

Import necessary libraries

```
In [1]: 1 import warnings
        2 warnings.filterwarnings('ignore')
```

```
In [3]: 1 from sklearn.tree import DecisionTreeClassifier
        2 import seaborn as sns
        3 from sklearn.model_selection import TimeSeriesSplit
        4 from scipy.sparse import *
        5 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
        6 from sklearn.preprocessing import StandardScaler
        7 from sklearn.metrics import *
        8 import pickle
        9 from tqdm import tqdm
       10 import numpy as np
       11 import matplotlib.pyplot as plt
       12 import pandas as pd
       13 from sklearn.model_selection import train_test_split
       14 import graphviz
       15 from sklearn.tree import export_graphviz
       16 from sklearn.externals import joblib
       17 #from prettytable import PrettyTable
       18 from wordcloud import WordCloud
```

Load preprocessed data

```
In [42]: 1 #Functions to save objects for later use and retrieve it
2 def savetofile(obj,filename):
3     pickle.dump(obj,open(filename+".pkl","wb"))
4 def openfromfile(filename):
5     temp = pickle.load(open(filename+".pkl","rb"))
6     return temp
7 y_train =openfromfile('y_train')
8 y_test =openfromfile('y_test')
9
10 count_vect =openfromfile('count_vect')
11 X_train_bigram = openfromfile('X_train_bigram')
12 X_test_bigram = openfromfile('X_test_bigram')
13
14 tf_idf_vect =openfromfile('tf_idf_vect')
15 X_train_tfidf =openfromfile('X_train_tfidf')
16 X_test_tfidf =openfromfile('X_test_tfidf')
17
18 avg_sent_vectors=openfromfile('avg_sent_vectors')
19 avg_sent_vectors_test=openfromfile('avg_sent_vectors_test')
20
21 tfidf_sent_vectors=openfromfile('tfidf_sent_vectors')
22 tfidf_sent_vectors_test=openfromfile('tfidf_sent_vectors_test')
```

Save and Load Model

```
In [4]: 1 def saveModeltofile(obj,filename):
2     joblib.dump(obj,open(filename+".pkl","wb"))
3 def openModelfromfile(filename):
4     temp = joblib.load(open(filename+".pkl","rb"))
5     return temp
```

Standardizing data

```
In [44]: 1 def std_data(train,test,mean):  
2         scaler=StandardScaler(with_mean=mean)  
3         std_train=scaler.fit_transform(train)  
4         std_test=scaler.transform(test)  
5         return std_train, std_test
```

Observation:

1. In Decision Tree based algorithm we are not dealing with distance at all.
2. So, Data Standardization is not required for DecisionTree.

Decision Trees

Function for hyperparameter tuning using corss validation and error plot using heatmap:

In [45]:

```

1 # find Optimal value of hyperparam by TimeSeriesSplit and 10_fold_cross_validation
2 # using RandomizedSearchCV and GridSearchCV.
3 def DT_Classifier(x_train,y_train,TBS,params,searchMethod,vect):
4     ''' FUNCTION FOR FINDING OPTIMAL VALUE OF HYPERPARAM AND DRAW HEATMAP WITH SCORE AND HYPERPARAM'''
5     #INITIALIZE DECISION-TREE CLASSIFIER
6     clf=DecisionTreeClassifier(class_weight='balanced',criterion='gini')
7     # APPLY RANDOM OR GRID SEARCH FOR HYPERPARAMETER TUNNING
8     if searchMethod=='grid':
9         model=GridSearchCV(clf,\
10                             cv=TBS,\
11                             n_jobs=-1,\
12                             param_grid=params,\
13                             return_train_score=True,\
14                             scoring=make_scorer(roc_auc_score,average='weighted'))
15         model.fit(x_train,y_train)
16     elif searchMethod=='random':
17         model=RandomizedSearchCV(clf,\
18                                 n_jobs=-1,\
19                                 cv=TBS,\
20                                 param_distributions=params,\
21                                 n_iter=len(params['max_depth']),\
22                                 return_train_score=True,\
23                                 scoring=make_scorer(roc_auc_score,average='weighted'))
24         model.fit(x_train,y_train)
25     #PLOT THE PERFORMANCE OF MODEL ON CROSSVALIDATION DATA FOR EACH HYPERPARAM VALUE
26     auc_results=[]
27     auc_results.append(model.cv_results_['mean_test_score'])
28     auc_results.append(model.cv_results_['mean_train_score'])
29     data=['CV','Train'];i=0;
30     plt.figure(figsize= (17,6))
31     for auc in auc_results:
32         auc=np.array(auc).reshape(len(params['min_samples_split']),len(params['max_depth']))
33         cv_auc_df=pd.DataFrame(auc, np.array(params['min_samples_split']),np.array(params['max_depth']))
34         cv_auc_df=cv_auc_df.round(4)
35         plt.subplot(int('12'+str(i+1)))
36         plt.title('Hyperparam Tunning %s-Data(%s)' %(data[i],vect))
37         sns.set(font_scale=1.4)#for label size
38         ax=sns.heatmap(cv_auc_df, annot=True,annot_kws={"size": 12}, fmt='g',)
39         ax.set(xlabel='Depth-Values', ylabel='Min #Samples to split')
40         i+=1
41     plt.show()
42     return model

```

Function which calculate performance on test data with optimal hyperparam :

```

In [175]: 1 # ===== Decision Tree with optimal depth and optimalmin_samples_split=====
2
3 def test_performance(x_train,y_train,x_test,y_test,optimal,vect,summarize):
4     '''FUNCTION FOR TEST PERFORMANCE(PLOT ROC CURVE FOR BOTH TRAIN AND TEST) WITH OPTIMAL HYPERPARAM'''
5     #INITIALIZE DECISION TREE WITH OPTIMAL VALUE OF HYPERPARAMS
6     clf=DecisionTreeClassifier(min_samples_split=optimal['min_samples_split'],\
7                               max_depth=optimal['max_depth'],\
8                               class_weight='balanced',\
9                               criterion='gini')
10    clf.fit(x_train,y_train)
11
12    train_prob=clf.predict_proba(x_train)[:,-1]
13    test_prob=clf.predict_proba(x_test)[:,-1]
14    y_pred=clf.predict(x_test)
15
16    fpr_test, tpr_test, threshold_test = roc_curve(y_test, test_prob,pos_label=1)
17    fpr_train, tpr_train, threshold_train = roc_curve(y_train, train_prob,pos_label=1)
18    auc_score_test=auc(fpr_test, tpr_test)
19    auc_score_train=auc(fpr_train, tpr_train)
20
21    f1=f1_score(y_test,y_pred,average='weighted')
22
23    #ADD RESULTS TO PRETTY TABLE
24    summarize.add_row([vect, optimal['max_depth'],optimal['min_samples_split'],\
25                      '%.3f' %auc_score_test,'%3f' %auc_score_train,'%3f' %f1])
26
27    plt.figure(1,figsize=(14,5))
28    plt.subplot(121)
29    plt.title('ROC Curve (%s)' %vect)
30    #IDEAL ROC CURVE
31    plt.plot([0,1],[0,1],'k--')
32    #ROC CURVE OF TEST DATA
33    plt.plot(fpr_test, tpr_test , 'b', label='Test_AUC= %.2f' %auc_score_test)
34    #ROC CURVE OF TRAIN DATA
35    plt.plot(fpr_train, tpr_train , 'g', label='Train_AUC= %.2f' %auc_score_train)
36    plt.xlim([-0.1,1.1])
37    plt.ylim([-0.1,1.1])
38    plt.xlabel('False Positive Rate')
39    plt.ylabel('True Positive Rate')
40    plt.grid(True)
41    plt.legend(loc='lower right')
42

```

```
43 #PLOT CONFUSION MATRIX USING HEATMAP
44 plt.subplot(122)
45 plt.title('Confusion-Matrix(Test Data)')
46 df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), ['Negative', 'Positive'], ['Negative', 'Positive'])
47 sns.set(font_scale=1.4)#for label size
48 sns.heatmap(df_cm, cmap='gist_earth', annot=True, annot_kws={"size": 16}, fmt='g')
49 plt.show()
50 return clf
51
```

Function which print top important features (feature importance):

```

In [5]: 1 def feature_importance(vectorizer,clf,n):
        2     '''FUNCTION FOR OVERALL FEATURE IMPORTANCE'''
        3     # FEATURE IMPORTANCES FROM DECISION TREE
        4     fe_importances = clf.feature_importances_
        5
        6     # INDEX OF SORTED IMPORTANT FEATURES
        7     indices = np.argsort(fe_importances)[::-1][:n]
        8
        9     # FEATURE NAMES
       10     names = vectorizer.get_feature_names()
       11     names=np.array(names)
       12
       13     #WORDCLOUD PLOT
       14     wordcloud = WordCloud(max_font_size=50, max_words=100,collocations=False).\
       15     generate(str(names[indices]))
       16     plt.figure(1,figsize=(14,13))
       17     plt.title("WordCloud(Important Feature)")
       18     plt.imshow(wordcloud, interpolation="bilinear")
       19     plt.axis("off")
       20
       21     #BAR CHART
       22     plt.figure(2,figsize=(13,8))
       23     sns.set(rc={'figure.figsize':(11.7,8.27)})
       24     plt.title("Feature Importance")
       25     # ADD BARS
       26     plt.bar(range(n), fe_importances[indices])
       27     # FEATURE NAMES ON X AXIS
       28     plt.xticks(range(n), names[indices], rotation=70)
       29     # Show plot
       30     plt.show()

```

Function for visualizing tree:

```

In [48]: 1 def visualizeDecisionTree(clf,vect,treeName):
        2     export_graphviz(clf, out_file=treeName+".dot",feature_names=vect.get_feature_names(),filled=True,proport

```

Initialization of common objects required for all vectorization:


```

In [102]: 1 #(best depth in range [1, 5, 10, 50, 100, 500], and the best min_samples_split in range [5, 10, 100, 500])
          2 #VECTORIZER
          3 vect=['Bow', 'TF-IDF', 'AVG-W2V', 'TFIDF-W2V']
          4 #OBJECT FOR TIMESERIES CROSS VALIDATION
          5 TBS=TimeSeriesSplit(n_splits=10)
          6 #METHOD USE FOR HYPER PARAMETER TUNNING
          7 searchMethod='grid'
          8 #RANGE OF VALUES FOR HYPERPARAM
          9 samples=[5, 10, 100, 500]
         10 depth=[1, 5, 10, 50, 100, 500]
         11 params={'max_depth':depth, 'min_samples_split':samples}
         12 #INITIALIZE PRETTY TABLE OBJECT
         13 summarize = PrettyTable()
         14 summarize.field_names = ['Vectorizer', 'Optimal-Depth', 'Optimal #Samples', 'Test(AUC)', 'Train(AUC)', 'Test(f

```

[1.1] Applying Decision Trees on BOW, SET 1

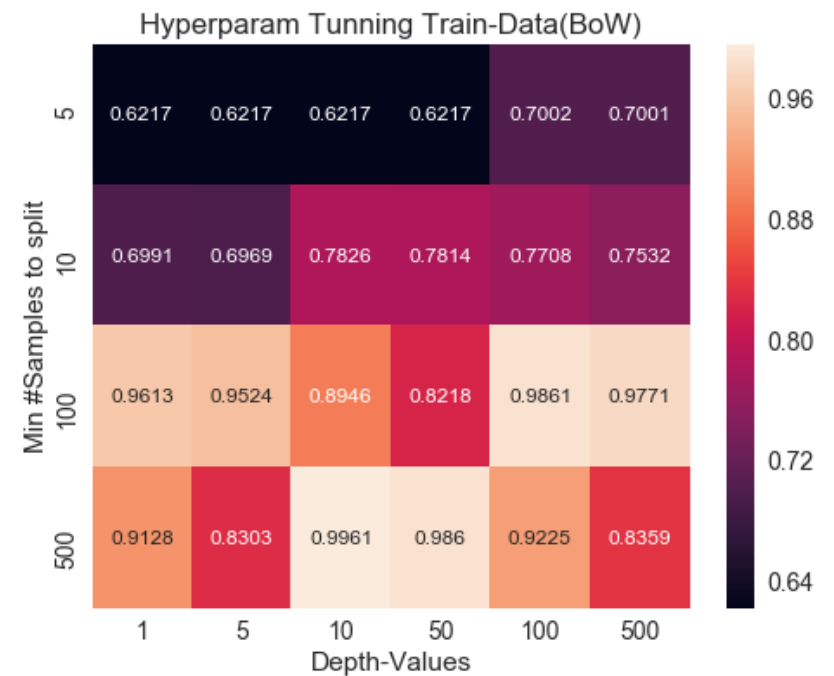
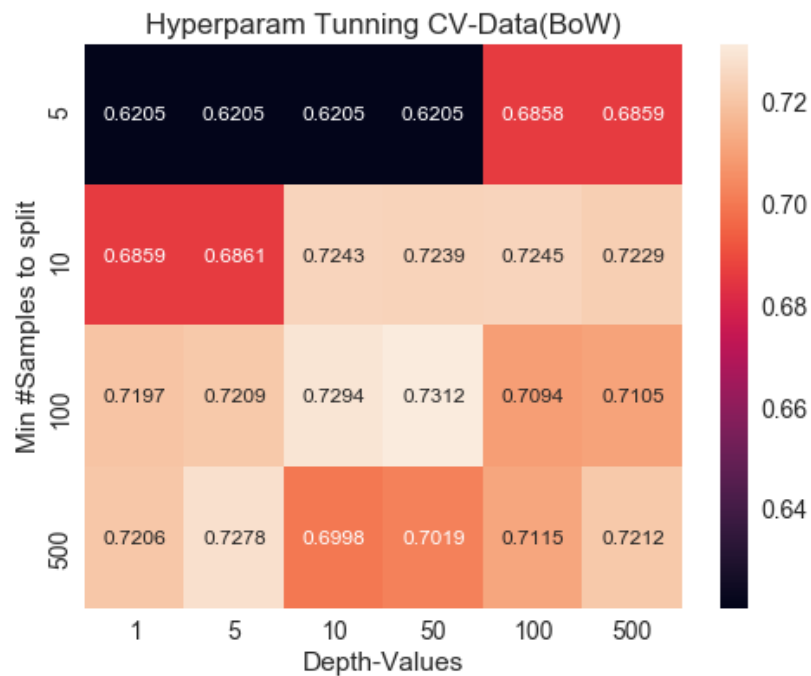
[1.1.1] Hyperparam tuning and plot Heatmap for hyperparam:

In [103]:

```

1 #TRAIN AND TEST DATA
2 train=X_train_bigram; test=X_test_bigram;
3 #HYPERPARAMS TUNNING
4 %time model=DT_Classifier(train,y_train,TBS,params,searchMethod,vect[0])
5 print(model.best_params_)
6 #SAVE CURRENT STATE OF MODEL
7 saveModeltofile(model,'model_bow_dt')

```



Wall time: 12min 8s

{ 'max_depth': 50, 'min_samples_split': 500 }

Observation:

1. From the above heatmaps we pick optimal value of hyperparam such that our model is not overfit and underfit.
2. We pick optimal value of hyperparams:

- a. max_depth=50
- b. min_smples_split=500

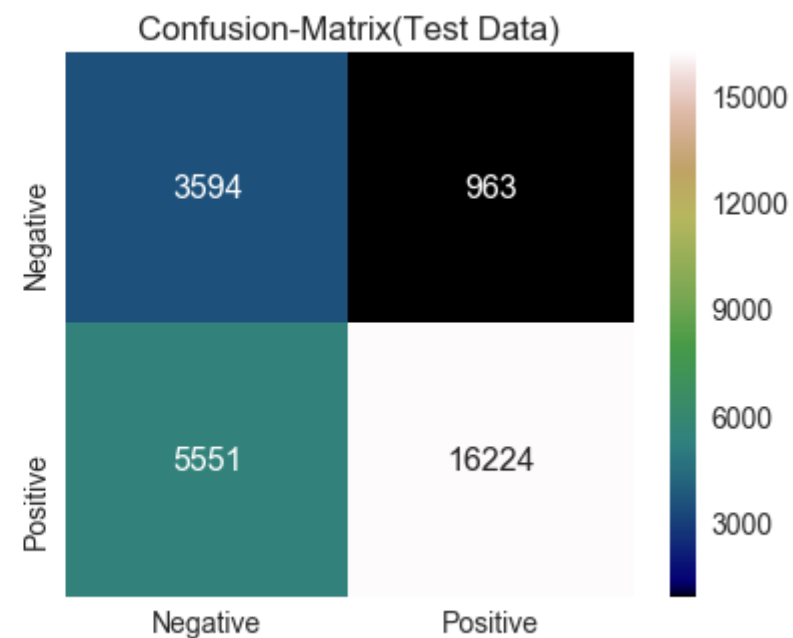
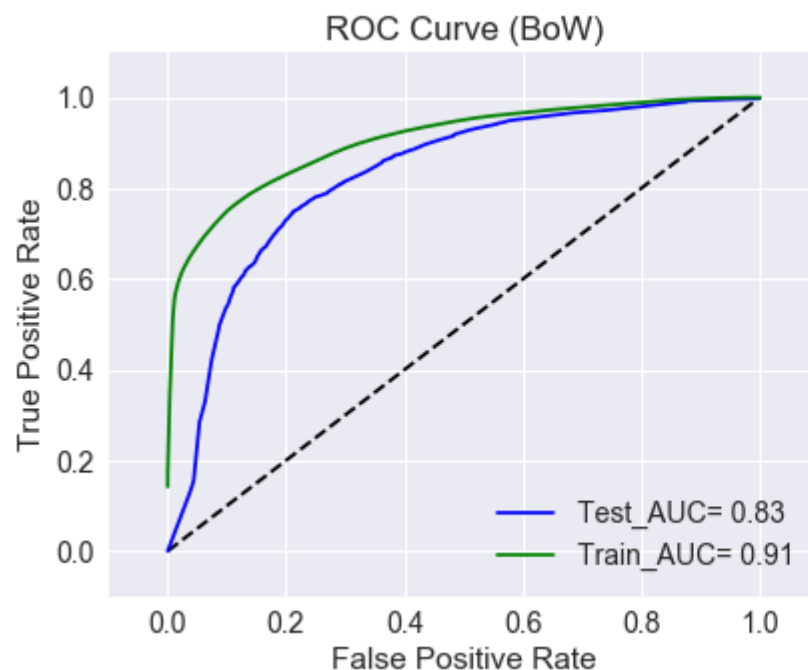
3. for those optimal hyperparam, model performance:

- a. auc for crossvalidation data is .7019 and
- b. auc for train data is .9860

```
In [140]: 1 best_params={'max_depth':50,'min_samples_split':500}
```

[1.1.2] Performance on test data with optimal value of hyperparam:

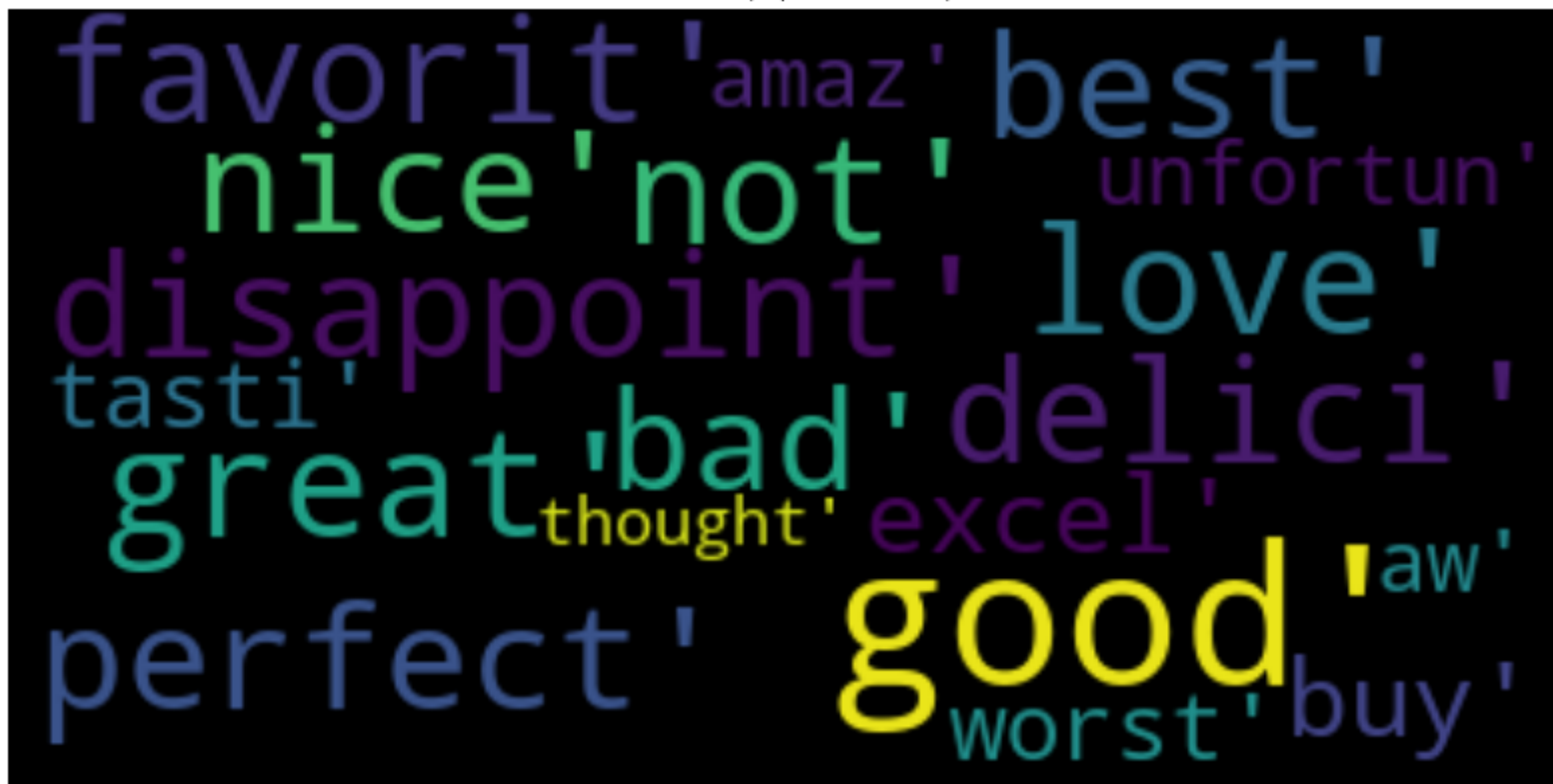
```
In [141]: 1 clf=test_performance(train,y_train,test,y_test,best_params,vect[0],summarize)
```

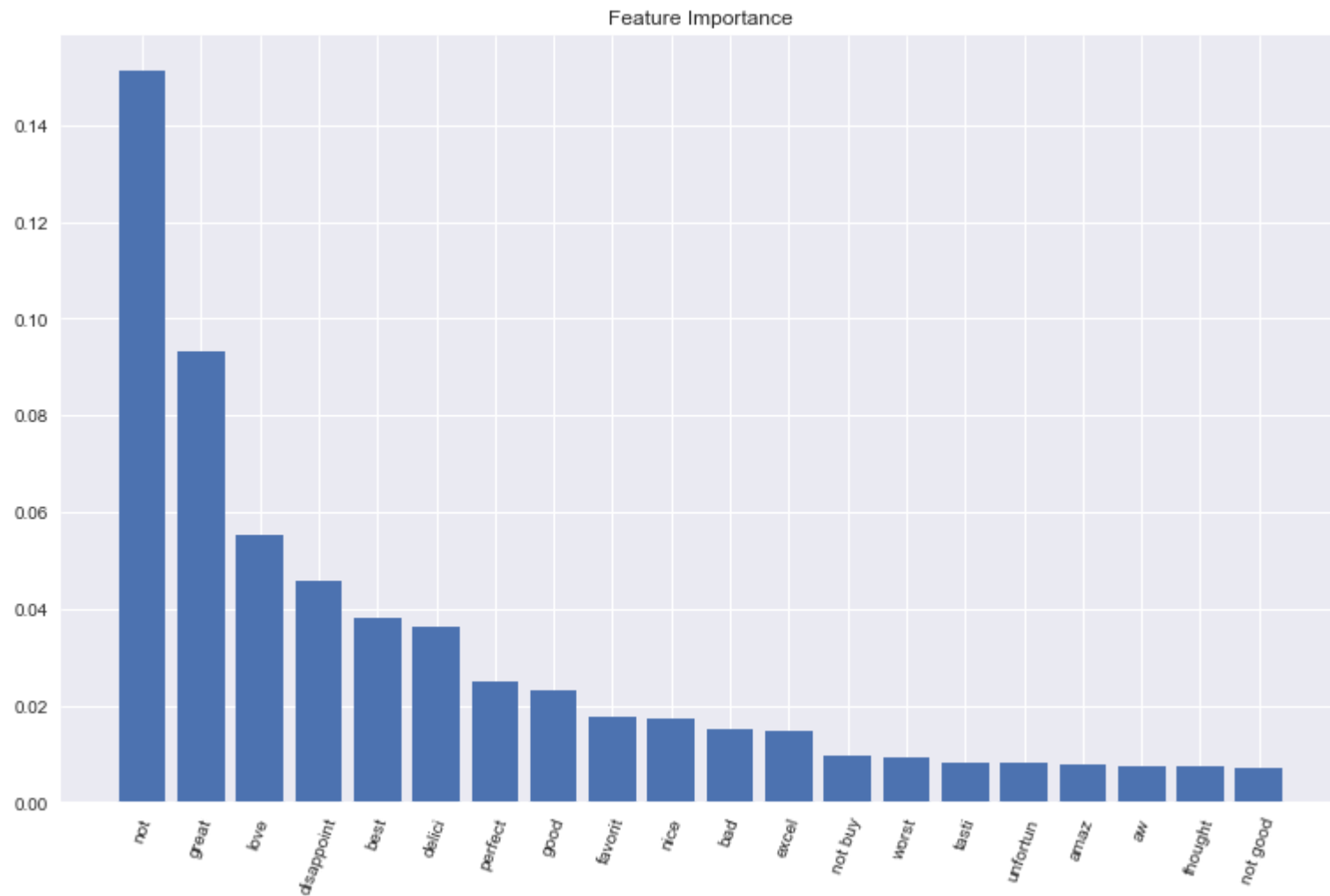


[1.1.3] Top 20 important features:

```
In [7]: 1 no_of_imp_features=20  
2 feature_importance(count_vect,clf,no_of_imp_features)
```

WordCloud(Important Feature)





Observation:

1. Decision Tree give feature importance based on reduction in entroy or gini impurity due to a feature in whole Decision Tree.
2. We can't get class based feature importance in Decision Tree.

[1.1.5]Graphviz visualization of Decision Tree:

```
In [144]: 1 treeName='decision_bow'
          2 visualizeDecisionTree(clf,count_vect,treeName)
          3 with open(treeName+'.dot') as f:
          4     dot_graph = f.read()
          5 graphviz.Source(dot_graph)
```

```
Out[144]: <graphviz.files.Source at 0x17862e10>
```

Note:

1. Please check the graph_bow.png to visualize the whole graph.

[2.1] Applying Decision Tree on TFIDF, SET 2

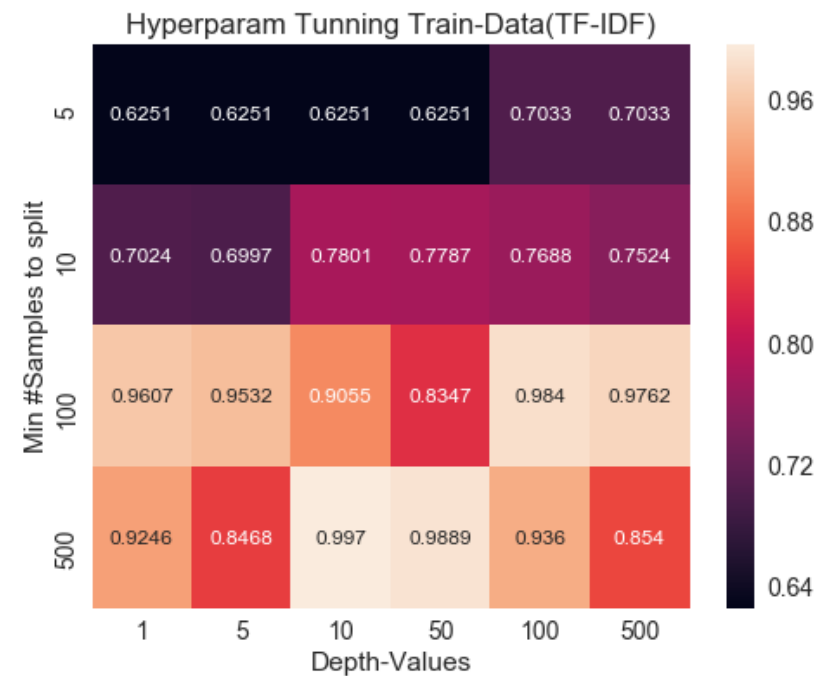
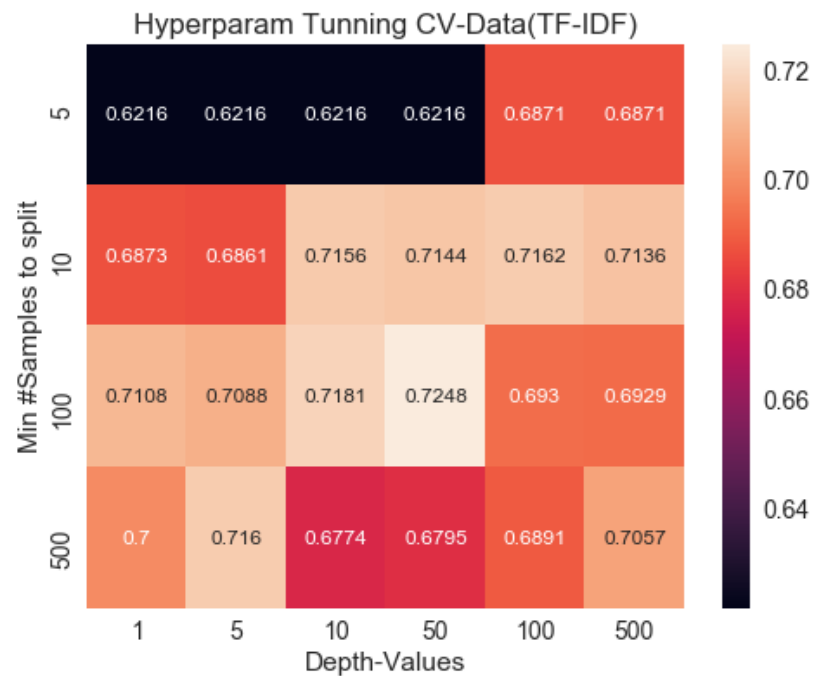
[2.1.1] Hyperparam tuning and plot Heatmap for hyperparam:

In [145]:

```

1 #TRAIN AND TEST DATA
2 train=X_train_tfidf; test=X_test_tfidf;
3 #HYPERPARAMS TUNNING
4 %time model=DT_Classifier(train,y_train,TBS,params,searchMethod,vect[1])
5 print(model.best_params_)
6 #SAVE CURRENT STATE OF MODEL
7 saveModeltofile(model,'model_tfidf_dt')

```



Wall time: 13min 25s

{ 'max_depth': 50, 'min_samples_split': 500 }

Observation:

1. From the above heatmaps we pick optimal value of hyperparam such that our model is not overfit and underfit.
2. We pick optimal value of hyperparams:

- a. max_depth=50
- b. min_smples_split=500

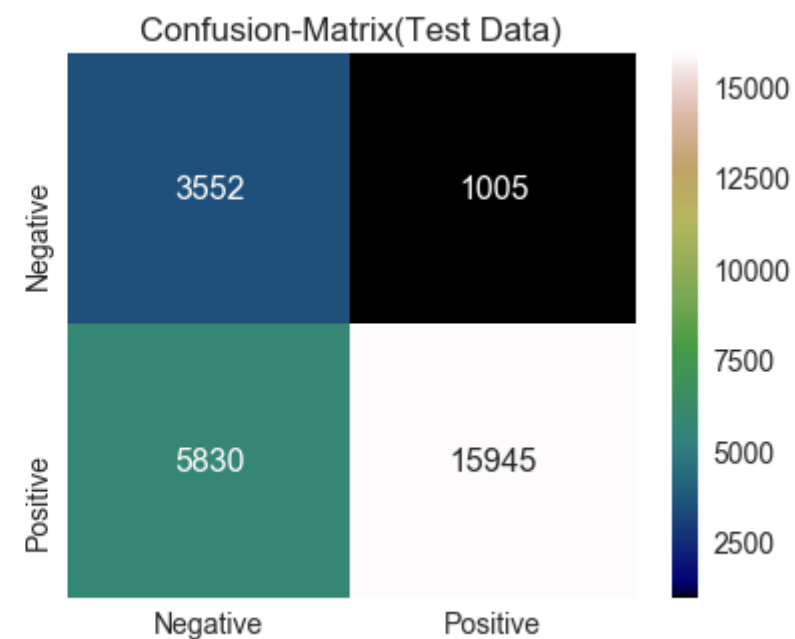
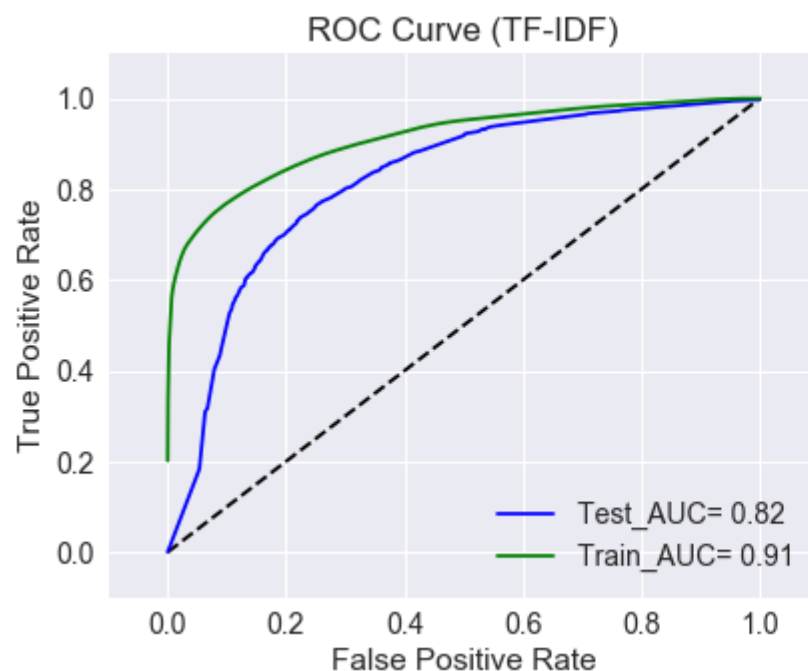
3. for those optimal hyperparam, model performance:

- a. auc for crossvalidation data is .6795 and
- b. auc for train data is .9889

```
In [161]: 1 best_params={'max_depth':50,'min_samples_split':500}
```

[2.1.2] Performance on test data with optimal value of hyperparam:

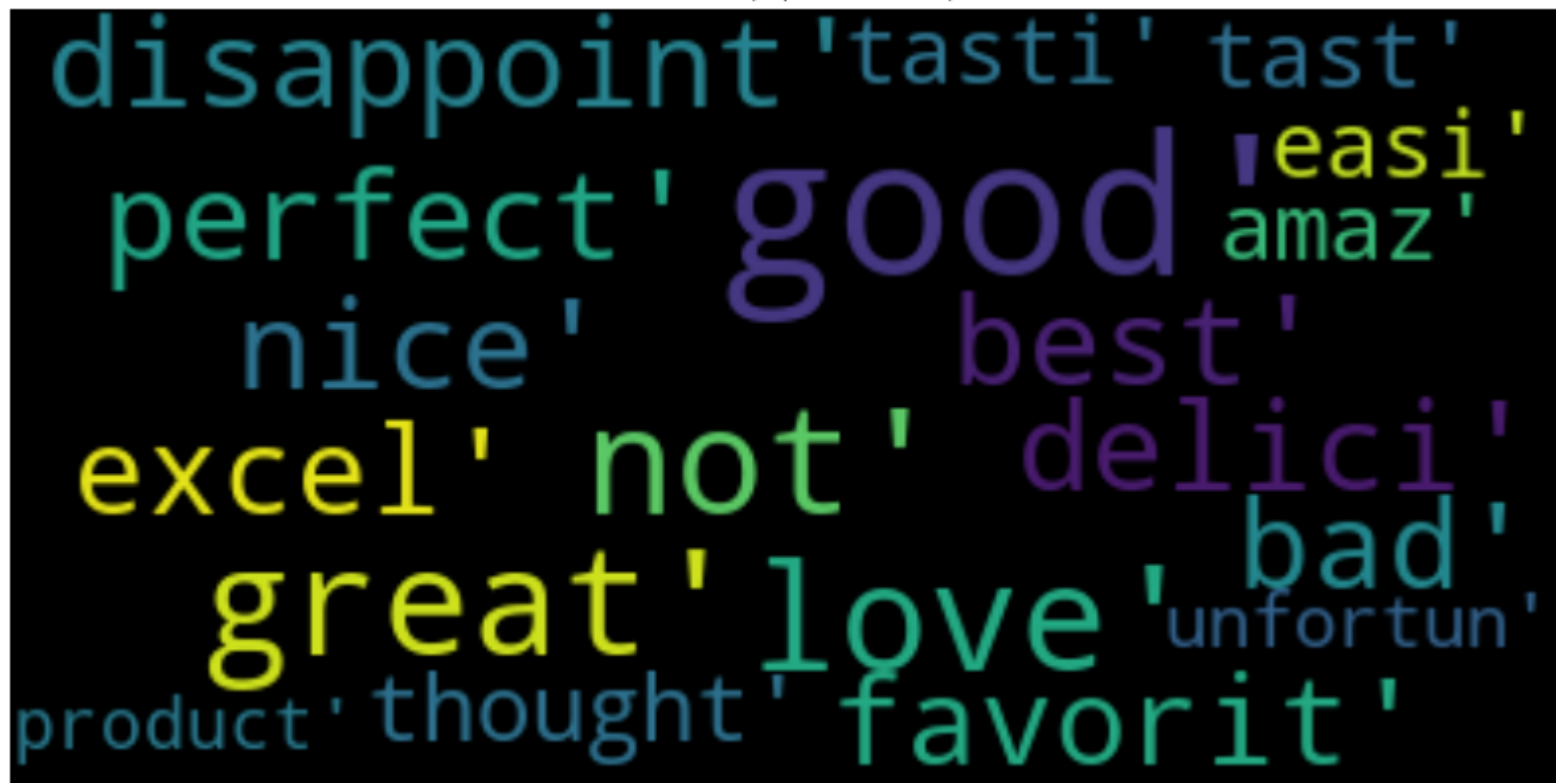
```
In [162]: 1 clf=test_performance(train,y_train,test,y_test,best_params,vect[1],summarize)
```

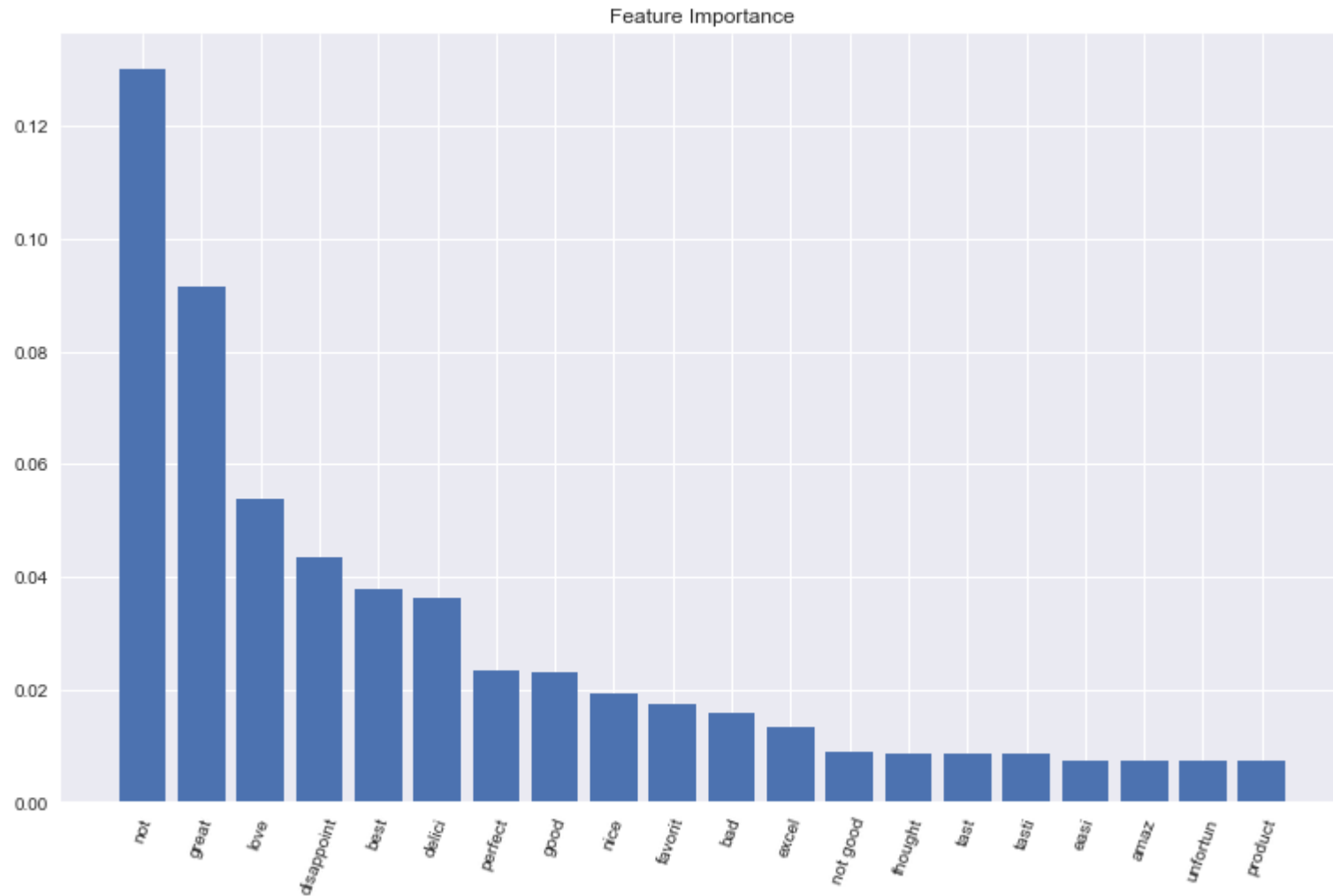


[2.1.3] Top 20 important features:


```
In [8]: 1 no_of_imp_features=20  
2 feature_importance(tf_idf_vect,clf,no_of_imp_features)
```

WordCloud(Important Feature)





Observation:

1. Decision Tree give feature importance based on reduction in entroy or gini impurity due to a feature in whole Decision Tree.
2. We can't get class based feature importance in Decision Tree.

[2.1.4] Graphviz visualization of Decision Tree:

```
In [163]: 1 treeName='decision_tfidf'
          2 visualizeDecisionTree(clf,tf_idf_vect,treeName)
          3 with open(treeName+'.dot') as f:
          4     dot_graph = f.read()
          5 graphviz.Source(dot_graph)
```

```
Out[163]: <graphviz.files.Source at 0x16fa9668>
```

Note:

1. Please check the graph_tfidf.png to visualize the whole graph.

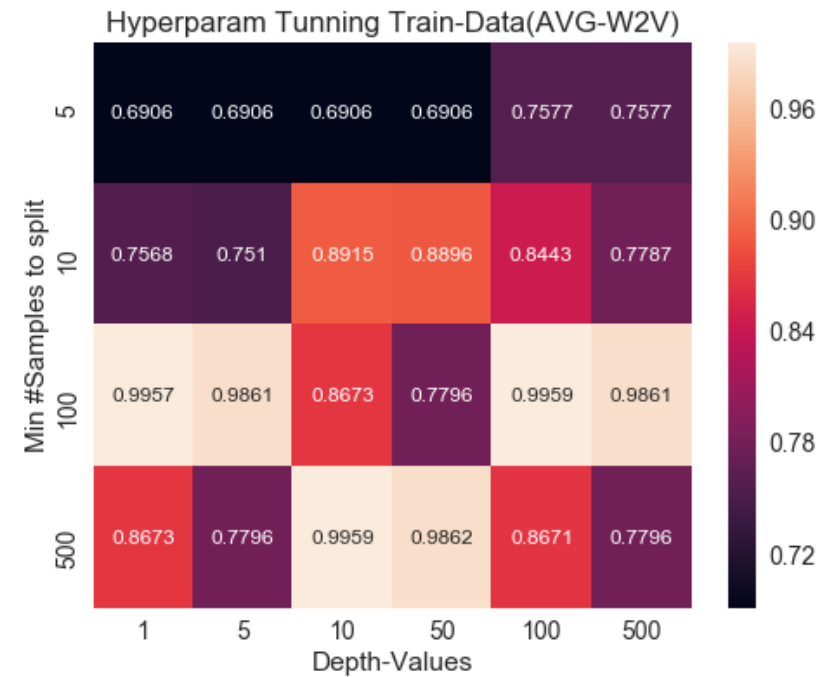
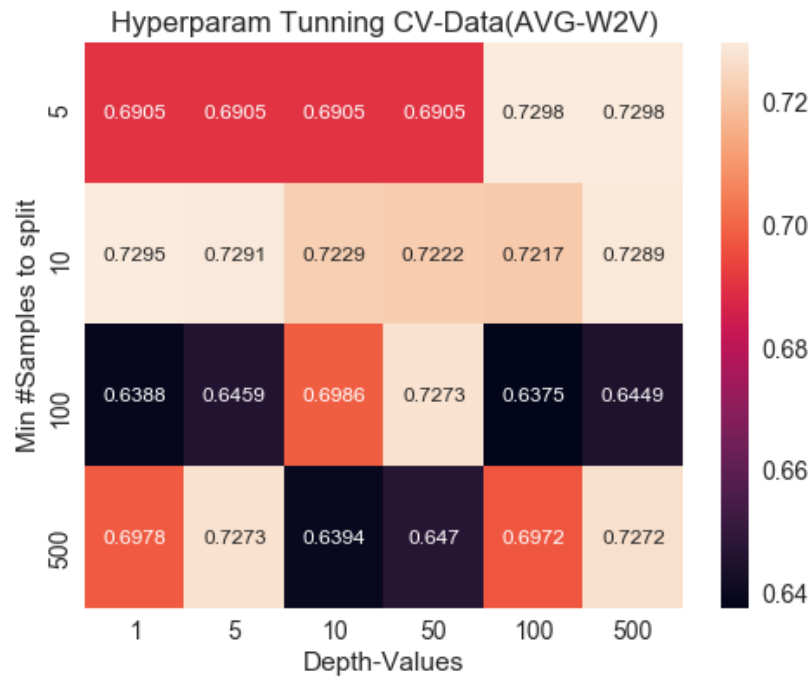
[3.1] Applying Decision Tree on AVG-W2V, SET 3**[3.1.1] Hyperparam tuning and plot Heatmap for hyperparam:**

In [164]:

```

1 #TRAIN AND TEST DATA
2 train=avg_sent_vectors; test=avg_sent_vectors_test;
3 #HYPERPARAMS TUNNING
4 %time model=DT_Classifier(train,y_train,TBS,params,searchMethod,vect[2])
5 print(model.best_params_)
6 #SAVE CURRENT STATE OF MODEL
7 saveModeltofile(model,'model_avgw2v_dt')

```



Wall time: 6min 11s

{ 'max_depth': 5, 'min_samples_split': 5 }

Observation:

1. From the above heatmaps we pick optimal value of hyperparam such that our model is not overfit and underfit.
2. We pick optimal value of hyperparams:

- a. max_depth=5
- b. min_smples_split=5

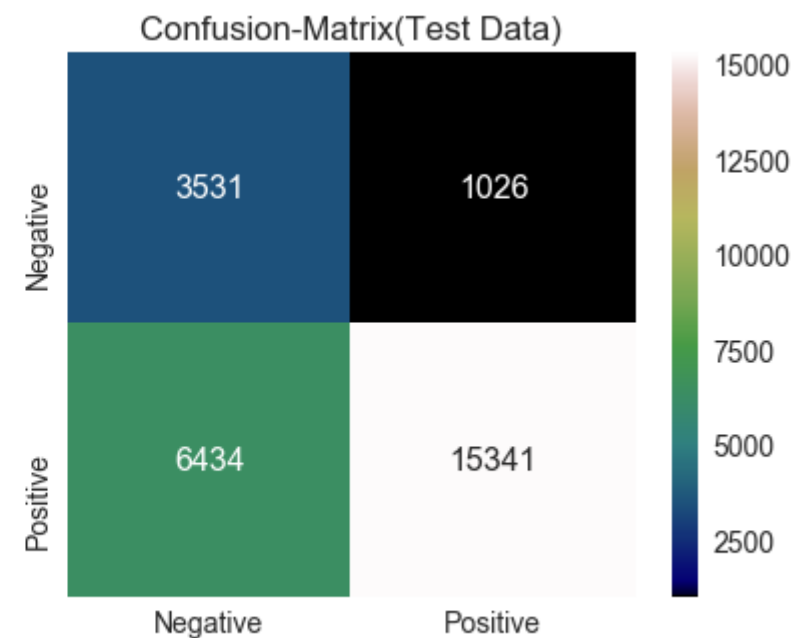
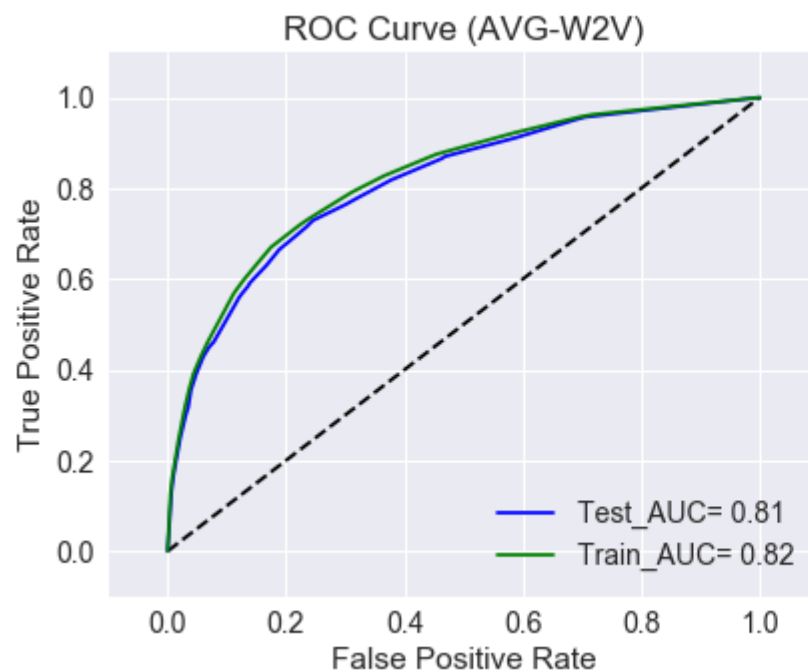
3. for those optimal hyperparam, model performance:

- a. auc for crossvalidation data is .6905 and
- b. auc for train data is .6906

```
In [173]: 1 best_params={'max_depth':5,'min_samples_split':5}
```

[3.1.2] Performance on test data with optimal value of hyperparam:

```
In [176]: 1 clf=test_performance(train,y_train,test,y_test,best_params,vect[2],summarize)
```



[4.1] Applying Decision Tree on TFIDF-W2V, SET 4

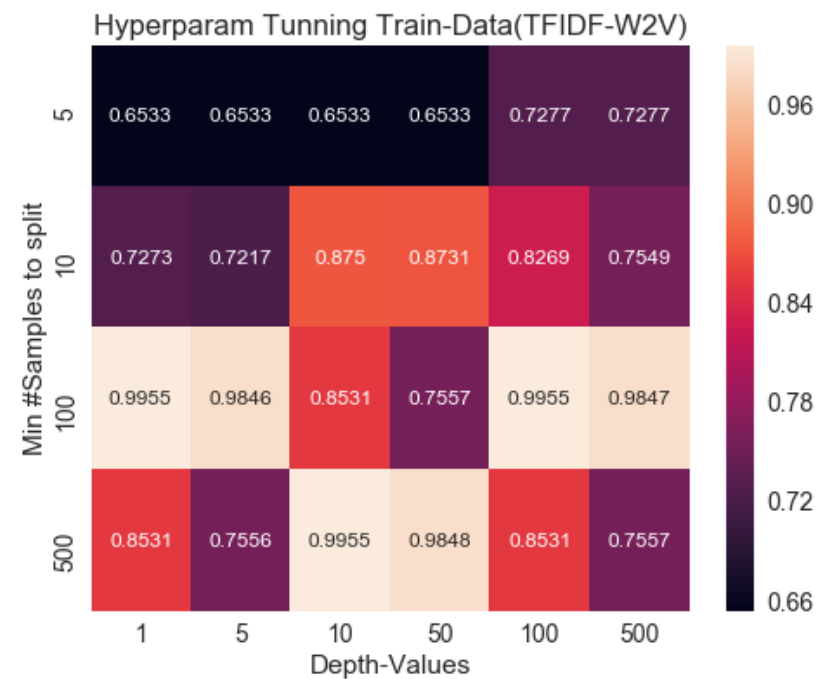
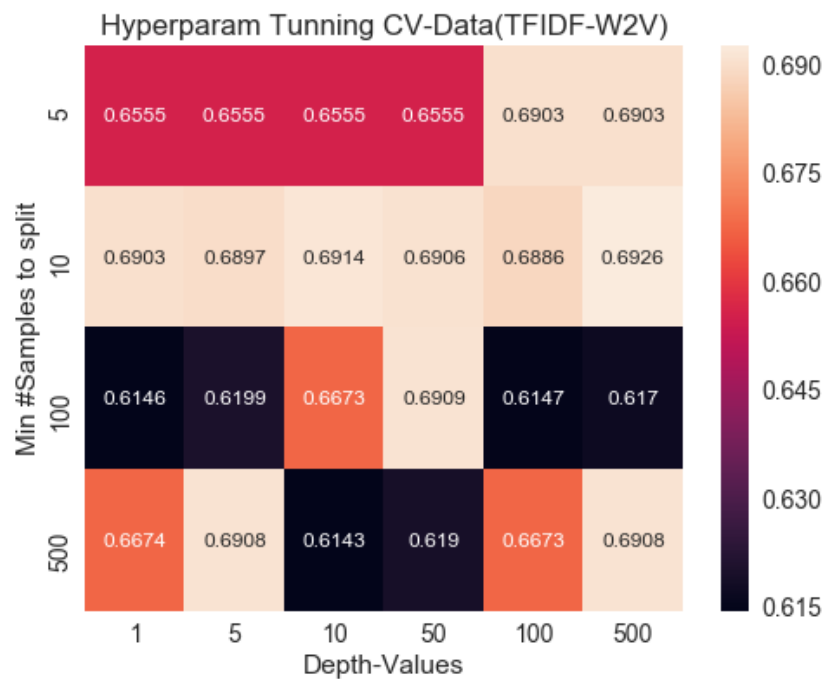
[4.1.1] Hyperparam tuning and plot Heatmap for hyperparam:

In [178]:

```

1 # TRAIN AND TEST DATA
2 train=tfidf_sent_vectors; test=tfidf_sent_vectors_test;
3 #HYPERPARAMS TUNNING
4 %time model=DT_Classifier(train,y_train,TBS,params,searchMethod,vect[3])
5 print(model.best_params_)
6 #SAVE CURRENT STATE OF MODEL
7 saveModeltofile(model,'model_tfidf2v_dt')

```



Wall time: 6min 10s

{ 'max_depth': 10, 'min_samples_split': 500 }

Observation:

1. From the above heatmaps we pick optimal value of hyperparam such that our model is not overfit and underfit.
2. We pick optimal value of hyperparams:

- a. max_depth=10
- b. min_smples_split=500

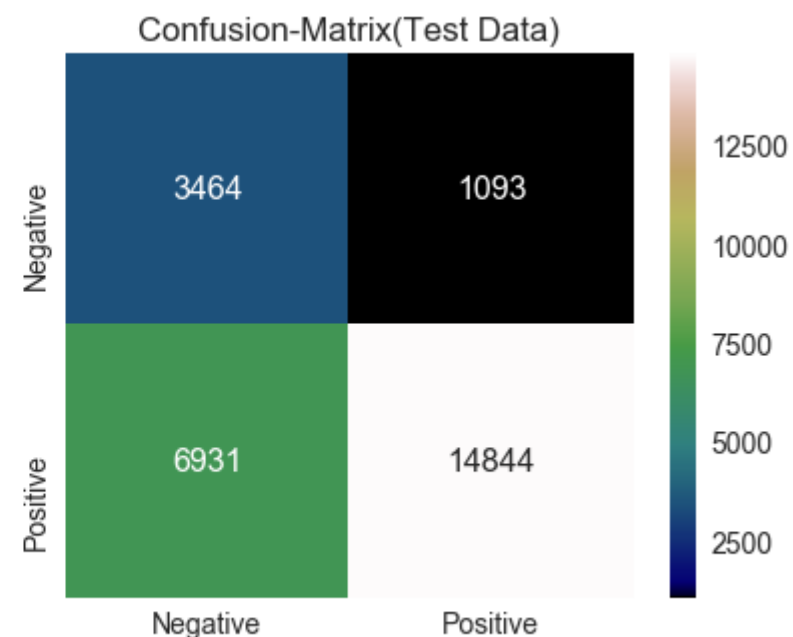
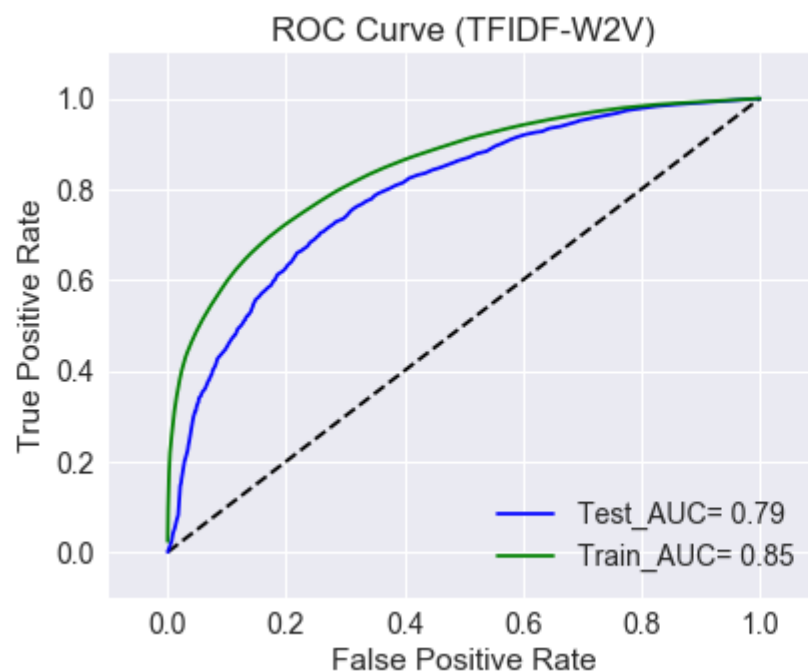
3. for those optimal hyperparam, model performance:

- auc for crossvalidation data is .6143 and
- auc for train data is .9955

```
In [179]: 1 best_params={'max_depth':10,'min_samples_split':500}
```

[4.1.2] Performance on test data with optimal value of hyperparam:

```
In [180]: 1 clf=test_performance(train,y_train,test,y_test,best_params,vect[3],summarize)
```



Observation:

1. We can't get feature importance in case of AVG-W2V and TFIDF-W2V vectorizer.
2. Because in case of AVG-W2V and TFIDF-W2V vectors the dimension of a vector is not represented by a feature.

Conclusion:

In [182]: 1 `print(summarize)`

Vectorizer	Optimal-Depth	Optimal #Samples	Test(AUC)	Train(AUC)	Test(f1-score)
BoW	50	500	0.828	0.905	0.779
TF-IDF	50	500	0.820	0.913	0.771
AVG-W2V	5	5	0.811	0.822	0.749
TFIDF-W2V	10	500	0.792	0.846	0.731

1. from the above table we can observe that the optimal performance is give by:

- a. BoW vectorizer
- b. f1-score=.779 and auc=.828

In []: 1