

Import necessary libraries

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

```
In [2]: 1 from sklearn.naive_bayes import MultinomialNB
        2 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
        3 from sklearn.preprocessing import StandardScaler
        4 from sklearn.metrics import *
        5 import pickle
        6 from tqdm import tqdm_notebook
        7 import seaborn as sns
        8 from sklearn.model_selection import TimeSeriesSplit
        9 from sklearn.model_selection import cross_val_score
       10 import numpy as np
       11 import matplotlib.pyplot as plt
       12 import pandas as pd
       13 from scipy.sparse import *
       14 from prettytable import PrettyTable
       15 from sklearn.externals import joblib
```

Load preprocessed data

```
In [3]: 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
8 y_train =openfromfile('y_train')
9 y_test =openfromfile('y_test')
10
11 count_vect =openfromfile('count_vect')
12 X_train_bigram = openfromfile('X_train_bigram')
13 X_test_bigram = openfromfile('X_test_bigram')
14
15 count_vect_fe =openfromfile('count_vect_fe')
16 X_train_bigram_fe = openfromfile('X_train_bigram_fe')
17 X_test_bigram_fe = openfromfile('X_test_bigram_fe')
18
19 tf_idf_vect =openfromfile('tf_idf_vect')
20 X_train_tfidf =openfromfile('X_train_tfidf')
21 X_test_tfidf =openfromfile('X_test_tfidf')
22
23 tf_idf_vect_fe =openfromfile('tf_idf_vect_fe')
24 X_train_tfidf_fe =openfromfile('X_train_tfidf_fe')
25 X_test_tfidf_fe =openfromfile('X_test_tfidf_fe')
```

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 [5]: 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
```

Naive Bayes

Function for finding optimal value of hyperparameter nd draw error plot :

```

In [6]: 1 def NB_Classifier(x_train,y_train,TBS,params,searchMethod,vect):
2         ''' FUNCTION FOR FINDING OPTIMAL VALUE OF HYPERPARAM AND DRAW ERROR PLOT'''
3         #INITIALIZE MULTINOMIAL_NB OBJECT
4         clf=MultinomialNB(fit_prior=True)
5
6         # APPLY RANDOM OR GRID SEARCH FOR HYPERPARAMETER TUNNING
7         if searchMethod=='grid':
8             model=GridSearchCV(clf,\
9                                 cv=TBS,\
10                                n_jobs=-1,\
11                                param_grid=params,\
12                                return_train_score=True,\
13                                scoring=make_scorer(roc_auc_score,average='weighted'))
14             model.fit(x_train,y_train)
15         elif searchMethod=='random':
16             model=RandomizedSearchCV(clf,\
17                                     n_jobs=-1,\
18                                     cv=TBS,\
19                                     param_distributions=params,\
20                                     n_iter=len(params['alpha']),\
21                                     return_train_score=True,\
22                                     scoring=make_scorer(roc_auc_score,average='weighted'))
23             model.fit(x_train,y_train)
24
25         #PLOT HYPERPARAM VS AUC VALUES(FOR BOTH CV AND TRAIN)
26         train_auc= model.cv_results_['mean_train_score']
27         train_auc_std= model.cv_results_['std_train_score']
28         cv_auc = model.cv_results_['mean_test_score']
29         cv_auc_std= model.cv_results_['std_test_score']
30         plt.figure(1,figsize=(10,6))
31         plt.plot(params['alpha'], train_auc, label='Train AUC')
32         # Reference Link: https://stackoverflow.com/a/48803361/4084039
33         # gca(): get current axis
34         plt.gca().fill_between(params['alpha'],train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkorange')
35         plt.plot(params['alpha'], cv_auc, label='CV AUC')
36         # Reference Link: https://stackoverflow.com/a/48803361/4084039
37         plt.gca().fill_between(params['alpha'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
38
39         plt.title('ERROR PLOT (%s)' %vect)
40         plt.xlabel('Alpha: Hyperparam')
41         plt.ylabel('AUC')
42         plt.grid(True)

```

```
43 plt.legend()  
44 plt.show()  
45 return model  
46
```

Function which calculate performance on test data with optimal hyperparam :

```

In [7]: 1 def test_performance(x_train,y_train,x_test,y_test,optimal_alpha,vect,summarize):
2         '''FUNCTION FOR TEST PERFORMANCE(PLOT ROC CURVE FOR BOTH TRAIN AND TEST) WITH OPTIMAL_K'''
3         #INITIALIZE MultinomialNB WITH OPTIMAL HYPERPARAM
4         clf=MultinomialNB(alpha=optimal_alpha,fit_prior=True)
5         clf.fit(x_train,y_train)
6
7         y_pred=clf.predict(x_test)
8         test_probability = clf.predict_proba(x_test)[:,-1]
9         train_probability = clf.predict_proba(x_train)[:,-1]
10        fpr_test, tpr_test, threshold_test = roc_curve(y_test, test_probability,pos_label=1)
11        fpr_train, tpr_train, threshold_train = roc_curve(y_train, train_probability,pos_label=1)
12        auc_score_test=auc(fpr_test, tpr_test)
13        auc_score_train=auc(fpr_train, tpr_train)
14        f1=f1_score(y_test,y_pred,average='weighted')
15
16        #ADD RESULTS TO PRETTY TABLE
17        summarize.add_row([vect, optimal_alpha, '%.3f' %auc_score_test,'%%.3f' %f1])
18
19        plt.figure(1,figsize=(14,5))
20        plt.subplot(121)
21        plt.title('ROC Curve (%s)' %vect)
22        #IDEAL ROC CURVE
23        plt.plot([0,1],[0,1],'k--')
24        #ROC CURVE OF TEST DATA
25        plt.plot(fpr_test, tpr_test , 'b', label='Test_AUC= %.2f' %auc_score_test)
26        #ROC CURVE OF TRAIN DATA
27        plt.plot(fpr_train, tpr_train , 'g', label='Train_AUC= %.2f' %auc_score_train)
28        plt.xlim([-0.1,1.1])
29        plt.ylim([-0.1,1.1])
30        plt.xlabel('False Positive Rate')
31        plt.ylabel('True Positive Rate')
32        plt.grid(True)
33        plt.legend(loc='lower right')
34        #PLOT CONFUSION MATRIX USING HEATMAP
35        plt.subplot(122)
36        plt.title('Confusion-Matrix(Test Data)')
37        df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), ['Negative','Positive'],['Negative','Positive'])
38        sns.set(font_scale=1.4)#for label size
39        sns.heatmap(df_cm,cmap='gist_earth', annot=True,annot_kws={"size": 16}, fmt='g')
40        plt.show()
41        return clf

```

Function which print top important fetures and plot them using Bar plot :

```

In [8]: 1 # FUNCTION CREATED BY SELF
2 def feature_importance(vectorizer,clf,n,top_features):
3     coef_n=[];coef_p=[];coef_np=[];i=1;
4     names_n=[];names_p=[];names=[];
5     feature_names=vectorizer.get_feature_names()
6
7     #METHOD-1
8     #PROBABILITY FEATURE BELONGS TO CLASS 0
9     prob_n=clf.feature_log_prob_[0,:]
10    #PROBABILITY FEATURE BELONGS TO CLASS 1
11    prob_p=clf.feature_log_prob_[1,:]
12    #SELECT FEATURES WHICH HAVE HIGHEST PROBABILITY AND BELONGS TO -VE CLASS
13    sorted_neg_coef_feat=sorted(zip(prob_n,feature_names),reverse=True)
14    #SELECT FEATURES WHICH HAVE HIGHEST PROBABILITY AND BELONGS TO +VE CLASS
15    sorted_pos_coef_feat=sorted(zip(prob_p,feature_names),reverse=True)
16
17    for (coef_neg , feat_neg), (coef_pos, feat_pos) in zip(sorted_neg_coef_feat,sorted_pos_coef_feat):
18        top_features.add_row([i, '%.4f' %coef_neg, feat_neg,'%%.4f' %coef_pos, feat_pos])
19        coef_n.append(coef_neg)
20        names_n.append(feat_neg)
21        coef_p.append(coef_pos)
22        names_p.append(feat_pos)
23        i+=1
24        if i==n+1:
25            break
26    #METHOD-2
27    '''print(np.take(feature_names,prob_n.argsort())[:-(n + 1):-1])
28    print(np.take(feature_names,prob_p.argsort())[:-(n + 1):-1])'''
29
30    names.extend(names_n)
31    names.extend(names_p)
32    coef_np.extend(coef_n)
33    coef_np.extend(coef_p)
34    names=np.array(names)
35    #BAR CHART
36    plt.figure(2,figsize=(13,6))
37    sns.set(rc={'figure.figsize':(11.7,8.27)})
38    plt.title("Feature Importance(top %d positive and negative class features)" % n)
39    # ADD BARS
40    plt.bar(range(2*n), coef_np)
41    # ADD FEATURE NAMES
42    plt.xticks(range(2*n), names, rotation=80)

```



```
43 plt.show()  
44
```

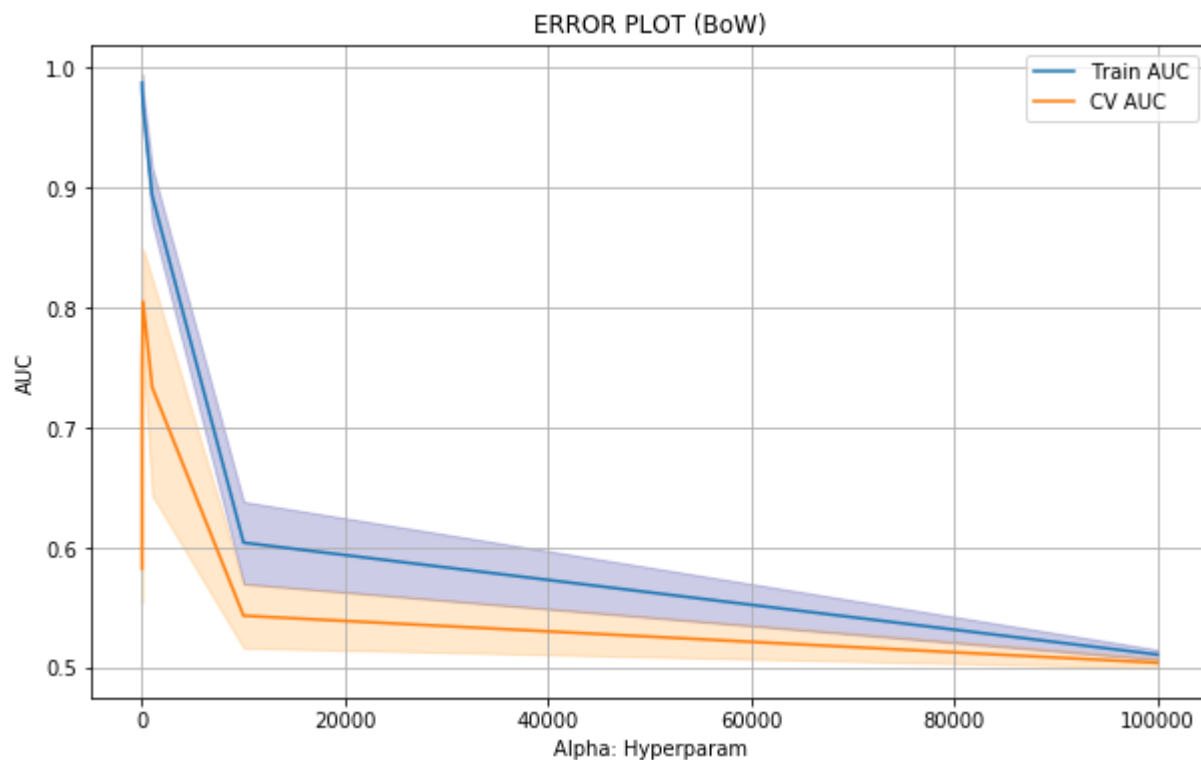
Initialization of common objects required for all vectorization:

```
In [9]: 1 #VECTORIZER  
2 vect=['Bow', 'Bow-FE', 'TFIDF', 'TFIDF-FE']  
3 #OBJECT FOR TIMESERIES CROSS VALIDATION  
4 TBS=TimeSeriesSplit(n_splits=10)  
5 #METHOD USE FOR HYPER PARAMETER TUNNING  
6 searchMethod='grid'  
7 #RANGE OF VALUES(HYPERPARAM)#[100000,10000,1000,500,100,50,10,5,1,.5,.1,.05,.01,.005,.001,.0005,.0001,.00001]  
8 alpha_ranges=[10**x for x in range(-6,6)]  
9 params={'alpha':alpha_ranges}  
10 #INITIALIZE PRETTY TABLE OBJECT  
11 summarize = PrettyTable()  
12 summarize.field_names = ['Vectorizer', 'Optimal-alpha', 'Test(AUC)', 'Test(f1-score)']
```

[1.] Naive Bayes on BOW, SET 1

[1.1] Hyperparam Tunning SET 1

```
In [10]: 1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=X_train_bigram,test=X_test_bigram,mean=False)
3 #HYPERPARAM TUNNING
4 %time model=NB_Classifier(train,y_train,TBS,params,searchMethod,vect[0])
5 #PRINT OPTIMAL VALUE OF HYPERPARAM
6 print('Optimal value of Alpha: ',model.best_params_)
7 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
8 saveModeltofile(model,'model_bow_mnb')
```



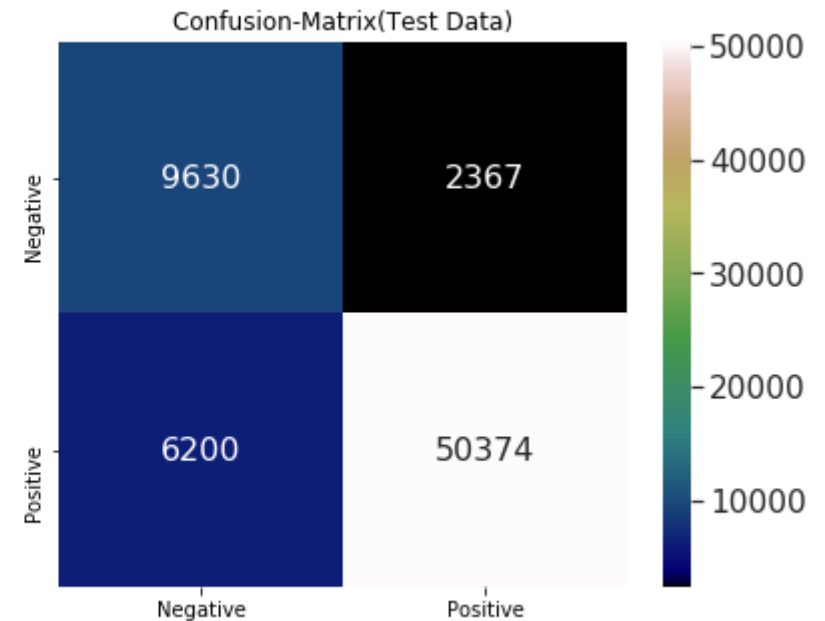
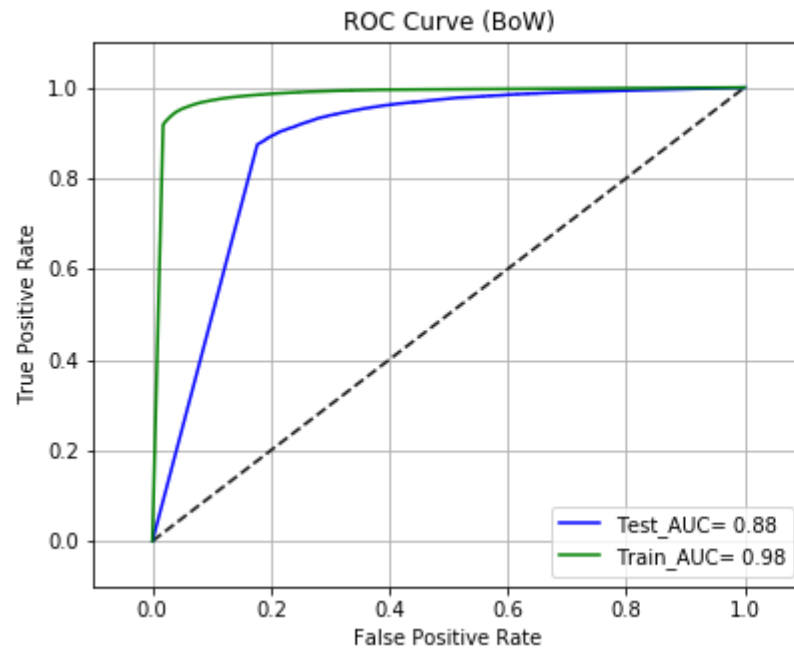
CPU times: user 29.2 s, sys: 552 ms, total: 29.8 s

Wall time: 28.1 s

Optimal value of Alpha: {'alpha': 100}

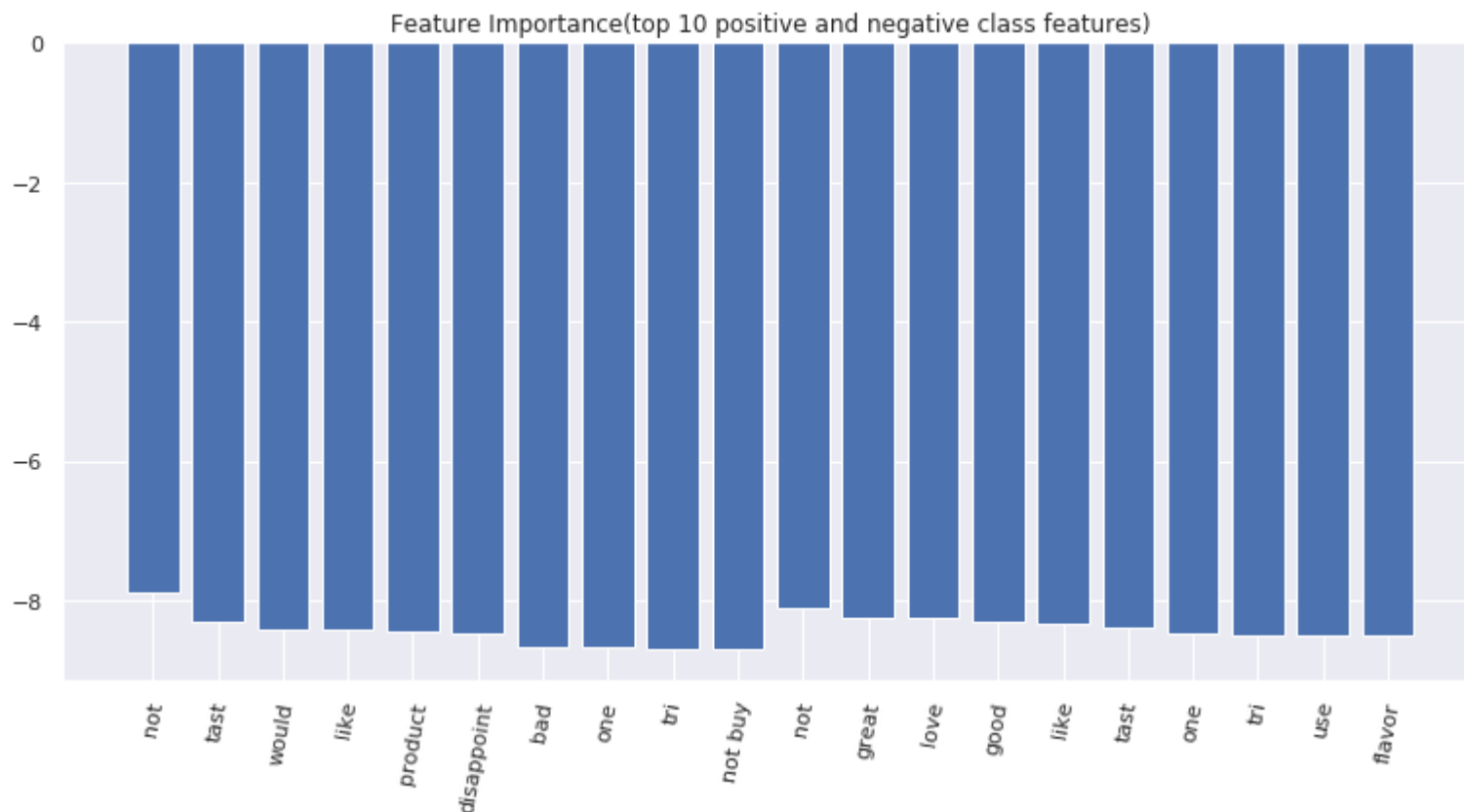
[1.2] Performance on test data with optimal value of hyperparam SET 1

```
In [11]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['alpha'],vect[0],summarize)
```



[1.3] Top 10 important features of positive and negative class from **SET 1**

```
In [12]: 1 #NO. OF IMPORTANT FEATURE TO DISPLAY
2 no_of_imp_features=10
3 #INITIALIZE PRETTYTABLE OBJECT
4 top_features=PrettyTable()
5 top_features.field_names=['S.No.', 'Log-Prob-Neg', 'Negative-Feature', 'Log-Prob-Pos', 'Positive-Feature']
6 feature_importance(count_vect,clf,no_of_imp_features,top_features)
7 print(top_features)
```



S.No.	Log-Prob-Neg	Negative-Feature	Log-Prob-Pos	Positive-Feature
1	-7.8959	not	-8.1201	not
2	-8.3074	tast	-8.2480	great
3	-8.4258	would	-8.2497	love
4	-8.4264	like	-8.3044	good
5	-8.4593	product	-8.3372	like
6	-8.4722	disappoint	-8.3892	tast

	7		-8.6593		bad		-8.4849		one	
	8		-8.6718		one		-8.4891		tri	
	9		-8.6868		tri		-8.5090		use	
	10		-8.6984		not buy		-8.5125		flavor	
+-----+-----+-----+-----+-----+										

[2.] Naive Bayes on BOW with Feature Engineering, SET 1

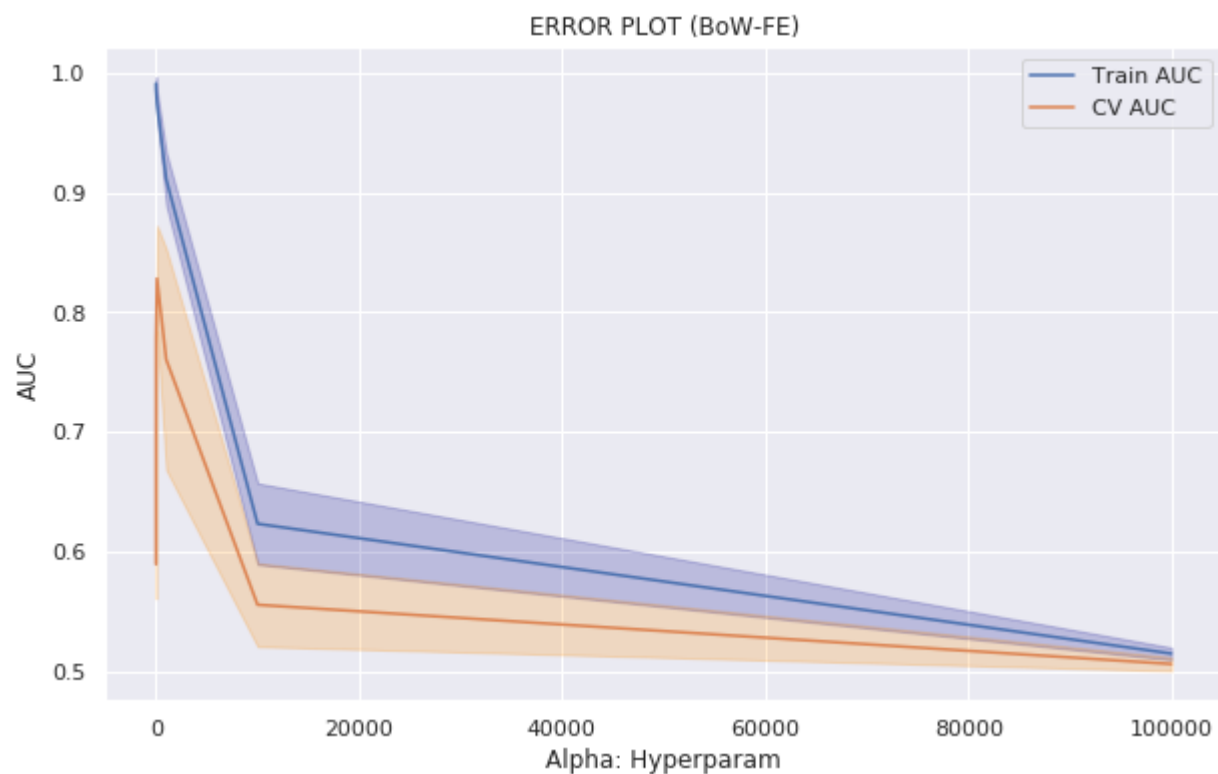
[2.1] Hyperparam Tunning SET 1

In [13]:

```

1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=X_train_bigram_fe,test=X_test_bigram_fe,mean=False)
3 #HYPERPARAM TUNNING
4 %time model=NB_Classifier(train,y_train,TBS,params,searchMethod,vect[1])
5 #PRINT OPTIMAL VALUE OF HYPERPARAM
6 print('Optimal value of Alpha: ',model.best_params_)
7 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
8 saveModeltofile(model,'model_bowfe_mnb')

```



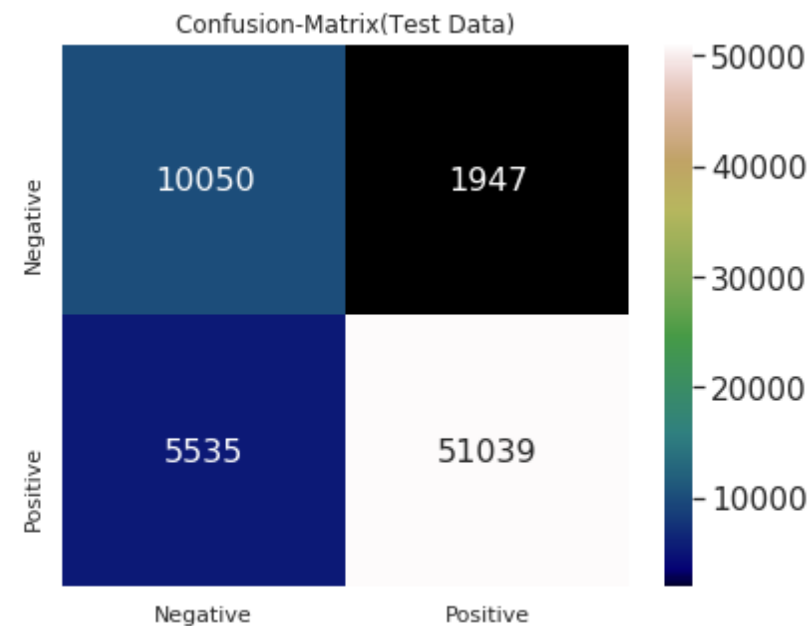
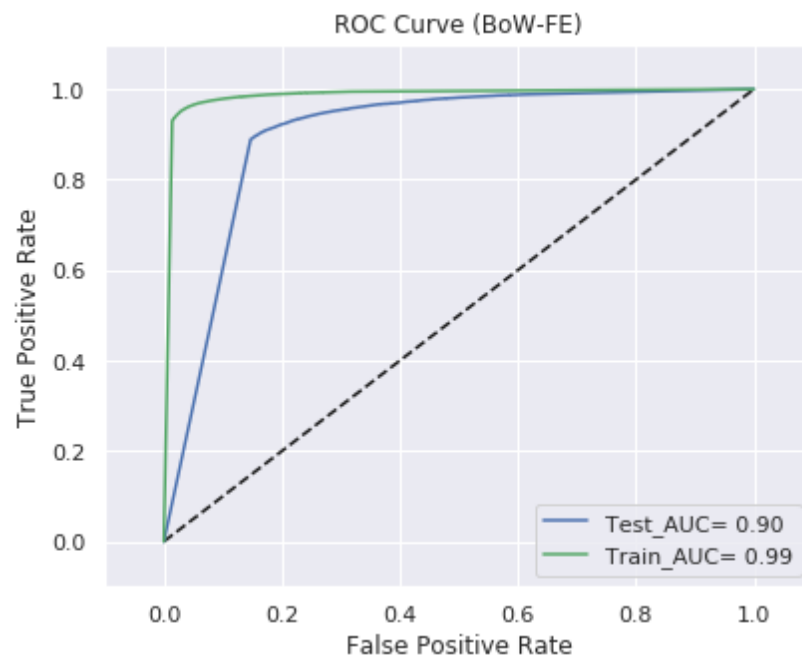
CPU times: user 30.9 s, sys: 520 ms, total: 31.4 s

Wall time: 29.6 s

Optimal value of Alpha: {'alpha': 100}

[2.2] Performance on test data with optimal value of hyperparam SET 1

```
In [14]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['alpha'],vect[1],summarize)
```



Observation:

1. from the above BOW vectors without feature engineering and BOW vectors with feature engineering we can observe that AUC score is improved from '.877' to '.898'.

[3.] Naive Bayes on TFIDF, SET 2

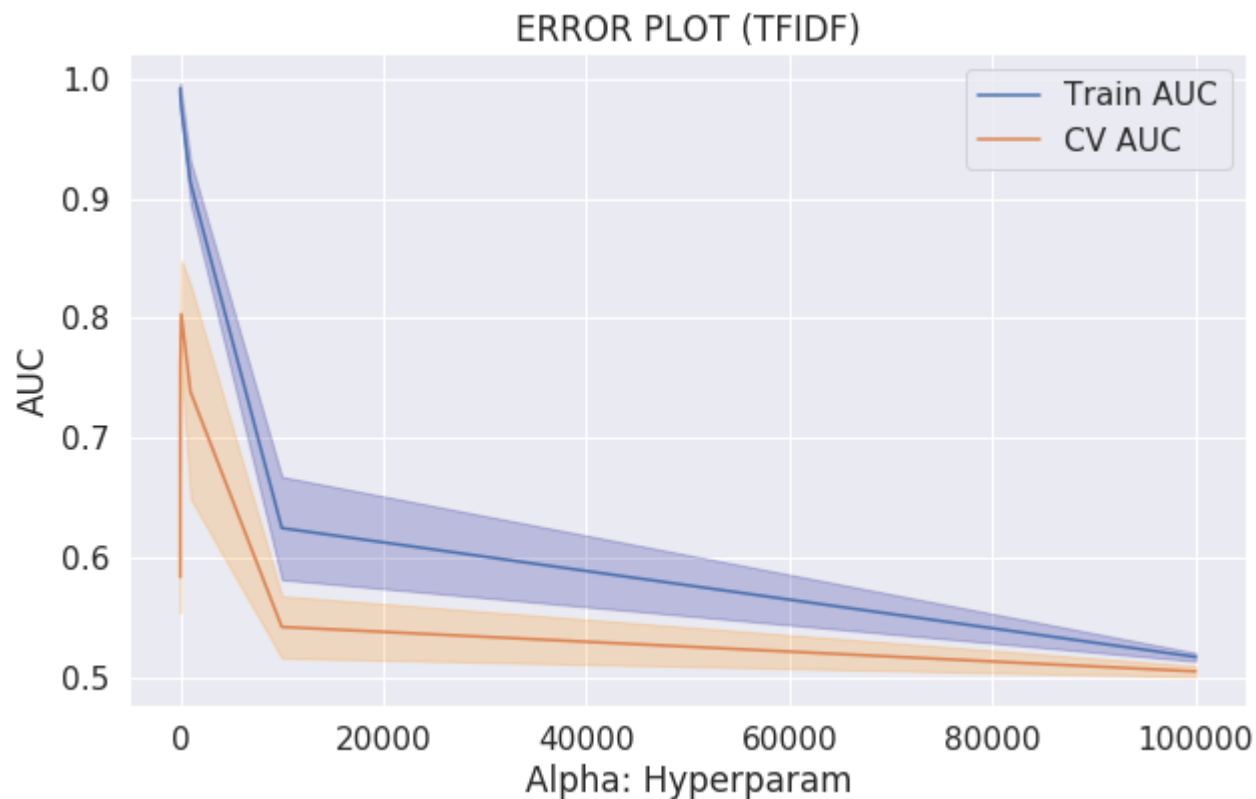
[3.1] Hyperparam Tunning SET 2

In [15]:

```

1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=X_train_tfidf,test=X_test_tfidf,mean=False)
3 #HYPERPARAM TUNNING
4 %time model=NB_Classifier(train,y_train,TBS,params,searchMethod,vect[2])
5 print('Optimal value of Alpha: ',model.best_params_)
6 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
7 saveModeltofile(model,'model_tfidf_mnb')

```



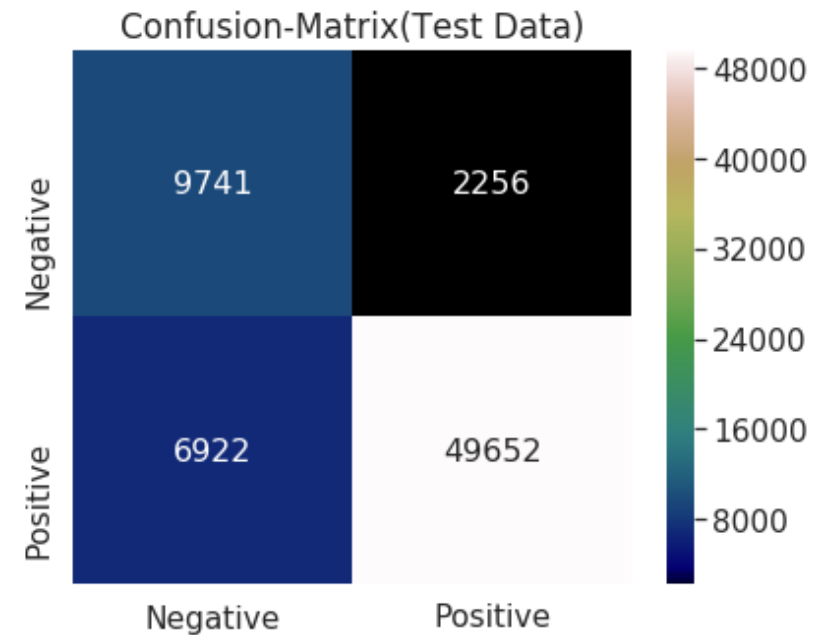
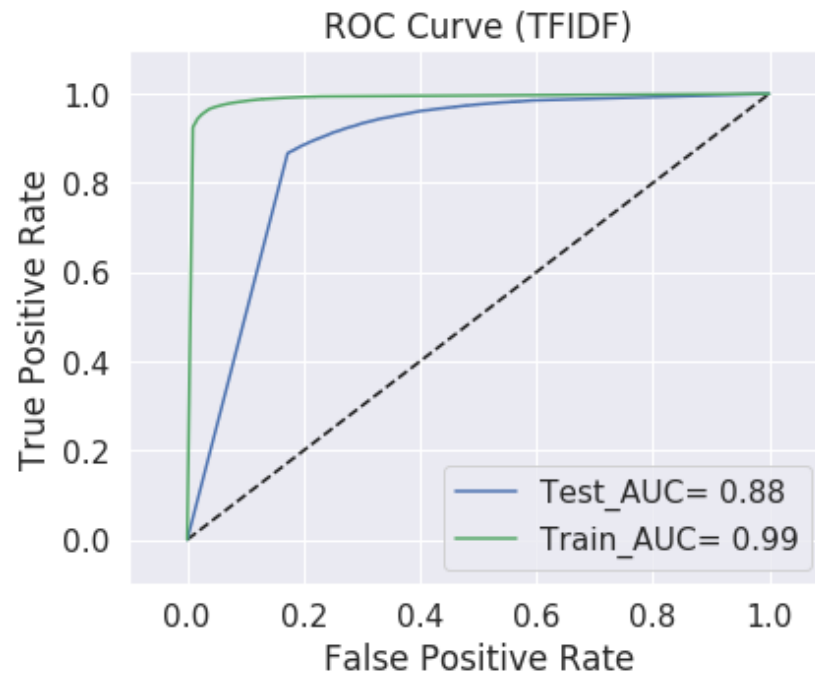
CPU times: user 28.9 s, sys: 604 ms, total: 29.5 s

Wall time: 28 s

Optimal value of Alpha: {'alpha': 100}

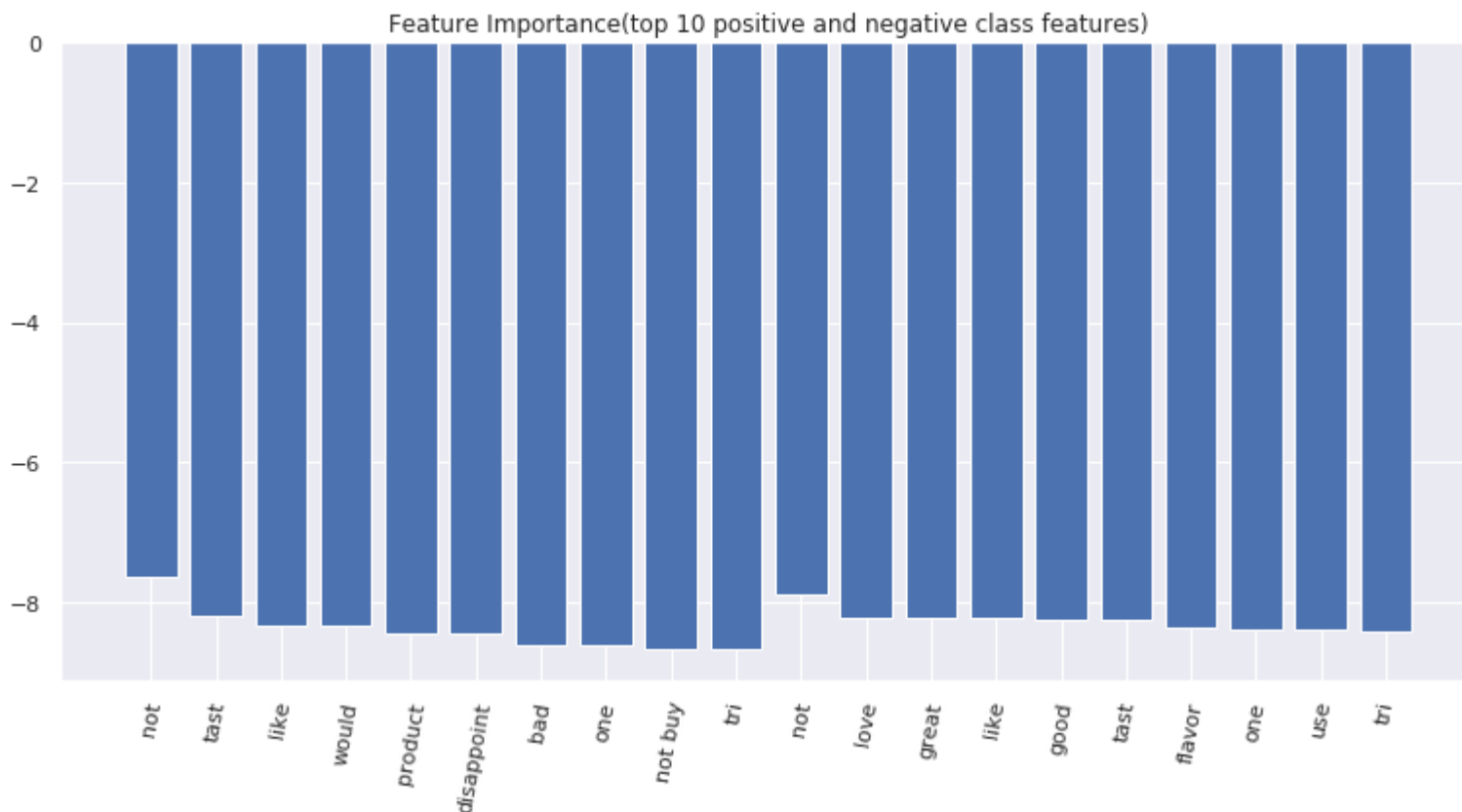
[3.2] Performance on test data with optimal value of hyperparam **SET 2**


```
In [16]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['alpha'],vect[2],summarize)
```



[3.3] Top 10 important features of positive and negative class from SET 2

```
In [17]: 1 no_of_imp_features=10
2 top_features_tf=PrettyTable()
3 top_features_tf.field_names=['S.No.', 'Log-Prob-Neg', 'Negative-Feature', 'Log-Prob-Pos', 'Positive-Feature']
4 feature_importance(tf_idf_vect, clf, no_of_imp_features, top_features_tf)
5 print(top_features_tf)
```



S.No.	Log-Prob-Neg	Negative-Feature	Log-Prob-Pos	Positive-Feature
1	-7.6394	not	-7.9056	not
2	-8.2048	tast	-8.2189	love
3	-8.3319	like	-8.2302	great
4	-8.3336	would	-8.2387	like
5	-8.4429	product	-8.2532	good
6	-8.4443	disappoint	-8.2669	tast
7	-8.6207	bad	-8.3704	flavor
8	-8.6223	one	-8.3885	one

	9		-8.6737		not buy		-8.3976		use	
	10		-8.6749		tri		-8.4302		tri	
+-----+-----+-----+-----+-----+										

[4.] Naive Bayes on TFIDF with Feature Engineering, SET 2

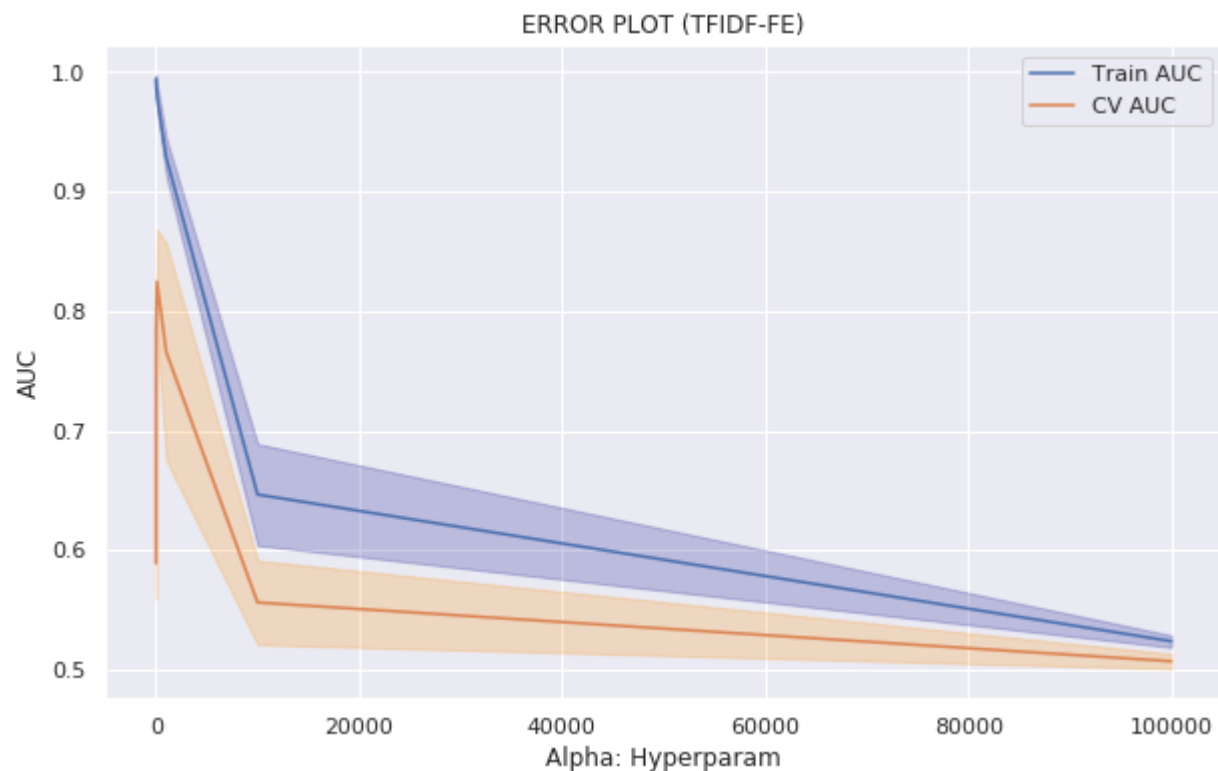
[4.1] Hyperparam Tunning SET 2

In [18]:

```

1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=X_train_tfidf_fe,test=X_test_tfidf_fe,mean=False)
3 #HYPERPARAM TUNNING
4 %time model=NB_Classifier(train,y_train,TBS,params,searchMethod,vect[3])
5 print('Optimal value of Alpha: ',model.best_params_)
6 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
7 saveModeltofile(model,'model_tfidf_fe_mnb')

```



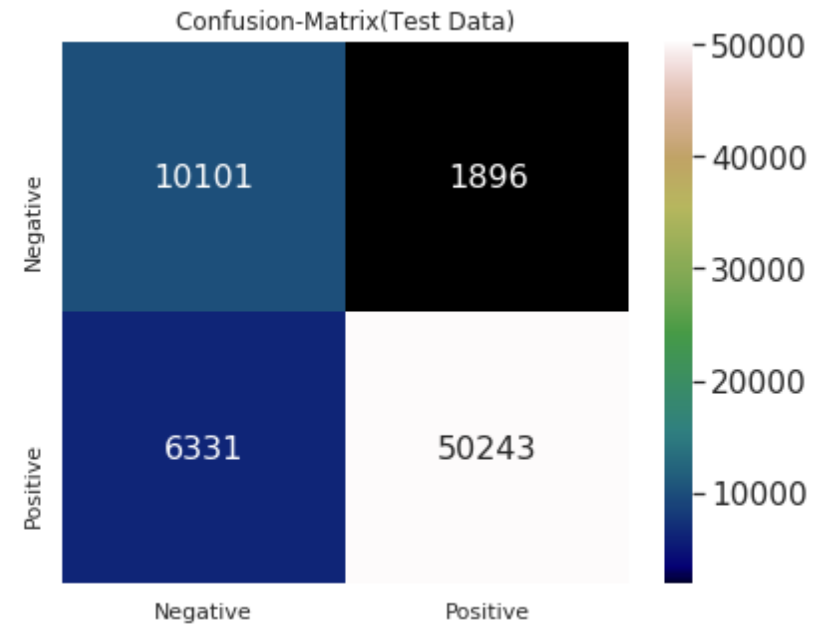
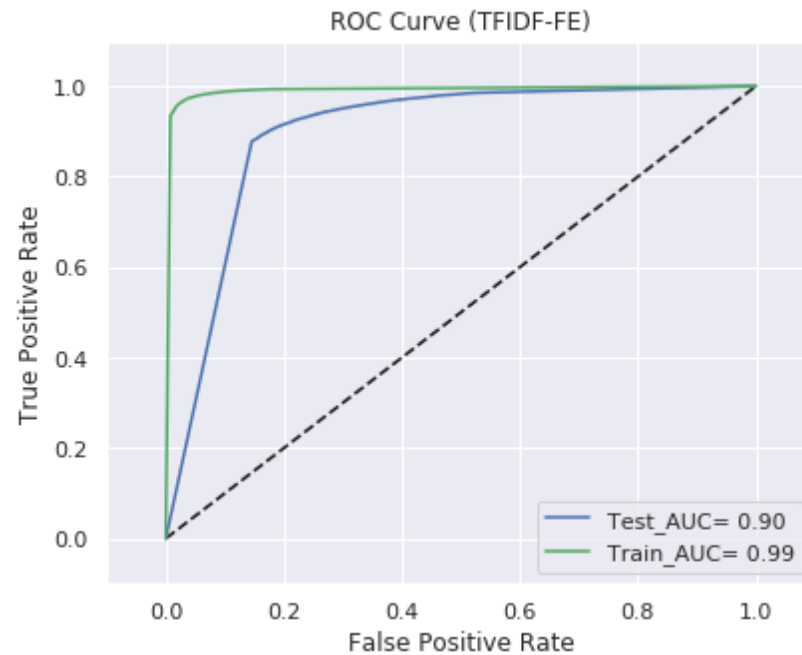
CPU times: user 30.5 s, sys: 576 ms, total: 31.1 s

Wall time: 29.3 s

Optimal value of Alpha: {'alpha': 100}

[4.2] Performance on test data with optimal value of hyperparam **SET 2**

```
In [19]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['alpha'],vect[3],summarize)
```



Observation:

1. from the above TFIDF vectors without feature engineering and TFIDF vectors with feature engineering we can observe that AUC score is improved from '.877' to '.896'.

Conclusions:

```
In [21]: 1 print(summarize)
```

Vectorizer	Optimal-alpha	Test(AUC)	Test(f1-score)
BoW	100	0.877	0.881
BoW-FE	100	0.898	0.896
TFIDF	100	0.877	0.874
TFIDF-FE	100	0.896	0.887

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

- Bag Of Word vectorizer with new features(such as 'length of reviews')
- f1-score=.896 and auc=.898

```
In [ ]: 1
```