# 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://nycdatascience.com/blog/student-works/amazon-fine-foods-visualization/ (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 - unqiue 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

# [1]. Reading Data

# Assignment 9: Random Forests

## [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 wil be set to "positive". Otherwise, it will be set to "negative".

# Mounting Google Drive locally

In [3]:
```python
from google.colab import drive
drive.mount('/content/gdrive')
```

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awg%3Aoauth%3A2.0%3Aoob&scope=email%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdocs.test%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fpeopleapi.readonly&response_type=code (https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awg%3Aoauth%3A2.0%3Aoob&scope=email%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdocs.test%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fpeopleapi.readonly&response_type=code)

Enter your authorization code:
..........
Mounted at /content/gdrive

In [0]:
```python
pip install paramiko
```

In [0]:
```python
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")


import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os
```

In [6]:
```python
# using SQLite Table to read data.
con = sqlite3.connect("/content/gdrive/My Drive/Dataset/database.sqlite")


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)
```

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

Out[6]:

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

In [0]:
```python
display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

In [0]:
```
print(display.shape)
display.head()
```

(80668, 7)

Out[6]:

| | UserId | ProductId | ProfileName | Time | Score | Text | COUNT(*) |
|---|---|---|---|---|---|---|---|
| 0 | #oc-R115TNMSPFT9I7 | B007Y59HVM | Breyton | 1331510400 | 2 | Overall its just OK when considering the price... | 2 |
| 1 | #oc-R11D9D7SHXIJB9 | B005HG9ET0 | Louis E. Emory "hoppy" | 1342396800 | 5 | My wife has recurring extreme muscle spasms, u... | 3 |
| 2 | #oc-R11DNU2NBKQ23Z | B007Y59HVM | Kim Cieszykowski | 1348531200 | 1 | This coffee is horrible and unfortunately not ... | 2 |
| 3 | #oc-R11O5J5ZVQE25C | B005HG9ET0 | Penguin Chick | 1346889600 | 5 | This will be the bottle that you grab from the... | 3 |
| 4 | #oc-R12KPBODL2B5ZD | B007OSBE1U | Christopher P. Presta | 1348617600 | 1 | I didnt like this coffee. Instead of telling y... | 2 |

In [0]:
```
display[display['UserId']=='AZY10LLTJ71NX']
```

Out[7]:

| | 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 [0]:
```
display['COUNT(*)'].sum()
```

Out[8]:   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 [0]:  display= pd.read_sql_query("""
         SELECT *
         FROM Reviews
         WHERE Score != 3 AND UserId="AR5J8UI46CURR"
         ORDER BY ProductID
         """, con)
         display.head()
```

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

```
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
          (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 calcualtions

In [0]:
```
display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)

display.head()
```

Out[9]:

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

In [11]:
```
final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
final.shape
```

Out[11]: (87773, 10)

In [12]:
```
#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[12]: 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 obsereved to be better than Porter Stemming)

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

In [0]:
```python
# 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 [0]:
```python
# 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', "you're", "you've",\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of',\
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"])
```

In [15]:
```python
# Combining all the above stundents
from bs4 import BeautifulSoup
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentance in tqdm(final['Text'].values):
    sentance = re.sub(r"http\S+", "", sentance)
    sentance = BeautifulSoup(sentance, 'lxml').get_text()
    sentance = decontracted(sentance)
    sentance = re.sub("\S*\d\S*", "", sentance).strip()
    sentance = re.sub('[^A-Za-z]+', ' ', sentance)
    # https://gist.github.com/sebleier/554280
    sentance = ' '.join(e.lower() for e in sentance.split() if e.lower() not in stopwords)
    preprocessed_reviews.append(sentance.strip())
```

100%|████████████████| 87773/87773 [00:34<00:00, 2559.25it/s]

In [0]:
```python
final["CleanText"] = [preprocessed_reviews[i] for i in range(len(final))]
```

In [0]:
```python
final.head(2)
```

Out[14]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenom |
|---|---|---|---|---|---|---|
| **22620** | 24750 | 2734888454 | A13ISQV0U9GZIC | Sandikaye | 1 | |
| **22621** | 24751 | 2734888454 | A1C298ITT645B6 | Hugh G. Pritchard | 0 | |

# [4] Featurization

In [0]:
```python
from sklearn.model_selection import train_test_split

from sklearn.metrics import accuracy_score

from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_auc_score
import seaborn as sns

from sklearn.metrics import confusion_matrix

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
```

In [0]:
```python
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
```

In [0]:
```python
Total_X = final['CleanText'].values
Total_y = final['Score'].values
```

In [0]:
```python
# split the data set into train and test
X_train, X_test, y_train, y_test = train_test_split(Total_X, Total_y, test_size=0.33)

# split the train data set into cross validation train and cross validation test
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33)
```

In [21]:
```python
print(f"Train Data : ({len(X_train)} , {len(y_train)})")
print(f"CV Data : ({len(X_cv)} , {len(y_cv)})")
print(f"Test Data : ({len(X_test)} , {len( y_test)})")
```

```
Train Data : (39400 , 39400)
CV Data : (19407 , 19407)
Test Data : (28966 , 28966)
```

# Testing the max_depth with Test datapoints and Confusion Matrix

## [4.1] BAG OF WORDS

## [5.1] Applying RF

In [0]:
```python
from sklearn.metrics import precision_score
from sklearn.metrics import f1_score
from sklearn.metrics import recall_score
```

In [0]:
```python
def RF(X_train_reg,X_cv_reg, y_train=y_train, y_cv=y_cv, y_test=y_test):

    max_depth = [2,3,4,5,6,7,8,9,10]
    min_samples_split = [5, 10, 100,200, 500,1000]
    tuned_parameters = [{'max_depth': max_depth, 'min_samples_split':min_samples_split}]

    #Using GridSearchCV
    model = GridSearchCV(RandomForestClassifier( class_weight = "balanced"), tuned_parameters, n_jobs=2 ,scor
    model.fit(X_train_reg, y_train)

    print(model.best_estimator_)
    print("_"*10)
    print("Best HyperParameter: ",model.best_params_)
    print(f"Best Accuracy: {model.best_score_*100}")

    tr_auc = model.cv_results_["mean_train_score"]
    cv_auc = model.cv_results_["mean_test_score"]

    reshape_tr_auc = tr_auc.reshape(len(max_depth),len( min_samples_split))
    reshape_cv_auc = cv_auc.reshape(len(max_depth),len( min_samples_split))

    plt.figure(figsize = (16,5))
    ax = sns.heatmap(reshape_tr_auc, annot=True, fmt="g", cmap='viridis')
    plt.xlabel(" min_samples_split")
    plt.ylabel("max_depth")
    plt.title("Train Set Score")
    plt.show()


    plt.figure(figsize = (16,5))
    sns.heatmap(reshape_cv_auc, annot=True, fmt="g", cmap='viridis')
    plt.xlabel(" min_samples_split")
    plt.ylabel("max_depth")
    plt.title("CV Set Score")
    plt.show()
```

In [0]:
```python
def testing_RF(X_train_reg,X_test_reg, max_d,min_s, y_train=y_train,  y_test=y_test):

    clf= RandomForestClassifier(max_depth =  max_d, min_samples_split =min_s, class_weight = "balanced")
    clf.fit(X_test_reg, y_test)

    train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(X_train_reg)[:,1])
    test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_test_reg)[:,1])

    plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
    plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
    plt.xscale(value = 'log')
    plt.legend()
    plt.xlabel("False Positive Rate")
    plt.ylabel("True Positive Rate")
    plt.title("ROC curves")
    plt.show()

    print(f" AUC score on test Data with max_depth = {max_d} and min_samples_split = {min_s } is :  {roc_auc_sco


    print(f"Precision on test data: {precision_score(y_test, clf.predict(X_test_reg))}")
    print(f"Recall on test data: {recall_score(y_test, clf.predict(X_test_reg))}")
    print(f"F1-Score on test data: {f1_score(y_test, clf.predict(X_test_reg))}")

    print("\nConfusion Matrix of Train and Test set:\n [ [TN  FP]\n   [FN TP] ]\n")
    confusionMatrix_train=confusion_matrix(y_train, clf.predict(X_train_reg))
    confusionMatrix_test=confusion_matrix(y_test, clf.predict(X_test_reg))
    df_cm_tr = pd.DataFrame(confusionMatrix_train, range(2),range(2))
    df_cm_te = pd.DataFrame(confusionMatrix_test, range(2),range(2))
    plt.figure(figsize = (7,5))
    sns.set(font_scale=1)#for label size
    plt.ylabel("Predicted label")
    plt.xlabel("Actual label")
    plt.title("Confusion Matrix of Train Set")

    sns.heatmap(df_cm_tr, annot=True,annot_kws={"size": 12},fmt="d")

    plt.figure(figsize = (7,6))
    plt.ylabel("Predicted label")
    plt.xlabel("Actual label")
    plt.title("Confusion Matrix of Test Set")
    sns.heatmap(df_cm_te, annot=True,annot_kws={"size": 12},fmt="d")
```

## [5.1.1] Applying Random Forests on BOW, SET 1

In [0]:
```python
#BoW
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(X_train)
print("some feature names ", count_vect.get_feature_names()[1000:1010])
print('='*50)


# we use the fitted CountVectorizer to convert the text to vector
X_train_bow = count_vect.transform(X_train)
X_cv_bow = count_vect.transform(X_cv)
X_test_bow = count_vect.transform(X_test)



print("After vectorizations")
print(X_train_bow.shape, y_train.shape)
print(X_cv_bow.shape, y_cv.shape)
print(X_test_bow.shape, y_test.shape)
print("="*100)
```

some feature names  ['alternitives', 'alterra', 'alters', 'altho', 'althogh', 'althoug', 'although', 'althought', 'altinb
as', 'altitude']
================================================
After vectorizations
(39400, 37579) (39400,)
(19407, 37579) (19407,)
(28966, 37579) (28966,)
================================================================================
===========

In [0]:
```python
#BoW
```

In [0]:  `RF(X_train_bow,X_cv_bow)`

RandomForestClassifier(bootstrap=True, class_weight='balanced',
            criterion='gini', max_depth=10, max_features='auto',
            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,
            n_estimators=10, n_jobs=None, oob_score=False,
            random_state=None, verbose=0, warm_start=False)
_____
Best HyperParameter:  {'max_depth': 10, 'min_samples_split': 500}
Best Accuracy: 79.97777629534181

Train Set Score

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.650819 | 0.66156 | 0.648499 | 0.679671 | 0.639097 | 0.631612 |
| 1 | 0.693313 | 0.682872 | 0.684816 | 0.706624 | 0.692592 | 0.696437 |
| 2 | 0.735848 | 0.725191 | 0.710414 | 0.728126 | 0.711389 | 0.715476 |
| 3 | 0.745809 | 0.758806 | 0.761403 | 0.739067 | 0.74449 | 0.73582 |
| 4 | 0.788669 | 0.772098 | 0.767335 | 0.75908 | 0.765541 | 0.767218 |
| 5 | 0.780025 | 0.774962 | 0.778549 | 0.778204 | 0.772192 | 0.77184 |
| 6 | 0.788287 | 0.802548 | 0.787008 | 0.792446 | 0.799343 | 0.796131 |
| 7 | 0.812091 | 0.815873 | 0.805423 | 0.805797 | 0.789631 | 0.801694 |
| 8 | 0.827449 | 0.823229 | 0.816419 | 0.816213 | 0.814546 | 0.803396 |

CV Set Score

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.645854 | 0.65858 | 0.644087 | 0.679935 | 0.6369 | 0.624378 |
| 1 | 0.684688 | 0.677174 | 0.677342 | 0.699037 | 0.687269 | 0.695126 |
| 2 | 0.722869 | 0.712498 | 0.699018 | 0.716196 | 0.705081 | 0.707889 |
| 3 | 0.732259 | 0.744873 | 0.746391 | 0.729979 | 0.73376 | 0.727909 |
| 4 | 0.772602 | 0.751499 | 0.749605 | 0.747423 | 0.758051 | 0.759367 |
| 5 | 0.759382 | 0.751569 | 0.760704 | 0.760763 | 0.758925 | 0.761357 |
| 6 | 0.751482 | 0.772711 | 0.763938 | 0.778643 | 0.778207 | 0.785212 |
| 7 | 0.777945 | 0.785006 | 0.779086 | 0.790377 | 0.775842 | 0.787521 |
| 8 | 0.791228 | 0.792932 | 0.784235 | 0.793899 | 0.799778 | 0.793498 |

In [0]:  testing_RF(X_train_bow,X_test_bow,10,500)

### ROC curves



 AUC score on test Data with max_depth = 10 and min_samples_split = 500 is :  75.3513742584806 %
Precision on test data: 0.9440500604722091
Recall on test data: 0.7373803754055941
F1-Score on test data: 0.8280140208467853

Confusion Matrix of Train and Test set:
 [ [TN  FP]
   [FN TP] ]

### Confusion Matrix of Train Set

Confusion Matrix of Test Set



## [5.1.2] Wordcloud of top 20 important features from <span style="color:red">SET 1</span>

**Code reference:** https://www.datacamp.com/community/tutorials/wordcloud-python (https://www.datacamp.com/community/tutorials/wordcloud-python)

https://python-graph-gallery.com/wordcloud/ (https://python-graph-gallery.com/wordcloud/)

https://stackoverflow.com/questions/11116697/how-to-get-most-informative-features-for-scikit-learn-classifiers (https://stackoverflow.com/questions/11116697/how-to-get-most-informative-features-for-scikit-learn-classifiers)

In [0]:
```python
from wordcloud import WordCloud, STOPWORDS
```

In [0]:
```python
def important_features(vect,max_depth,min_samples_split,X_train_reg, n):
    clf = RandomForestClassifier(max_depth = max_depth, min_samples_split=min_samples_split)
    clf.fit(X_train_reg, y_train)

    features =vect.get_feature_names()
    coef = clf.feature_importances_
    coef_df = pd.DataFrame({'word': features, 'coeficient': coef}, index = None)
    df = coef_df.sort_values("coeficient", ascending = False)[:n]
    cloud = " ".join(word for word in df.word)
    stopwords = set(STOPWORDS)
    wordcloud = WordCloud(width = 1000, height = 600, background_color ='white', stopwords = stopwords).gene

    # plot the WordCloud image
    plt.figure(figsize = (10, 8))
    plt.imshow(wordcloud, interpolation = 'bilinear')
    plt.axis("off")
    plt.title(f"Top {n} most important features")
    plt.tight_layout(pad = 0)

    plt.show()
```

In [0]:
```python
important_features(count_vect,10,500,X_train_bow,20)
```



Top 20 most important features

## [5.1.3] Applying Random Forests on TFIDF, SET 2

In [0]:
```python
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(X_train)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names()[0:10])
print('='*50)


# we use the fitted CountVectorizer to convert the text to vector
X_train_tf_idf = tf_idf_vect.transform(X_train)
X_cv_tf_idf = tf_idf_vect.transform(X_cv)
X_test_tf_idf = tf_idf_vect.transform(X_test)


print("After vectorizations")
print(X_train_tf_idf.shape, y_train.shape)
print(X_cv_tf_idf.shape, y_cv.shape)
print(X_test_tf_idf.shape, y_test.shape)
print("="*100)
```

```
some sample features(unique words in the corpus) ['ability', 'able', 'able add', 'able buy', 'able chew', 'able dri
nk', 'able eat', 'able enjoy', 'able find', 'able finish']
==================================================
After vectorizations
(39400, 23390) (39400,)
(19407, 23390) (19407,)
(28966, 23390) (28966,)
================================================================================
===========
```
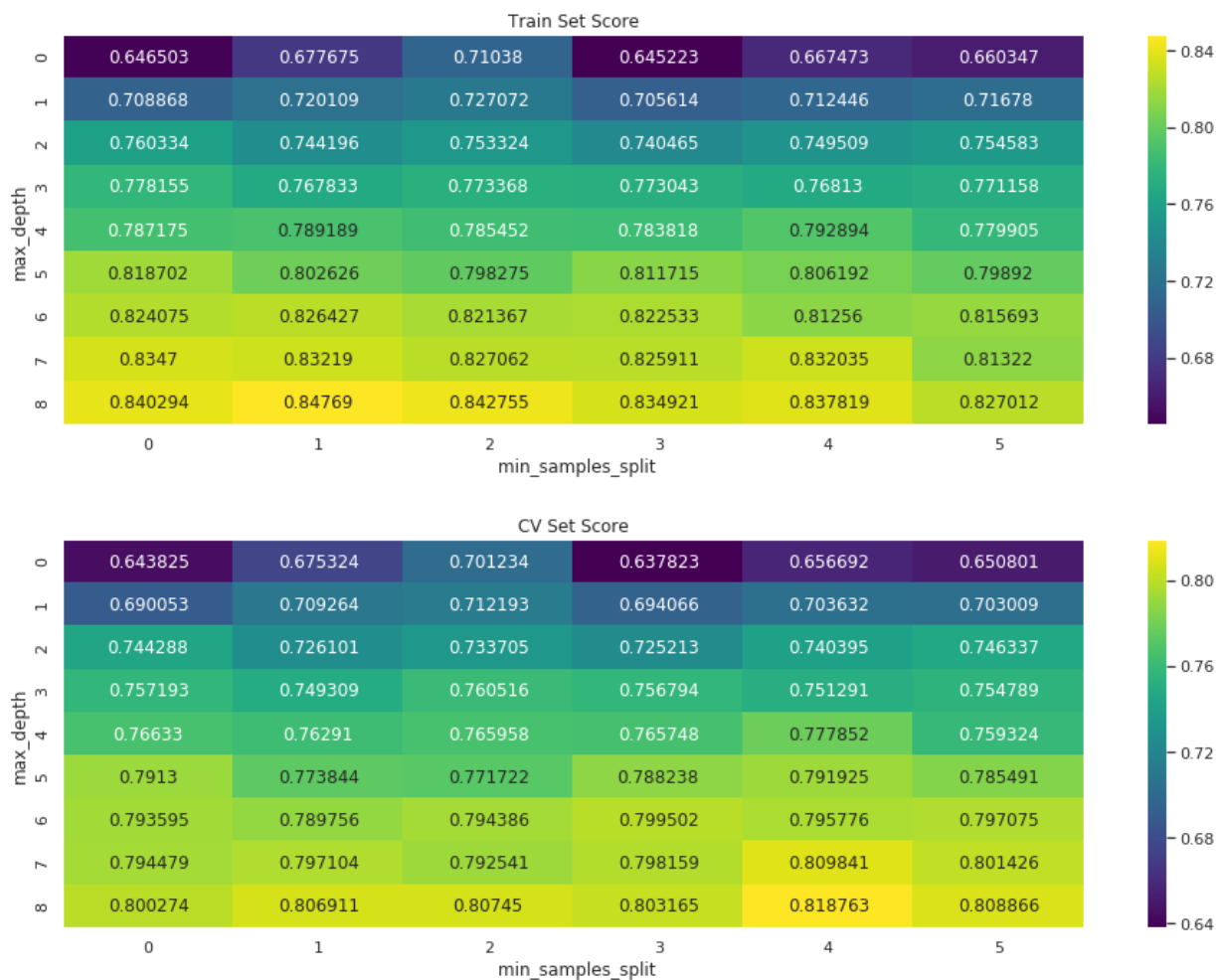
In [0]:   RF(X_train_tf_idf,X_cv_tf_idf)

RandomForestClassifier(bootstrap=True, class_weight='balanced',
        criterion='gini', max_depth=10, max_features='auto',
        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,
        n_estimators=10, n_jobs=None, oob_score=False,
        random_state=None, verbose=0, warm_start=False)
_____
Best HyperParameter:  {'max_depth': 10, 'min_samples_split': 500}
Best Accuracy: 81.87625778659117

Train Set Score

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.646503 | 0.677675 | 0.71038 | 0.645223 | 0.667473 | 0.660347 |
| 1 | 0.708868 | 0.720109 | 0.727072 | 0.705614 | 0.712446 | 0.71678 |
| 2 | 0.760334 | 0.744196 | 0.753324 | 0.740465 | 0.749509 | 0.754583 |
| 3 | 0.778155 | 0.767833 | 0.773368 | 0.773043 | 0.76813 | 0.771158 |
| 4 | 0.787175 | 0.789189 | 0.785452 | 0.783818 | 0.792894 | 0.779905 |
| 5 | 0.818702 | 0.802626 | 0.798275 | 0.811715 | 0.806192 | 0.79892 |
| 6 | 0.824075 | 0.826427 | 0.821367 | 0.822533 | 0.81256 | 0.815693 |
| 7 | 0.8347 | 0.83219 | 0.827062 | 0.825911 | 0.832035 | 0.81322 |
| 8 | 0.840294 | 0.84769 | 0.842755 | 0.834921 | 0.837819 | 0.827012 |

CV Set Score

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.643825 | 0.675324 | 0.701234 | 0.637823 | 0.656692 | 0.650801 |
| 1 | 0.690053 | 0.709264 | 0.712193 | 0.694066 | 0.703632 | 0.703009 |
| 2 | 0.744288 | 0.726101 | 0.733705 | 0.725213 | 0.740395 | 0.746337 |
| 3 | 0.757193 | 0.749309 | 0.760516 | 0.756794 | 0.751291 | 0.754789 |
| 4 | 0.76633 | 0.76291 | 0.765958 | 0.765748 | 0.777852 | 0.759324 |
| 5 | 0.7913 | 0.773844 | 0.771722 | 0.788238 | 0.791925 | 0.785491 |
| 6 | 0.793595 | 0.789756 | 0.794386 | 0.799502 | 0.795776 | 0.797075 |
| 7 | 0.794479 | 0.797104 | 0.792541 | 0.798159 | 0.809841 | 0.801426 |
| 8 | 0.800274 | 0.806911 | 0.80745 | 0.803165 | 0.818763 | 0.808866 |

In [0]:    testing_RF(X_train_tf_idf,X_test_tf_idf,10,500)

### ROC curves



 AUC score on test Data with max_depth = 10 and min_samples_split = 500 is :  73.27537931272717 %
Precision on test data: 0.9209900265349071
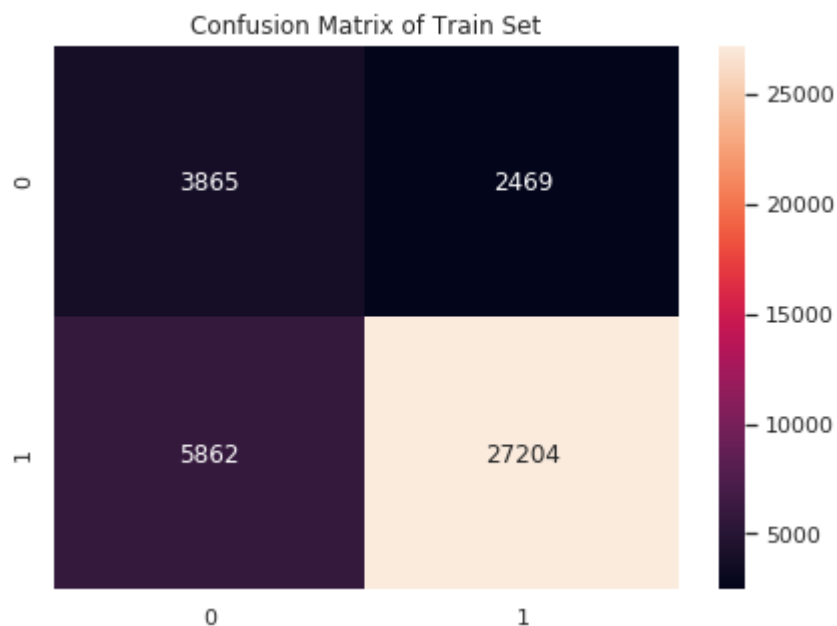Recall on test data: 0.830623865324311
F1-Score on test data: 0.8734759404694754

Confusion Matrix of Train and Test set:
 [ [TN  FP]
   [FN TP] ]

### Confusion Matrix of Train Set

Confusion Matrix of Test Set



**[5.1.4] Wordcloud of top 20 important features from SET 2**

In [0]:    important_features(tf_idf_vect,50,500,X_train_tf_idf,20)

Top 20 most important features



# [4.4] Word2Vec

In [0]:
```python
i=0

w2v_train=[]
w2v_cv=[]
w2v_test=[]

for sentance in X_train:
    w2v_train.append(sentance.split())

for sentance in X_cv:
    w2v_cv.append(sentance.split())


for sentance in X_test:
    w2v_test.append(sentance.split())
```

In [0]:
```python
want_to_train_w2v = True
if want_to_train_w2v:
# min_count = 5 considers only words that occured atleast 5 times
#w2v_model=Word2Vec(list_of_sentance,min_count=5,size=50, workers=4)
    w2v_model_train = Word2Vec(w2v_train,min_count=5,size=50, workers=4)
    print(w2v_model_train.wv.most_similar('great'))
    print('='*50)
else:
    pass
```

[('awesome', 0.8310557007789612), ('fantastic', 0.8176830410957336), ('good', 0.8043232560157776), ('exc
ellent', 0.7917436361312866), ('wonderful', 0.770425021648407), ('amazing', 0.7472847700119019), ('perfe
ct', 0.7148235440254211), ('terrific', 0.6971834897994995), ('decent', 0.682517945766449), ('ideal', 0.6448
290348052979)]
==================================================

In [0]:
```python
w2v_words_train = list(w2v_model_train.wv.vocab)

print("number of words that occured minimum 5 times ",len(w2v_words_train ))
print("sample words ", w2v_words_train[0:50])
```

number of words that occured minimum 5 times  12107
sample words  ['kind', 'opposite', 'effect', 'friend', 'recommended', 'try', 'hour', 'energy', 'use', 'replace', 'morni
ng', 'cup', 'coffee', 'anything', 'actually', 'made', 'tired', 'groggy', 'guess', 'different', 'influence', 'individual', 'wo
uld', 'suggest', 'trying', 'one', 'first', 'buying', 'bulk', 'best', 'gluten', 'free', 'bread', 'ever', 'not', 'break', 'apart', 'li
ke', 'others', 'almost', 'cake', 'texture', 'wonderful', 'flavor', 'amazon', 'approximately', 'cheaper', 'per', 'loaf', 'm
arket']

#Converting text into vectors using Avg W2V, TFIDF-W2V

## [5.1.5] Applying Random Forests on AVG W2V, <span style="color:red">SET 3</span>

In [0]:
```python
train_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(w2v_train): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change this to 300 if you us
    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_train:
            vec = w2v_model_train.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    train_vectors.append(sent_vec)
print()
print(len(train_vectors))
print(len(train_vectors[0]))
```

```
100%|████████████████| 39400/39400 [01:04<00:00, 614.48it/s]
```

```
39400
50
```

In [0]:
```python
# compute average word2vec for each review.
cv_vectors = [] # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(w2v_cv): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change this to 300 if you us
    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_train:
            vec = w2v_model_train.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    cv_vectors.append(sent_vec)
print()
print(len(cv_vectors))
print(len(cv_vectors[0]))
```

```
100%|████████████████| 19407/19407 [00:32<00:00, 606.30it/s]
```

```
19407
50
```

In [0]:
```python
test_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(w2v_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change this to 300 if you us
    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_train:
            vec =  w2v_model_train.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    test_vectors.append(sent_vec)
print()
print(len(test_vectors))
print(len(test_vectors[0]))
```

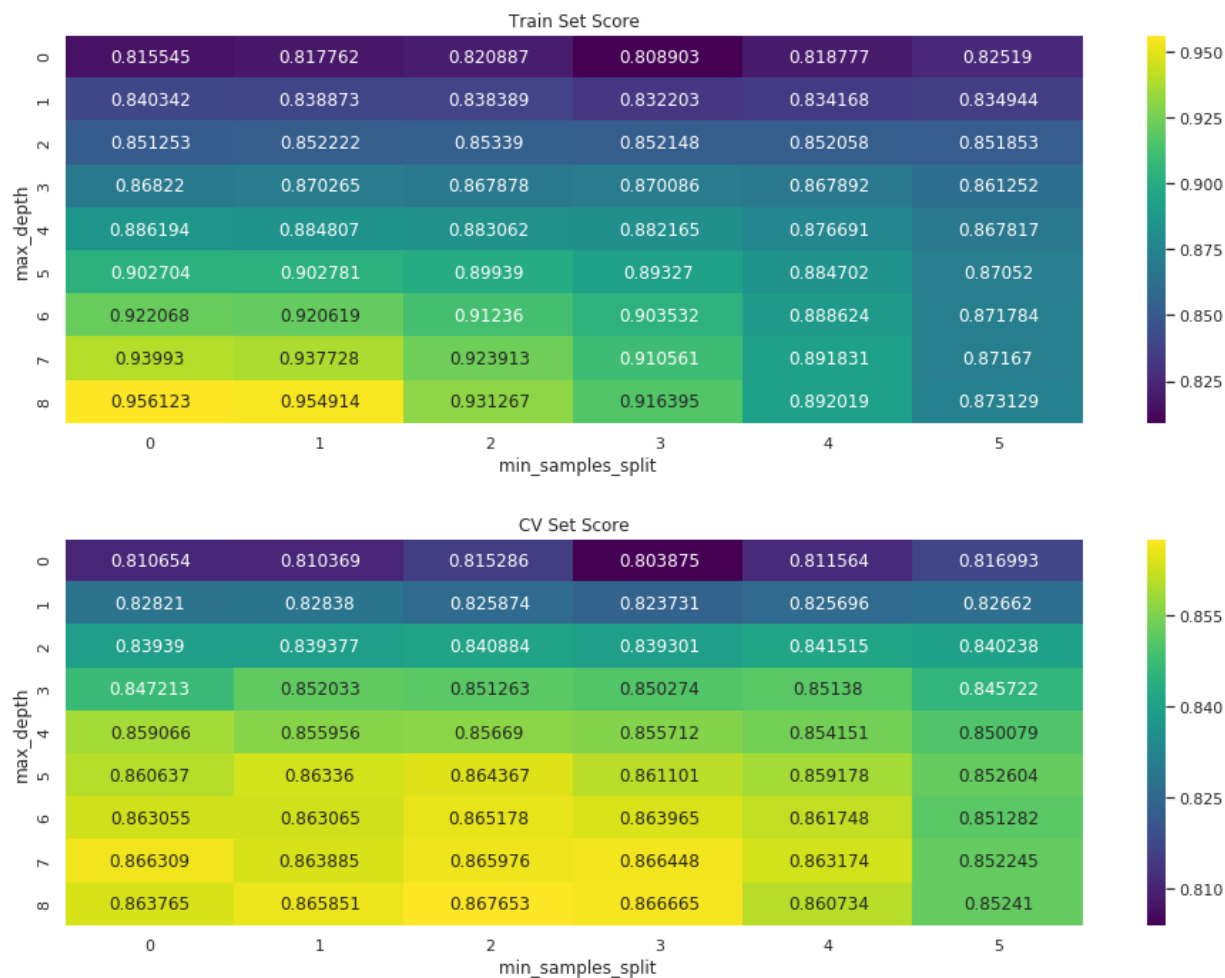100%|████████████████| 28966/28966 [00:46<00:00, 616.71it/s]


28966
50

In [0]:

In [0]: `RF(train_vectors,cv_vectors)`

RandomForestClassifier(bootstrap=True, class_weight='balanced',
        criterion='gini', max_depth=10, max_features='auto',
        max_leaf_nodes=None, min_impurity_decrease=0.0,
        min_impurity_split=None, min_samples_leaf=1,
        min_samples_split=100, min_weight_fraction_leaf=0.0,
        n_estimators=10, n_jobs=None, oob_score=False,
        random_state=None, verbose=0, warm_start=False)
_____

Best HyperParameter:  {'max_depth': 10, 'min_samples_split': 100}
Best Accuracy: 86.76530885027215

**Train Set Score**

| max_depth | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.815545 | 0.817762 | 0.820887 | 0.808903 | 0.818777 | 0.82519 |
| 1 | 0.840342 | 0.838873 | 0.838389 | 0.832203 | 0.834168 | 0.834944 |
| 2 | 0.851253 | 0.852222 | 0.85339 | 0.852148 | 0.852058 | 0.851853 |
| 3 | 0.86822 | 0.870265 | 0.867878 | 0.870086 | 0.867892 | 0.861252 |
| 4 | 0.886194 | 0.884807 | 0.883062 | 0.882165 | 0.876691 | 0.867817 |
| 5 | 0.902704 | 0.902781 | 0.89939 | 0.89327 | 0.884702 | 0.87052 |
| 6 | 0.922068 | 0.920619 | 0.91236 | 0.903532 | 0.888624 | 0.871784 |
| 7 | 0.93993 | 0.937728 | 0.923913 | 0.910561 | 0.891831 | 0.87167 |
| 8 | 0.956123 | 0.954914 | 0.931267 | 0.916395 | 0.892019 | 0.873129 |

min_samples_split

**CV Set Score**

| max_depth | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.810654 | 0.810369 | 0.815286 | 0.803875 | 0.811564 | 0.816993 |
| 1 | 0.82821 | 0.82838 | 0.825874 | 0.823731 | 0.825696 | 0.82662 |
| 2 | 0.83939 | 0.839377 | 0.840884 | 0.839301 | 0.841515 | 0.840238 |
| 3 | 0.847213 | 0.852033 | 0.851263 | 0.850274 | 0.85138 | 0.845722 |
| 4 | 0.859066 | 0.855956 | 0.85669 | 0.855712 | 0.854151 | 0.850079 |
| 5 | 0.860637 | 0.86336 | 0.864367 | 0.861101 | 0.859178 | 0.852604 |
| 6 | 0.863055 | 0.863065 | 0.865178 | 0.863965 | 0.861748 | 0.851282 |
| 7 | 0.866309 | 0.863885 | 0.865976 | 0.866448 | 0.863174 | 0.852245 |
| 8 | 0.863765 | 0.865851 | 0.867653 | 0.866665 | 0.860734 | 0.85241 |

min_samples_split

In [0]:   testing_RF(train_vectors, test_vectors,10,100)

### ROC curves



 AUC score on test Data with max_depth = 10 and min_samples_split = 100 is :  85.82035634792057 %
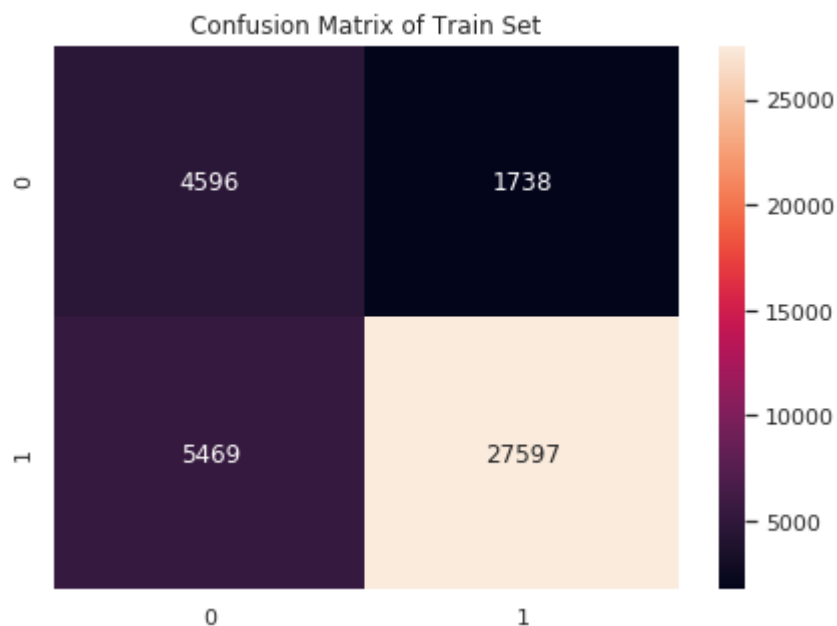Precision on test data: 0.9696600234466588
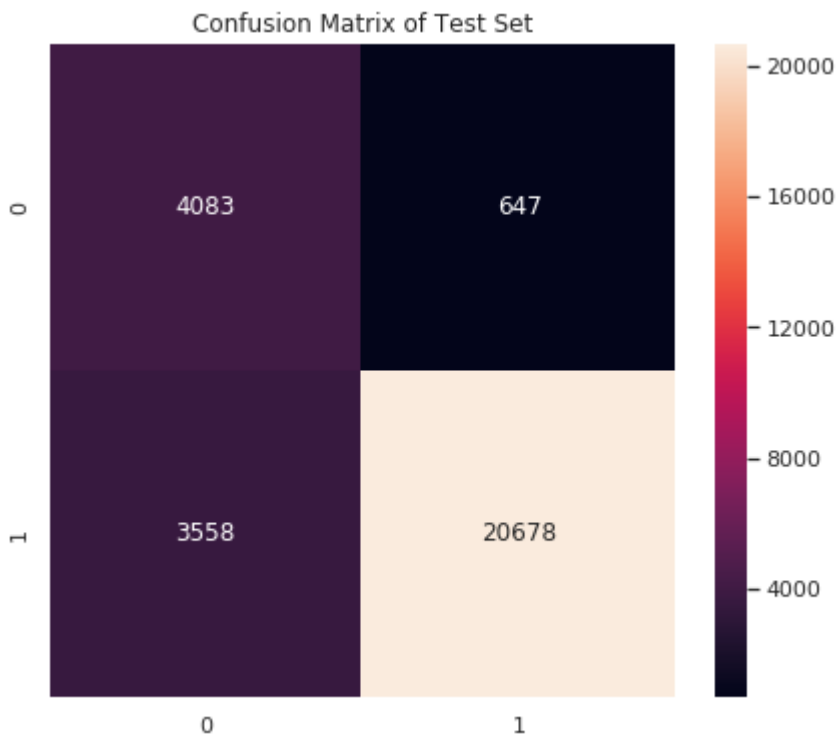Recall on test data: 0.8531935963030203
F1-Score on test data: 0.9077061521915674

Confusion Matrix of Train and Test set:
 [ [TN  FP]
   [FN TP] ]

### Confusion Matrix of Train Set

Confusion Matrix of Test Set



## [5.1.6] Applying Random Forests on TFIDF W2V, SET 4

In [0]:
```
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(X_train)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

In [0]:
```
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

train_tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(w2v_train): # 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_words_train and word in tfidf_feat:
            vec = w2v_model_train.wv[word]
            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
    train_tfidf_sent_vectors.append(np.array(sent_vec))
    row += 1
```

100%|████████████████| 39400/39400 [14:09<00:00, 46.35it/s]

In [0]:
```python
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

cv_tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(w2v_cv): # 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_words_train and word in tfidf_feat:
            vec = w2v_model_train.wv[word]
#            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    cv_tfidf_sent_vectors.append(np.array(sent_vec))
    row += 1
```

100%|████████████████| 19407/19407 [06:45<00:00, 47.85it/s]

In [0]:
```python
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

test_tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(w2v_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words_train and word in tfidf_feat:
            vec = w2v_model_train.wv[word]
#            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    test_tfidf_sent_vectors.append(np.array(sent_vec))
    row += 1
```
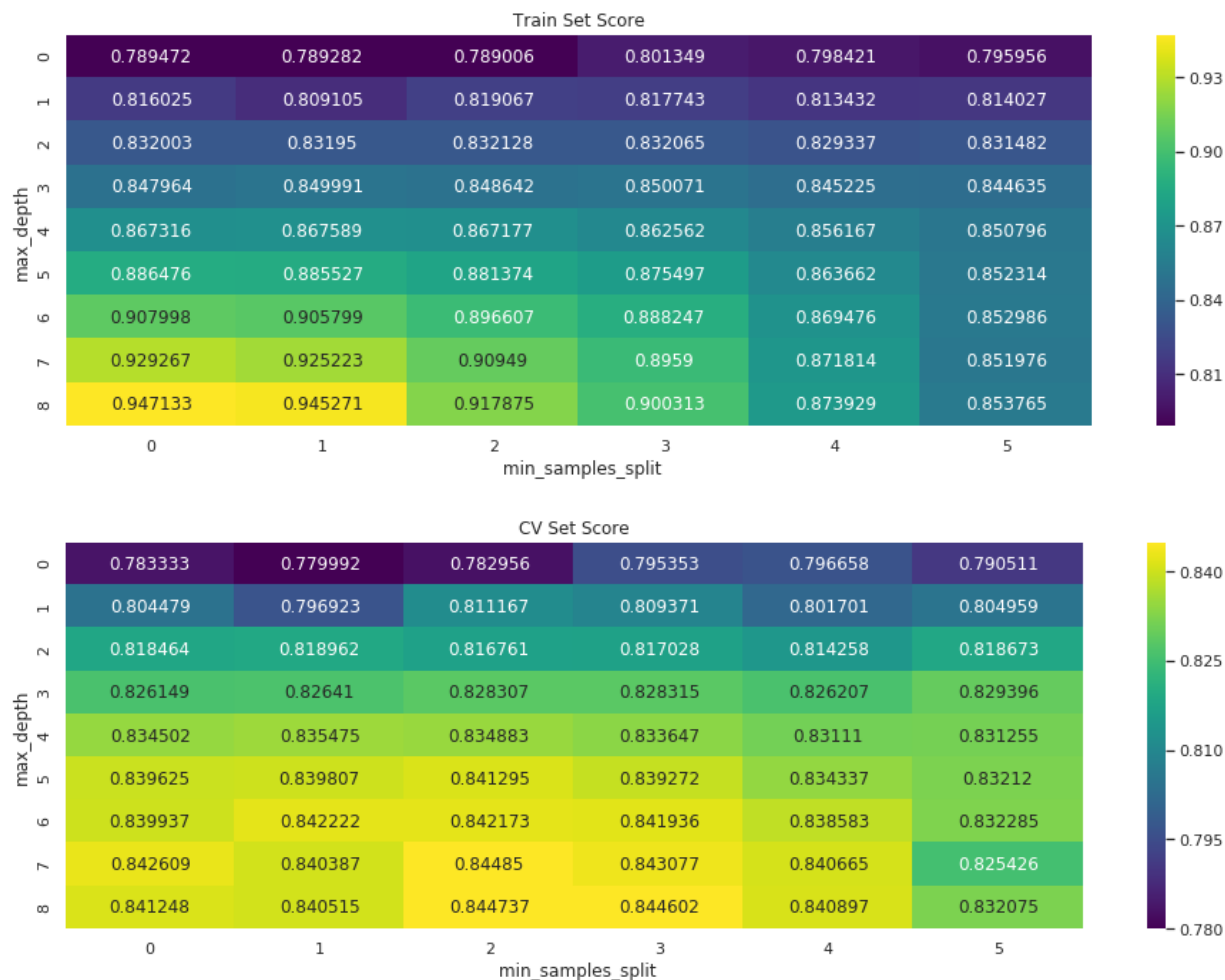
100%|████████████████| 28966/28966 [10:11<00:00, 47.37it/s]

In [0]: `RF(train_tfidf_sent_vectors, cv_tfidf_sent_vectors)`

RandomForestClassifier(bootstrap=True, class_weight='balanced',
            criterion='gini', max_depth=9, max_features='auto',
            max_leaf_nodes=None, min_impurity_decrease=0.0,
            min_impurity_split=None, min_samples_leaf=1,
            min_samples_split=100, min_weight_fraction_leaf=0.0,
            n_estimators=10, n_jobs=None, oob_score=False,
            random_state=None, verbose=0, warm_start=False)

—————————

Best HyperParameter: {'max_depth': 9, 'min_samples_split': 100}
Best Accuracy: 84.48501925072983

Train Set Score

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.789472 | 0.789282 | 0.789006 | 0.801349 | 0.798421 | 0.795956 |
| 1 | 0.816025 | 0.809105 | 0.819067 | 0.817743 | 0.813432 | 0.814027 |
| 2 | 0.832003 | 0.83195 | 0.832128 | 0.832065 | 0.829337 | 0.831482 |
| 3 | 0.847964 | 0.849991 | 0.848642 | 0.850071 | 0.845225 | 0.844635 |
| 4 | 0.867316 | 0.867589 | 0.867177 | 0.862562 | 0.856167 | 0.850796 |
| 5 | 0.886476 | 0.885527 | 0.881374 | 0.875497 | 0.863662 | 0.852314 |
| 6 | 0.907998 | 0.905799 | 0.896607 | 0.888247 | 0.869476 | 0.852986 |
| 7 | 0.929267 | 0.925223 | 0.90949 | 0.8959 | 0.871814 | 0.851976 |
| 8 | 0.947133 | 0.945271 | 0.917875 | 0.900313 | 0.873929 | 0.853765 |

CV Set Score

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.783333 | 0.779992 | 0.782956 | 0.795353 | 0.796658 | 0.790511 |
| 1 | 0.804479 | 0.796923 | 0.811167 | 0.809371 | 0.801701 | 0.804959 |
| 2 | 0.818464 | 0.818962 | 0.816761 | 0.817028 | 0.814258 | 0.818673 |
| 3 | 0.826149 | 0.82641 | 0.828307 | 0.828315 | 0.826207 | 0.829396 |
| 4 | 0.834502 | 0.835475 | 0.834883 | 0.833647 | 0.83111 | 0.831255 |
| 5 | 0.839625 | 0.839807 | 0.841295 | 0.839272 | 0.834337 | 0.83212 |
| 6 | 0.839937 | 0.842222 | 0.842173 | 0.841936 | 0.838583 | 0.832285 |
| 7 | 0.842609 | 0.840387 | 0.84485 | 0.843077 | 0.840665 | 0.825426 |
| 8 | 0.841248 | 0.840515 | 0.844737 | 0.844602 | 0.840897 | 0.832075 |

In [0]: `testing_RF(train_tfidf_sent_vectors, test_tfidf_sent_vectors,9,200)`

### ROC curves



 AUC score on test Data with max_depth = 9 and min_samples_split = 200 is :  81.70441591440336 %
Precision on test data: 0.9564035296923444
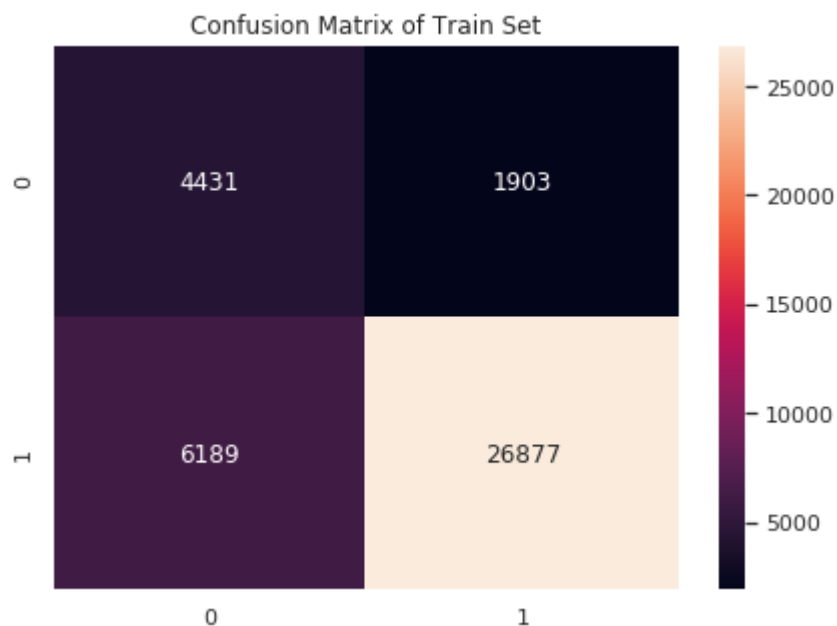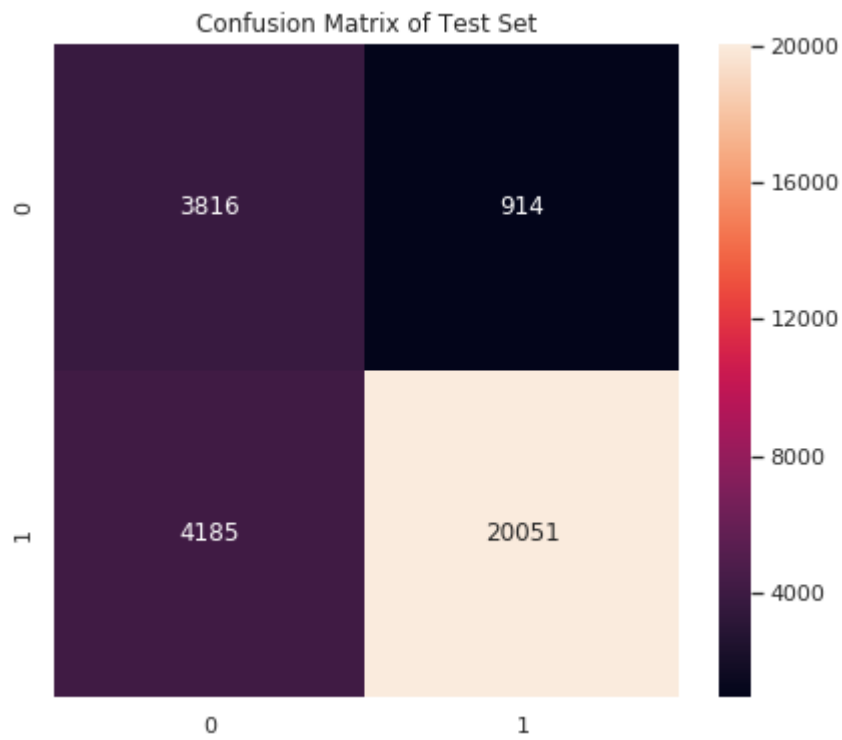Recall on test data: 0.827322990592507
F1-Score on test data: 0.8871927612220969

Confusion Matrix of Train and Test set:
 [ [TN  FP]
   [FN TP] ]

### Confusion Matrix of Train Set

Confusion Matrix of Test Set



## [5.2] Applying GBDT using XGBOOST

```python
In [0]: def XGB(X_train_reg,X_cv_reg, y_train=y_train, y_cv=y_cv, y_test=y_test):

            max_depth = [2,3,4,5,6,7,8,9,10]
            min_samples_split = [5, 10, 100,200, 500,1000]
            tuned_parameters = [{'max_depth': max_depth, 'min_samples_split':min_samples_split}]

            #Using GridSearchCV
            model = GridSearchCV(XGBClassifier(), tuned_parameters, n_jobs=1 ,scoring = 'roc_auc', cv=5, return_train_sc
            model.fit(X_train_reg, y_train)

            print(model.best_estimator_)
            print("_"*10)
            print("Best HyperParameter: ",model.best_params_)
            print(f"Best Accuracy: {model.best_score_*100}")

            tr_auc = model.cv_results_["mean_train_score"]
            cv_auc = model.cv_results_["mean_test_score"]

            reshape_tr_auc = tr_auc.reshape(len(max_depth),len( min_samples_split))
            reshape_cv_auc = cv_auc.reshape(len(max_depth),len( min_samples_split))

            plt.figure(figsize = (15,5))
            ax = sns.heatmap(reshape_tr_auc, annot=True, fmt="g", cmap='viridis')
            plt.xlabel(" min_samples_split")
            plt.ylabel("max_depth")
            plt.title("Train Set Score")
            plt.show()


            plt.figure(figsize = (15,5))
            sns.heatmap(reshape_cv_auc, annot=True, fmt="g", cmap='viridis')
            plt.xlabel(" min_samples_split")
            plt.ylabel("max_depth")
            plt.title("CV Set Score")
            plt.show()
```

In [0]:

```python
def testing_XGB(X_train_reg,X_test_reg, max_d,min_s, y_train=y_train,  y_test=y_test):

    clf= XGBClassifier(max_depth =  max_d, min_samples_split =min_s, class_weight = "balanced")
    clf.fit(X_test_reg, y_test)

    train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(X_train_reg)[:,1])
    test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_test_reg)[:,1])

    plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
    plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
    plt.xscale(value = 'log')
    plt.legend()
    plt.xlabel("False Positive Rate")
    plt.ylabel("True Positive Rate")
    plt.title("ROC curves")
    plt.show()

    print(f" AUC score on test Data with max_depth = {max_d} and min_samples_split = {min_s } is :  {roc_auc_sco


    print(f"Precision on test data: {precision_score(y_test, clf.predict(X_test_reg))}")
    print(f"Recall on test data: {recall_score(y_test, clf.predict(X_test_reg))}")
    print(f"F1-Score on test data: {f1_score(y_test, clf.predict(X_test_reg))}")

    print("\nConfusion Matrix of Train and Test set:\n [ [TN  FP]\n   [FN TP] ]\n")
    confusionMatrix_train=confusion_matrix(y_train, clf.predict(X_train_reg))
    confusionMatrix_test=confusion_matrix(y_test, clf.predict(X_test_reg))
    df_cm_tr = pd.DataFrame(confusionMatrix_train, range(2),range(2))
    df_cm_te = pd.DataFrame(confusionMatrix_test, range(2),range(2))
    plt.figure(figsize = (7,5))
    sns.set(font_scale=1)#for label size
    plt.ylabel("Predicted label")
    plt.xlabel("Actual label")
    plt.title("Confusion Matrix of Train Set")

    sns.heatmap(df_cm_tr, annot=True,annot_kws={"size": 12},fmt="d")

    plt.figure(figsize = (7,6))
    plt.ylabel("Predicted label")
    plt.xlabel("Actual label")
    plt.title("Confusion Matrix of Test Set")
    sns.heatmap(df_cm_te, annot=True,annot_kws={"size": 12},fmt="d")
```
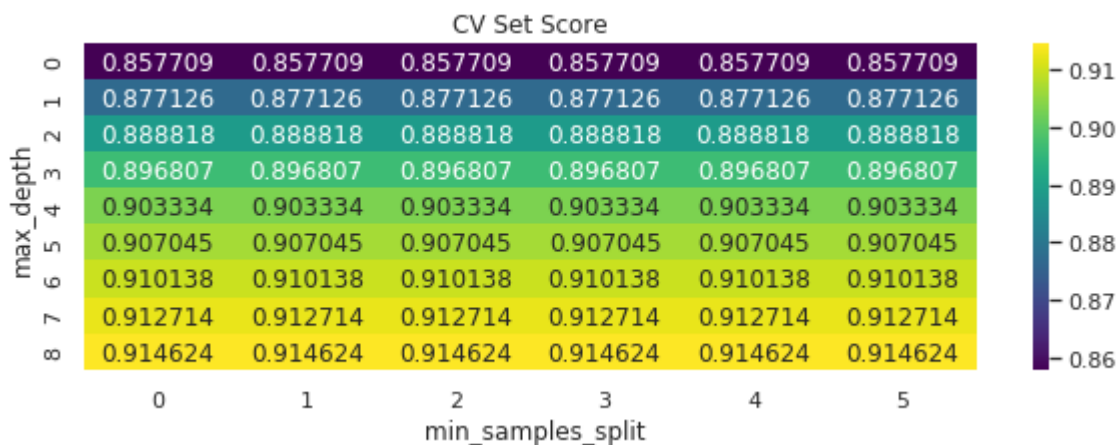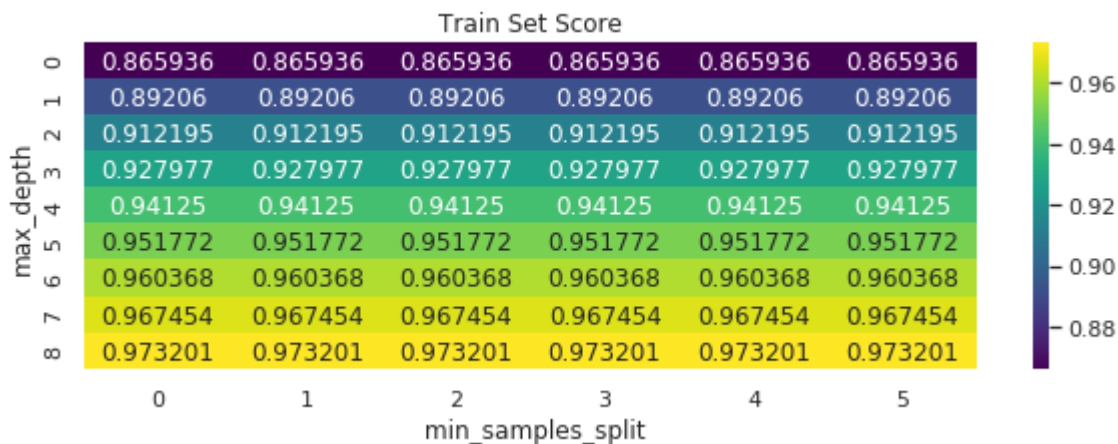
## [5.2.1] Applying XGBOOST on BOW, SET 1
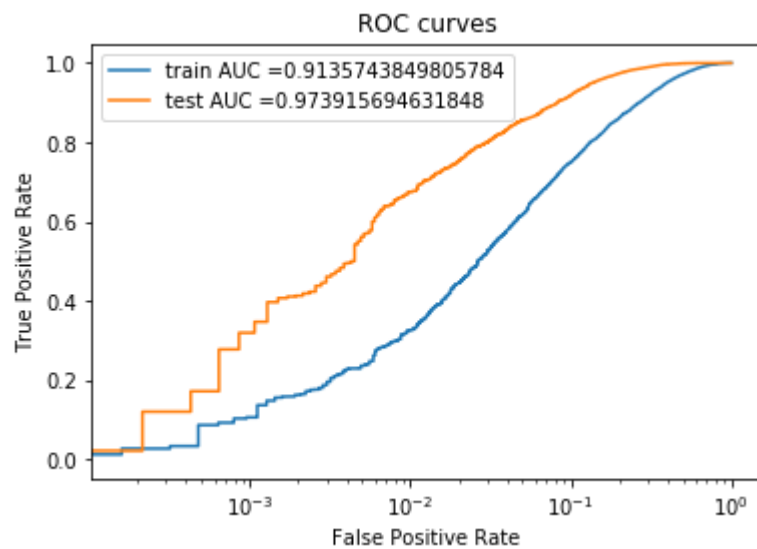
In [0]:    XGB(X_train_bow,X_cv_bow)

XGBClassifier(base_score=0.5, booster='gbtree', class_weight='balanced',
       colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1,
       gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=10,
       min_child_weight=1, min_samples_split=5, missing=None,
       n_estimators=100, n_jobs=1, nthread=None,
       objective='binary:logistic', random_state=0, reg_alpha=0,
       reg_lambda=1, scale_pos_weight=1, seed=None, silent=None,
       subsample=1, verbosity=1)
_____
Best HyperParameter: {'max_depth': 10, 'min_samples_split': 5}
Best Accuracy: 91.46241015786842

**Train Set Score**

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.865936 | 0.865936 | 0.865936 | 0.865936 | 0.865936 | 0.865936 |
| 1 | 0.89206 | 0.89206 | 0.89206 | 0.89206 | 0.89206 | 0.89206 |
| 2 | 0.912195 | 0.912195 | 0.912195 | 0.912195 | 0.912195 | 0.912195 |
| 3 | 0.927977 | 0.927977 | 0.927977 | 0.927977 | 0.927977 | 0.927977 |
| 4 | 0.94125 | 0.94125 | 0.94125 | 0.94125 | 0.94125 | 0.94125 |
| 5 | 0.951772 | 0.951772 | 0.951772 | 0.951772 | 0.951772 | 0.951772 |
| 6 | 0.960368 | 0.960368 | 0.960368 | 0.960368 | 0.960368 | 0.960368 |
| 7 | 0.967454 | 0.967454 | 0.967454 | 0.967454 | 0.967454 | 0.967454 |
| 8 | 0.973201 | 0.973201 | 0.973201 | 0.973201 | 0.973201 | 0.973201 |

**CV Set Score**

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.857709 | 0.857709 | 0.857709 | 0.857709 | 0.857709 | 0.857709 |
| 1 | 0.877126 | 0.877126 | 0.877126 | 0.877126 | 0.877126 | 0.877126 |
| 2 | 0.888818 | 0.888818 | 0.888818 | 0.888818 | 0.888818 | 0.888818 |
| 3 | 0.896807 | 0.896807 | 0.896807 | 0.896807 | 0.896807 | 0.896807 |
| 4 | 0.903334 | 0.903334 | 0.903334 | 0.903334 | 0.903334 | 0.903334 |
| 5 | 0.907045 | 0.907045 | 0.907045 | 0.907045 | 0.907045 | 0.907045 |
| 6 | 0.910138 | 0.910138 | 0.910138 | 0.910138 | 0.910138 | 0.910138 |
| 7 | 0.912714 | 0.912714 | 0.912714 | 0.912714 | 0.912714 | 0.912714 |
| 8 | 0.914624 | 0.914624 | 0.914624 | 0.914624 | 0.914624 | 0.914624 |

In [0]: testing_XGB(X_train_bow,X_test_bow,10,5)

### ROC curves



 AUC score on test Data with max_depth = 10 and min_samples_split = 5 is :  80.47504242534302 %
Precision on test data: 0.9303980003845415
Recall on test data: 0.996252985258997
F1-Score on test data: 0.9621999960231453

Confusion Matrix of Train and Test set:
 [ [TN  FP]
   [FN TP] ]

### Confusion Matrix of Train Set

Confusion Matrix of Test Set



# Observation

**When tested model on unseen data(test data) the auc score is 80%.In a nutshell we can say the generalization error is low means this model works well with unseen data.**

## [5.2.2] Applying XGBOOST on TFIDF, SET 2

In [0]:    XGB(X_train_tf_idf,X_cv_tf_idf)

XGBClassifier(base_score=0.5, booster='gbtree', class_weight='balanced',
       colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1,
       gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=10,
       min_child_weight=1, min_samples_split=5, missing=None,
       n_estimators=100, n_jobs=1, nthread=None,
       objective='binary:logistic', random_state=0, reg_alpha=0,
       reg_lambda=1, scale_pos_weight=1, seed=None, silent=None,
       subsample=1, verbosity=1)

_____

Best HyperParameter: {'max_depth': 10, 'min_samples_split': 5}
Best Accuracy: 92.17994529638467

In [0]:     testing_XGB(X_train_tf_idf,X_test_tf_idf,10,5)

ROC curves



 AUC score on test Data with max_depth = 10 and min_samples_split = 5 is :  82.57095941825007 %
Precision on test data: 0.9383440080253106
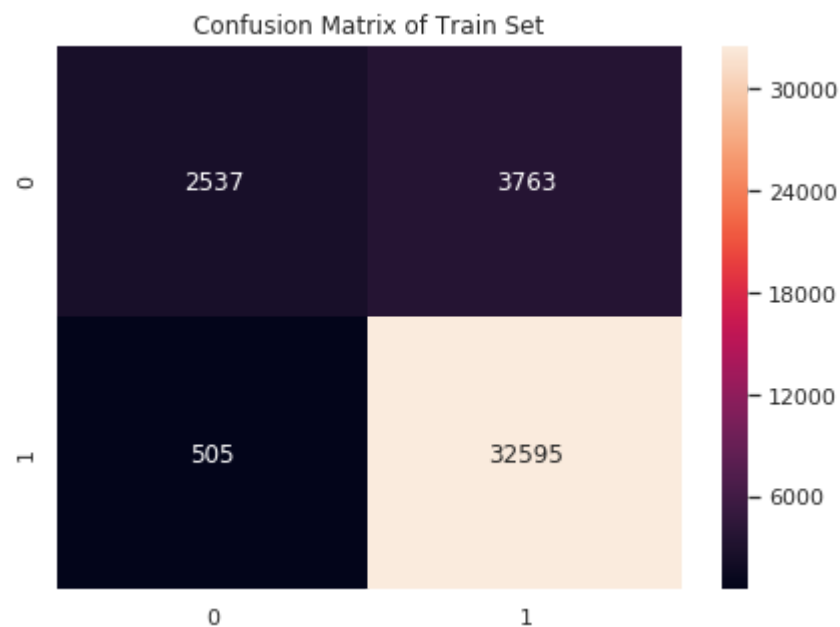Recall on test data: 0.9983579638752053
F1-Score on test data: 0.9674211384701062

Confusion Matrix of Train and Test set:
 [ [TN  FP]
   [FN TP] ]

Confusion Matrix of Train Set

Confusion Matrix of Test Set

|   | 0 | 1 |
|---|---|---|
| 0 | 3008 | 1598 |
| 1 | 40 | 24320 |

### [5.2.3] Applying XGBOOST on AVG W2V, SET 3

```
In [0]:   trv = np.array(train_vectors)
          cvv = np.array(cv_vectors)
          tev = np.array(test_vectors)
```

In [0]:   XGB(trv,cvv)
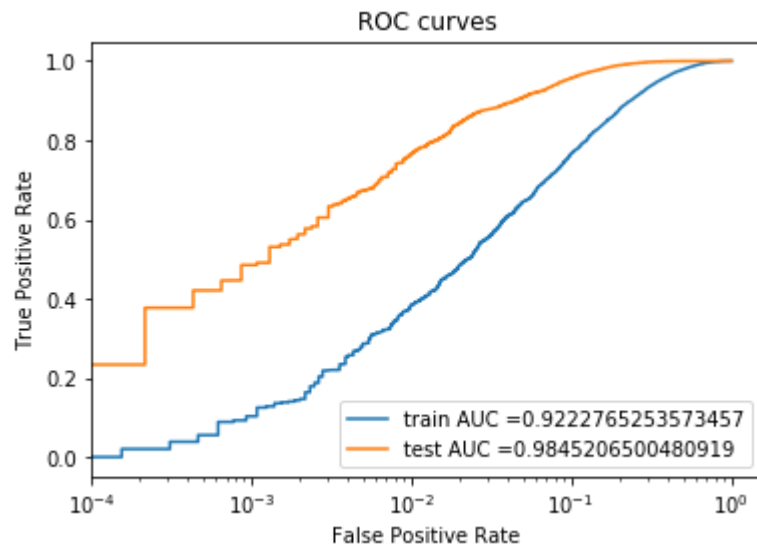
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
           colsample_bynode=1, colsample_bytree=1, gamma=0,
           learning_rate=0.1, max_delta_step=0, max_depth=7,
           min_child_weight=1, min_samples_split=5, missing=None,
           n_estimators=100, n_jobs=1, nthread=None,
           objective='binary:logistic', random_state=0, reg_alpha=0,
           reg_lambda=1, scale_pos_weight=1, seed=None, silent=None,
           subsample=1, verbosity=1)

_____

Best HyperParameter:  {'max_depth': 7, 'min_samples_split': 5}
Best Accuracy: 89.64293624581944

**Train Set Score**

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.892303 | 0.892303 | 0.892303 | 0.892303 | 0.892303 | 0.892303 |
| 1 | 0.910131 | 0.910131 | 0.910131 | 0.910131 | 0.910131 | 0.910131 |
| 2 | 0.928263 | 0.928263 | 0.928263 | 0.928263 | 0.928263 | 0.928263 |
| 3 | 0.948664 | 0.948664 | 0.948664 | 0.948664 | 0.948664 | 0.948664 |
| 4 | 0.969257 | 0.969257 | 0.969257 | 0.969257 | 0.969257 | 0.969257 |
| 5 | 0.985675 | 0.985675 | 0.985675 | 0.985675 | 0.985675 | 0.985675 |
| 6 | 0.995261 | 0.995261 | 0.995261 | 0.995261 | 0.995261 | 0.995261 |
| 7 | 0.998991 | 0.998991 | 0.998991 | 0.998991 | 0.998991 | 0.998991 |
| 8 | 0.999854 | 0.999854 | 0.999854 | 0.999854 | 0.999854 | 0.999854 |

**CV Set Score**

| max_depth \ min_samples_split | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.879851 | 0.879851 | 0.879851 | 0.879851 | 0.879851 | 0.879851 |
| 1 | 0.888751 | 0.888751 | 0.888751 | 0.888751 | 0.888751 | 0.888751 |
| 2 | 0.893553 | 0.893553 | 0.893553 | 0.893553 | 0.893553 | 0.893553 |
| 3 | 0.896322 | 0.896322 | 0.896322 | 0.896322 | 0.896322 | 0.896322 |
| 4 | 0.896155 | 0.896155 | 0.896155 | 0.896155 | 0.896155 | 0.896155 |
| 5 | 0.896429 | 0.896429 | 0.896429 | 0.896429 | 0.896429 | 0.896429 |
| 6 | 0.895839 | 0.895839 | 0.895839 | 0.895839 | 0.895839 | 0.895839 |
| 7 | 0.89554 | 0.89554 | 0.89554 | 0.89554 | 0.89554 | 0.89554 |
| 8 | 0.895487 | 0.895487 | 0.895487 | 0.895487 | 0.895487 | 0.895487 |

In [0]: testing_XGB(trv, tev,7,5)

ROC curves



 AUC score on test Data with max_depth = 7 and min_samples_split = 5 is :  88.38223855471843 %
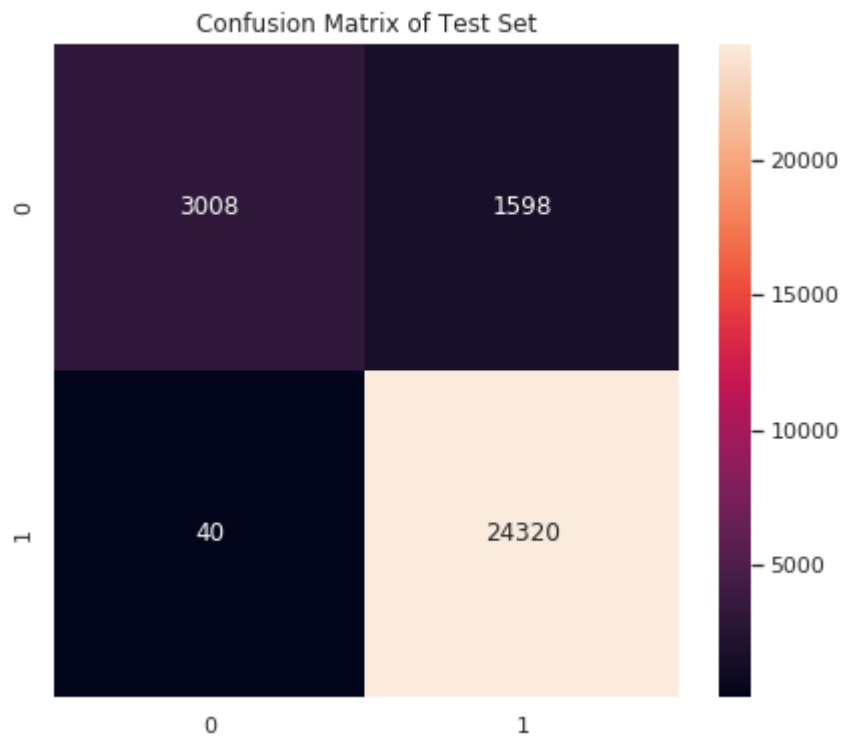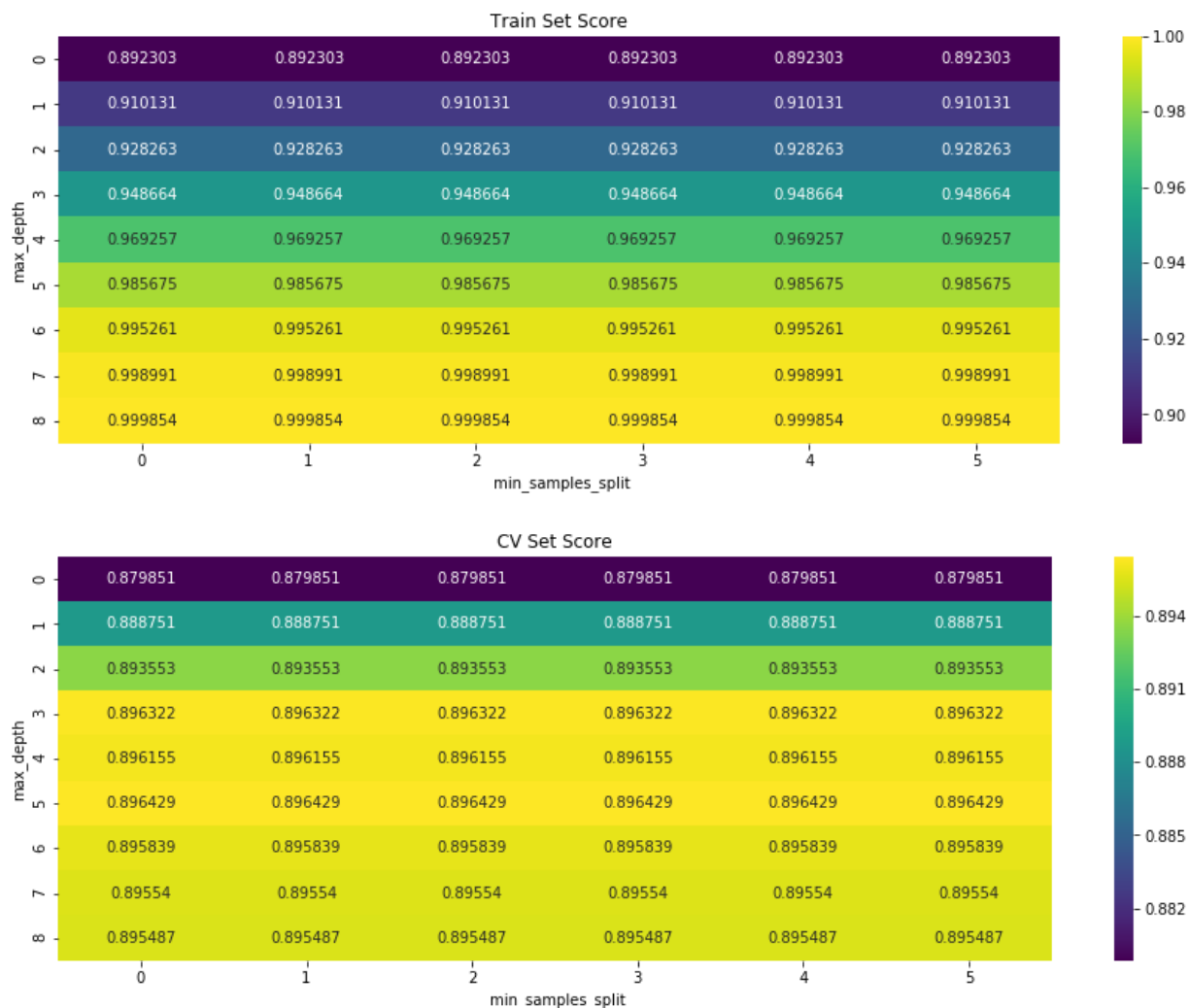Precision on test data: 0.9589768147172673
Recall on test data: 0.988482579902102
F1-Score on test data: 0.9735061778407941

Confusion Matrix of Train and Test set:
 [ [TN  FP]
   [FN TP] ]

Confusion Matrix of Train Set

Confusion Matrix of Test Set



## [5.2.4] Applying XGBOOST on TFIDF W2V, SET 4

In [0]: 
```
XGB(np.array(train_tfidf_sent_vectors), np.array(cv_tfidf_sent_vectors))
```
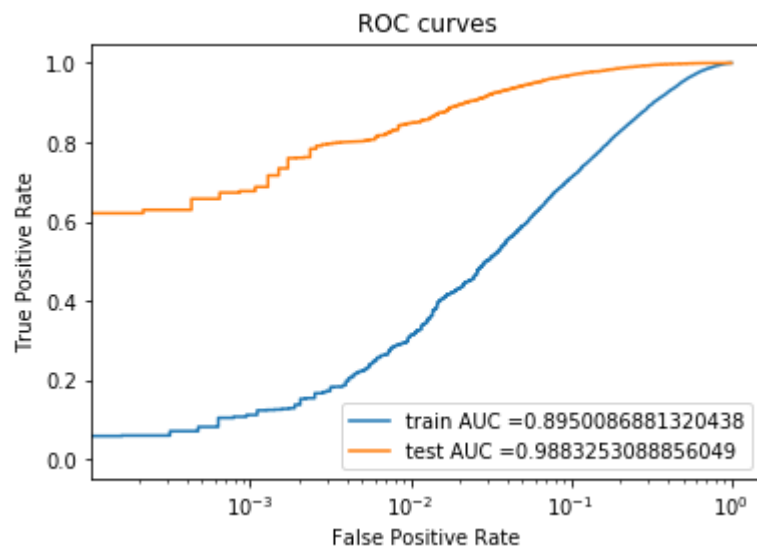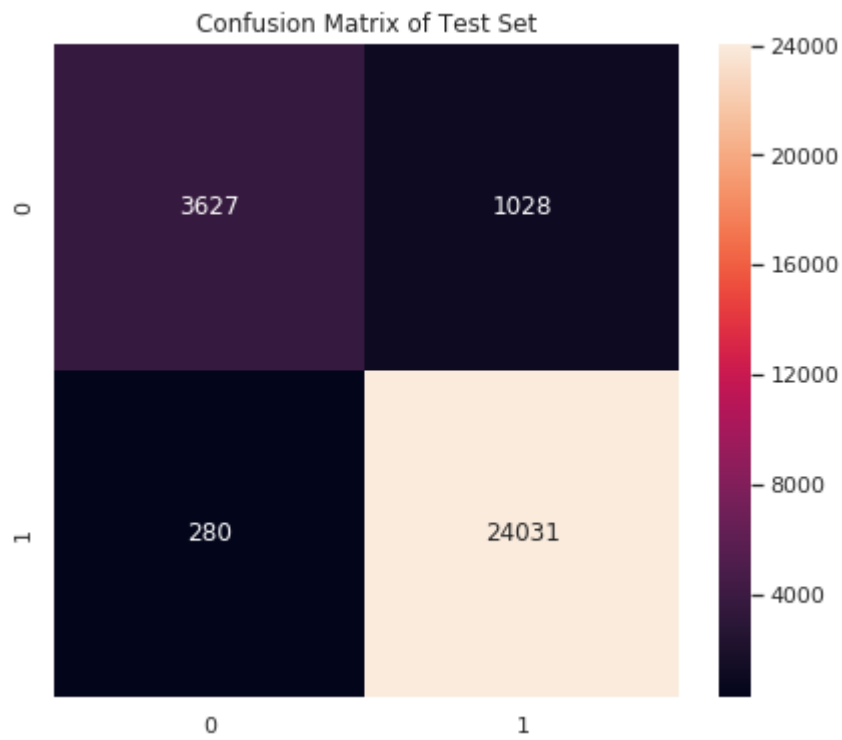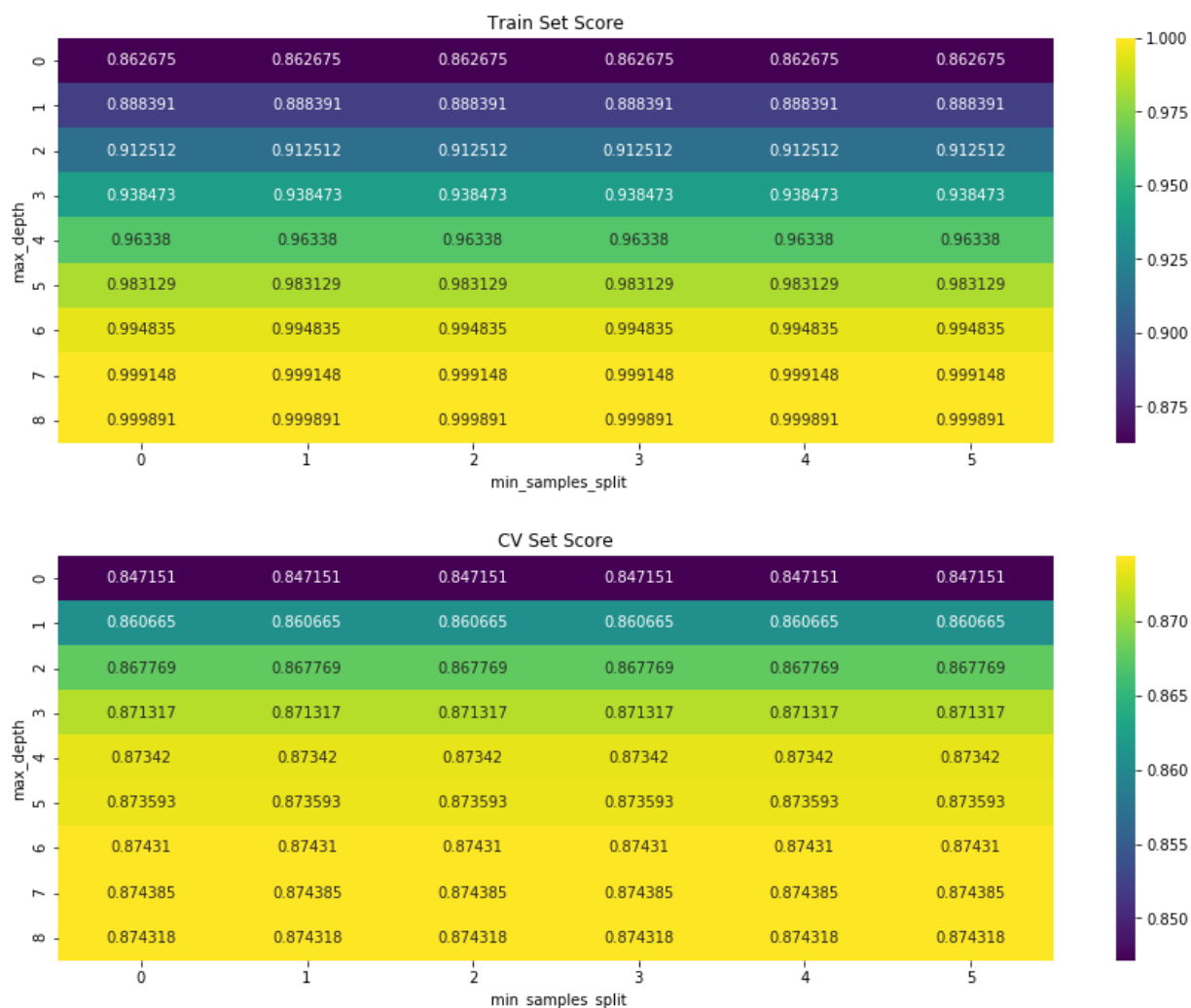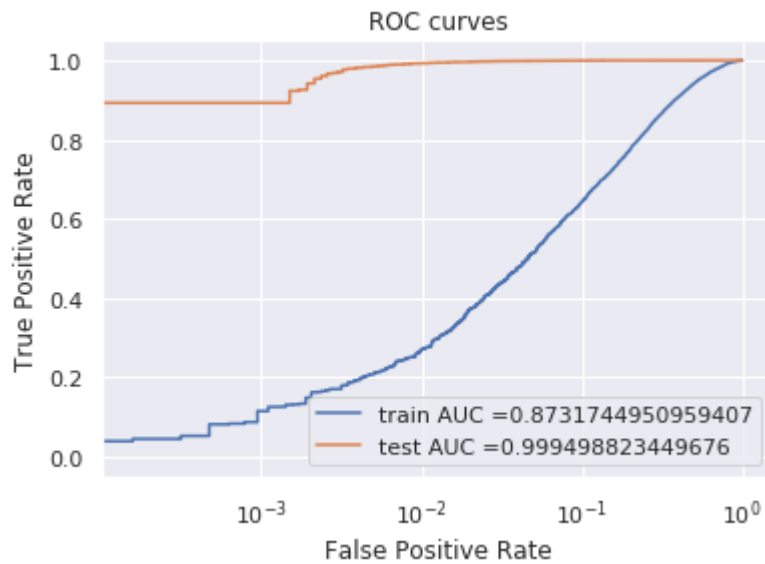
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
        colsample_bynode=1, colsample_bytree=1, gamma=0,
        learning_rate=0.1, max_delta_step=0, max_depth=9,
        min_child_weight=1, min_samples_split=5, missing=None,
        n_estimators=100, n_jobs=1, nthread=None,
        objective='binary:logistic', random_state=0, reg_alpha=0,
        reg_lambda=1, scale_pos_weight=1, seed=None, silent=None,
        subsample=1, verbosity=1)
_____
Best HyperParameter:  {'max_depth': 9, 'min_samples_split': 5}
Best Accuracy: 87.43851381680938

Train Set Score

| max_depth | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.862675 | 0.862675 | 0.862675 | 0.862675 | 0.862675 | 0.862675 |
| 1 | 0.888391 | 0.888391 | 0.888391 | 0.888391 | 0.888391 | 0.888391 |
| 2 | 0.912512 | 0.912512 | 0.912512 | 0.912512 | 0.912512 | 0.912512 |
| 3 | 0.938473 | 0.938473 | 0.938473 | 0.938473 | 0.938473 | 0.938473 |
| 4 | 0.96338 | 0.96338 | 0.96338 | 0.96338 | 0.96338 | 0.96338 |
| 5 | 0.983129 | 0.983129 | 0.983129 | 0.983129 | 0.983129 | 0.983129 |
| 6 | 0.994835 | 0.994835 | 0.994835 | 0.994835 | 0.994835 | 0.994835 |
| 7 | 0.999148 | 0.999148 | 0.999148 | 0.999148 | 0.999148 | 0.999148 |
| 8 | 0.999891 | 0.999891 | 0.999891 | 0.999891 | 0.999891 | 0.999891 |

min_samples_split

CV Set Score

| max_depth | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 0.847151 | 0.847151 | 0.847151 | 0.847151 | 0.847151 | 0.847151 |
| 1 | 0.860665 | 0.860665 | 0.860665 | 0.860665 | 0.860665 | 0.860665 |
| 2 | 0.867769 | 0.867769 | 0.867769 | 0.867769 | 0.867769 | 0.867769 |
| 3 | 0.871317 | 0.871317 | 0.871317 | 0.871317 | 0.871317 | 0.871317 |
| 4 | 0.87342 | 0.87342 | 0.87342 | 0.87342 | 0.87342 | 0.87342 |
| 5 | 0.873593 | 0.873593 | 0.873593 | 0.873593 | 0.873593 | 0.873593 |
| 6 | 0.87431 | 0.87431 | 0.87431 | 0.87431 | 0.87431 | 0.87431 |
| 7 | 0.874385 | 0.874385 | 0.874385 | 0.874385 | 0.874385 | 0.874385 |
| 8 | 0.874318 | 0.874318 | 0.874318 | 0.874318 | 0.874318 | 0.874318 |

min_samples_split

In [0]: `testing_XGB(np.array(train_tfidf_sent_vectors), np.array(test_tfidf_sent_vectors),9,5)`

### ROC curves



AUC score on test Data with max_depth = 9 and min_samples_split = 5 is : 97.69385930376515 %
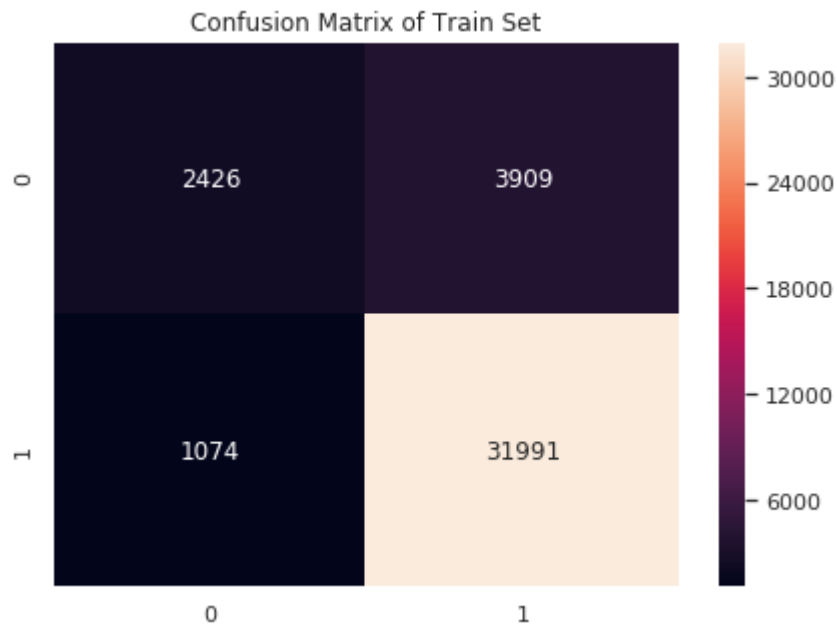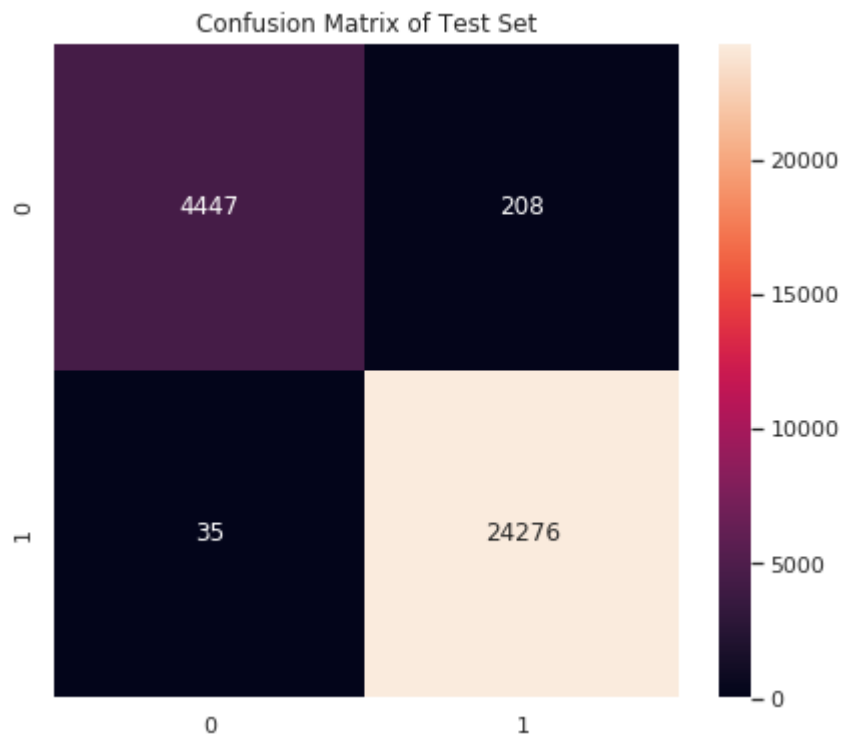Precision on test data: 0.9915046561019442
Recall on test data: 0.9985603224877627
F1-Score on test data: 0.9950199815554872

Confusion Matrix of Train and Test set:
 [ [TN  FP]
   [FN TP] ]

### Confusion Matrix of Train Set

Confusion Matrix of Test Set

| | 0 | 1 |
|---|---|---|
| **0** | 4447 | 208 |
| **1** | 35 | 24276 |

# [6] Conclusions

In [0]:
```
from prettytable import PrettyTable
```

In [66]:

```python
x = PrettyTable(["Vectorizer", "Model", "Max_Depth", "Min_samples_split", "Train_AUC","Test_AUC"])
y = PrettyTable(["Vectorizer", "Model", "Max_Depth", "Min_samples_split", "Train_AUC","Test_AUC"])

print("Random Forest:")
x.add_row(["BoW","Random Forest" , "10", "500", 79.87, 75.35])
x.add_row(["Tf-Idf", "Random Forest" , "10", "500", 81.87, 73.27])
x.add_row(["AVG_W2V", "Random Forest" , "10", "100", 86.76, 85.82])
x.add_row(["TFIDF_W2V", "Random Forest" , "7", "100", 84.48, 81.70])

print(x)
print("\n")

print("GBDT using XGBOOST:")
y.add_row(["BoW", "GBDT " , "10", "5", 91.46, 80.47])
y.add_row(["Tf-Idf", "GBDT " , "10", "5", 92.17, 82.57])
y.add_row(["AVG_W2V", "GBDT " , "7", "5", 89.64, 88.38])
y.add_row(["TFIDF_W2V", "GBDT " , "9", "5", 87.43, 97.69])

print(y)
```

Random Forest:

```
+------------+---------------+-----------+-------------------+-----------+----------+
| Vectorizer |     Model     | Max_Depth | Min_samples_split | Train_AUC | Test_AUC |
+------------+---------------+-----------+-------------------+-----------+----------+
|    BoW     | Random Forest |    10     |        500        |   79.87   |  75.35   |
|   Tf-Idf   | Random Forest |    10     |        500        |   81.87   |  73.27   |
|  AVG_W2V   | Random Forest |    10     |        100        |   86.76   |  85.82   |
| TFIDF_W2V  | Random Forest |     7     |        100        |   84.48   |   81.7   |
+------------+---------------+-----------+-------------------+-----------+----------+
```

GBDT using XGBOOST:

```
+------------+-------+-----------+-------------------+-----------+----------+
| Vectorizer | Model | Max_Depth | Min_samples_split | Train_AUC | Test_AUC |
+------------+-------+-----------+-------------------+-----------+----------+
|    BoW     | GBDT  |    10     |         5         |   91.46   |  80.47   |
|   Tf-Idf   | GBDT  |    10     |         5         |   92.17   |  82.57   |
|  AVG_W2V   | GBDT  |     7     |         5         |   89.64   |  88.38   |
| TFIDF_W2V  | GBDT  |     9     |         5         |   87.43   |  97.69   |
+------------+-------+-----------+-------------------+-----------+----------+
```

# Observation:

**1) Random Forest with TF-IDF and GBDT using XGBOOST with BOW has given the lowest Test_Error.**

**2) Random Forest of TFIDF_W2V with depth = 9 and GBDT using XGBOOST of AVG_W2V with depth = 7 is the lowest depth.**