

## Import necessary libraries

```
In [3]: 1 import warnings
        2 warnings.filterwarnings('ignore')
```

```
In [4]: 1 import xgboost as xgb
        2 #IMPORTING XGBOOST WRAPPER OF SKLEARN
        3 from xgboost.sklearn import XGBClassifier
        4 import seaborn as sns
        5 from sklearn.ensemble import RandomForestClassifier
        6 from sklearn.model_selection import TimeSeriesSplit
        7 from scipy.sparse import *
        8 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
        9 from sklearn.preprocessing import StandardScaler
       10 from sklearn.metrics import *
       11 import pickle
       12 from tqdm import tqdm
       13 import numpy as np
       14 import matplotlib.pyplot as plt
       15 import pandas as pd
       16 from sklearn.model_selection import train_test_split
       17 from prettytable import PrettyTable
       18 from wordcloud import WordCloud
```

## Load preprocessed data

In [37]:

```

1  #Functions to save objects for later use and retireve 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')

```

## Standardizing data

In [6]:

```

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, GBDT and Random-Forest.

## Ensemble Models

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

In [7]:

```

1  # find Optimal value of hyperparam by TimeSeriesSplit and 10_fold_cross_validation
2  # using RandomizedSearchCV and GridSearchCV.
3  def Ensemble_Classifier(x_train,y_train,TBS,params,searchMethod,vect,classifier):
4      if classifier=='random_forest':
5          #INITIALIZE RANDOM-FOREST CLASSIFIER
6          clf=RandomForestClassifier(class_weight='balanced',\
7                                     criterion='gini',\
8                                     oob_score=True)
9
10     elif classifier == 'gbdt':
11         #INITIALIZE GBDT CLASSIFIER
12         clf=xgb.XGBClassifier(nthread=8,\
13                               learning_rate=.1,\
14                               gamma=0,\
15                               subsample=.8,\
16                               colsample_bytree=.8,\
17                               booster='gbtree',\
18                               objective='binary:logistic')
19     # APPLY RANDOM OR GRID SEARCH FOR HYPERPARAMETER TUNNING
20     if searchMethod=='grid':
21         model=GridSearchCV(clf,\
22                             cv=TBS,\
23                             n_jobs=-1,\
24                             param_grid=params,\
25                             return_train_score=False,\
26                             scoring=make_scorer(roc_auc_score,average='weighted'))
27         model.fit(x_train,y_train)
28     elif searchMethod=='random':
29         model=RandomizedSearchCV(clf,\
30                                 n_jobs=-1,\
31                                 cv=TBS,\
32                                 param_distributions=params,\
33                                 n_iter=len(params['max_depth']),\
34                                 return_train_score=False,\
35                                 scoring=make_scorer(roc_auc_score,average='weighted'))
36         model.fit(x_train,y_train)
37     #PLOT THE PERFORMANCE OF MODEL ON CROSSVALIDATION DATA FOR EACH HYPERPARAM VALUE
38     cv_auc=model.cv_results_['mean_test_score']
39     cv_auc=np.array(cv_auc).reshape(len(params['max_depth']),len(params['n_estimators']))
40     cv_auc_df=pd.DataFrame(cv_auc, np.array(params['max_depth']),np.array(params['n_estimators']))
41     cv_auc_df=cv_auc_df.round(4)
42     plt.figure(figsize= (10,8))

```

```
42 plt.title('Hyperparam Tunning Results(%)' %vect)
43 sns.set(font_scale=1.4)#for label size
44 ax=sns.heatmap(cv_auc_df, annot=True,annot_kws={"size": 15}, fmt='g',)
45 ax.set(xlabel='No. of Estimators', ylabel='Depth-Values')
46 plt.show()
47 return model
```

**Function which calculate performance on test data with optimal hyperparam :**

```

In [8]: 1 def test_performance(x_train,y_train,x_test,y_test,optimal,vect,summarize,classifier):
2         '''FUNCTION FOR TEST PERFORMANCE(PLOT ROC CURVE FOR BOTH TRAIN AND TEST) WITH OPTIMAL HYPERPARAM'''
3         if classifier=='random_forest':
4             #INITIALIZE RANDOM FOREST WITH OPTIMAL VALUE OF HYPERPARAMS
5             clf=RandomForestClassifier(n_estimators=optimal['n_estimators'],\
6                                       max_depth=optimal['max_depth'],\
7                                       class_weight='balanced',\
8                                       criterion='gini',\
9                                       oob_score=True,\
10                                      n_jobs=-1)
11         elif classifier=='gbdt':
12             #INITIALIZE GBDT WITH OPTIMAL VALUE OF HYPERPARAMS
13             clf=xgb.XGBClassifier(n_estimators=optimal['n_estimators'],\
14                                  max_depth=optimal['max_depth'],\
15                                  nthread=8,\
16                                  learning_rate=.1,\
17                                  gamma=0,\
18                                  subsample=.8,\
19                                  colsample_bytree=.8,\
20                                  booster='gbtree',\
21                                  objective='binary:logistic',\
22                                  n_jobs=-1)
23         clf.fit(x_train,y_train)
24         train_prob=clf.predict_proba(x_train)[:,-1]
25         test_prob=clf.predict_proba(x_test)[:,-1]
26         fpr_test, tpr_test, threshold_test = roc_curve(y_test, test_prob,pos_label=1)
27         fpr_train, tpr_train, threshold_train = roc_curve(y_train, train_prob,pos_label=1)
28         auc_score_test=auc(fpr_test, tpr_test)
29         auc_score_train=auc(fpr_train, tpr_train)
30         y_pred=clf.predict(x_test)
31         f1=f1_score(y_test,y_pred,average='weighted')
32
33         #ADD RESULTS TO PRETTY TABLE
34         summarize.add_row([vect, optimal['max_depth'],optimal['n_estimators'], '%.3f' %auc_score_test,'%f1'])
35
36         plt.figure(1,figsize=(14,5))
37         plt.subplot(121)
38         plt.title('ROC Curve (%s)' %vect)
39         #IDEAL ROC CURVE
40         plt.plot([0,1],[0,1],'k--')
41         #ROC CURVE OF TEST DATA

```

```
42 plt.plot(fpr_test, tpr_test , 'b', label='Test_AUC= %.2f' %auc_score_test)
43 #ROC CURVE OF TRAIN DATA
44 plt.plot(fpr_train, tpr_train , 'g', label='Train_AUC= %.2f' %auc_score_train)
45 plt.xlim([-0.1,1.1])
46 plt.ylim([-0.1,1.1])
47 plt.xlabel('False Positive Rate')
48 plt.ylabel('True Positive Rate')
49 plt.grid(True)
50 plt.legend(loc='lower right')
51
52 #PLOT CONFUSION MATRIX USING HEATMAP
53 plt.subplot(122)
54 plt.title('Confusion-Matrix(Test Data)')
55 df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), ['Negative', 'Positive'], ['Negative', 'Positive'])
56 sns.set(font_scale=1.4)#for label size
57 sns.heatmap(df_cm,cmap='gist_earth', annot=True,annot_kws={"size": 16}, fmt='g')
58 plt.show()
59 return clf
```

**Function which print top important features (feature importance):**

```

In [9]: 1 def feature_importance(vectorizer,clf,n):
2         '''FUNCTION FOR FEATURE IMPORTANCE AND PLOT CORRESPONDING IMPORTANT FEATURES IN WORDCLOUD AND BARCHART'''
3         #CALCULATE FEATURE IMPORTANCES FROM ENSEMBLE MODEL
4         importances = clf.feature_importances_
5
6         # SORT FEATURE IMPORTANCES IN DECENDING ORDER
7         indices = np.argsort(importances)[::-1][:n]
8
9         # Rearrange feature names so they match the sorted feature importances
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        # Create plot title
25        plt.title("Feature Importance")
26        # Add bars
27        plt.bar(range(n), importances[indices])
28        # Add feature names as x-axis labels
29        plt.xticks(range(n), names[indices], rotation=70)
30        # Show plot
31        plt.show()

```

**Initialization of common objects required for all vectorization:**



```
In [10]: 1 #ENSEMBLE MODEL TO USE
2 classifier=['random_forest','gbdt']
3 #VECTORIZER
4 vect=['BoW','TF-IDF','AVG-W2V','TFIDF-W2V']
5 #OBJECT FOR TIMESERIES CROSS VALIDATION
6 TBS=TimeSeriesSplit(n_splits=10)
7 #METHOD USE FOR HYPER PARAMETER TUNNING
8 searchMethod='grid'
9 #RANGE OF VALUES FOR HYPERPARAM
10 estimators=[16,32,64,128,256,512]
11 depth=[9,16,18,20]
12 params={'max_depth':depth,'n_estimators':estimators}
13 #INITIALIZE PRETTY TABLE OBJECT
14 summarize = PrettyTable()
15 summarize.field_names = ['Vectorizer', 'Optimal-Depth', 'Optimal #Estimators', 'Test(AUC)', 'Test(f1-score)']
```

## Observation:

1. In case of Random Forest, base learners are high variance(low train error) models.
2. Here the base learners are Decision Tree and Decision Tree is having high variance when the tree is of significant depth(high depth).
3. So here we are taking high depth values in our hyperparam as well.

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

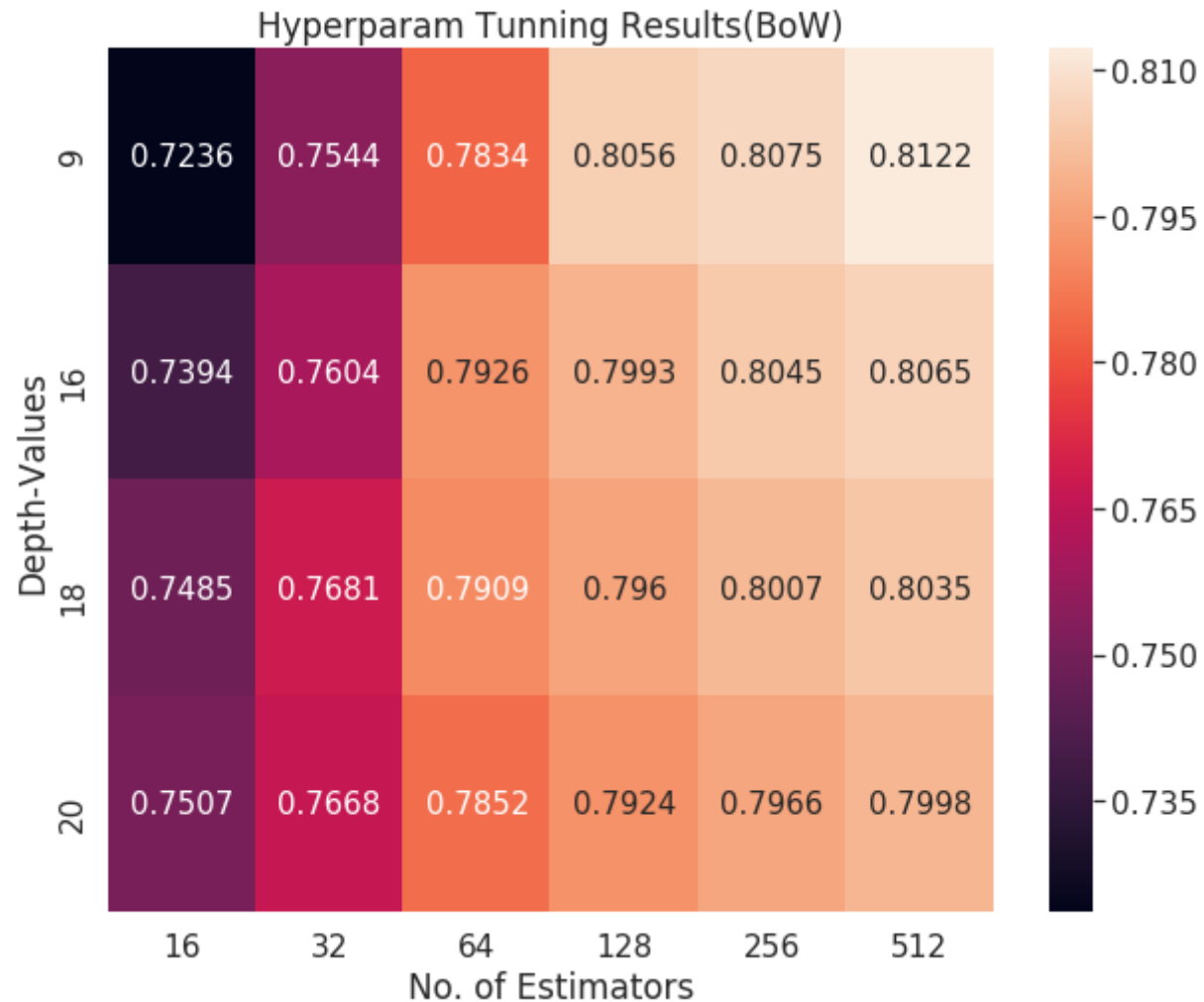
### [1.1.1] Hyperparam tuning and plot Heatmap for hyperparam:

In [82]:

```

1 #TRAIN AND TEST DATA
2 train=X_train_bigram;test=X_test_bigram;
3 #default search tech='GridSearch'
4 %time model=Ensemble_Classifier(train,y_train,TBS,params,searchMethod,vect[0],classifier[0])
5 print(model.best_params_)

```



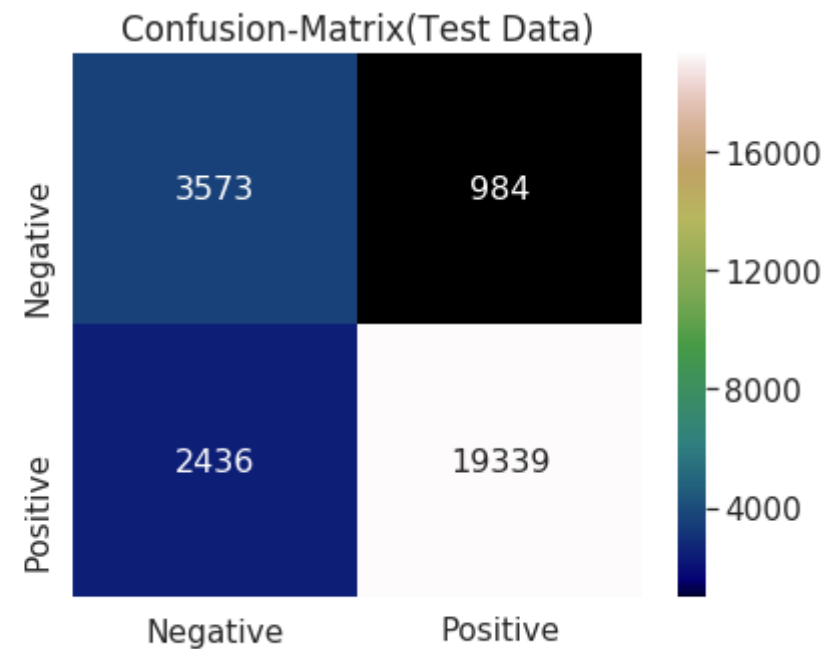
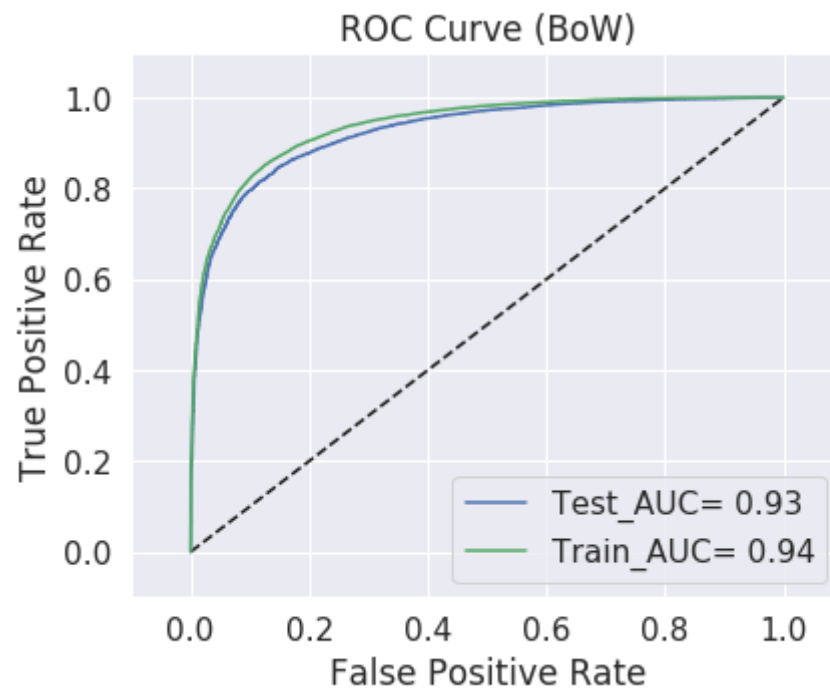
CPU times: user 49.3 s, sys: 788 ms, total: 50 s

Wall time: 5min 47s

{'n\_estimators': 512, 'max\_depth': 9}

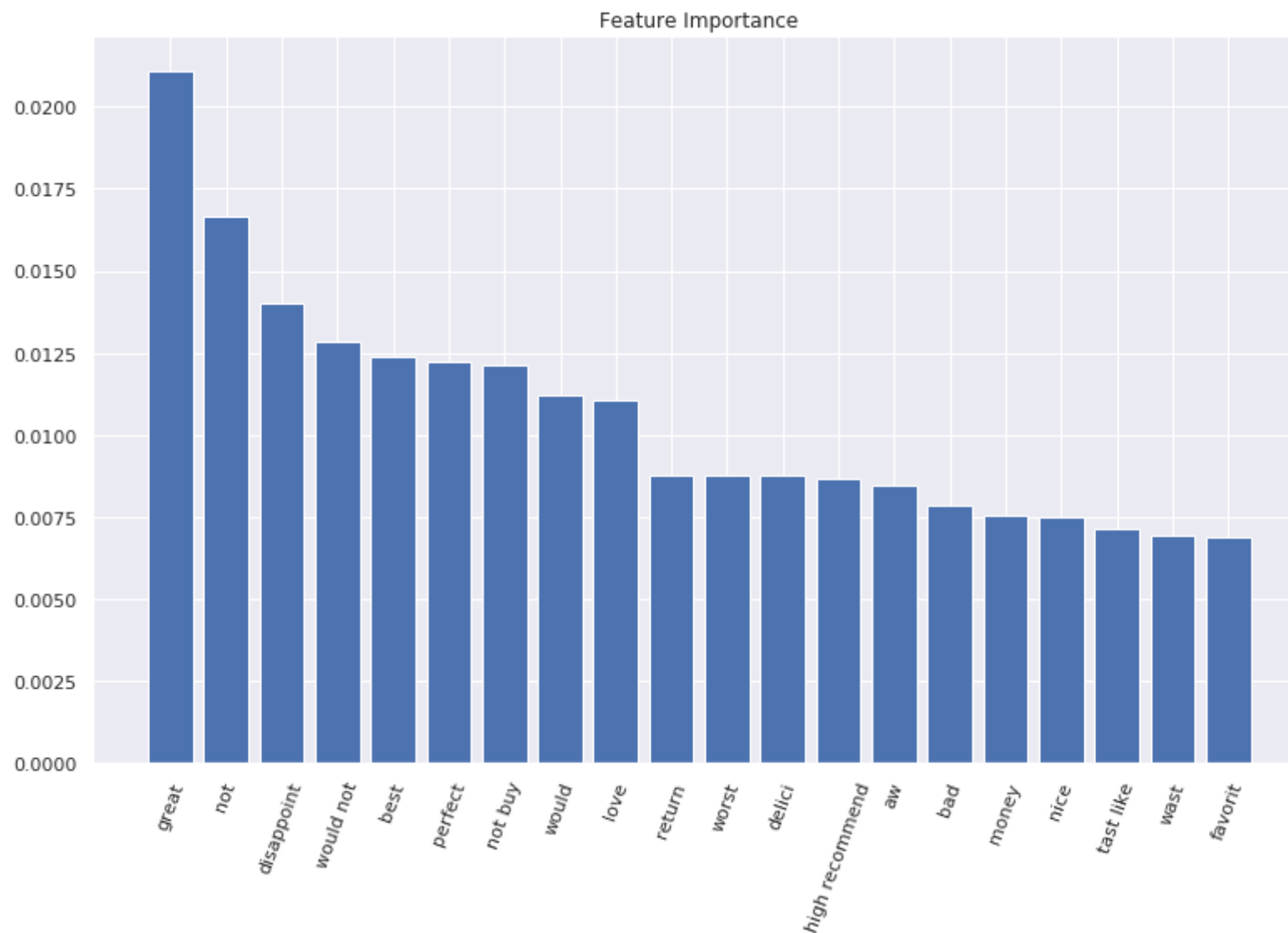
**[1.1.2] Performance on test data with optimal value of hyperparam:**

```
In [83]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_,vect[0],summarize,classifier[0])
```

**[1.1.3] Top 20 important features:**

```
In [84]: 1 no_of_imp_features=20  
2 feature_importance(count_vect,clf,no_of_imp_features)
```





### Observation:

1. Random Forest is good at overall interpretation.
2. `feature_importances_` provides the overall important feature.
3. We can't get class based feature importance in Random Forest.

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

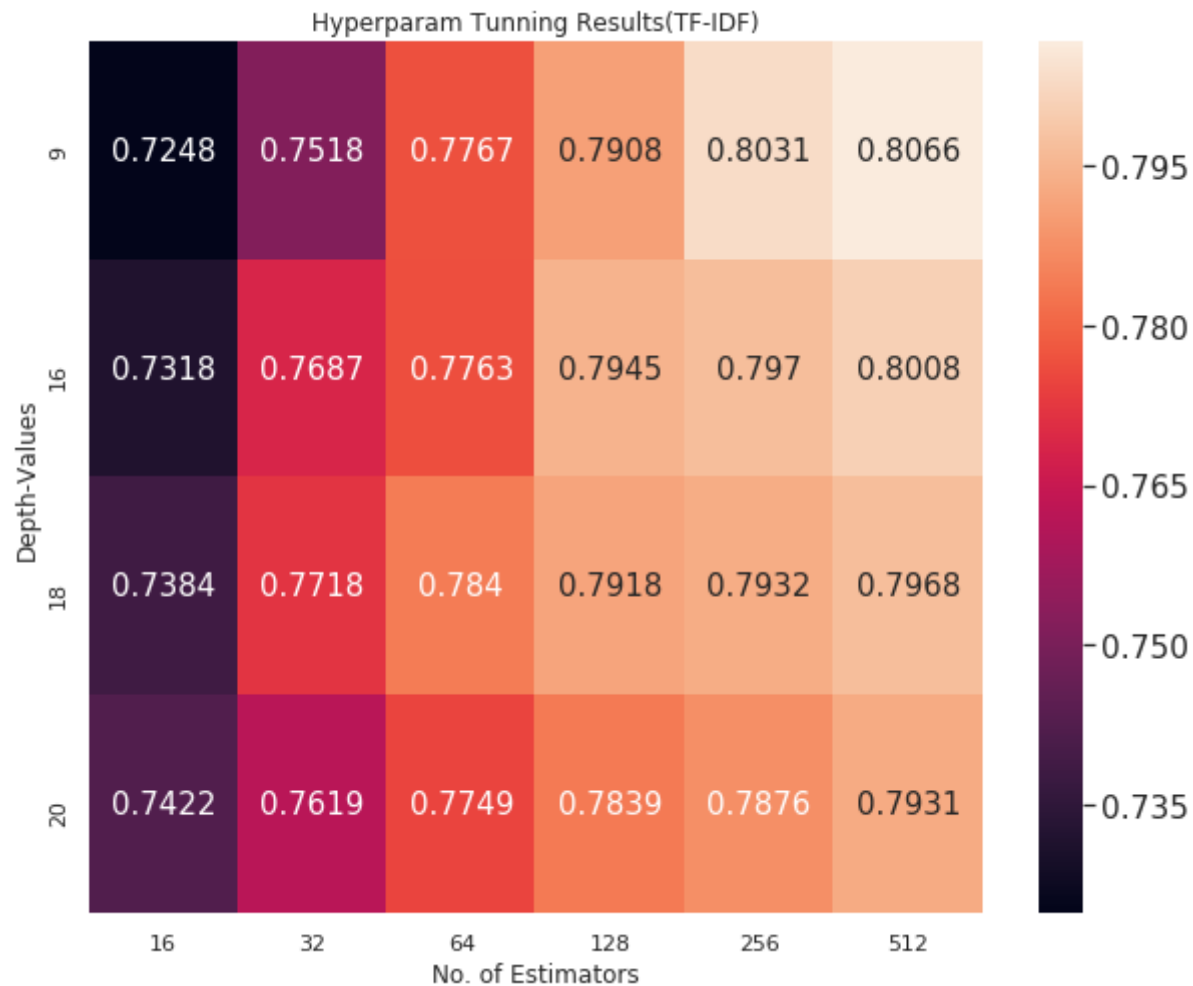
[2.1.1] Hyperparam tuning and plot Heatmap for hyperparam:

In [86]:

```

1 #TRAIN AND TEST DATA
2 #train, test=std_data(train=X_train_tfidf,test=X_test_tfidf,mean=False)
3 train=X_train_tfidf;test=X_test_tfidf;
4 #HYPERPARAM TUNNING
5 %time model=Ensemble_Classifier(train,y_train,TBS,params,searchMethod,vect[1],classifier[0])
6 #PRINT OPTIMAL VALUE OF HYPERPARAM
7 print('Optimal value of hyperparam: ',model.best_params_)

```



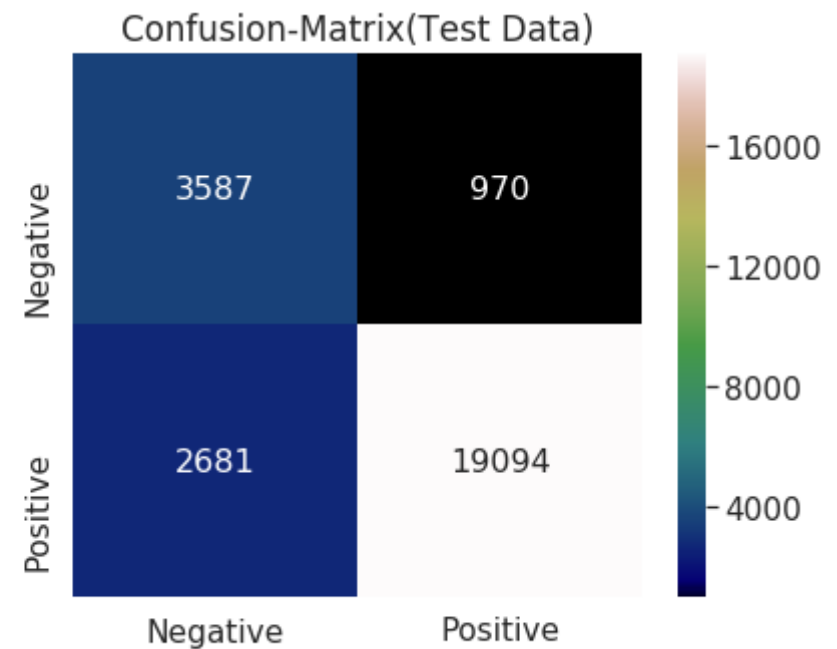
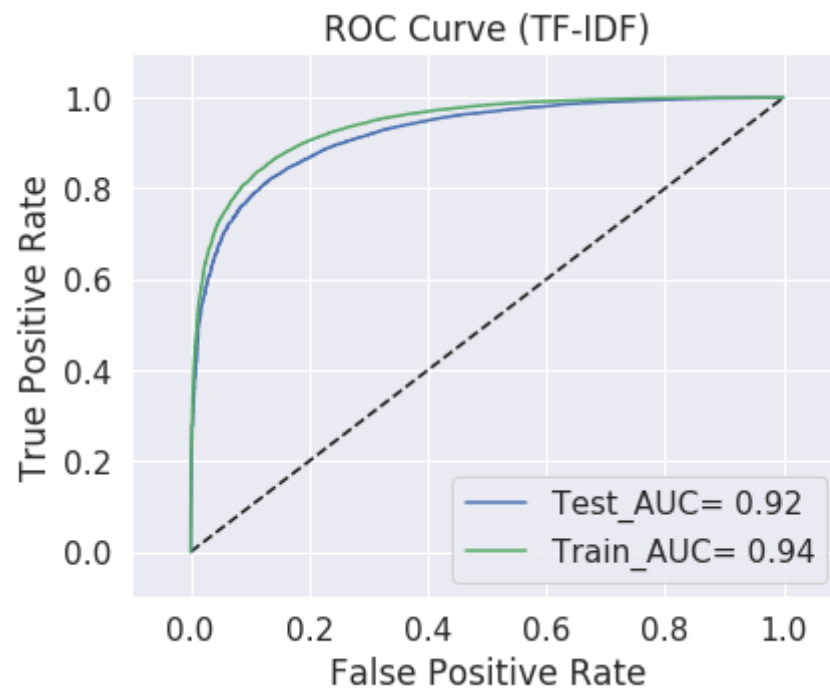
CPU times: user 50.7 s, sys: 884 ms, total: 51.6 s

Wall time: 5min 59s

Optimal value of hyperparam: {'n\_estimators': 512, 'max\_depth': 9}

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

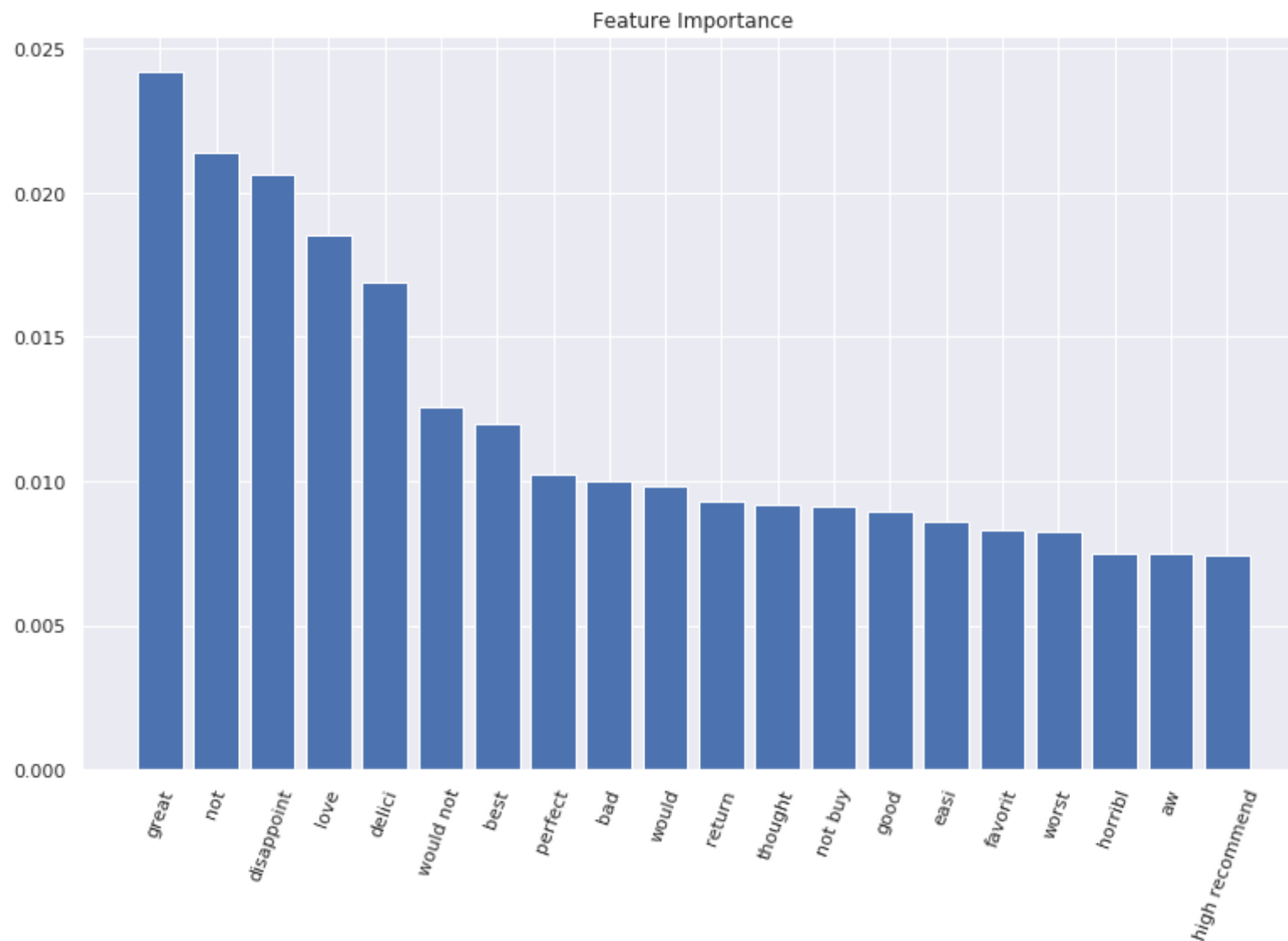
```
In [87]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_,vect[1],summarize,classifier[0])
```

**[2.1.3] Top 20 important features:**



```
In [88]: 1 no_of_imp_features=20  
2 feature_importance(tf_idf_vect,clf,no_of_imp_features)
```





### Observation:

1. Random Forest is good at overall interpretation.
2. `feature_importances_` provides the overall important feature.
3. We can't get class based feature importance in Random Forest.

## **[3.1] Applying Random Forests on AVG-W2V, SET 3**

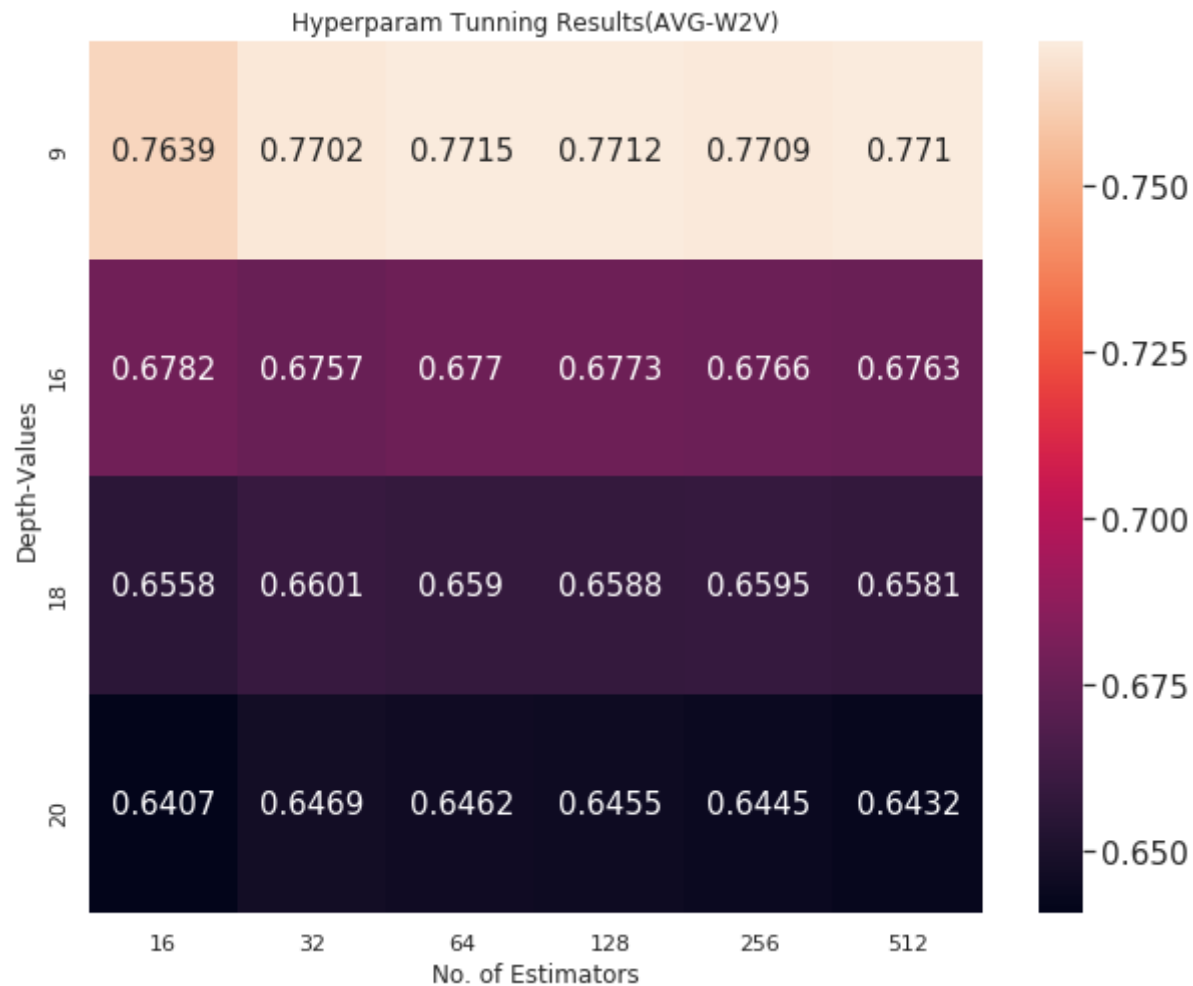
**[3.1.1] Hyperparam tuning and plot Heatmap for hyperparam:**

In [89]:

```

1 #TRAIN AND TEST DATA
2 #train, test=std_data(train=avg_sent_vectors,test=avg_sent_vectors_test,mean=True)
3 train=avg_sent_vectors;test=avg_sent_vectors_test;
4 #HYPERPARAM TUNNING
5 %time model=Ensemble_Classifier(train,y_train,TBS,params,searchMethod,vect[2],classifier[0])
6 #PRINT OPTIMAL VALUE OF HYPERPARAM
7 print('Optimal value of hyperparam: ',model.best_params_)

```



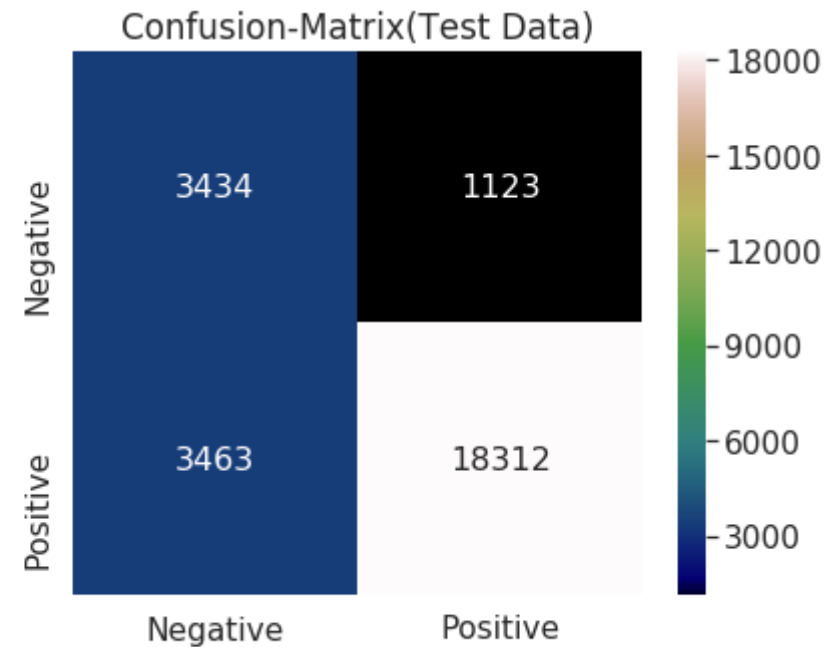
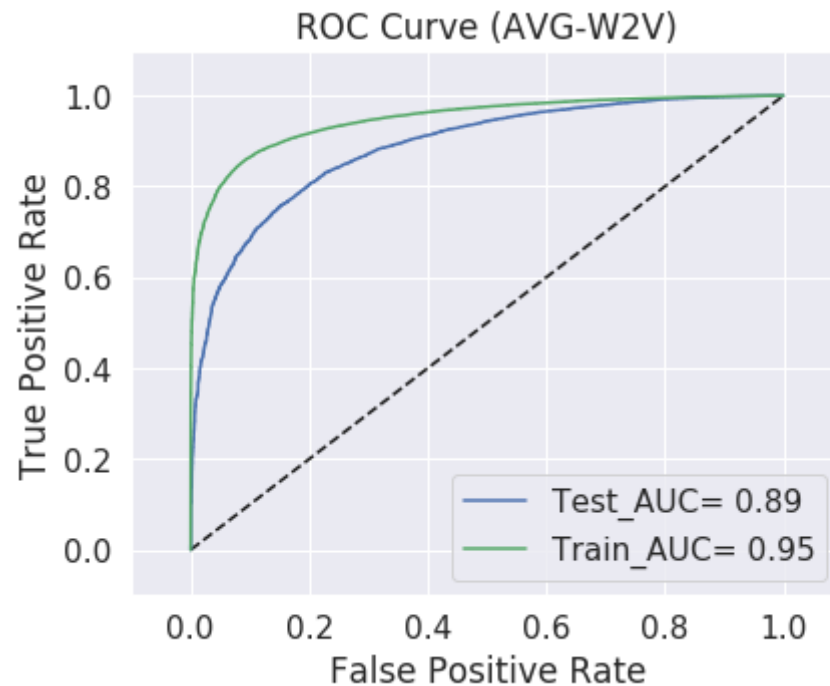
CPU times: user 6min 7s, sys: 29.9 s, total: 6min 37s

Wall time: 22min 2s

Optimal value of hyperparam: {'n\_estimators': 64, 'max\_depth': 9}

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

In [90]: `1 clf=test_performance(train,y_train,test,y_test,model.best_params_,vect[2],summarize,classifier[0])`



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

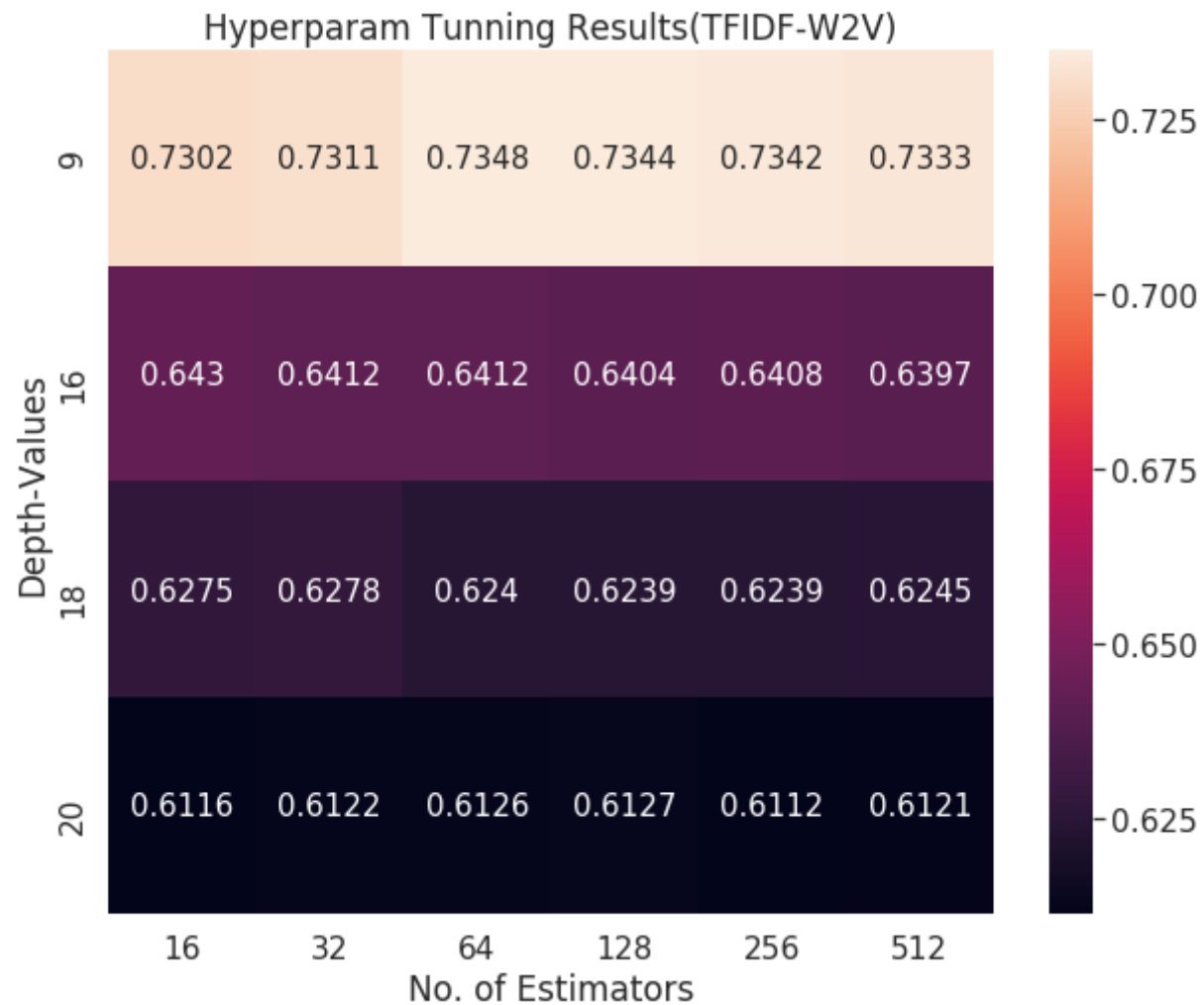
#### [4.1.1] Hyperparam tuning and plot Heatmap for hyperparam:

In [95]:

```

1 #TRAIN AND TEST DATA
2 #train, test=std_data(train=avg_sent_vectors,test=avg_sent_vectors_test,mean=True)
3 train=tfidf_sent_vectors;test=tfidf_sent_vectors_test;
4 #HYPERPARAM TUNNING
5 %time model=Ensemble_Classifier(train,y_train,TBS,params,searchMethod,vect[3],classifier[0])
6 #PRINT OPTIMAL VALUE OF HYPERPARAM
7 print('Optimal value of hyperparam: ',model.best_params_)

```



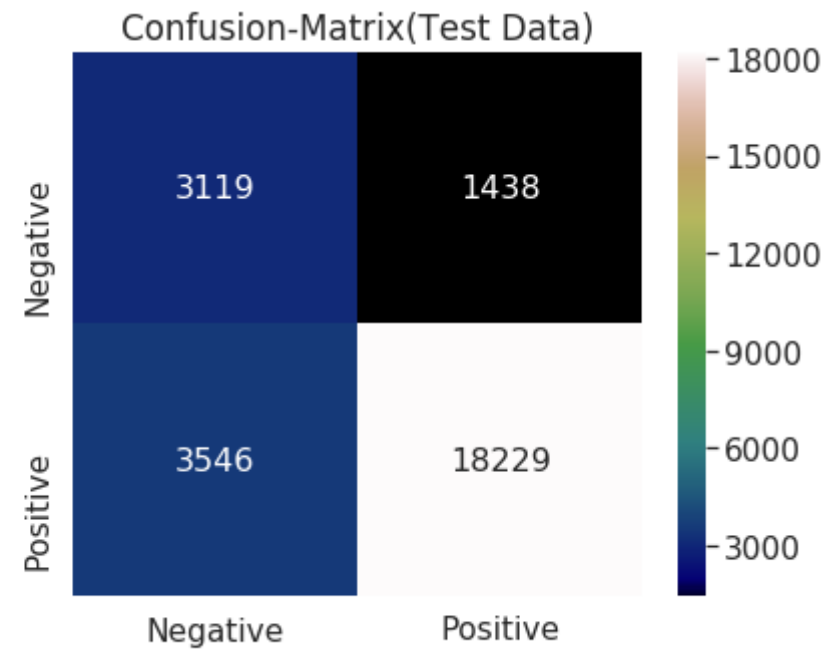
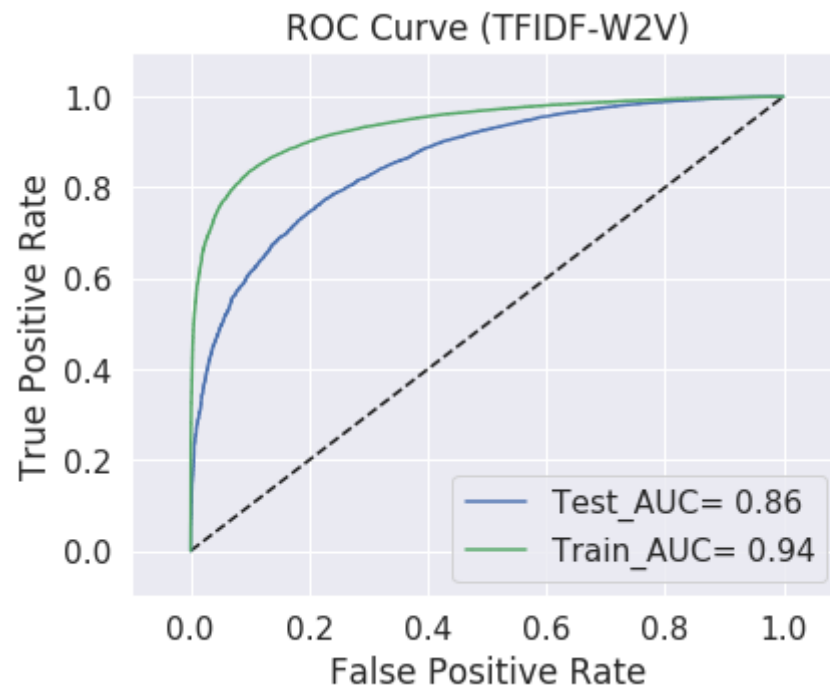
CPU times: user 6min 12s, sys: 30.6 s, total: 6min 42s

Wall time: 22min 19s

Optimal value of hyperparam: {'n\_estimators': 64, 'max\_depth': 9}

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

```
In [96]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_,vect[3],summarize,classifier[0])
```



**Conclusion:**

```
In [97]: 1 #summarize.del_row(2)
        2 print(summarize)
```

Vectorizer	Optimal-Depth	Optimal #Estimators	Test(AUC)	Test(f1-score)
Bow	9	512	0.926	0.877
TF-IDF	9	512	0.921	0.869
AVG-W2V	9	64	0.886	0.839
TFIDF-W2V	9	64	0.858	0.824

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

- Bag of word vectorizer
- f1-score=.877 and auc=.926

-----  
-----  
-----

## [2] GBDT(Gradient Booted Decision Tree)

```
In [49]: 1 #INITIALIZE PRETTY TABLE OBJECT
        2 summarize_gbdtdt = PrettyTable()
        3 summarize_gbdtdt.field_names = ['Vectorizer', 'Optimal-Depth', 'Optimal #Estimators', 'Test(AUC)', 'Test(f1-score)']
        4 #IN GBDT THE BASE LEARNERS ARE HIGH BIAS MODELS SO WE TAKE LOWER VALUES OF DEPTH IN HYPERPARAM
        5 params['max_depth']=[3,5,7,9,12]
```

### Observation:

- In case of GBDT, base learners are high bias(high train error) models.
- Here the base learners are Decision Tree and Decision Tree is highly biased when the tree is shallow(low depth).



3. So here we are taking low depth values in our hyperparam as well.

## **[1.1] Applying GBDT on BOW, SET 1**

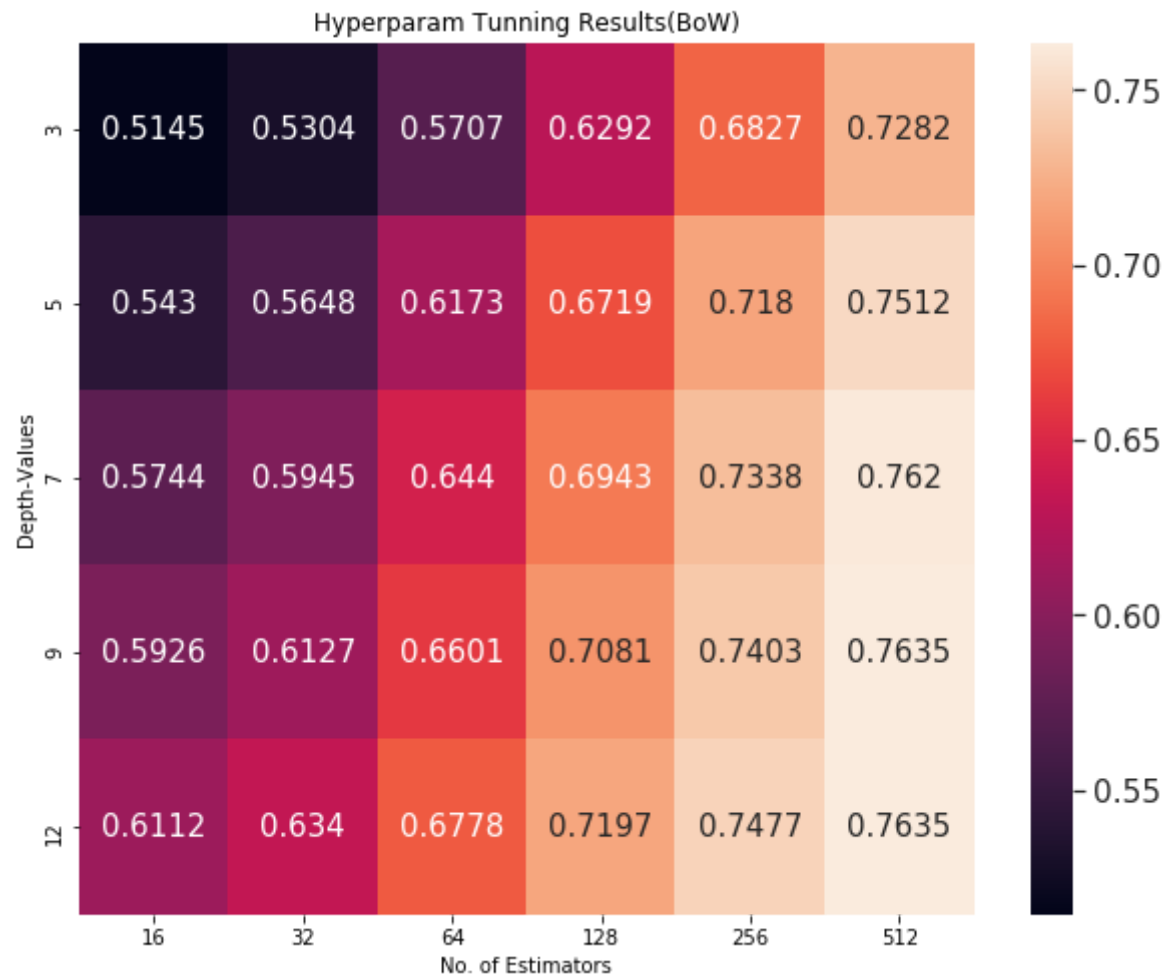
**[1.1.1] Hyperparam tuning and plot Heatmap for hyperparam:**

In [20]:

```

1 #TRAIN AND TEST DATA
2 #train, test=std_data(train=X_train_bigram,test=X_test_bigram,mean=False)
3 train=X_train_bigram;test=X_test_bigram;
4 #HYPERPARAM TUNNING
5 %time model=Ensemble_Classifier(train,y_train,TBS,params,searchMethod,vect[0],classifier[1])
6 #PRINT OPTIMAL VALUE OF HYPERPARAM
7 print('Optimal value of hyperparam: ',model.best_params_)

```



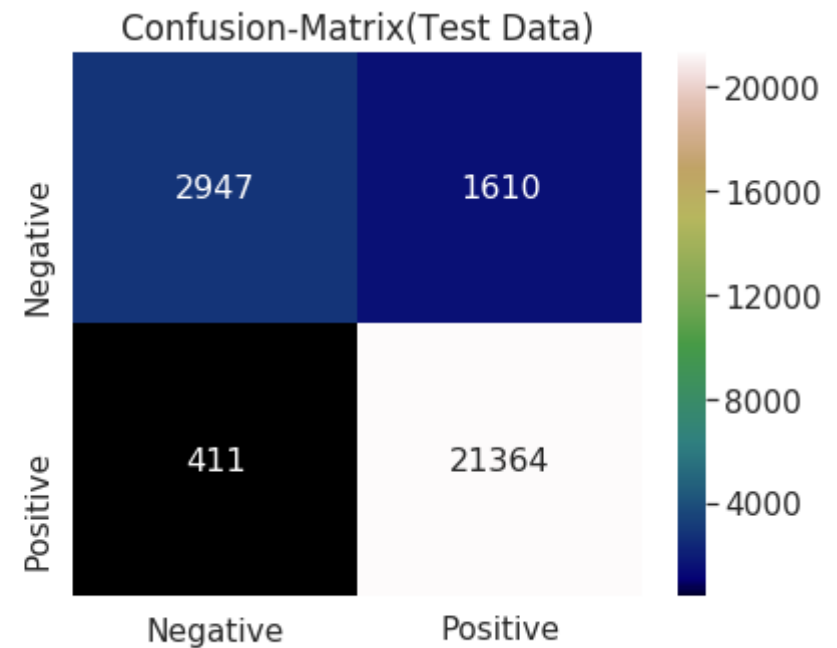
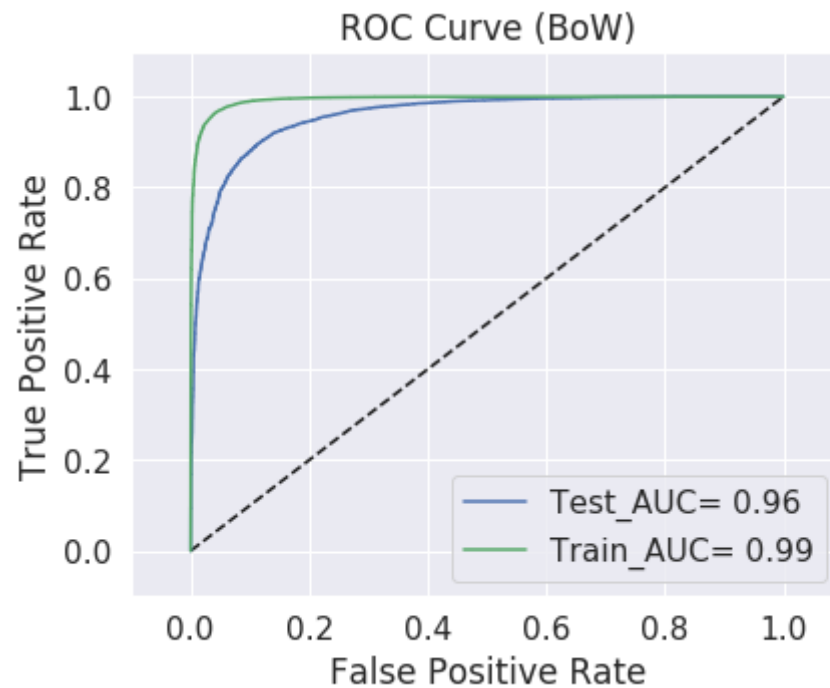
CPU times: user 7min 43s, sys: 2 s, total: 7min 45s

Wall time: 1h 3min 50s

Optimal value of hyperparam: {'n\_estimators': 512, 'max\_depth': 9}

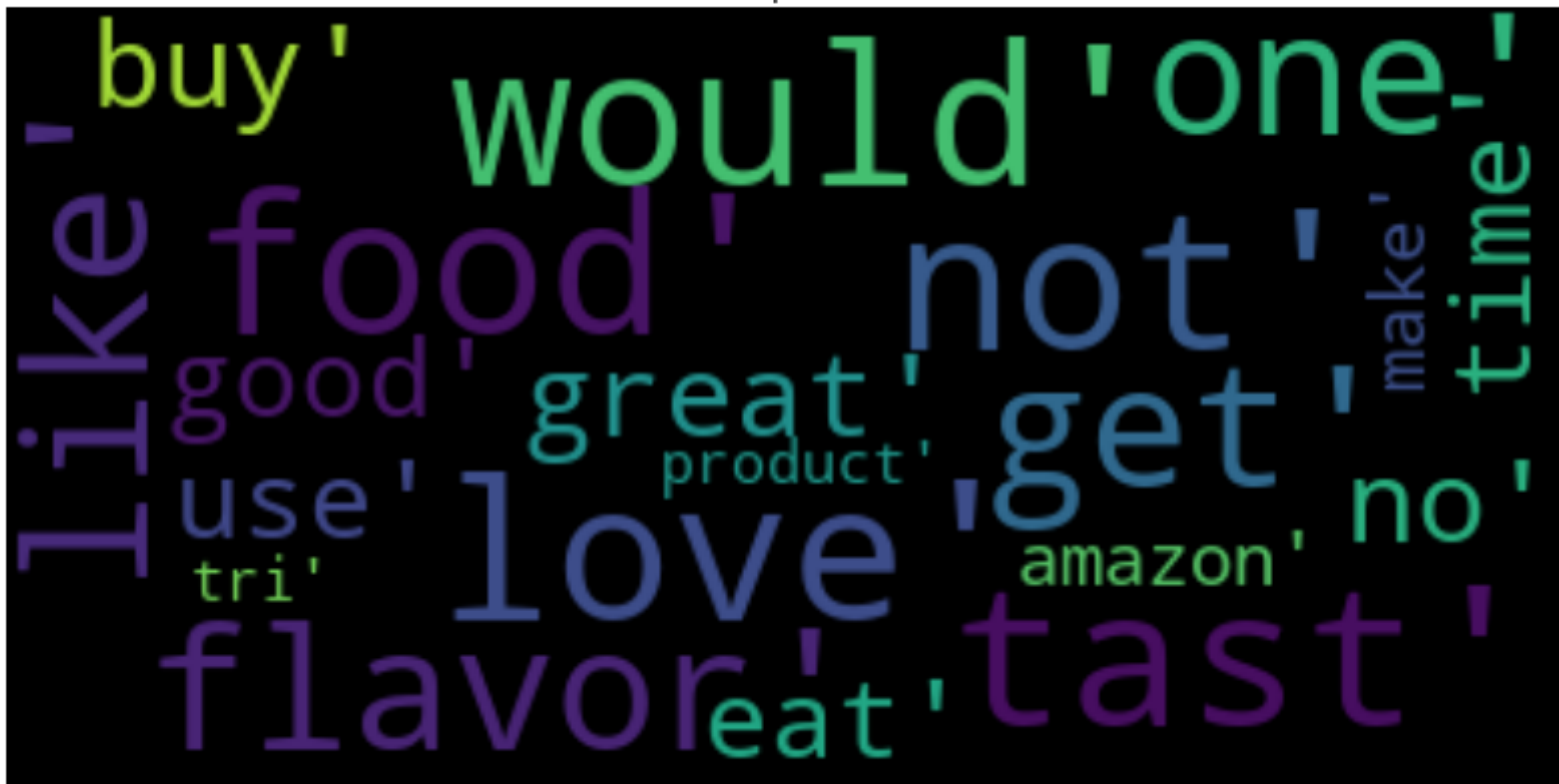
**[1.1.2] Performance on test data with optimal value of hyperparam:**

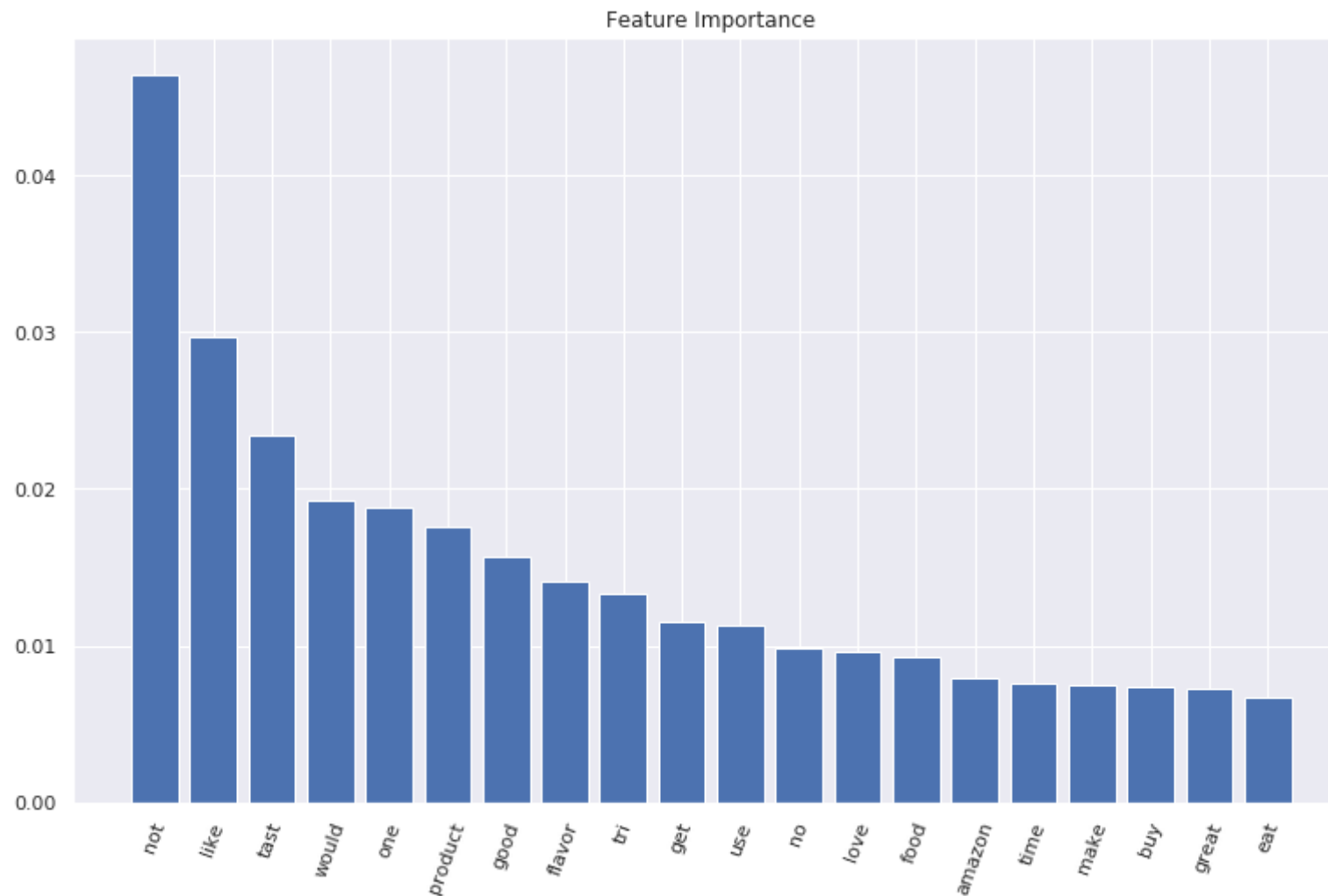
```
In [21]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_,vect[0],summarize_gbd,classifier[1])
```

**[1.1.3] Top 20 important features:**

```
In [22]: 1 no_of_imp_features=20  
2 feature_importance(count_vect,clf,no_of_imp_features)
```

WordCloud(Important Feature)





### Observation:

1. GBDT is good at overall interpretation.
2. `feature_importances_` provides the overall important feature.
3. We can't get class based feature importance in GBDT.

## [2.1] Applying GBDT on TFIDF, SET 2

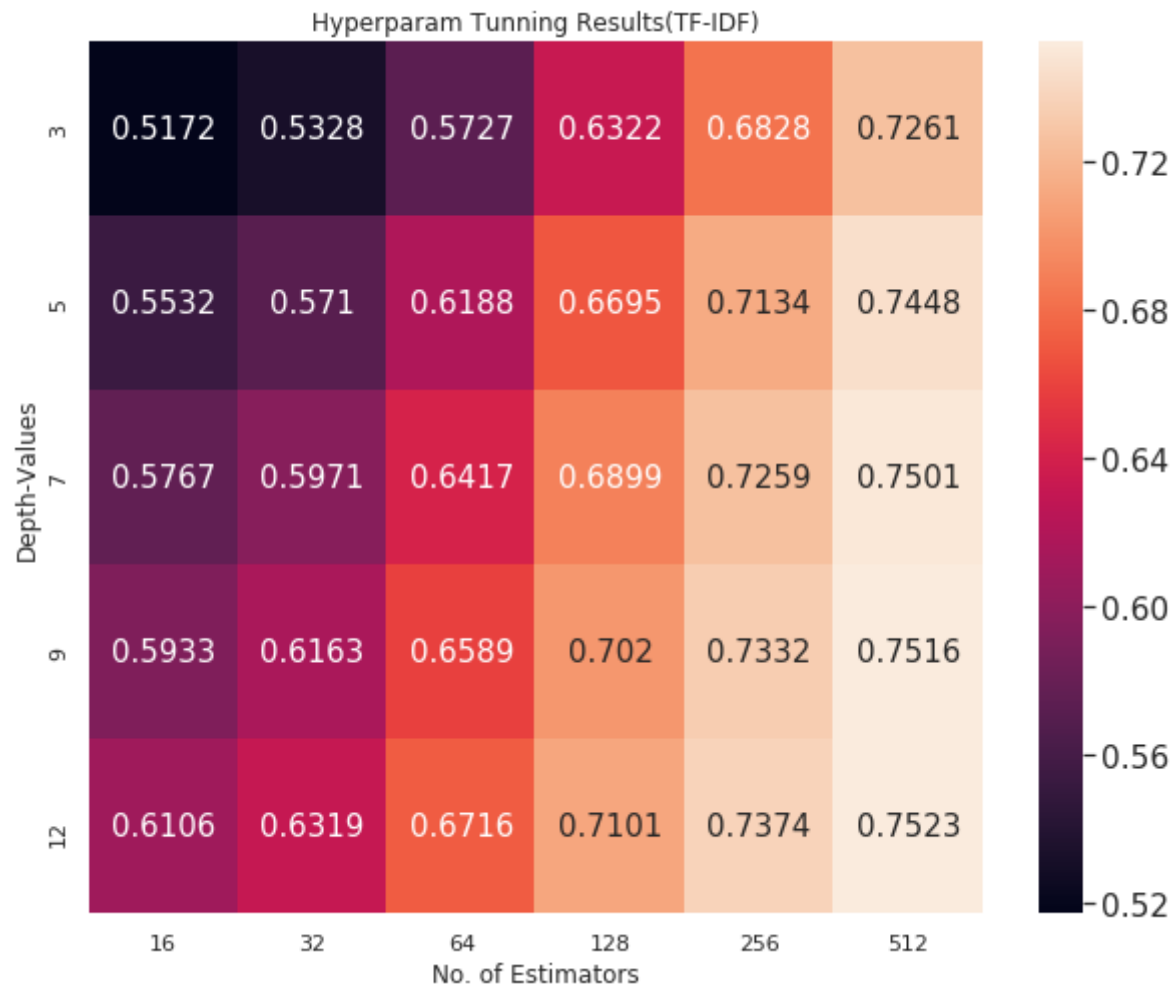
### [2.1.1] Hyperparam tuning and plot Heatmap for hyperparam:

In [23]:

```

1 # TRAIN AND TEST DATA
2 #train, test=std_data(train=X_train_tfidf,test=X_test_tfidf,mean=False)
3 train=X_train_tfidf;test=X_test_tfidf;
4 #HYPERPARAM TUNNING
5 %time model=Ensemble_Classifier(train,y_train,TBS,params,searchMethod,vect[1],classifier[1])
6 #PRINT OPTIMAL VALUE OF HYPERPARAM
7 print('Optimal value of hyperparam: ',model.best_params_)

```



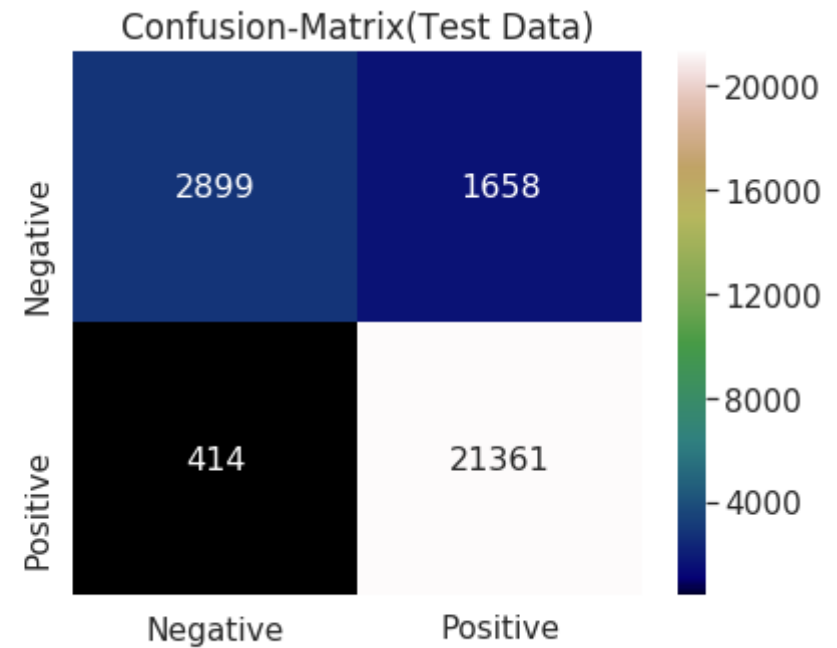
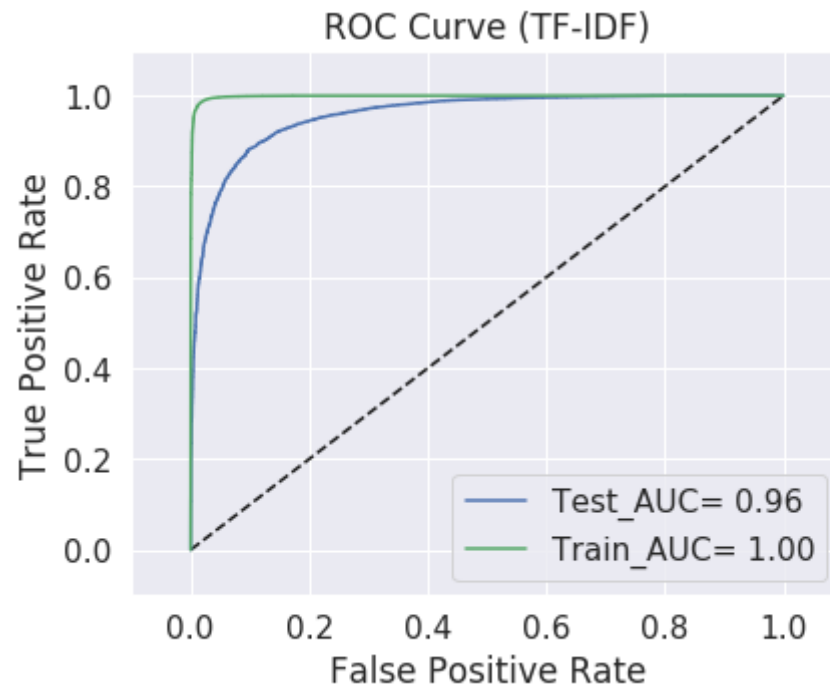
CPU times: user 14min 26s, sys: 2.77 s, total: 14min 29s

Wall time: 1h 21min 23s

Optimal value of hyperparam: {'n\_estimators': 512, 'max\_depth': 12}

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

```
In [24]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_,vect[1],summarize_gbd,classifier[1])
```

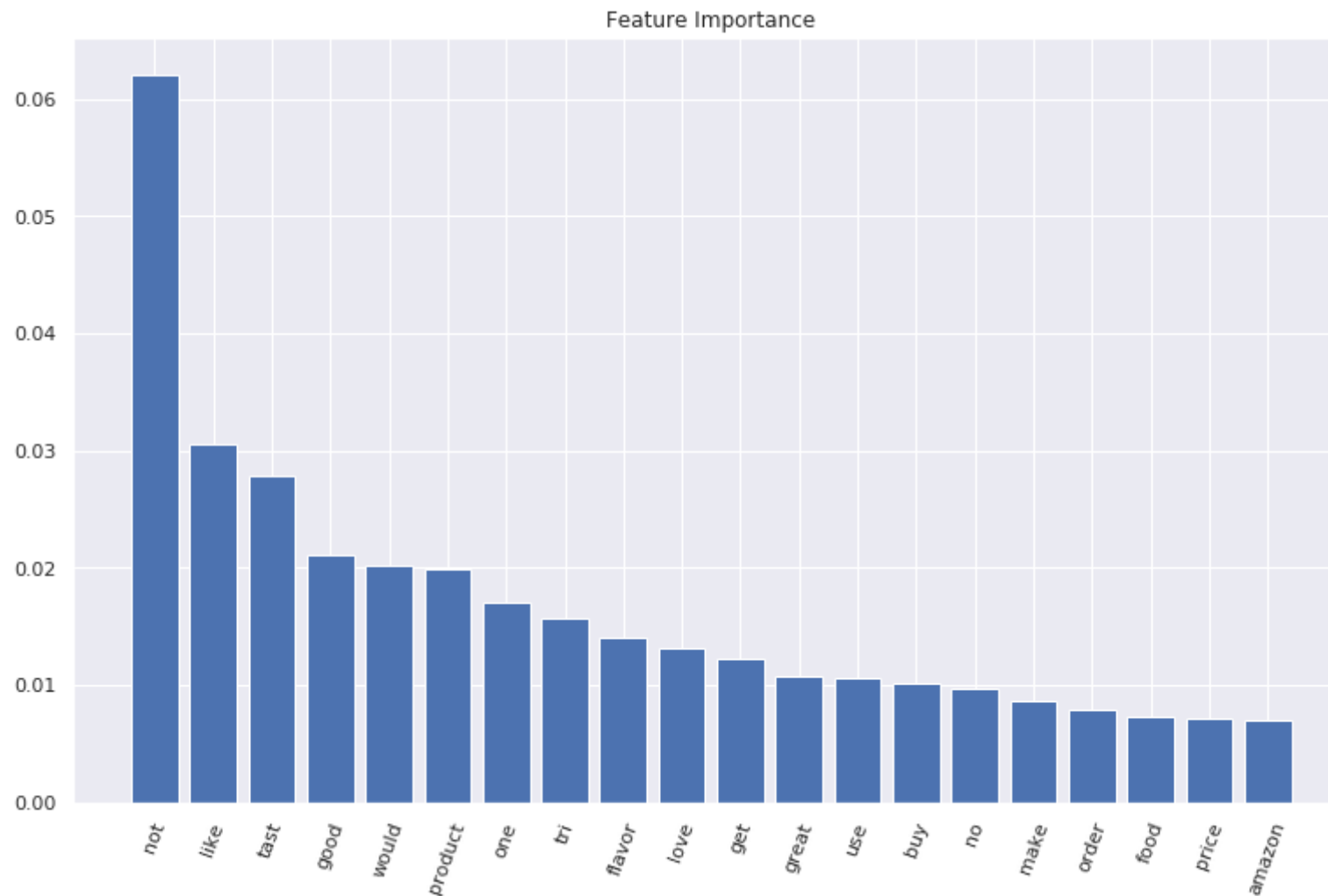
**[2.1.3] Top 20 important features:**



```
In [25]: 1 no_of_imp_features=20  
2 feature_importance(tf_idf_vect,clf,no_of_imp_features)
```

WordCloud(Important Feature)





### Observation:

1. GBDT is good at overall interpretation.
2. `feature_importances_` provides the overall important feature.
3. We can't get class based feature importance in GBDT.

## **[3.1] Applying GBDT on AVG W2V, SET 3**

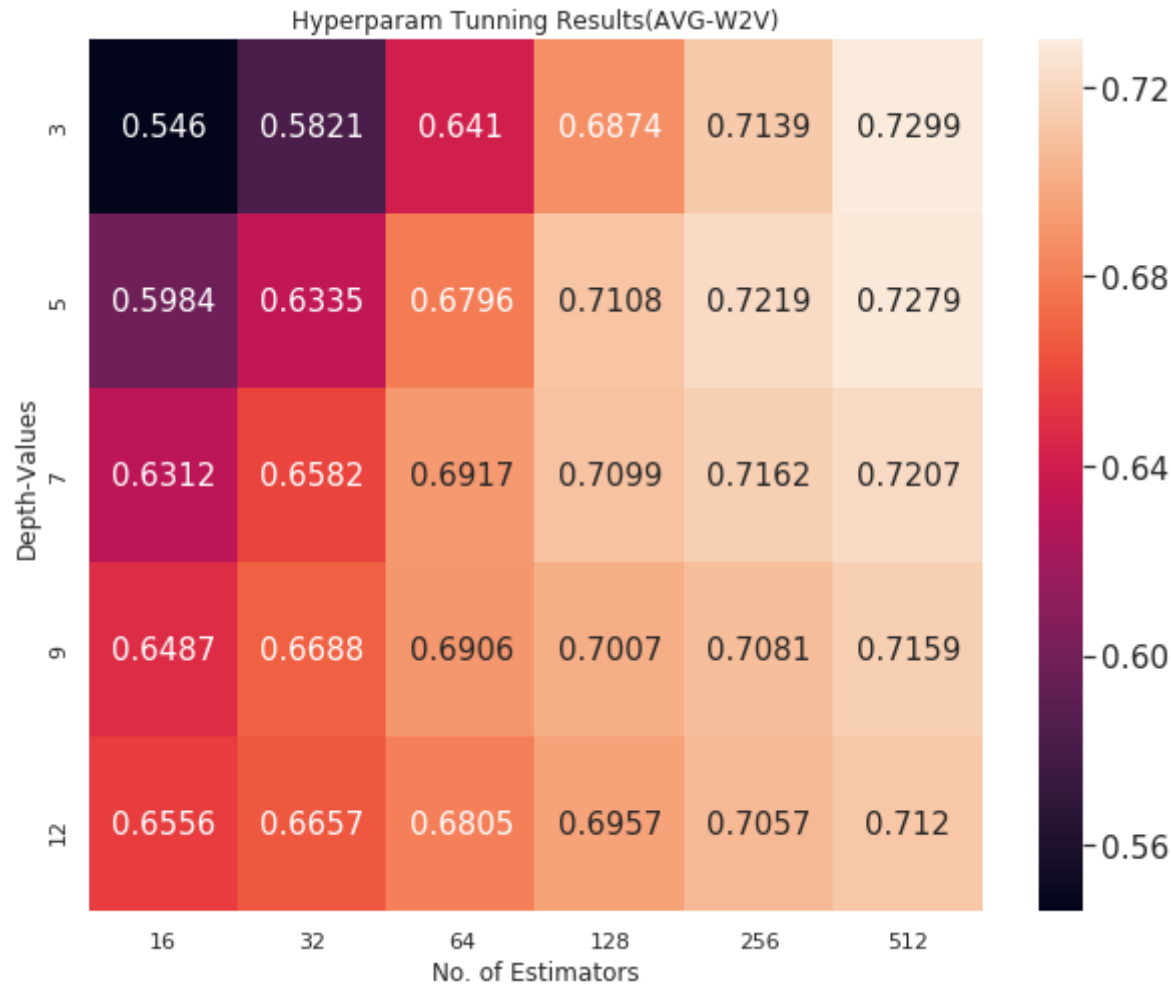
**[3.1.1] Hyperparam tuning and plot Heatmap for hyperparam:**

In [38]:

```

1 #TRAIN AND TEST DATA(convert the avgw2v to array)
2 train=np.array(avg_sent_vectors);test=np.array(avg_sent_vectors_test);
3 #HYPERPARAM TUNNING
4 %time model=Ensemble_Classifier(train,y_train,TBS,params,searchMethod,vect[2],classifier[1])
5 #PRINT OPTIMAL VALUE OF HYPERPARAM
6 print('Optimal value of hyperparam: ',model.best_params_)

```



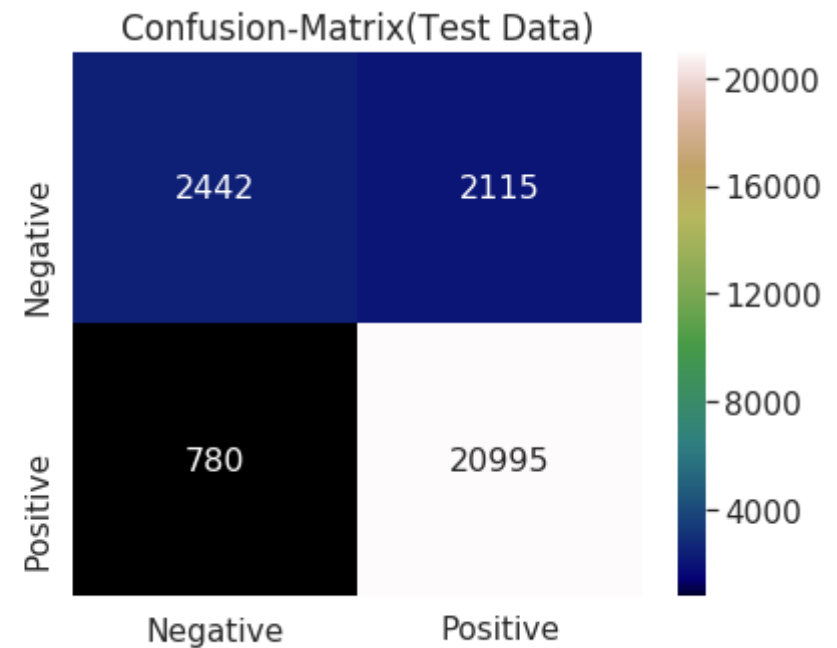
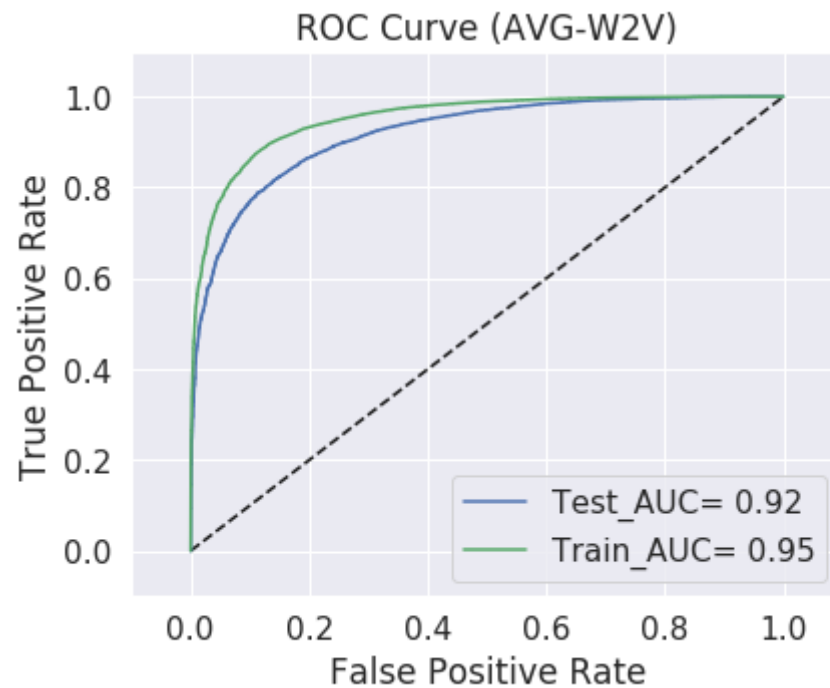
CPU times: user 3min 24s, sys: 1.69 s, total: 3min 25s

Wall time: 1h 2min 19s

Optimal value of hyperparam: {'n\_estimators': 512, 'max\_depth': 3}

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

```
In [39]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_,vect[2],summarize_gbd,classifier[1])
```



### [4.1] Applying GBDT on TFIDF W2V, SET 4

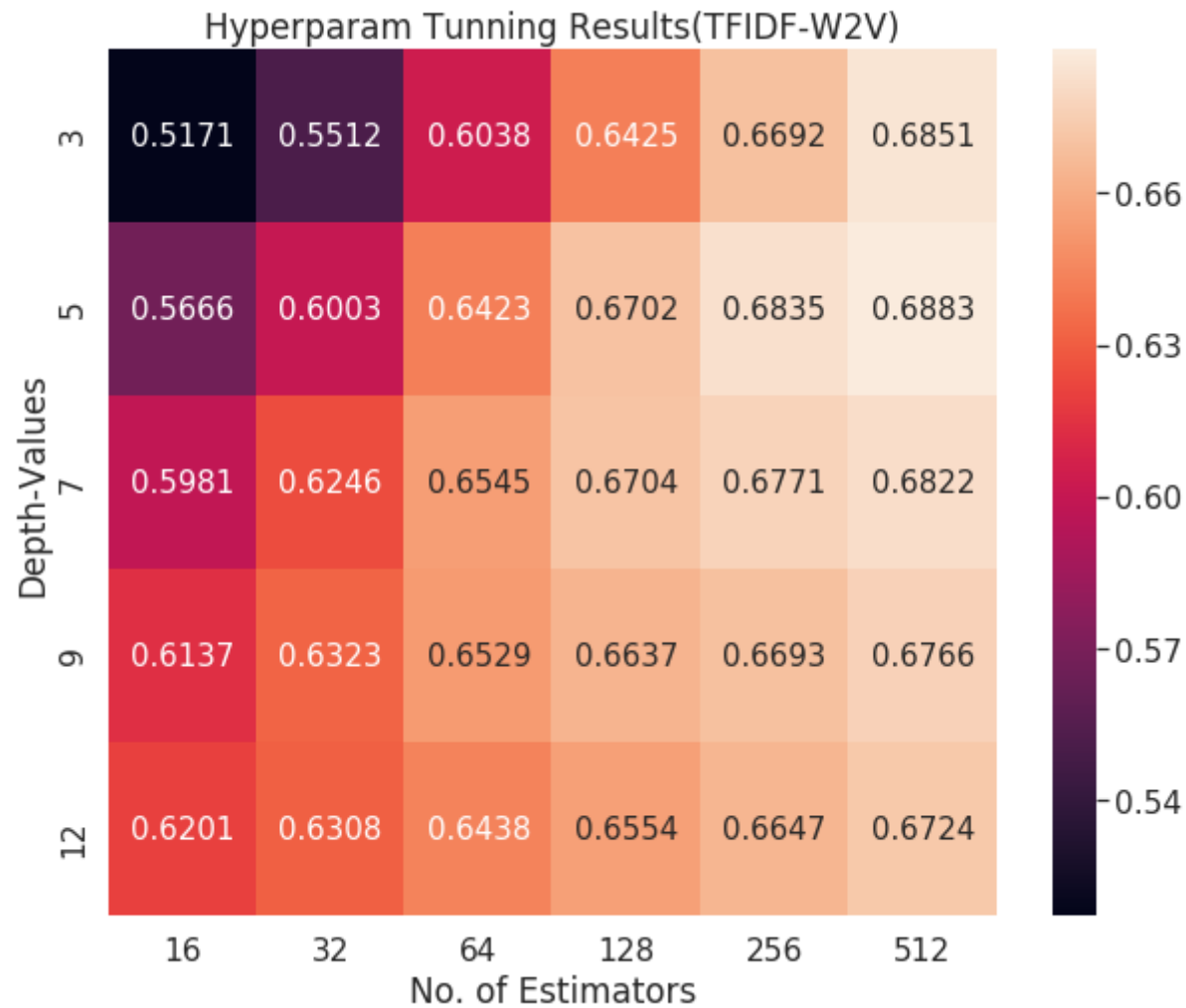
#### [4.1.1] Hyperparam tuning and plot Heatmap for hyperparam:

In [41]:

```

1 #TRAIN AND TEST DATA(convert tfidf2v to array)
2 train=np.array(tfidf_sent_vectors);test=np.array(tfidf_sent_vectors_test);
3 #HYPERPARAM TUNNING
4 %time model=Ensemble_Classifier(train,y_train,TBS,params,searchMethod,vect[3],classifier[1])
5 #PRINT OPTIMAL VALUE OF HYPERPARAM
6 print('Optimal value of hyperparam: ',model.best_params_)

```



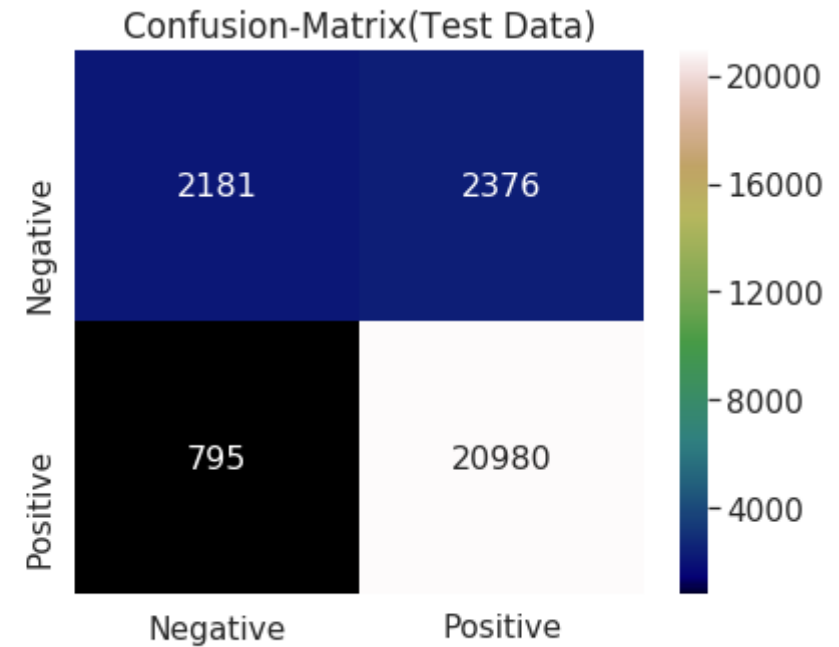
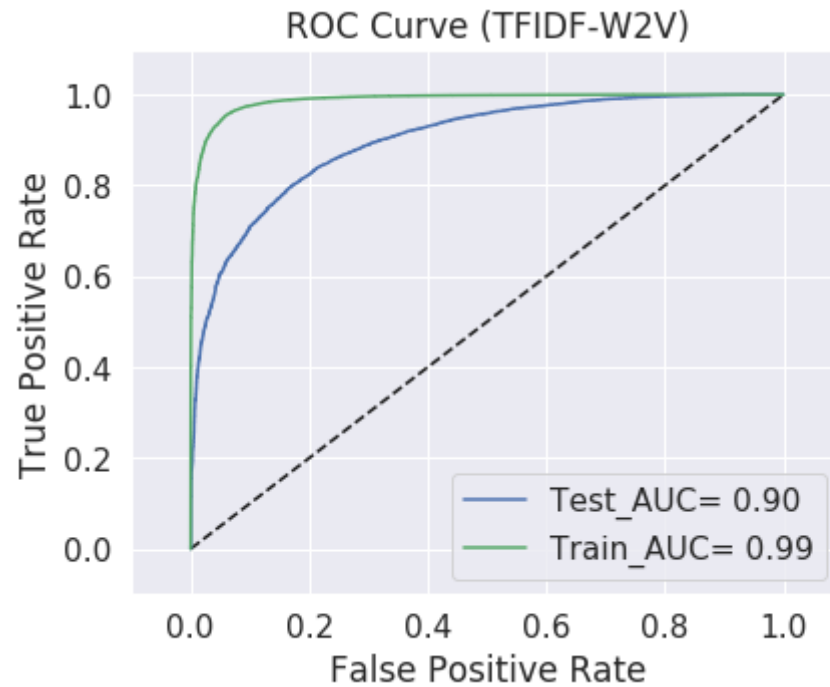
CPU times: user 5min 32s, sys: 2.06 s, total: 5min 34s

Wall time: 1h 4min 9s

Optimal value of hyperparam: {'n\_estimators': 512, 'max\_depth': 5}

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

```
In [42]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_,vect[3],summarize_gbd,classifier[1])
```

**Conclusion:**

```
In [47]: 1 print("Summary of GBDT model with optimal Hyperparam and Scores:" )
          2 print(summarize_gbdtd)
```

Summary of GBDT model with optimal Hyperparam and Scores:

Vectorizer	Optimal-Depth	Optimal #Estimators	Test(AUC)	Test(f1-score)
BoW	9	512	0.958	0.918
TF-IDF	12	512	0.957	0.916
AVG-W2V	3	512	0.919	0.882
TFIDF-W2V	5	512	0.899	0.869

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

- a. Bag of word vectorizer
- b. f1-score=.918 and auc=.958

In [ ]:

1