
('awful', 0.01744950993510839)
('waste', 0.01686007552218982)
('perfect', 0.015630476695974057)
('loves', 0.015271687168840937)
('disappointing', 0.015263760669904461)
('threw', 0.014663220813299996)
('nice', 0.01335400790815546)
('bad', 0.012334416821108925)
```

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

```
In [149]: 1 import graphviz
2 import pydotplus
3 import collections
4 from sklearn import tree
5 clf = DecisionTreeClassifier(max_depth=2)
6 clf.fit(final_counts, y_1)
7 dot_data = tree.export_graphviz(clf, feature_names = count_vect.get_feature_n
8 graph = pydotplus.graph_from_dot_data(dot_data)
9 colors = ('turquoise', 'orange')
10 edges = collections.defaultdict(list)
11 for edge in graph.get_edge_list():
12     edges[edge.get_source()].append(int(edge.get_destination()))
13
14 for edge in edges:
15     edges[edge].sort()
16     for i in range(2):
17         dest = graph.get_node(str(edges[edge][i]))[0]
18         dest.set_fillcolor(colors[i])
19
20 graph.write_png('tree.png')
```

Out[149]: True

[5.2] Applying Decision Trees on TFIDF, SET 2

```
In [30]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev
```

```

In [31]: 1 tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
2 tf_idf_vect.fit(X_1)
3 final_tf_idf = tf_idf_vect.transform(X_1)
4 final_test_count = tf_idf_vect.transform(X_test)
5
6 # split the train data set into cross validation train and cross validation test
7 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.2)
8
9 final_counts_tr_cv = tf_idf_vect.transform(X_tr)
10 final_test_count_cv = tf_idf_vect.transform(X_cv)
11
12 tuned_parameters = [{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples': [1, 5, 10, 50, 100, 500, 1000]}]
13
14 #Using GridSearchCV
15 model = GridSearchCV(DecisionTreeClassifier(), tuned_parameters, scoring = 'roc_auc')
16 model.fit(final_counts_tr_cv, y_tr)
17
18 print(model.best_estimator_)
19 print(model.score(final_test_count_cv, y_cv))
20
21 check_trade_off(final_counts_tr_cv, final_test_count_cv, y_tr, y_cv)

```

```

DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=50,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=500,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')

```

```
0.8175195450851837
```

```
AUC: 0.628
```

```
AUC: 0.703
```

```
AUC: 0.734
```

```
AUC: 0.695
```

```
AUC: 0.693
```

```
AUC: 0.703
```

```
AUC: 0.700
```

```
#####
```

```
AUC from train data #####
```

```
AUC: 0.626
```

```
AUC: 0.722
```

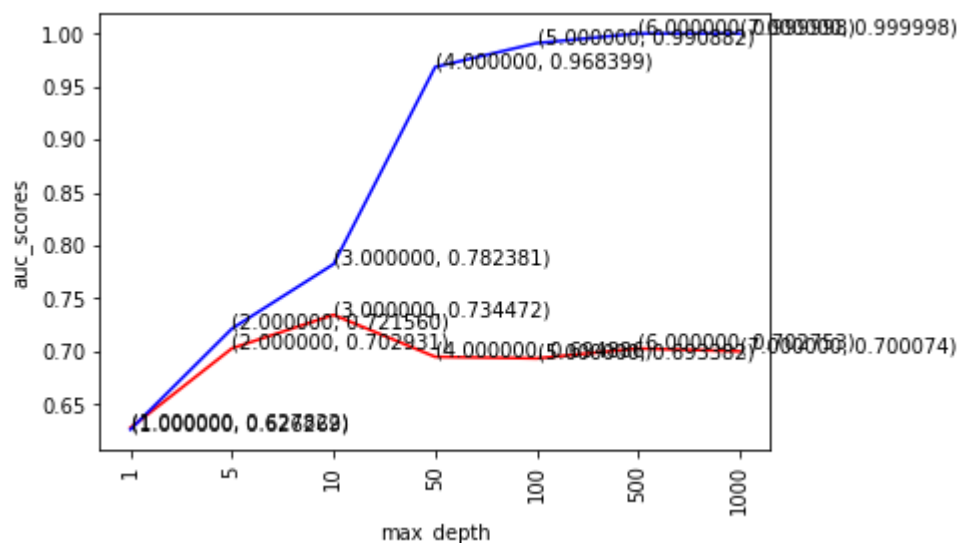
```
AUC: 0.782
```

```
AUC: 0.968
```

```
AUC: 0.991
```

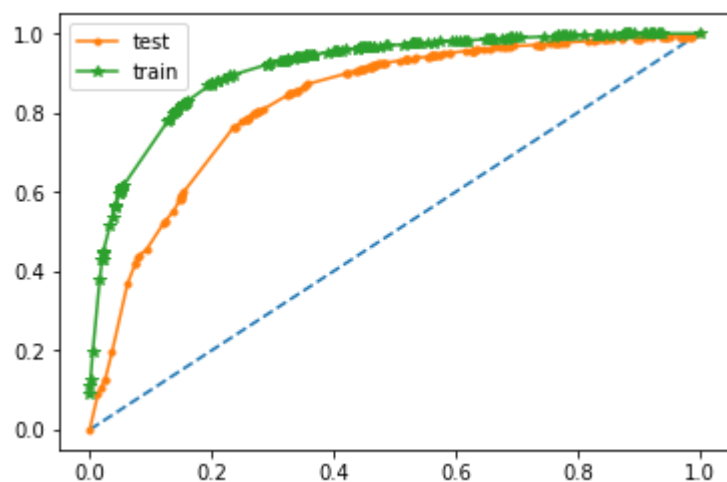
```
AUC: 1.000
```

```
AUC: 1.000
```

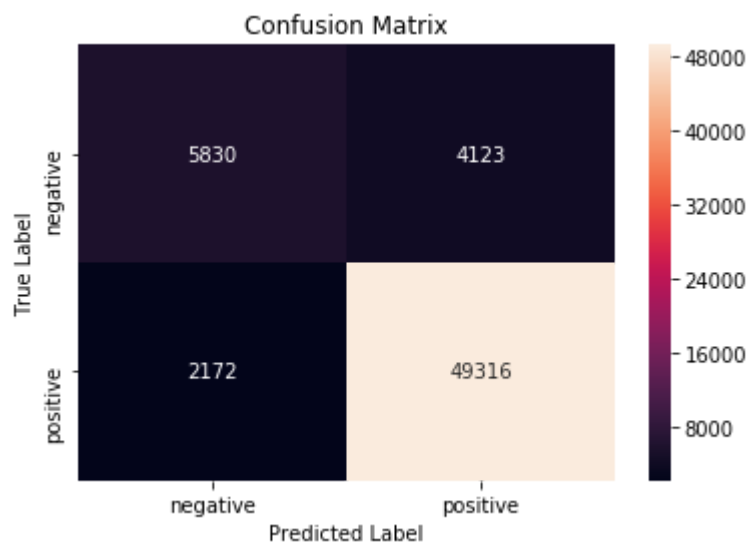
```
In [32]: 1 con_mat,con_mat_train,clf = dt_results(50,500,final_tf_idf,final_test_count,y
```

The accuracy of the DT classifier for maxDepth = 50 and min split = 500 is 86.180313%
 AUC: 0.825
 AUC: 0.911



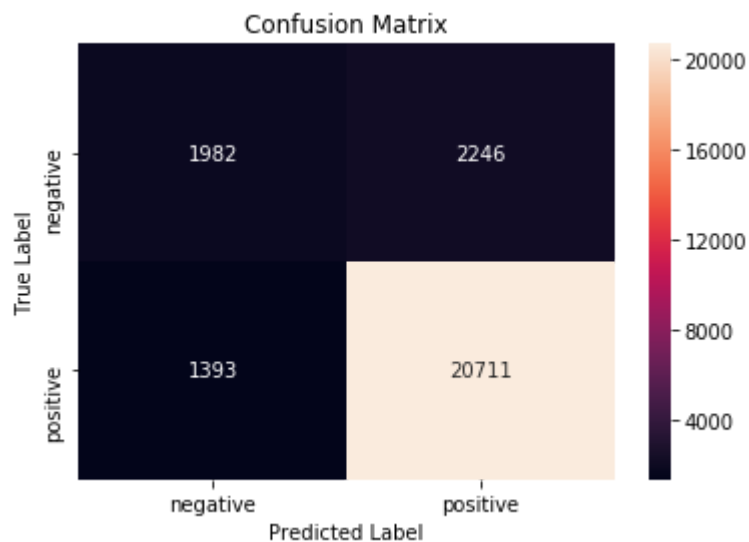
Observation: My model has predicted 86% accuracy with AUC: 0.825

```
In [33]: 1 showHeatMap(con_mat_train)
```



Observation: My model predicted 2172 + 4123 points wrongly in train data

```
In [35]: 1 showHeatMap(con_mat)
```



Observation: My model predicted 1393 + 2246 points wrongly

[5.2.1] Top 20 important features from SET 2

```
In [36]: 1 feature_names = np.array(tf_idf_vect.get_feature_names())
2 featureDict = dict(zip(feature_names, clf.feature_importances_))
3 sortedFeatures = sorted(featureDict.items(), key=lambda x: x[1], reverse=True)
4 print(type(sortedFeatures))
5 for i in range(0,20):
6     print(sortedFeatures[i])
```

```
<class 'list'>
('not', 0.10516319649370136)
('great', 0.05574250687918252)
('disappointed', 0.039630778078183465)
('worst', 0.03771103774223818)
('horrible', 0.026077555904695973)
('not buy', 0.02573570932893479)
('not worth', 0.025698879081758955)
('return', 0.022609828389066555)
('terrible', 0.022198814641220536)
('bad', 0.021413273008521163)
('awful', 0.018711373537716715)
('waste', 0.018038406237556696)
('love', 0.015707544506488345)
('delicious', 0.015009217862943312)
('refund', 0.014709302315080703)
('not recommend', 0.014437736293673174)
('disappointing', 0.014042244798650627)
('waste money', 0.013672608221397712)
('not good', 0.013005455660537103)
```

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

```
In [39]: 1 import graphviz
2 import pydotplus
3 import collections
4 from sklearn import tree
5 clf = DecisionTreeClassifier(max_depth=2)
6 clf.fit(final_tf_idf, y_1)
7 dot_data = tree.export_graphviz(clf, feature_names = tf_idf_vect.get_feature_
8 graph = pydotplus.graph_from_dot_data(dot_data)
9 colors = ('turquoise', 'orange')
10 edges = collections.defaultdict(list)
11 for edge in graph.get_edge_list():
12     edges[edge.get_source()].append(int(edge.get_destination()))
13
14 for edge in edges:
15     edges[edge].sort()
16     for i in range(2):
17         dest = graph.get_node(str(edges[edge][i]))[0]
18         dest.set_fillcolor(colors[i])
19
20 graph.write_png('treetfidf.png')
```

Out[39]: True

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

```
In [40]: 1 X_train, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed
```

In [41]:

```

1 i=0
2 list_of_sentence=[]
3 for sentence in X_train:
4     list_of_sentence.append(sentence.split())
5
6 w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
7 w2v_words = list(w2v_model.wv.vocab)
8
9 # average Word2Vec
10 # compute average word2vec for each review.
11 sent_vectors = []; # the avg-w2v for each sentence/review is stored in this L
12 for sent in tqdm(list_of_sentence): # for each review/sentence
13     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might
14     cnt_words = 0; # num of words with a valid vector in the sentence/review
15     for word in sent: # for each word in a review/sentence
16         if word in w2v_words:
17             vec = w2v_model.wv[word]
18             sent_vec += vec
19             cnt_words += 1
20     if cnt_words != 0:
21         sent_vec /= cnt_words
22     sent_vectors.append(sent_vec)
23 print(sent_vectors[0])
24
25
26
27 i=0
28 list_of_test_sentence=[]
29 for sentence in X_test:
30     list_of_test_sentence.append(sentence.split())
31
32 test_sent_vectors = [];
33
34 for sent in tqdm(list_of_test_sentence): # for each review/sentence
35     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might
36     cnt_words = 0; # num of words with a valid vector in the sentence/review
37     for word in sent: # for each word in a review/sentence
38         if word in w2v_words:
39             vec = w2v_model.wv[word]
40             sent_vec += vec
41             cnt_words += 1
42     if cnt_words != 0:
43         sent_vec /= cnt_words
44     test_sent_vectors.append(sent_vec)
45 print(test_sent_vectors[0])
46

```

100%|██████████| 61441/61441 [02:47<00:00, 366.80it/s]

```

[-1.10754768e+00  2.00193054e-01 -4.16910136e-01  3.16457859e-01
 1.67205180e-01  2.26236710e-01  1.63417535e-01  7.77784129e-04
 -2.26487562e-01 -6.12884563e-01 -3.77290405e-01  6.69056482e-02
 8.02721659e-01  3.09768991e-02  1.54656387e+00  1.16882547e-01
 5.74342138e-01 -3.87559963e-01  9.58900229e-01 -6.66402416e-01
 7.56670846e-02  3.80231652e-01 -1.59714155e-01  1.76354673e-01
 1.56241445e-01 -2.10104528e-01 -7.18050218e-01 -7.24558968e-01

```

```

-3.62788136e-01 -1.63233715e-01 1.31064441e-01 -1.40604076e-01
1.01911919e-01 2.49025648e-01 8.28877521e-01 3.88748020e-01
2.57363992e-01 1.51706172e-01 1.89670206e-01 -3.30035001e-02
3.07788125e-01 1.34390374e-01 -1.03560609e+00 -4.67261180e-01
-9.94279975e-02 3.30757251e-01 5.00301523e-01 -4.65017455e-01
1.55648759e-01 -5.28523285e-01]

```

100%|██████████| 26332/26332 [01:13<00:00, 360.05it/s]

```

[-0.74806017 -0.48290271 -0.47126988 0.4013922 0.29665439 -0.45408076
1.00896888 -0.37358878 -0.78329729 -0.19524113 0.06300695 -0.8424012
0.0484644 1.09750664 0.04071722 -0.30163987 0.01318479 0.61584562
-0.12420691 -0.39806303 0.19882826 -0.36472114 0.78072151 -0.86098378
0.0171463 1.47464345 -0.01129889 -0.1263671 -0.42119131 -0.42079127
0.62090134 0.06183675 -0.6788543 0.67373388 0.35000022 0.01159972
-0.26829883 0.48789187 0.06580197 0.46019822 -0.29895352 0.44290327
0.57227689 0.28952171 -0.47447127 0.20441469 0.1476062 -1.02437518
0.49800761 -0.89132388]

```

In [42]:

```

1  # split the train data set into cross validation train and cross validation test
2  X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_train, y_1, test_size=0.2)
3
4  i=0
5  list_of_cv_sentence=[]
6  for sentence in X_tr:
7      list_of_cv_sentence.append(sentence.split())
8
9  cv_train_sent_vectors = [];
10
11 for sent in tqdm(list_of_cv_sentence): # for each review/sentence
12     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might want to initialize them to zero
13     cnt_words = 0; # num of words with a valid vector in the sentence/review
14     for word in sent: # for each word in a review/sentence
15         if word in w2v_words:
16             vec = w2v_model.wv[word]
17             sent_vec += vec
18             cnt_words += 1
19     if cnt_words != 0:
20         sent_vec /= cnt_words
21     cv_train_sent_vectors.append(sent_vec)
22 print(cv_train_sent_vectors[0])
23
24 i=0
25 list_of_cv_test_sentence=[]
26 for sentence in X_cv:
27     list_of_cv_test_sentence.append(sentence.split())
28
29 cv_test_sent_vectors = [];
30
31 for sent in tqdm(list_of_cv_test_sentence): # for each review/sentence
32     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might want to initialize them to zero
33     cnt_words = 0; # num of words with a valid vector in the sentence/review
34     for word in sent: # for each word in a review/sentence
35         if word in w2v_words:
36             vec = w2v_model.wv[word]
37             sent_vec += vec
38             cnt_words += 1
39     if cnt_words != 0:
40         sent_vec /= cnt_words
41     cv_test_sent_vectors.append(sent_vec)
42 print(cv_test_sent_vectors[0])
43
44 tuned_parameters = [{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples': [1, 5, 10, 50, 100, 500, 1000]}]
45
46 #Using GridSearchCV
47 model = GridSearchCV(DecisionTreeClassifier(), tuned_parameters, scoring = 'r2')
48 model.fit(cv_train_sent_vectors, y_tr)
49
50 print(model.best_estimator_)
51 print(model.score(cv_test_sent_vectors, y_cv))
52
53 check_trade_off(cv_train_sent_vectors, cv_test_sent_vectors, y_tr, y_cv)

```

100% |██████████| 43008/43008 [02:04<00:00, 346.43it/s]

```
[ -1.17079823e+00  6.30106942e-01 -3.98037020e-01  1.26815554e-01
  4.77530123e-01 -4.38024840e-01  3.82910454e-01 -7.14504091e-02
 -1.97967817e-01 -2.57848058e-01  2.28154914e-01 -7.28499447e-01
  5.58125364e-01  4.65385314e-01  1.19031763e+00 -2.57445815e-01
 -2.23603821e-02 -6.10887574e-02  6.34071672e-01 -8.80125334e-01
 -5.31466887e-01  2.93856806e-01  2.20533359e-01 -1.97766191e-01
 -2.19944531e-01  2.29905581e-02 -2.82737630e-01 -6.82443065e-01
  1.40747673e-01 -1.31457557e-01  1.03656121e+00  8.39704917e-02
 -1.07011952e-03  1.03453724e+00  8.46687205e-01  4.14343814e-01
  6.55215979e-03 -6.73246777e-01 -8.67709714e-02 -1.68713270e-01
 -3.82535082e-01 -4.62584463e-01 -3.31906426e-01 -1.80009013e-01
 -9.42691050e-02  7.90953606e-02 -6.51220036e-02 -2.22707379e-01
  4.61115037e-01 -1.56566691e-01]
```

100%|██████████| 18433/18433 [00:51<00:00, 359.95it/s]

```
[ -0.65424046 -0.15814827 -0.54246559  0.41331746 -0.15781516  0.28864999
  0.18132091  0.29827008 -0.07347252 -0.01320163  0.11105867  0.3110009
  0.3764931 -0.14698873  0.9797962 -0.1694225 -0.13996538 -0.04780426
  0.18079759 -0.49656016 -0.24204404  0.03810901 -0.09690808 -0.25461577
 -0.1482041  0.70717171 -0.13982589  0.10306405 -0.55370422  0.03788182
  0.01893827 -0.03094728 -0.26235353  0.26219302  0.43014521  0.06087685
 -0.19330437 -0.36747148  0.34109956 -0.18614172  0.23745384 -0.01328977
 -0.53631501  0.02791026 -0.35983008  0.3108568  0.96507663 -0.31063059
 -0.25054761 -0.21689876]
```

```
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=10,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=500,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')
```

0.8177623575622466

AUC: 0.629

AUC: 0.799

AUC: 0.784

AUC: 0.664

AUC: 0.664

AUC: 0.663

AUC: 0.662

#####

AUC from train data #####

AUC: 0.636

AUC: 0.816

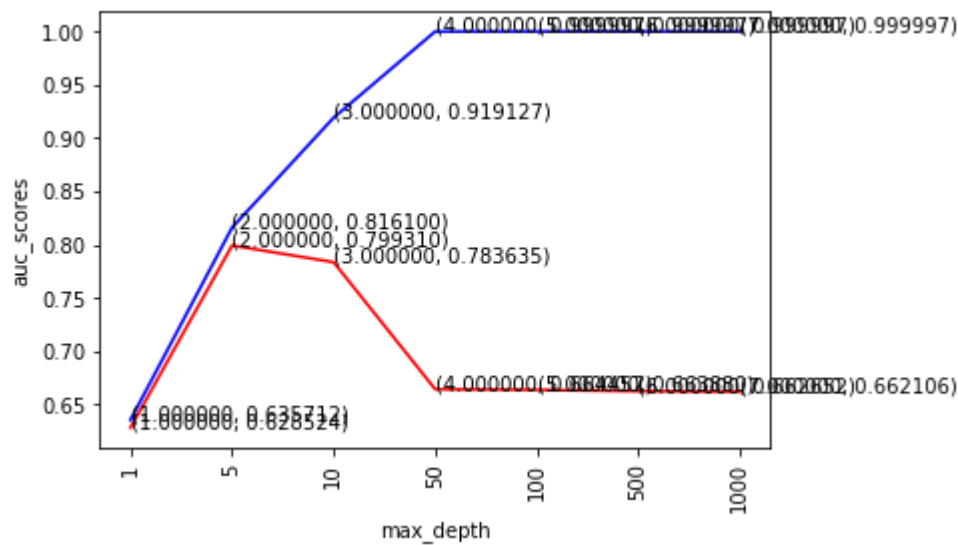
AUC: 0.919

AUC: 1.000

AUC: 1.000

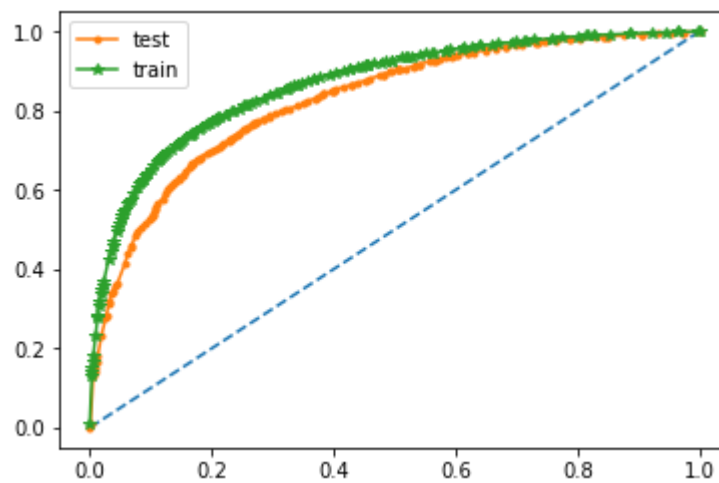
AUC: 1.000

AUC: 1.000



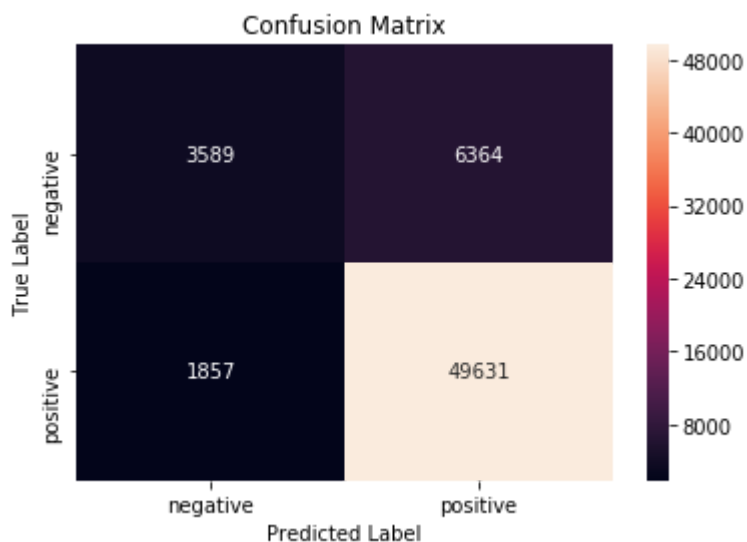
```
In [43]: 1 con_mat,con_mat_train,clf = dt_results(10,500,sent_vectors,test_sent_vectors,
```

The accuracy of the DT classifier for maxDepth = 10 and min split = 500 is 85.656236%
 AUC: 0.826
 AUC: 0.867



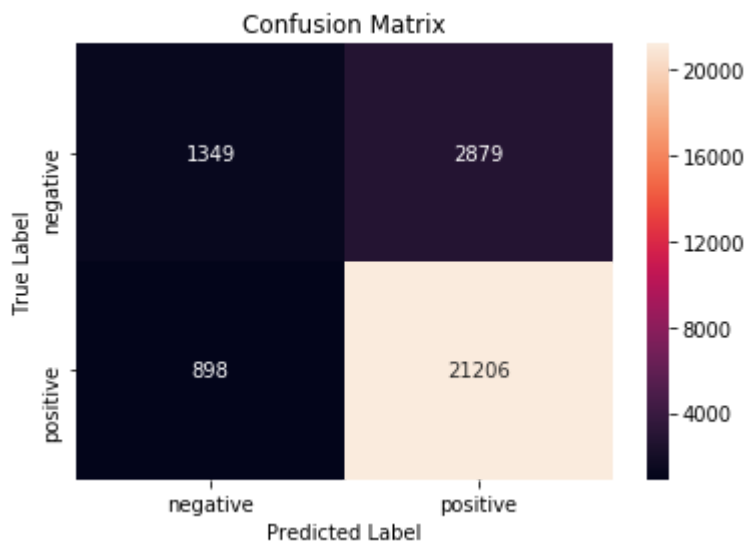
Observation: My model predicted with 85% accuracy with AUC: 0.832. Not much difference between Train ROC and Test ROC.

In [44]: 1 showHeatMap(con_mat_train)



Observation: My model predicted 1857 + 6364 points wrongly for train data set .

In [45]: 1 showHeatMap(con_mat)



Observation: My model predicted 898 + 2879 points wrongly

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

In [29]: 1 X_train, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed

In [30]: 1 model = TfidfVectorizer()
 2 X_train_transformed = model.fit_transform(X_train)
 3
 4 dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))

```
In [31]: 1 # Train your own Word2Vec model using your own text corpus
2 i=0
3 list_of_sentence=[]
4 for sentence in X_train:
5     list_of_sentence.append(sentence.split())
```

```
In [32]: 1 w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
2 w2v_words = list(w2v_model.wv.vocab)
```

```
In [33]: 1 # TF-IDF weighted Word2Vec
2 tfidf_feat = model.get_feature_names() # tfidf words/col-names
3 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
4
5 tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored i
6 row=0;
7 for sent in tqdm(list_of_sentence): # for each review/sentence
8     sent_vec = np.zeros(50) # as word vectors are of zero length
9     weight_sum =0; # num of words with a valid vector in the sentence/review
10    for word in sent: # for each word in a review/sentence
11        if word in w2v_words and word in tfidf_feat:
12            vec = w2v_model.wv[word]
13            # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
14            # to reduce the computation we are
15            # dictionary[word] = idf value of word in whole courpus
16            # sent.count(word) = tf valeus of word in this review
17            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
18            sent_vec += (vec * tf_idf)
19            weight_sum += tf_idf
20    if weight_sum != 0:
21        sent_vec /= weight_sum
22    tfidf_sent_vectors.append(sent_vec)
23    row += 1
```

100%|██████████| 61441/61441 [46:16<00:00, 32.22it/s]

```
In [34]: 1 i=0
2 list_of_test_sentence=[]
3 for sentence in X_test:
4     list_of_test_sentence.append(sentence.split())
```

```

In [35]: 1 # TF-IDF weighted Word2Vec
2 tfidf_feat = model.get_feature_names() # tfidf words/col-names
3 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
4
5 tfidf_test_sent_vectors = []; # the tfidf-w2v for each sentence/review is sto
6 row=0;
7 for sent in tqdm(list_of_test_sentence): # for each review/sentence
8     sent_vec = np.zeros(50) # as word vectors are of zero length
9     weight_sum =0; # num of words with a valid vector in the sentence/review
10    for word in sent: # for each word in a review/sentence
11        if word in w2v_words and word in tfidf_feat:
12            vec = w2v_model.wv[word]
13            # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
14            # to reduce the computation we are
15            # dictionary[word] = idf value of word in whole corpus
16            # sent.count(word) = tf value of word in this review
17            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
18            sent_vec += (vec * tf_idf)
19            weight_sum += tf_idf
20    if weight_sum != 0:
21        sent_vec /= weight_sum
22    tfidf_test_sent_vectors.append(sent_vec)
23    row += 1

```

100%|██████████| 26332/26332 [1:03:08<00:00, 6.95it/s]

```

In [36]: 1 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_train, y_1, test
2
3 i=0
4 list_of_cv_sentence=[]
5 for sentence in X_tr:
6     list_of_cv_sentence.append(sentence.split())
7
8
9 tfidf_feat = model.get_feature_names() # tfidf words/col-names
10 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
11
12 tfidf_cv_sent_vectors = []; # the tfidf-w2v for each sentence/review is store
13 row=0;
14 for sent in tqdm(list_of_cv_sentence): # for each review/sentence
15     sent_vec = np.zeros(50) # as word vectors are of zero length
16     weight_sum =0; # num of words with a valid vector in the sentence/review
17     for word in sent: # for each word in a review/sentence
18         if word in w2v_words and word in tfidf_feat:
19             vec = w2v_model.wv[word]
20             # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
21             # to reduce the computation we are
22             # dictionary[word] = idf value of word in whole courpus
23             # sent.count(word) = tf valeus of word in this review
24             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
25             sent_vec += (vec * tf_idf)
26             weight_sum += tf_idf
27         if weight_sum != 0:
28             sent_vec /= weight_sum
29         tfidf_cv_sent_vectors.append(sent_vec)
30         row += 1
31
32 i=0
33 list_of_cv_test_sentence=[]
34 for sentence in X_cv:
35     list_of_cv_test_sentence.append(sentence.split())
36
37
38 tfidf_cv_test_sent_vectors = []; # the tfidf-w2v for each sentence/review is
39 row=0;
40 for sent in tqdm(list_of_cv_test_sentence): # for each review/sentence
41     sent_vec = np.zeros(50) # as word vectors are of zero length
42     weight_sum =0; # num of words with a valid vector in the sentence/review
43     for word in sent: # for each word in a review/sentence
44         if word in w2v_words and word in tfidf_feat:
45             vec = w2v_model.wv[word]
46             # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
47             # to reduce the computation we are
48             # dictionary[word] = idf value of word in whole courpus
49             # sent.count(word) = tf valeus of word in this review
50             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
51             sent_vec += (vec * tf_idf)
52             weight_sum += tf_idf
53         if weight_sum != 0:
54             sent_vec /= weight_sum
55         tfidf_cv_test_sent_vectors.append(sent_vec)
56         row += 1

```

```

57
58
59 tuned_parameters = [{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples
60
61 #Using GridSearchCV
62 model = GridSearchCV(DecisionTreeClassifier(), tuned_parameters, scoring = 'r
63 model.fit(tfidf_cv_sent_vectors, y_tr)
64
65 print(model.best_estimator_)
66 print(model.score(tfidf_cv_test_sent_vectors, y_cv))
67
68 check_trade_off(tfidf_cv_sent_vectors,tfidf_cv_test_sent_vectors,y_tr,y_cv)
69

```

```
100%|██████████| 43008/43008 [34:39<00:00, 20.68it/s]
```

```
100%|██████████| 18433/18433 [21:01<00:00, 14.62it/s]
```

```

DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=1000,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=500,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')

```

```
0.7829388087147336
```

```
AUC: 0.604
```

```
AUC: 0.761
```

```
AUC: 0.746
```

```
AUC: 0.641
```

```
AUC: 0.641
```

```
AUC: 0.646
```

```
AUC: 0.643
```

```
#####
```

```
AUC from train data #####
```

```
AUC: 0.604
```

```
AUC: 0.775
```

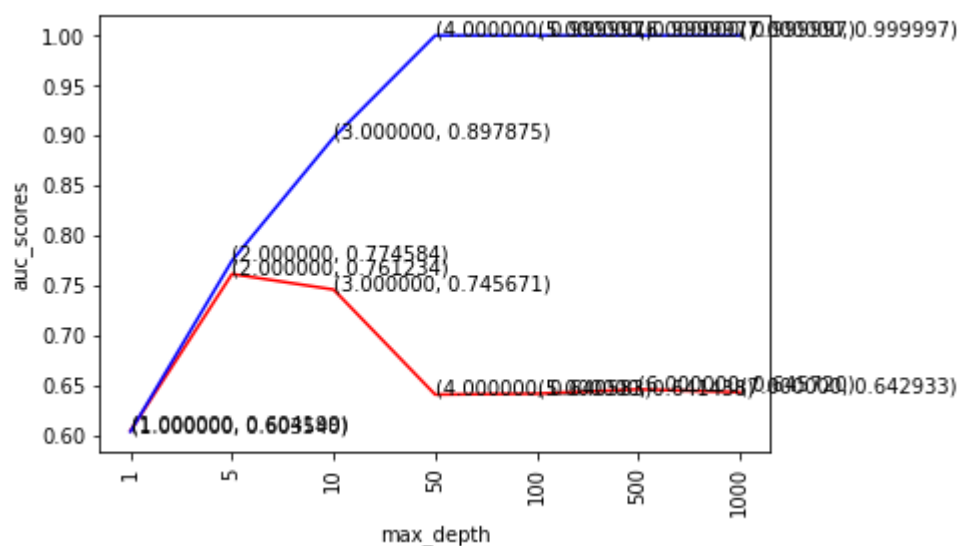
```
AUC: 0.898
```

```
AUC: 1.000
```

```
AUC: 1.000
```

```
AUC: 1.000
```

```
AUC: 1.000
```

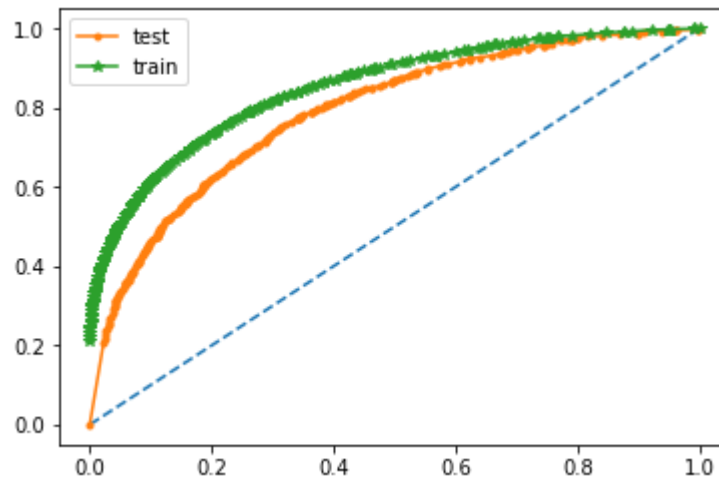


```
In [37]: 1 con_mat,con_mat_trai,clf = dt_results(1000,500,tfidf_sent_vectors,tfidf_test_
```

The accuracy of the DT classifier for maxDepth = 1000 and min split = 500 is 84.968859%

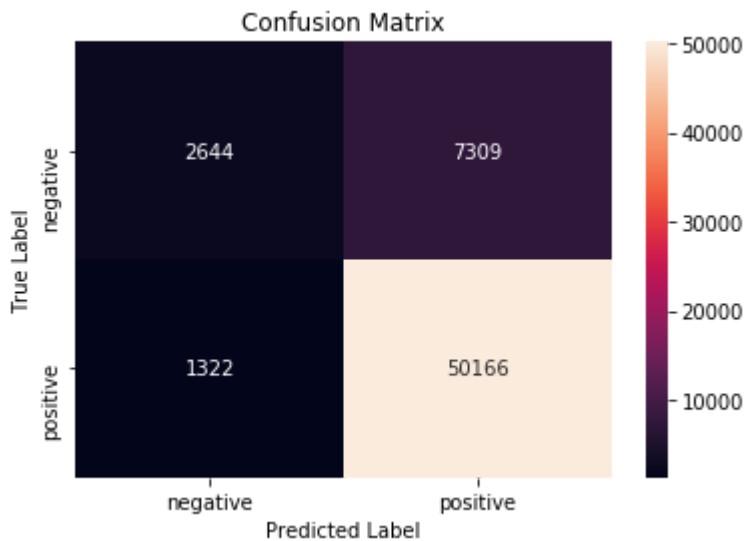
AUC: 0.791

AUC: 0.853



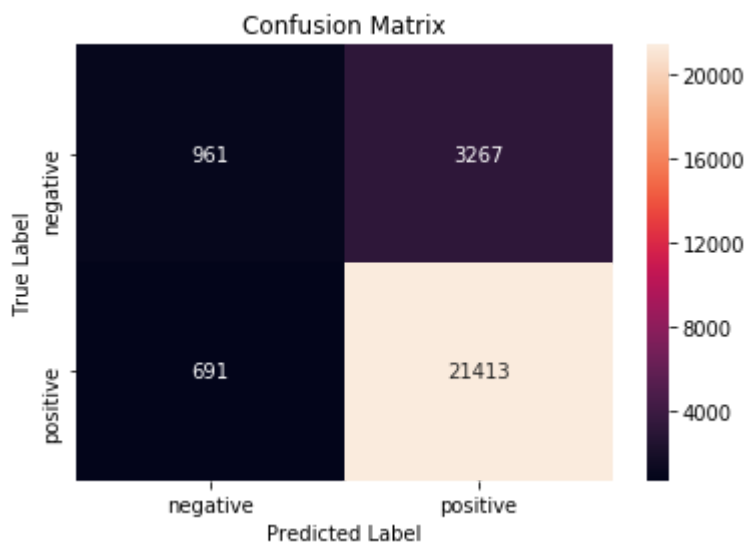
Observation: My model predicted with 84 % accuracy with AUC: 0.791

```
In [39]: 1 showHeatMap(con_mat_trai)
```



Observation : My model predicted 1322 + 7309 points wrongly for train data.

```
In [40]: 1 showHeatMap(con_mat)
```



Observation : My model predicted 691 + 3267 points wrongly

Repeat with extra features

```
In [41]: 1 mylen = np.vectorize(len)
2 newarr = mylen(preprocessed_summary)
```

```
In [42]: 1 newproce_reviews = np.asarray(preprocessed_reviews)
```

```
In [43]: 1 newproce_summary = np.asarray(preprocessed_summary)
```

```
In [44]: 1 df = pd.DataFrame({'desc':newproce_reviews, 'summary':newproce_summary, 'len':
```

```
In [45]: 1 df.head()
```

Out[45]:

	desc	summary	len
0	dogs loves chicken product china wont buying a...	made china	10
1	dogs love saw pet store tag attached regarding...	dog lover delites	17
2	infestation fruitflies literally everywhere fl...	one fruitfly stuck	18
3	worst product gotten long time would rate no s...	not work not waste money	24
4	wish would read reviews making purchase basica...	big rip	7

```
In [46]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(df, final['Score
```


In [47]:

```

1 import scipy
2 count_vect = CountVectorizer()
3 final_counts = count_vect.fit_transform(X_1['desc'])
4 final_test_count = count_vect.transform(X_test['desc'])
5
6 # split the train data set into cross validation train and cross validation test
7 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.2)
8
9 final_counts_tr_cv = count_vect.transform(X_tr['desc'])
10 final_test_count_cv = count_vect.transform(X_cv['desc'])
11
12 from scipy.sparse import csr_matrix, issparse
13
14 #####Adding len as feature#####
15 #if issparse(final_counts_tr_cv):
16     #print('sparse matrix')
17 len_sparse = scipy.sparse.coo_matrix(X_tr['len'])
18 len_sparse = len_sparse.transpose()
19
20 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, len_sparse])
21 print(final_counts_tr_cv.shape)
22
23 len_test_sparse = scipy.sparse.coo_matrix(X_cv['len'])
24 len_test_sparse = len_test_sparse.transpose()
25 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, len_test_sparse])
26 print("final_counts_tr_cv.shape after length = ", final_counts_tr_cv.shape)
27
28 #####Adding summary as feature#####
29 final_summary_count = count_vect.transform(X_tr['summary'])
30 final_test_summary_count_cv = count_vect.transform(X_cv['summary'])
31 columns=count_vect.get_feature_names()
32
33 print("sujet", final_summary_count[:,12].shape)
34 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, final_summary_count])
35 print("final_counts_tr_cv.shape after f1= ", final_counts_tr_cv.shape)
36
37 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, final_test_summary_count_cv])
38
39
40 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, final_summary_count])
41 print("final_counts_tr_cv.shape after f2= ", final_counts_tr_cv.shape)
42
43
44 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, final_test_summary_count_cv])
45
46 #####finding the new C #####
47
48 tuned_parameters = [{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples': [1, 5, 10, 50, 100, 500, 1000]}]
49
50 #Using GridSearchCV
51 model = GridSearchCV(DecisionTreeClassifier(), tuned_parameters, scoring = 'r2')
52 model.fit(final_counts_tr_cv, y_tr)
53
54 print(model.best_estimator_)
55 print(model.score(final_test_count_cv, y_cv))
56

```

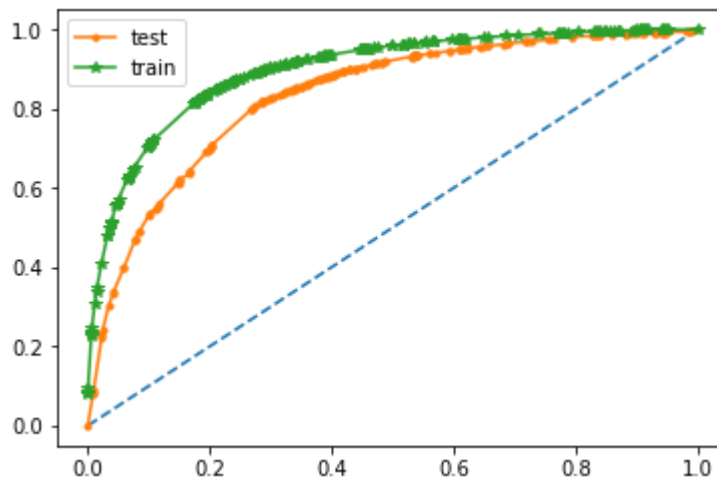
```
(43008, 46447)
final_counts_tr_cv.shape after length = (43008, 46447)
sujet (43008, 1)
final_counts_tr_cv.shape after f1= (43008, 46448)
final_counts_tr_cv.shape after f2= (43008, 46449)
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=50,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=500,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')
0.8124306837028364
```

```
In [48]: 1 con_mat,con_mat_train,clf = dt_results(50,500,final_counts,final_test_count,y
```

The accuracy of the DT classifier for maxDepth = 50 and min split = 500 is 85.861309%

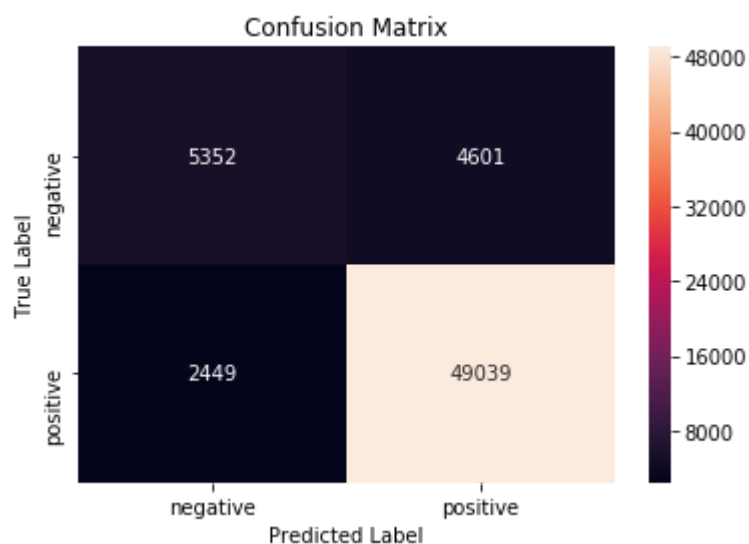
AUC: 0.833

AUC: 0.898



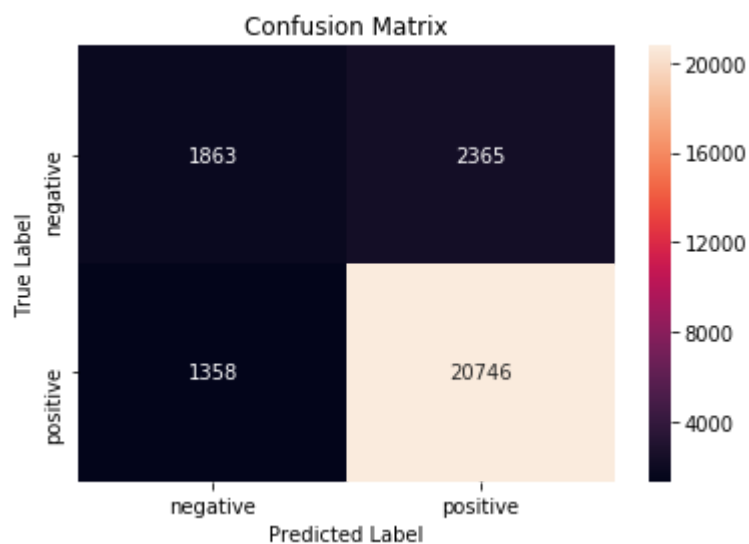
Observation: My model predicted with 85% accuracy with AUC: 0.833

In [49]: 1 showHeatMap(con_mat_train)



Observation: My model predicted 2449 + 4601 points wrongly even with train data

In [50]: 1 showHeatMap(con_mat)



Observation: My model predicted 1358 + 2365 points wrongly

[6] Conclusions

Method	No of samples	depth	split	accuracy	AUC Score
BOW	100000	50	500	85	0.833
TFIDF	100000	50	500	86	0.826
AVG W2VE	100000	10	500	85	0.832
TFIDF W2VE	100000	1000	500	85	0.817

Method	No of samples	depth	split	accuracy	AUC Score
BOW1	100000	50	500	85	0.833

In []:

1