

Load necessary libraries

```
In [5]: 1 import warnings
        2 warnings.filterwarnings('ignore')
```

```
In [6]: 1 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
        2 from sklearn.preprocessing import StandardScaler
        3 from sklearn.metrics import *
        4 import pickle
        5 from tqdm import tqdm_notebook
        6 from sklearn.neighbors import KNeighborsClassifier
        7 import seaborn as sns
        8 from sklearn.model_selection import TimeSeriesSplit
        9 from sklearn.model_selection import cross_val_score
       10 from sklearn.metrics import accuracy_score
       11 from sklearn.metrics import confusion_matrix
       12 from sklearn.metrics import f1_score
       13 from sklearn.metrics import precision_score
       14 import numpy as np
       15 import matplotlib.pyplot as plt
       16 import pandas as pd
       17 from prettytable import PrettyTable
       18 from sklearn.externals import joblib
       19 from imblearn.over_sampling import SMOTE
```

Load current state of object

In [3]:

```
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
8
9  y_train =openfromfile('y_train')
10 y_test =openfromfile('y_test')
11
12 count_vect =openfromfile('count_vect')
13 X_train_bigram = openfromfile('X_train_bigram')
14 X_test_bigram = openfromfile('X_test_bigram')
15
16 tf_idf_vect =openfromfile('tf_idf_vect')
17 X_train_tfidf =openfromfile('X_train_tfidf')
18 X_test_tfidf =openfromfile('X_test_tfidf')
19
20 avg_sent_vectors=openfromfile('avg_sent_vectors')
21 avg_sent_vectors_test=openfromfile('avg_sent_vectors_test')
22
23 tfidf_sent_vectors=openfromfile('tfidf_sent_vectors')
24 tfidf_sent_vectors_test=openfromfile('tfidf_sent_vectors_test')
25
26 y_train_kd =openfromfile('y_train_kd')
27 y_test_kd =openfromfile('y_test_kd')
28
29 count_vect_kd =openfromfile('count_vect_kd')
30 X_train_bigram_kd=openfromfile('X_train_bigram_kd')
31 X_test_bigram_kd=openfromfile('X_test_bigram_kd')
32
33 tf_idf_vect_kd =openfromfile('tf_idf_vect_kd')
34 X_train_tfidf_kd=openfromfile('X_train_tfidf_kd')
35 X_test_tfidf_kd=openfromfile('X_test_tfidf_kd')
36
37 avg_sent_vectors_kd=openfromfile('avg_sent_vectors_kd')
38 avg_sent_vectors_test_kd=openfromfile('avg_sent_vectors_test_kd')
39
40 #tfidf_sent_vectors_rbf
41 tfidf_sent_vectors_kd=openfromfile('tfidf_sent_vectors_kd')
```

```
42 tfidf_sent_vectors_test_kd=openfromfile('tfidf_sent_vectors_test_kd')
```

```
In [6]: 1 print('shape of train data used for brute force KNN model: ',X_train_bigram.shape)
2 print('shape of test data used for brute force KNN model: ',X_test_bigram.shape)
3
4 print('shape of train data used for kd-tree KNN model: ',X_train_bigram_kd.shape)
5 print('shape of test data used for kd-tree KNN model: ',X_test_bigram_kd.shape)
```

shape of train data used for brute force KNN model: (49000, 30357)
 shape of test data used for brute force KNN model: (21000, 30357)
 shape of train data used for kd-tree KNN model: (14000, 500)
 shape of test data used for kd-tree KNN model: (6000, 500)

Save and Load Model:

```
In [8]: 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 [9]: 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
```

KNN

Function for finding optimal value of hyperparameter nd plot missclassification error vs k :

```

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

```

```
42     plt.ylabel('AUC')
43     plt.grid(True)
44     plt.legend()
45     plt.show()
46     return model
47
```

Function which calculate performance on test data with optimal K :

```

In [11]: 1 def test_performance(x_train,y_train,x_test,y_test,optimal_k,algo,vect,summarize):
2         '''FUNCTION FOR TEST PERFORMANCE(PLOT ROC CURVE FOR BOTH TRAIN AND TEST) WITH OPTIMAL_K'''
3         clf=KNeighborsClassifier(algorithm=algo,n_neighbors=optimal_k,n_jobs=-1)
4         clf.fit(x_train,y_train)
5         data_used=['Test-Data', 'Train-Data']
6
7         test_probability = clf.predict_proba(x_test)[:,-1]
8         train_probability = clf.predict_proba(x_train)[:,-1]
9         fpr_test, tpr_test, threshold_test = roc_curve(y_test, test_probability)
10        fpr_train, tpr_train, threshold_train = roc_curve(y_train, train_probability)
11        auc_score_test=auc(fpr_test, tpr_test)
12        auc_score_train=auc(fpr_train, tpr_train)
13        y_pred={} ; y_act={};
14        y_pred[data_used[0]]=clf.predict(x_test)
15        y_pred[data_used[1]]=clf.predict(x_train)
16        y_act[data_used[0]]=y_test
17        y_act[data_used[1]]=y_train
18        saveModeltofile(auc_score_train,'auc_score_train'+algo+vect)
19        saveModeltofile(auc_score_test,'auc_score_test'+algo+vect)
20
21        f1=f1_score(y_test,y_pred[data_used[0]],average='weighted')
22
23        #ADD RESULTS TO PRETTY TABLE
24        summarize.add_row([vect, algo, optimal_k, '%.3f'%auc_score_train, '%.3f'%auc_score_test, '%.3f'%f1])
25
26        plt.figure(1,figsize=(14,5))
27        sns.set_style('darkgrid')
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     #PLOT CONFUSION MATRIX USING HEATMAP
43     plt.figure(2,figsize=(16,6))
44     sns.set_style('darkgrid')
45     for k in range(2):
46         #PLOT CONFUSION MATRIX USING HEATMAP
47         plt.subplot(int('12'+str(k+1)))
48         plt.title('Confusion-Matrix (%s)' %data_used[k])
49         df_cm = pd.DataFrame(confusion_matrix(y_act[data_used[k]],y_pred[data_used[k]]),\
50                             ['Negative','Positive'],['Negative','Positive'])
51         sns.set(font_scale=1.4)#for label size
52         sns.heatmap(df_cm,cmap='gist_earth', annot=True,annot_kws={"size": 16}, fmt='g')
53     plt.show()

```

Initialization of common objects required for all vectorization:

In [12]:

```

1  #ALGO USED
2  algo=['brute','kd_tree']
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 K VALUES(HYPERPARAM)
10 k={'n_neighbors':[x for x in range(3,30,2)]}
11 #INITIALIZE PRETTY TABLE OBJECT
12 summarize = PrettyTable()
13 summarize.field_names = ['Vectorizer', 'Algorithm', 'Optimal-K', 'Train(AUC)', 'Test(AUC)', 'Test(F1-Score)']

```

Apply KNN for BoW vectorizer

1. Brute force implementation

```
In [12]: 1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=X_train_bigram,test=X_test_bigram,mean=False)
3 #train, y_train= sm.fit_sample(train, y_train)
4 print('shape of oversampled data y_train:',y_train.shape,'X_train:',train.shape)
5 #HYPERPARAM TUNNING
6 %time model=KNN_Classifier(train,y_train,TBS,k,searchMethod,algo[0],vect[0])
7 print('Optimal value of K: ',model.best_params_)
8 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
9 saveModeltofile(model,'model_bow_knn')
```

shape of oversampled data y_train: (49000,) X_train: (49000, 30357)

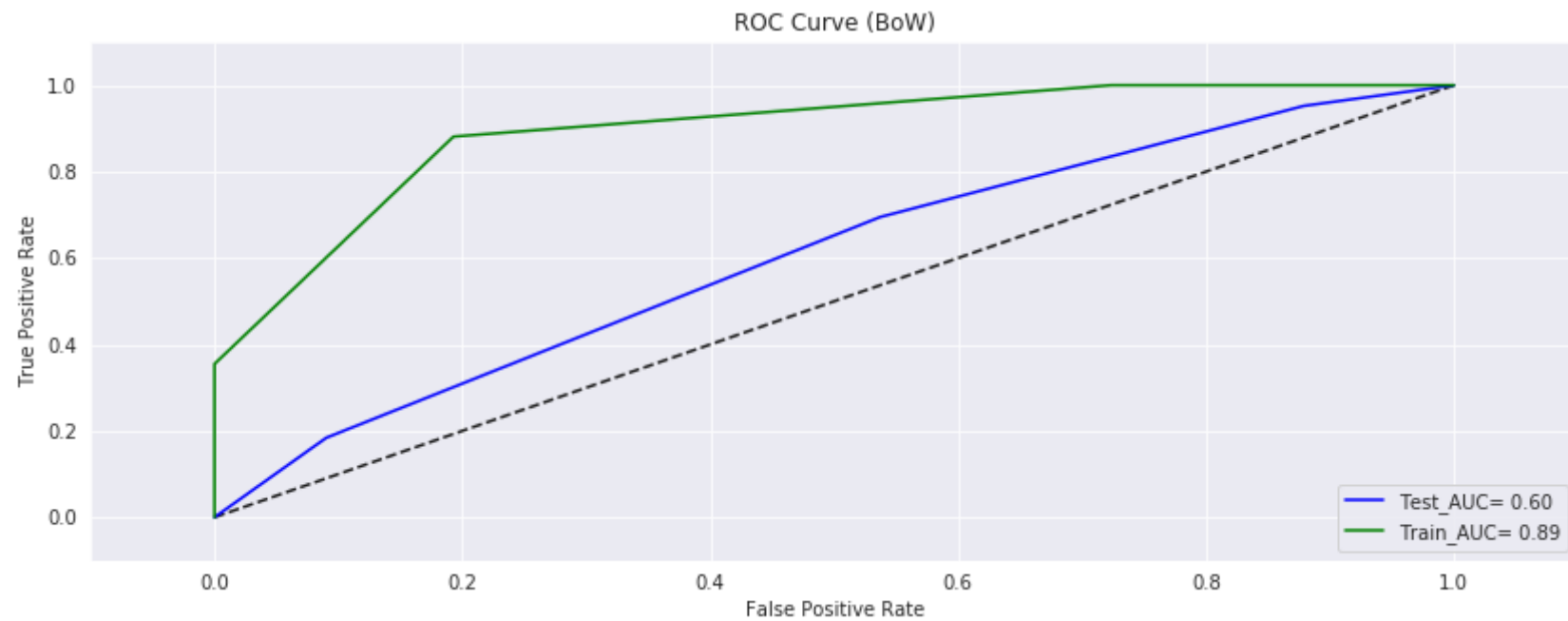


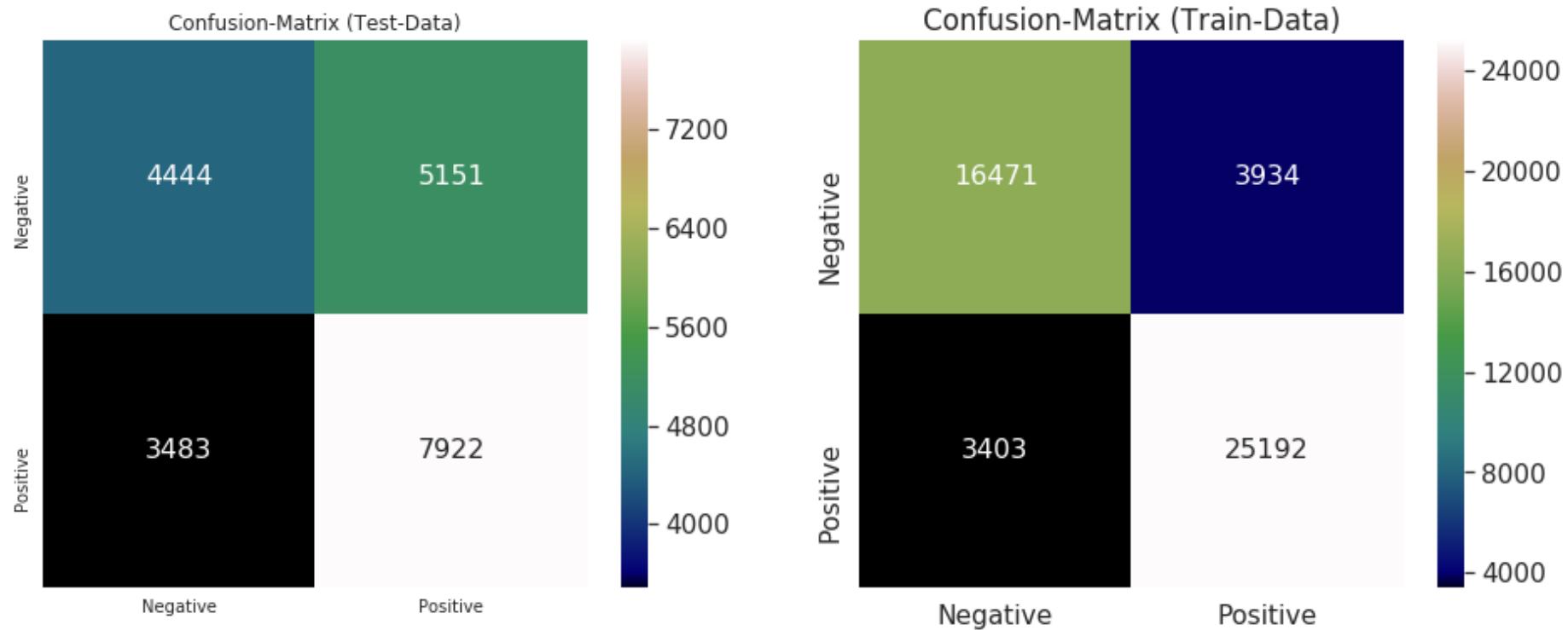
CPU times: user 1.55 s, sys: 2.04 s, total: 3.6 s

Wall time: 8min 29s

Optimal value of K: {'n_neighbors': 3}


```
In [13]: 1 test_performance(train,y_train,test,y_test,model.best_params_['n_neighbors'],algo[0],vect[0],summarize)
```





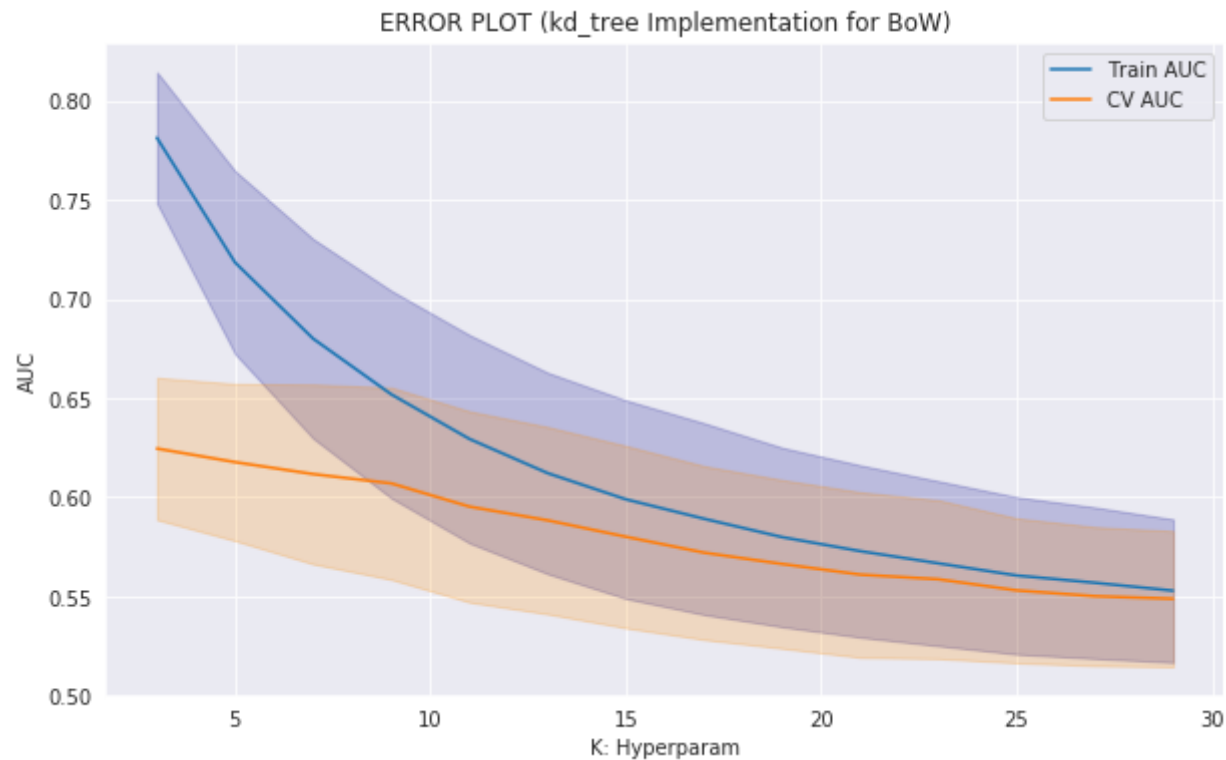
2. kd-tree implementation

In [12]:

```

1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=X_train_bigram_kd,test=X_test_bigram_kd,mean=False)
3 train_dense,test_dense=train.toarray(), test.toarray()
4 #HYPERPARAM TUNNING
5 %time model=KNN_Classifier(train_dense,y_train_kd,TBS,k,searchMethod,algo[1],vect[0])
6 print('Optimal value of K: ',model.best_params_)
7 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
8 saveModeltofile(model,'model_bow_kd_knn')

```

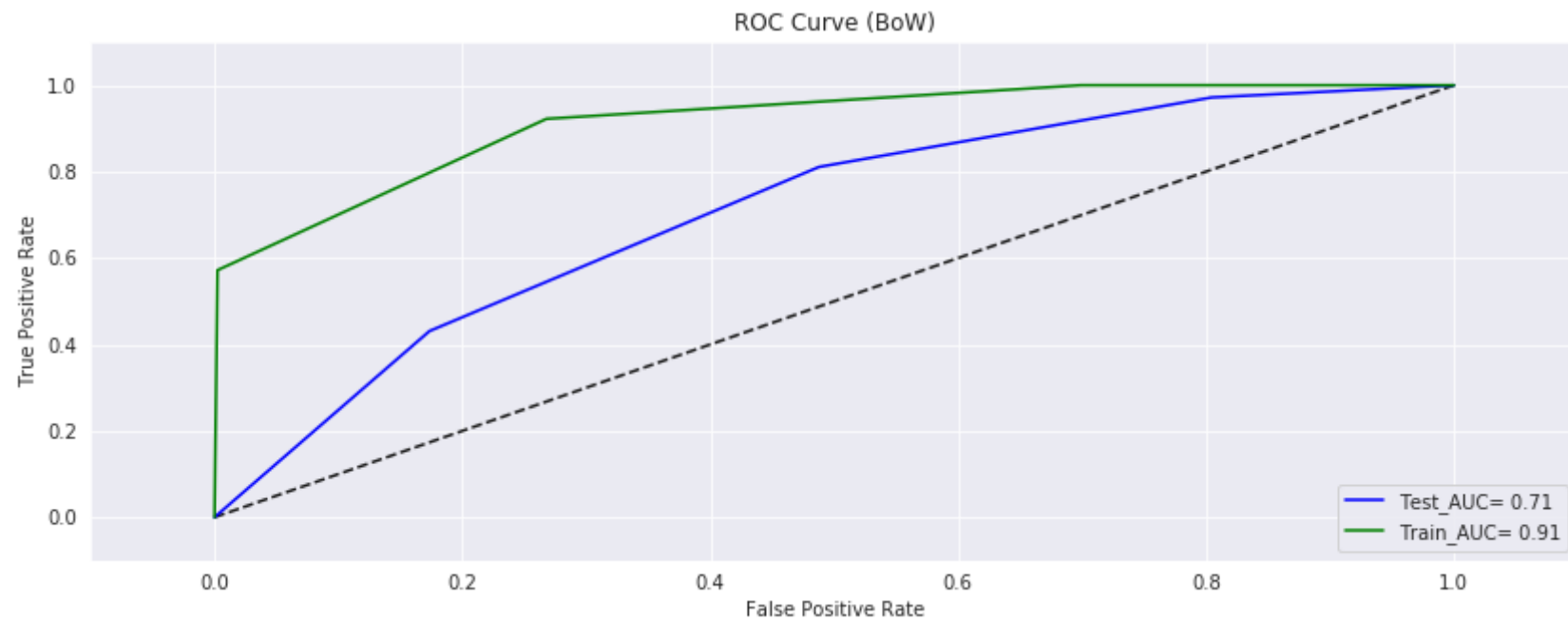


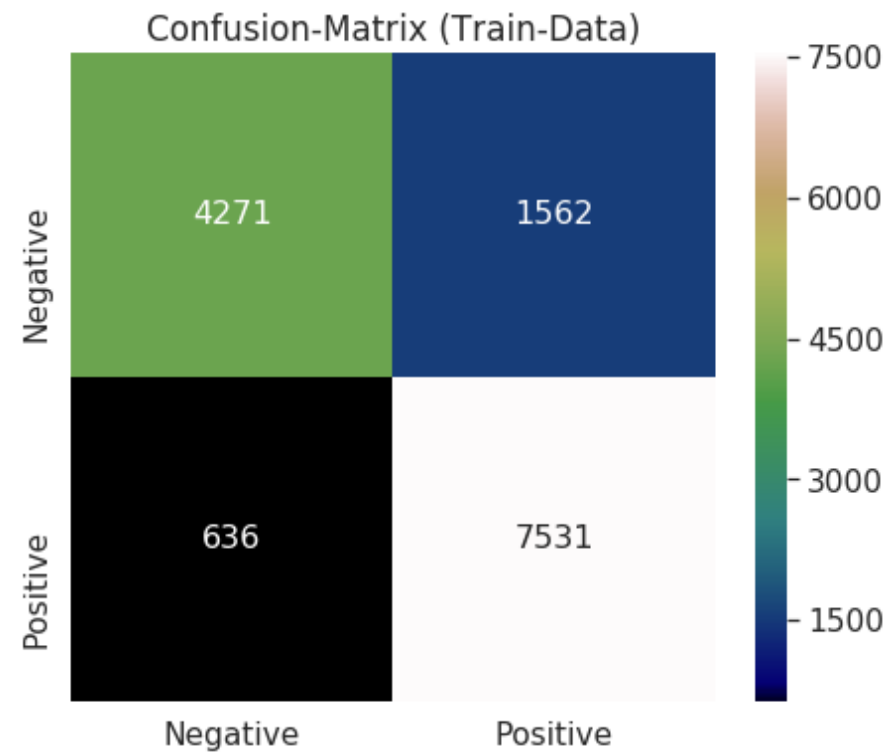
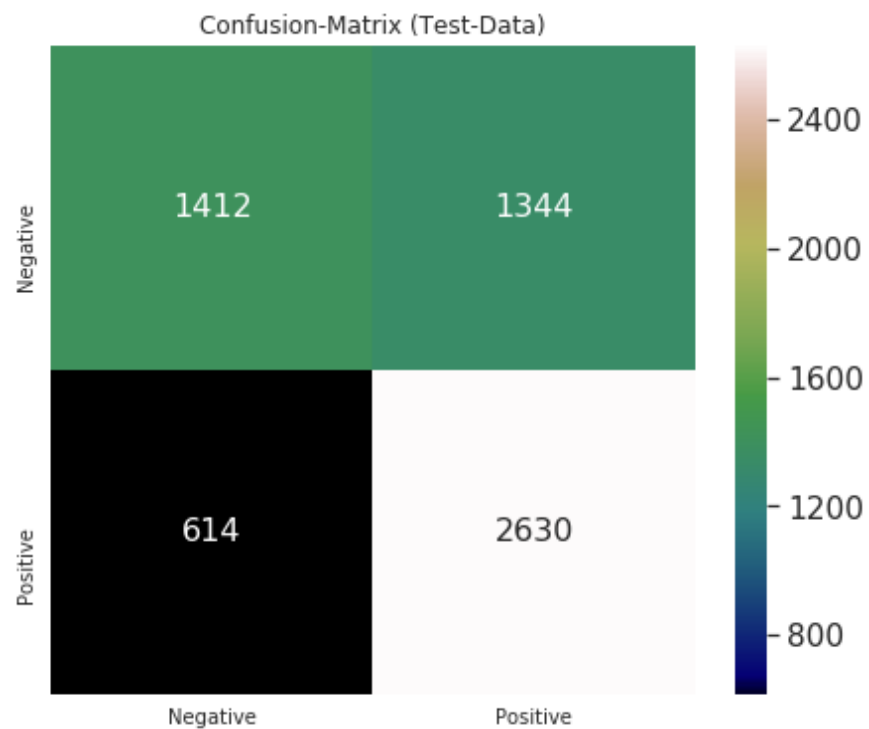
CPU times: user 2.99 s, sys: 2.11 s, total: 5.1 s

Wall time: 14min 52s

Optimal value of K: {'n_neighbors': 3}

```
In [13]: 1 test_performance(train_dense,y_train_kd,test_dense,y_test_kd,model.best_params_['n_neighbors'],algo[1],vect[0],summa
```





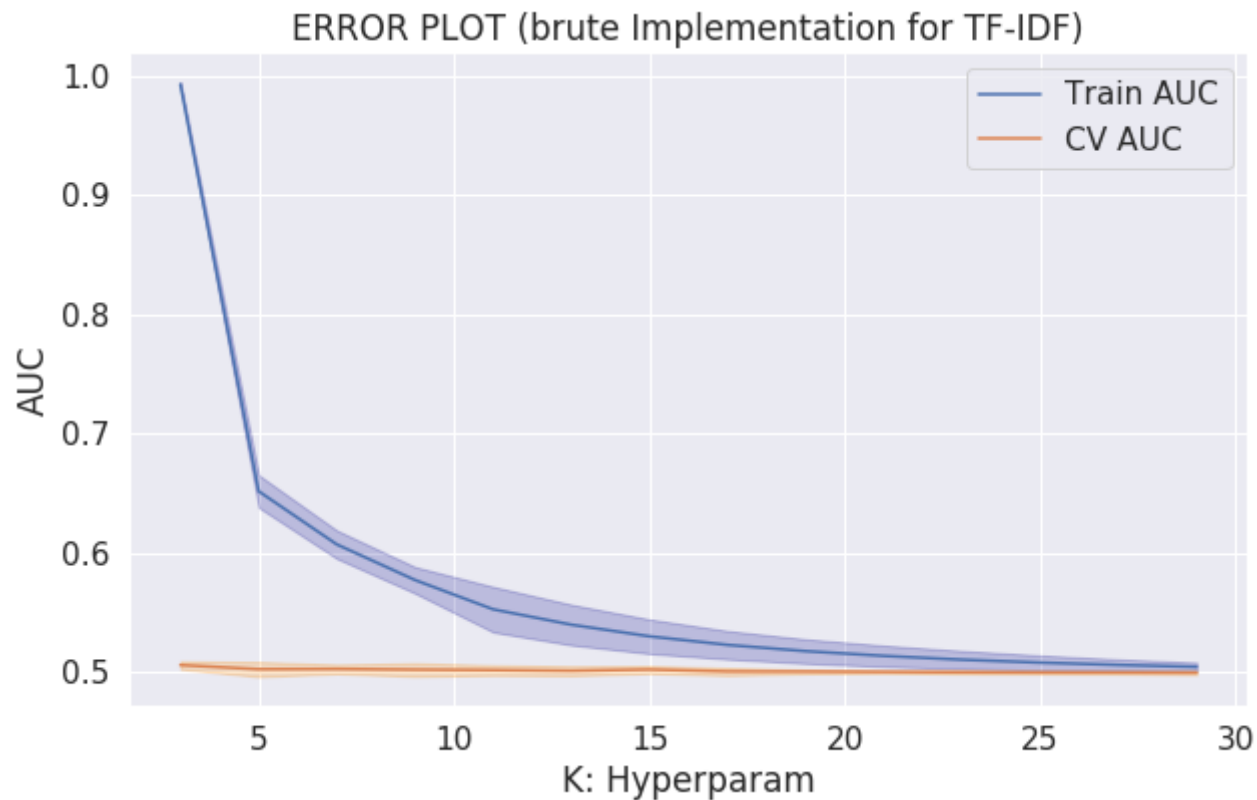
Apply KNN for TF-IDF vectorizer

1. Brute force implementation

```

In [15]: 1 #X_train, y_train =sm.fit_sample(X_train, y_train)
2 #STANDARDIZE TRAIN AND TEST DATA
3 train, test=std_data(train=X_train_tfidf,test=X_test_tfidf,mean=False)
4 #OVERSAMPLE TRAIN DATA
5 #train, y_train =sm.fit_sample(train, y_train)
6 #HYPERPARAM TUNNING
7 %time model=KNN_Classifier(train,y_train,TBS,k,searchMethod,algo[0],vect[1])
8 print('Optimal value of K: ',model.best_params_)
9 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
10 saveModeltofile(model,'model_tfidf_knn')

```

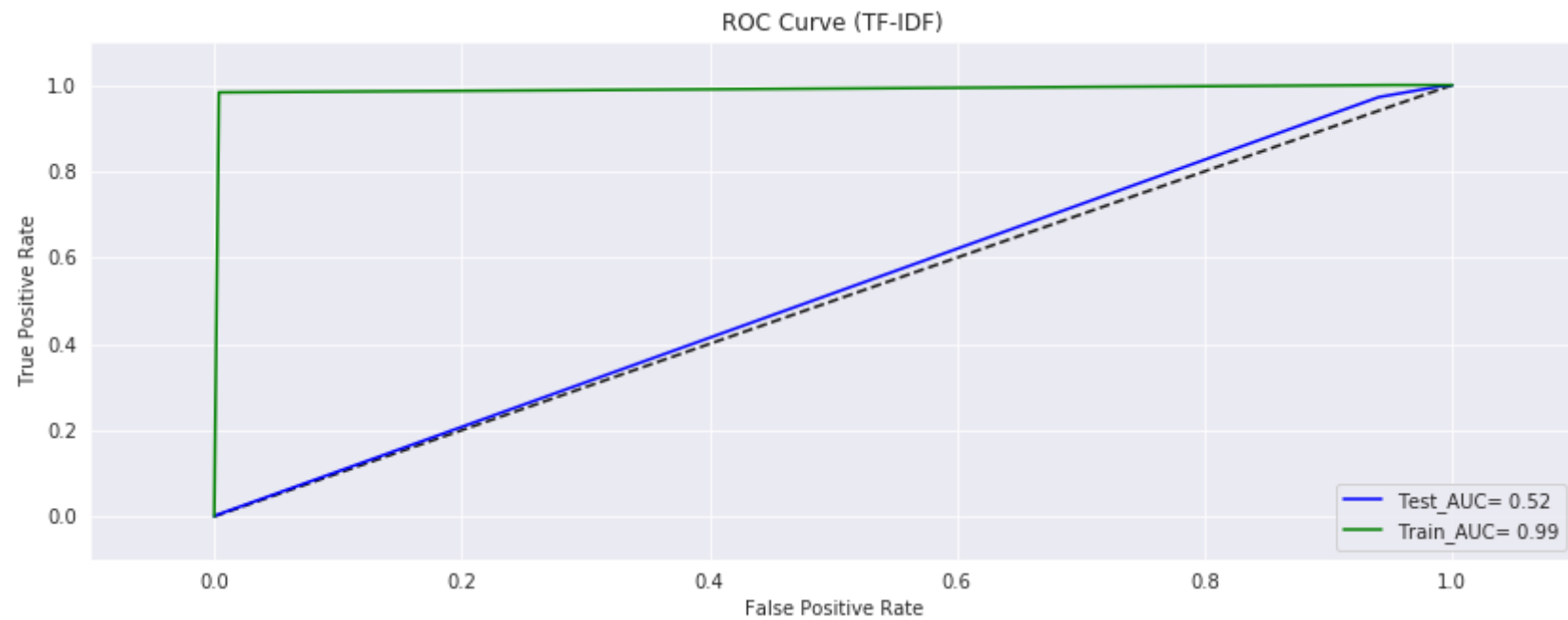


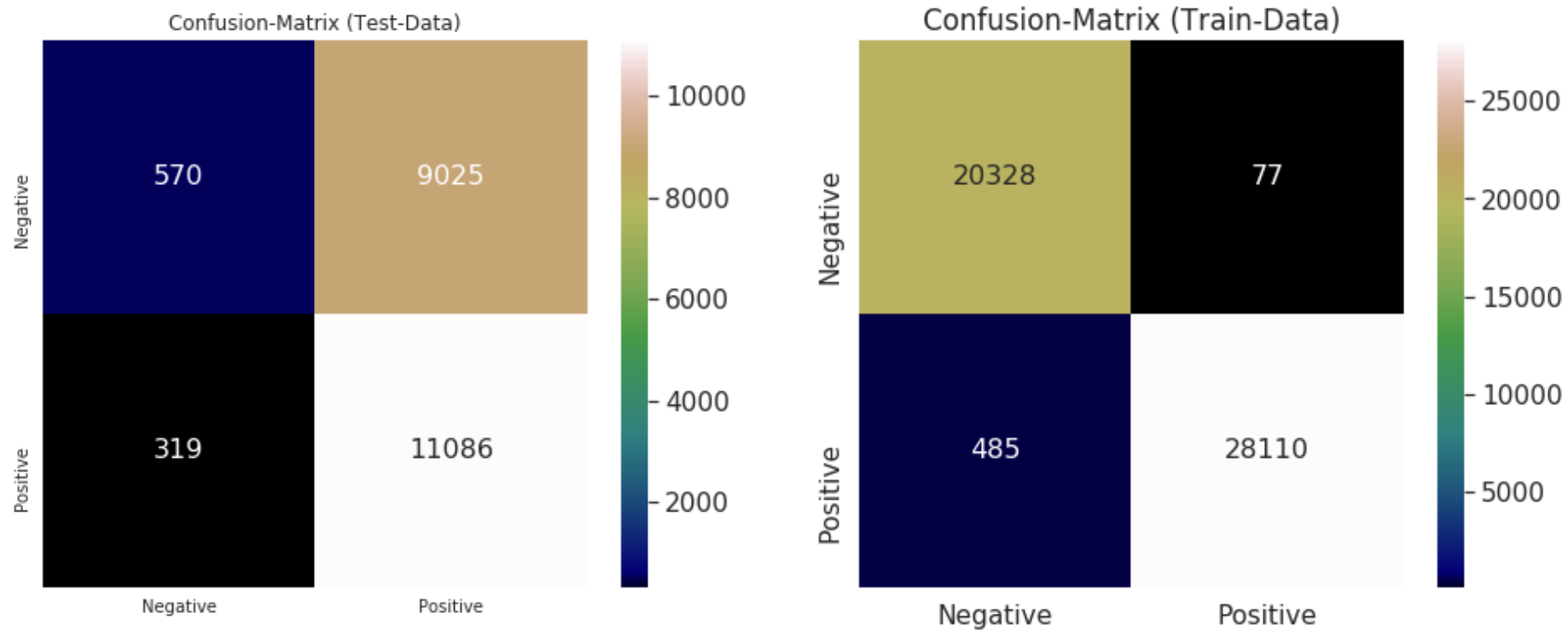
CPU times: user 1.17 s, sys: 820 ms, total: 1.99 s

Wall time: 9min 47s

Optimal value of K: {'n_neighbors': 3}

```
In [15]: 1 test_performance(train,y_train,test,y_test,model.best_params_['n_neighbors'],algo[0],vect[1],summarize)
```





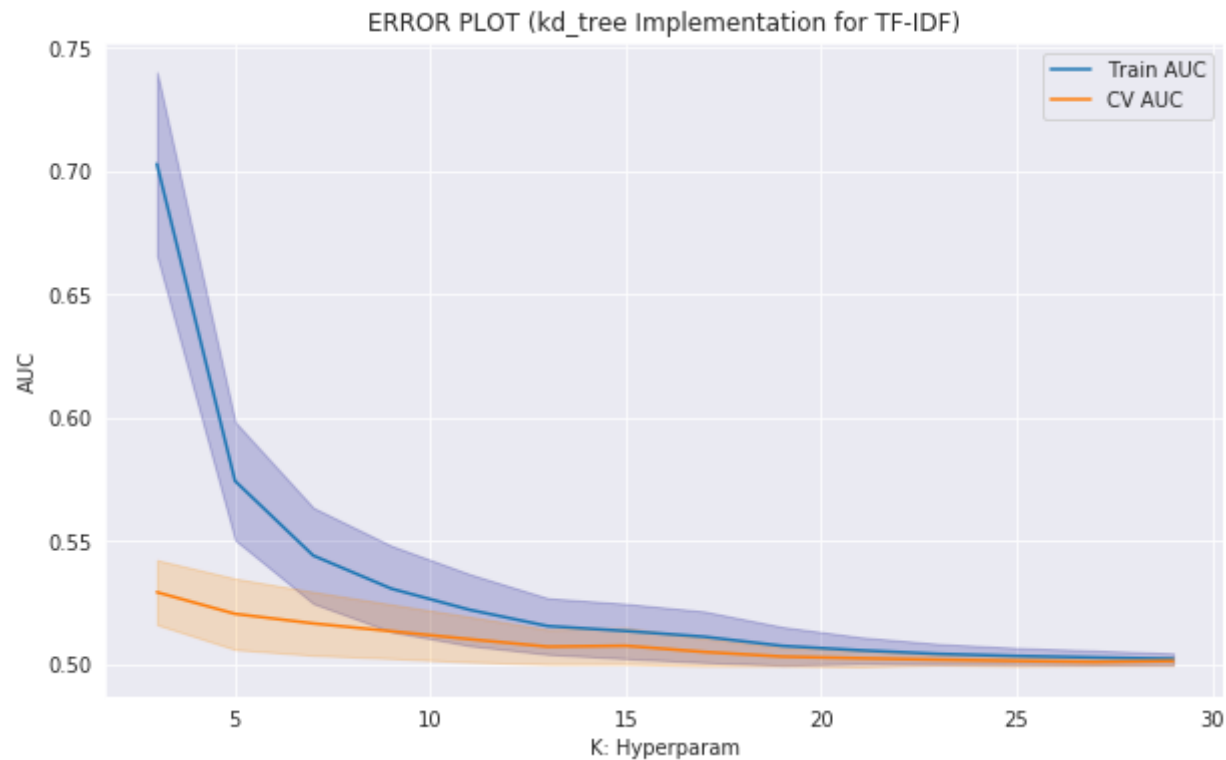
2. kd-tree implementation

In [15]:

```

1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=X_train_tfidf_kd,test=X_test_tfidf_kd,mean=False)
3 train_dense,test_dense=train.toarray(), test.toarray()
4 #HYPERPARAM TUNNING
5 %time model=KNN_Classifier(train_dense,y_train_kd,TBS,k,searchMethod,algo[1],vect[1])
6 print('Optimal value of K: ',model.best_params_)
7 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
8 saveModeltofile(model,'model_tfidf_kd_knn')

```

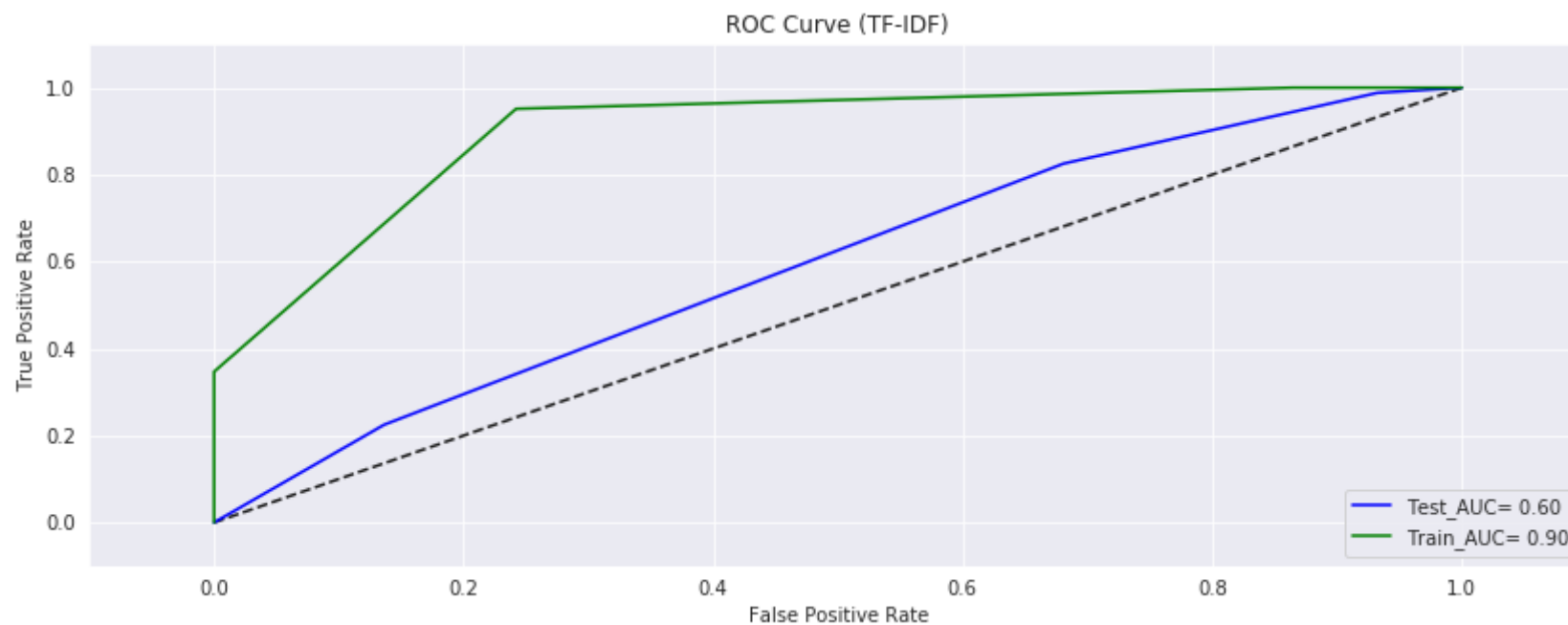


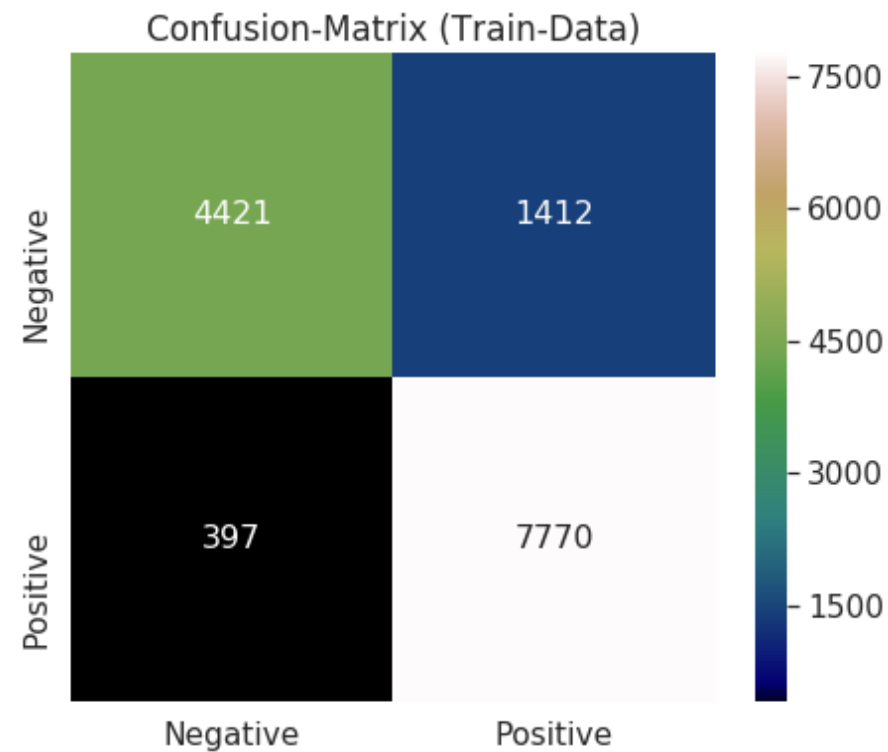
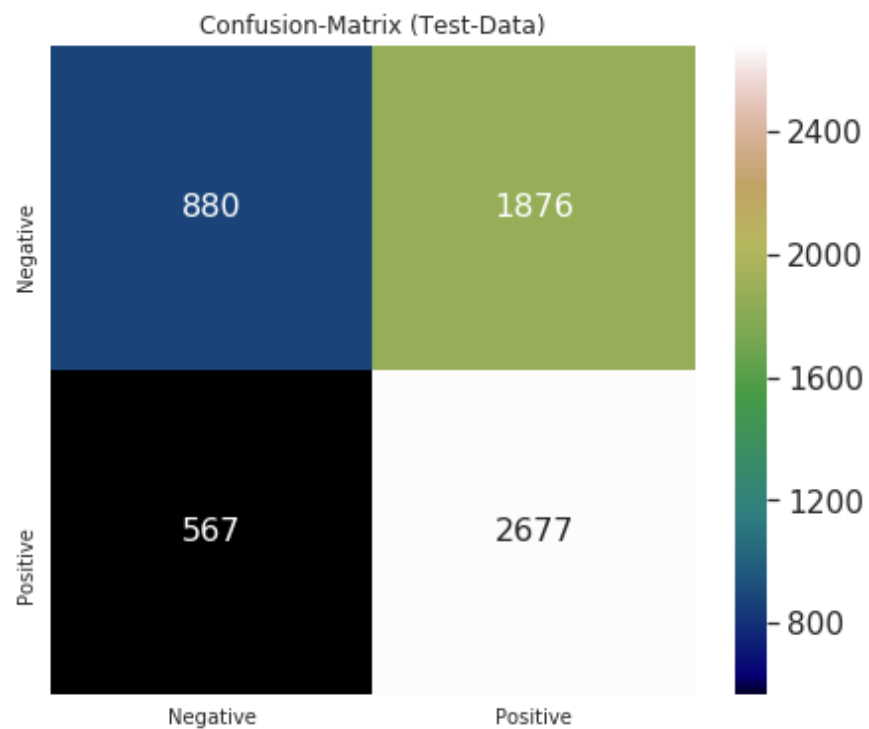
CPU times: user 2.56 s, sys: 1.18 s, total: 3.74 s

Wall time: 14min 55s

Optimal value of K: {'n_neighbors': 3}

```
In [16]: 1 test_performance(train_dense,y_train_kd,test_dense,y_test_kd,model.best_params_['n_neighbors'],algo[1],vect[1],summa
```





Apply KNN for AVG-W2V vectorizer

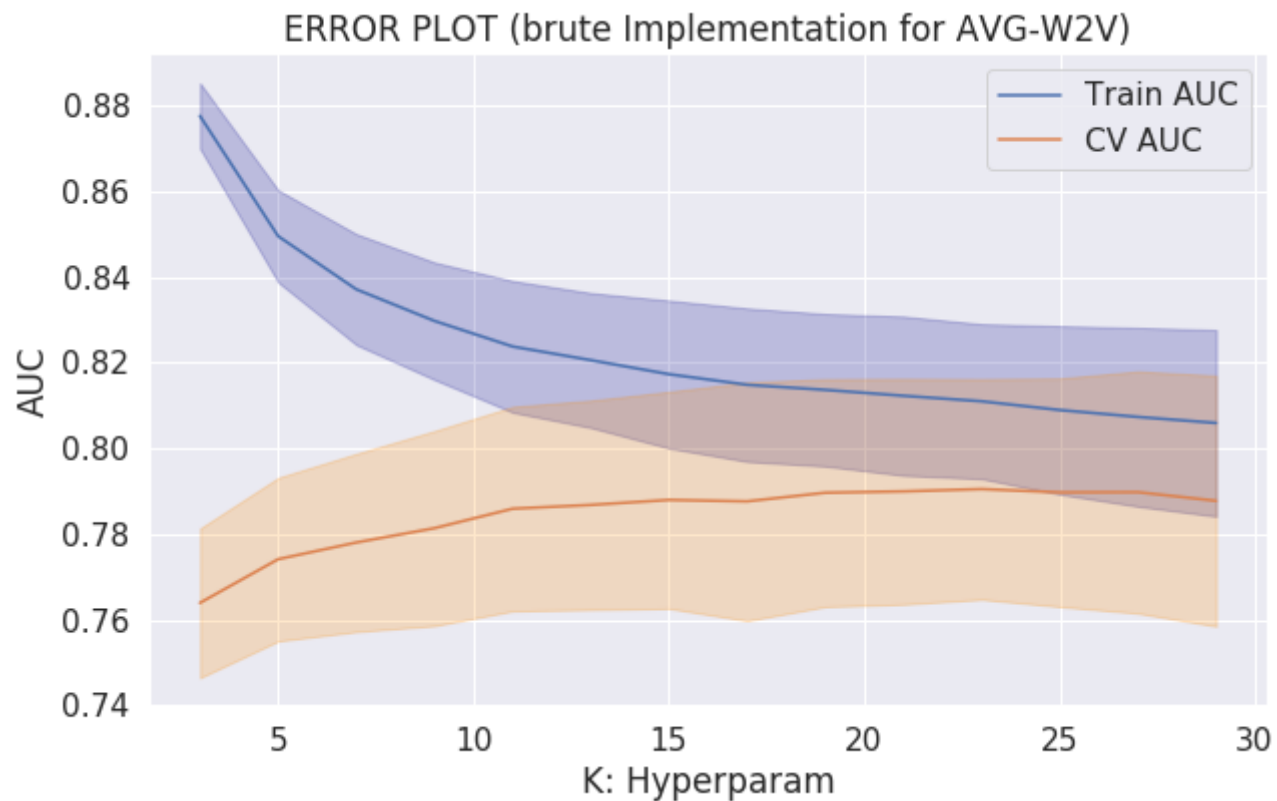
1. Brute force implementation

In [23]:

```

1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=avg_sent_vectors,test=avg_sent_vectors_test,mean=True)
3 #OVERSAMPLE TRAIN DATA
4 #train, y_train = oversampling(train, y_train)
5 #HYPERPARAM TUNNING
6 %time model=KNN_Classifier(train,y_train,TBS,k,searchMethod,algo[0],vect[2])
7 print('Optimal value of K: ',model.best_params_)
8 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
9 saveModeltofile(model,'model_avgw2v_knn')

```

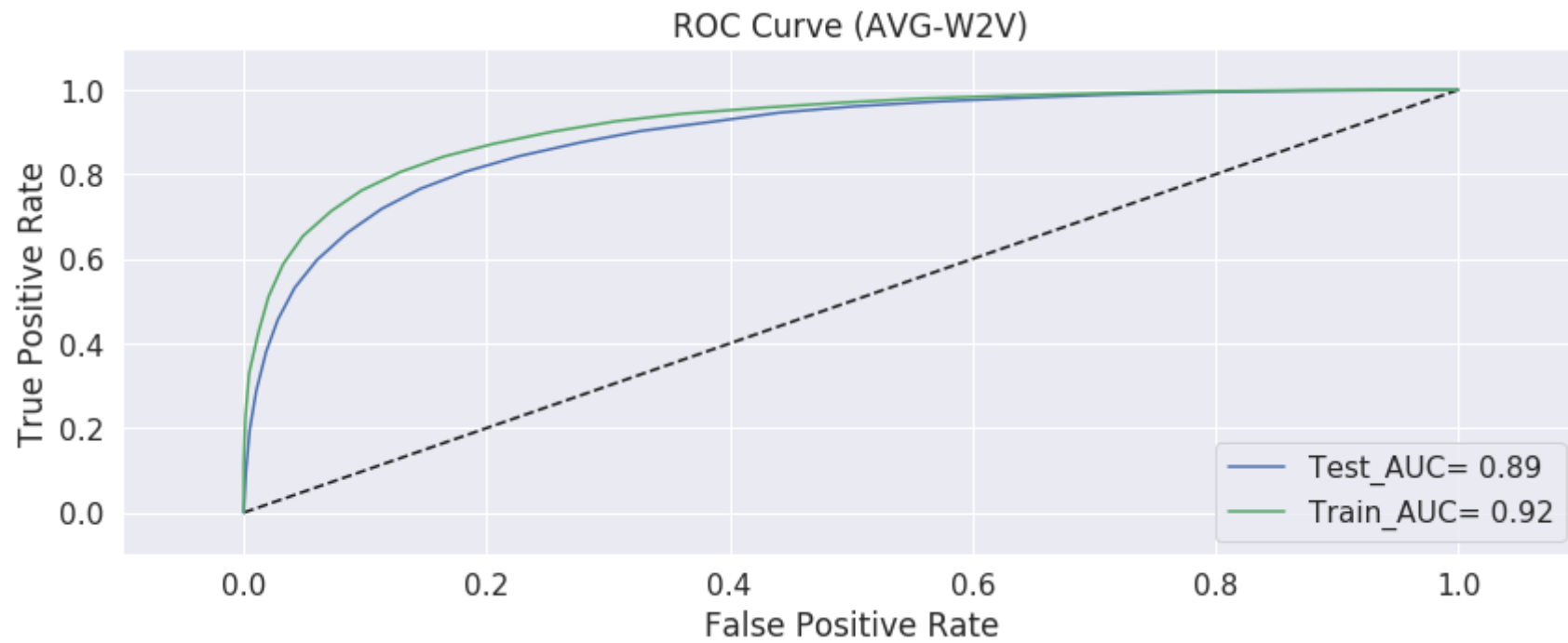


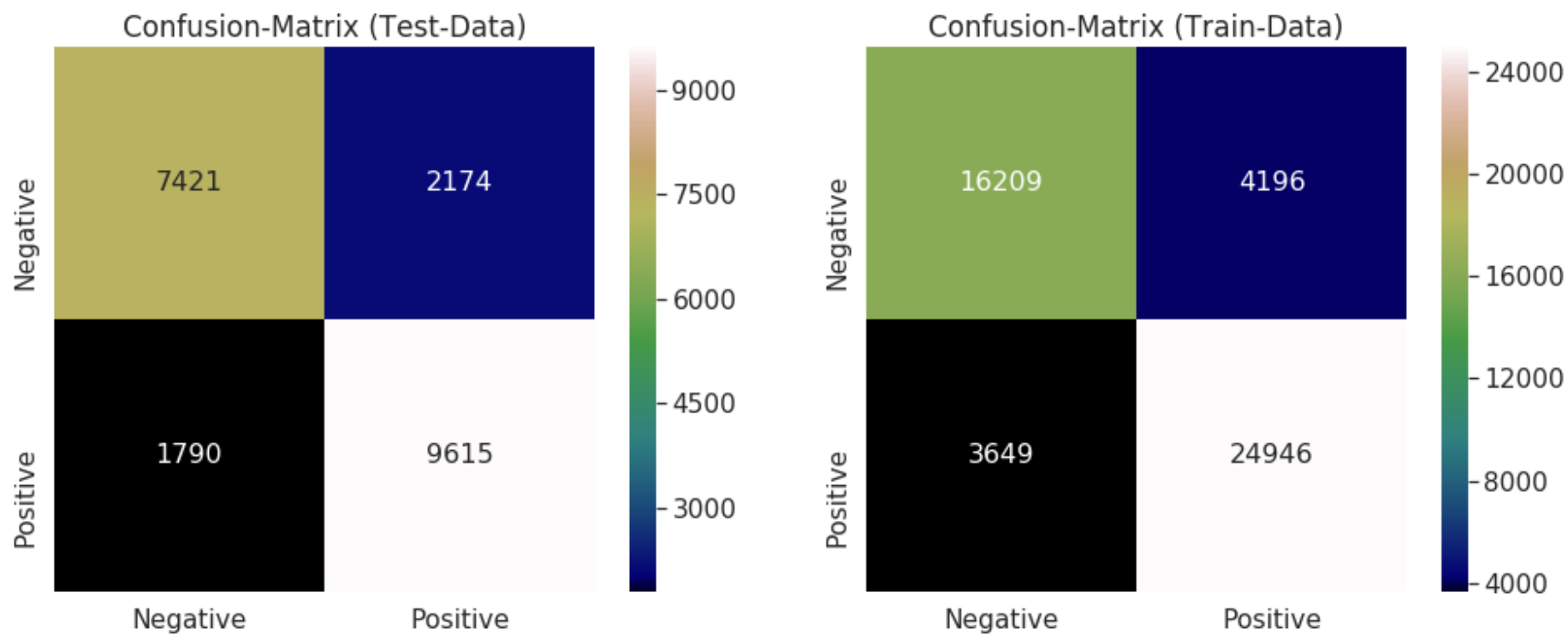
CPU times: user 1.54 s, sys: 1.7 s, total: 3.24 s

Wall time: 6min 14s

Optimal value of K: {'n_neighbors': 23}

```
In [15]: 1 test_performance(train,y_train,test,y_test,model.best_params_['n_neighbors'],algo[0],vect[2],summarize)
```



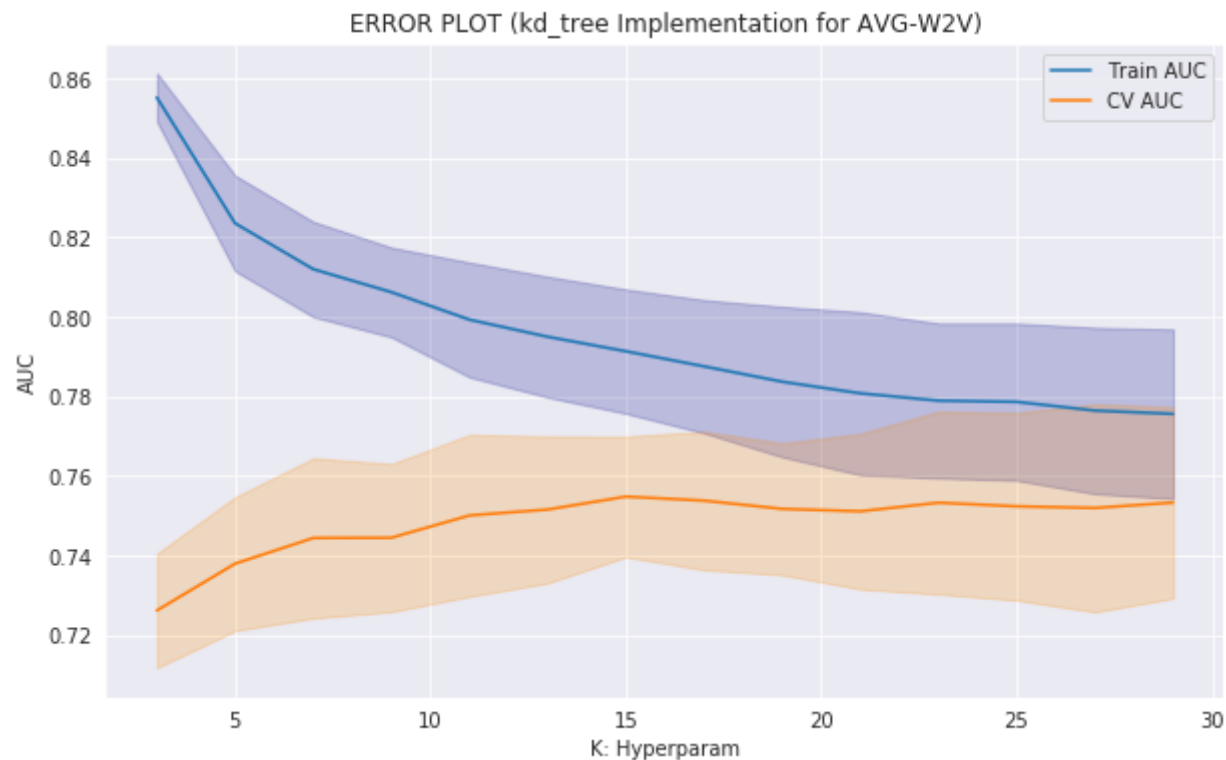


2. kd-tree implementation

```

In [11]: 1 train, test=std_data(train=avg_sent_vectors_kd,test=avg_sent_vectors_test_kd,mean=True)
          2 #HYPERPARAM TUNNING
          3 %time model=KNN_Classifier(train,y_train_kd,TBS,k,searchMethod,algo[1],vect[2])
          4 print('Optimal value of K: ',model.best_params_)
          5 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
          6 saveModeltofile(model,'model_avgw2v_kd_knn')

```

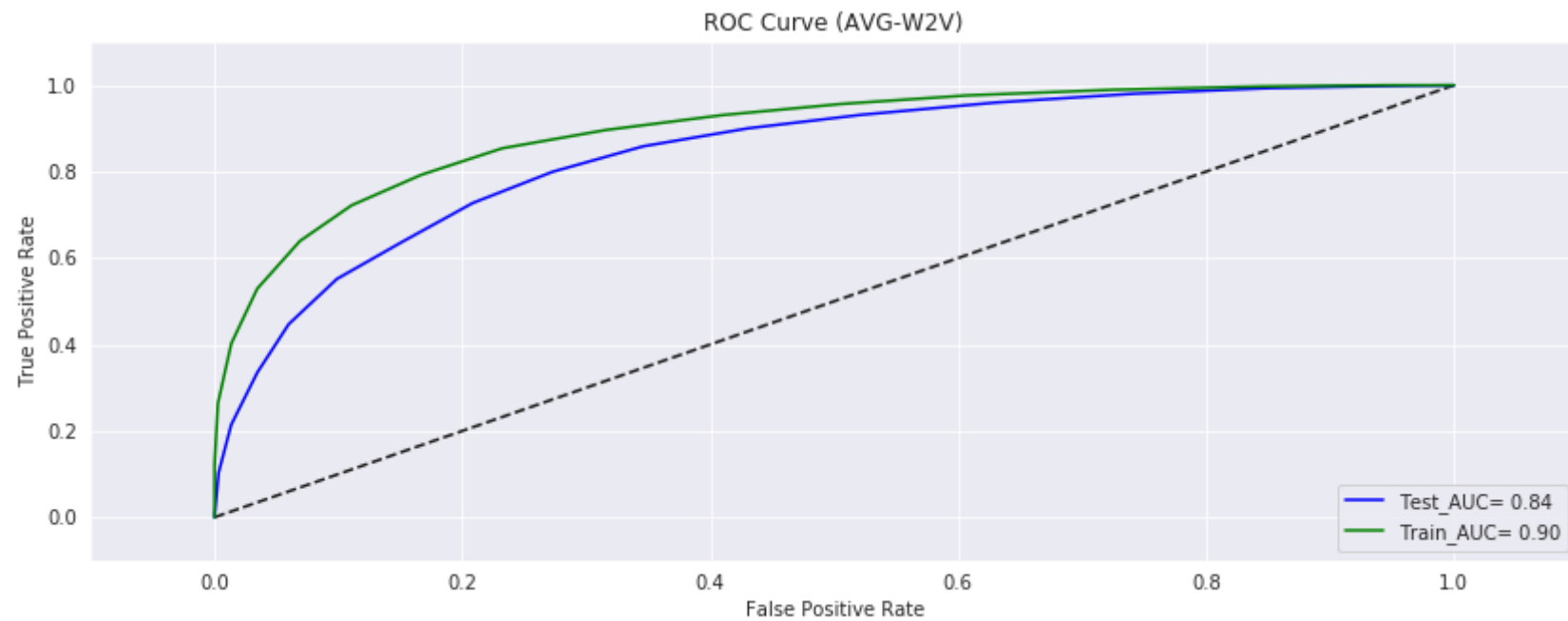


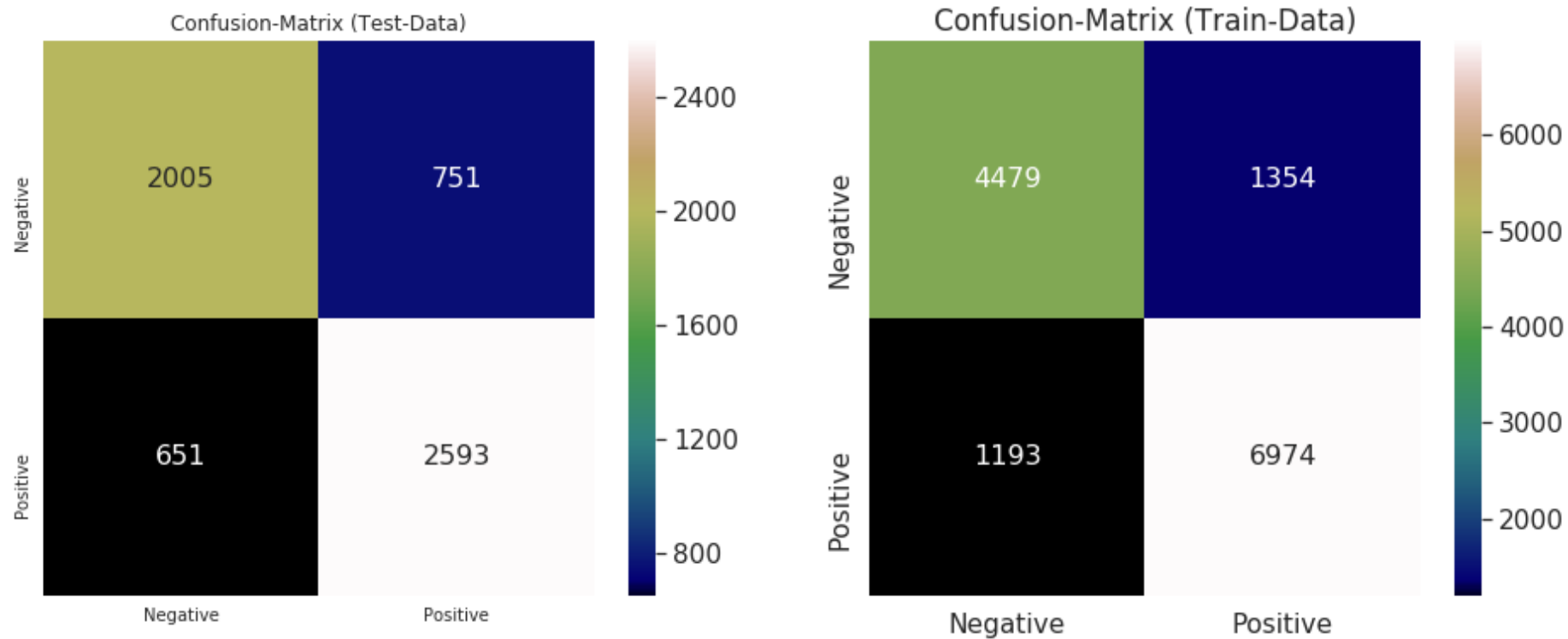
CPU times: user 1.37 s, sys: 1.73 s, total: 3.1 s

Wall time: 5min 20s

Optimal value of K: {'n_neighbors': 15}

```
In [12]: 1 test_performance(train,y_train_kd,test,y_test_kd,model.best_params_['n_neighbors'],algo[1],vect[2],summarize)
```





Apply KNN for TFIDF-W2V vectorizer

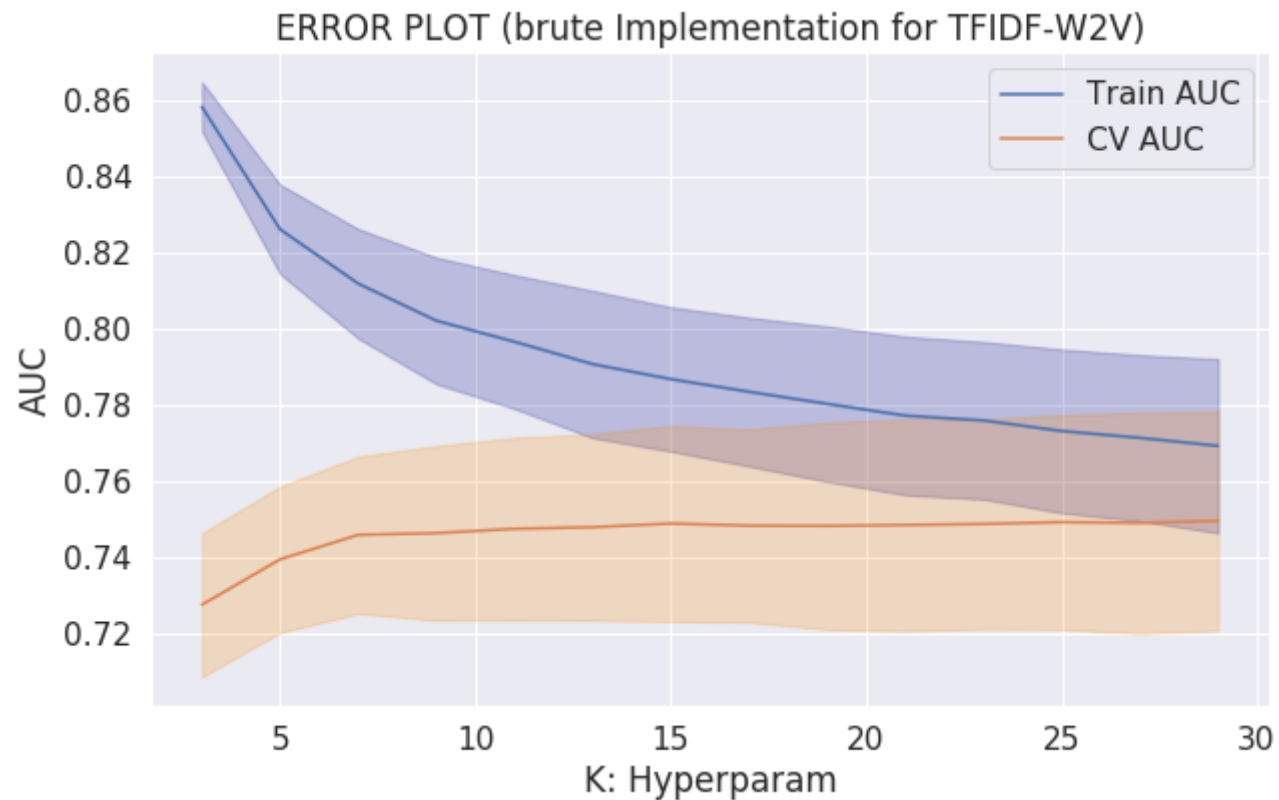
1.Brute force implementation

In [17]:

```

1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=tfidf_sent_vectors,test=tfidf_sent_vectors_test,mean=True)
3 #OVERSAMPLE TRAIN DATA
4 #train, y_train = oversampling(train, y_train)
5 #HYPERPARAM TUNNING
6 %time model=KNN_Classifier(train,y_train,TBS,k,searchMethod,algo[0],vect[3])
7 print('Optimal value of K: ',model.best_params_)
8 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
9 saveModeltofile(model,'model_tfw2v_knn')

```



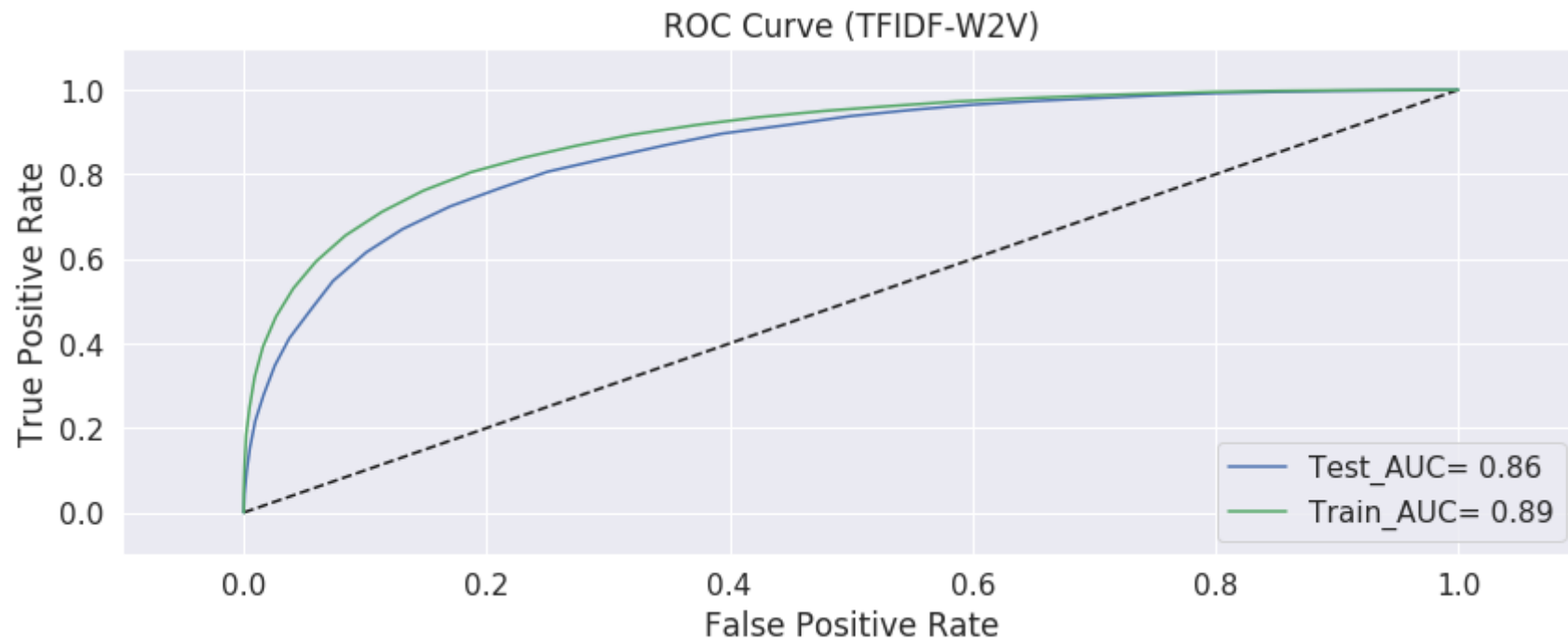
CPU times: user 1 s, sys: 1.12 s, total: 2.12 s

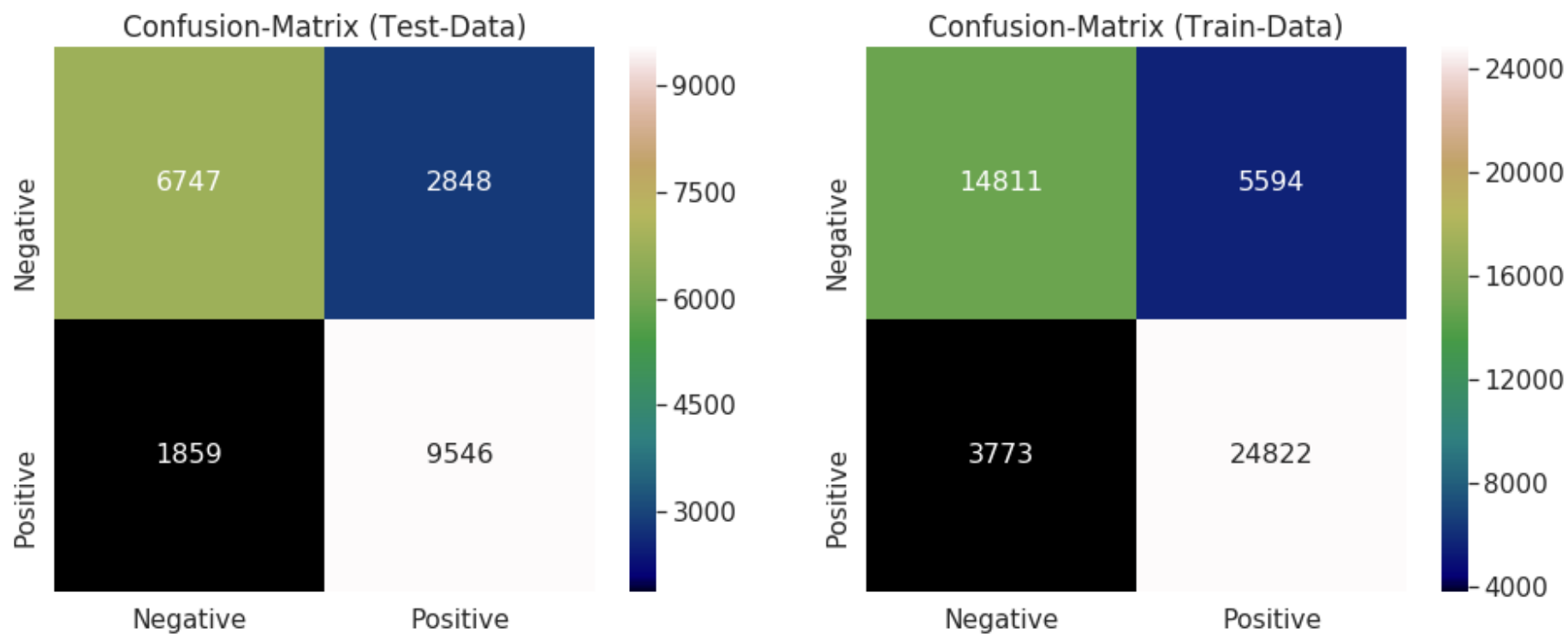
Wall time: 6min 39s

Optimal value of K: {'n_neighbors': 29}

[1.2] Test Performance:

```
In [18]: 1 test_performance(train,y_train,test,y_test,model.best_params_['n_neighbors'],algo[0],vect[3],summarize)
```





2. kd-tree implementation

In [13]:

```

1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=tfidf_sent_vectors_kd,test=tfidf_sent_vectors_test_kd,mean=True)
3 #HYPERPARAM TUNNING
4 %time model=KNN_Classifier(train,y_train_kd,TBS,k,searchMethod,algo[1],vect[3])
5 print('Optimal value of K: ',model.best_params_)
6 #SAVE CURRENT STATE OF ML-MODEL FOR FUTURE USE
7 saveModeltofile(model,'model_tfw2v_kd_knn')

```

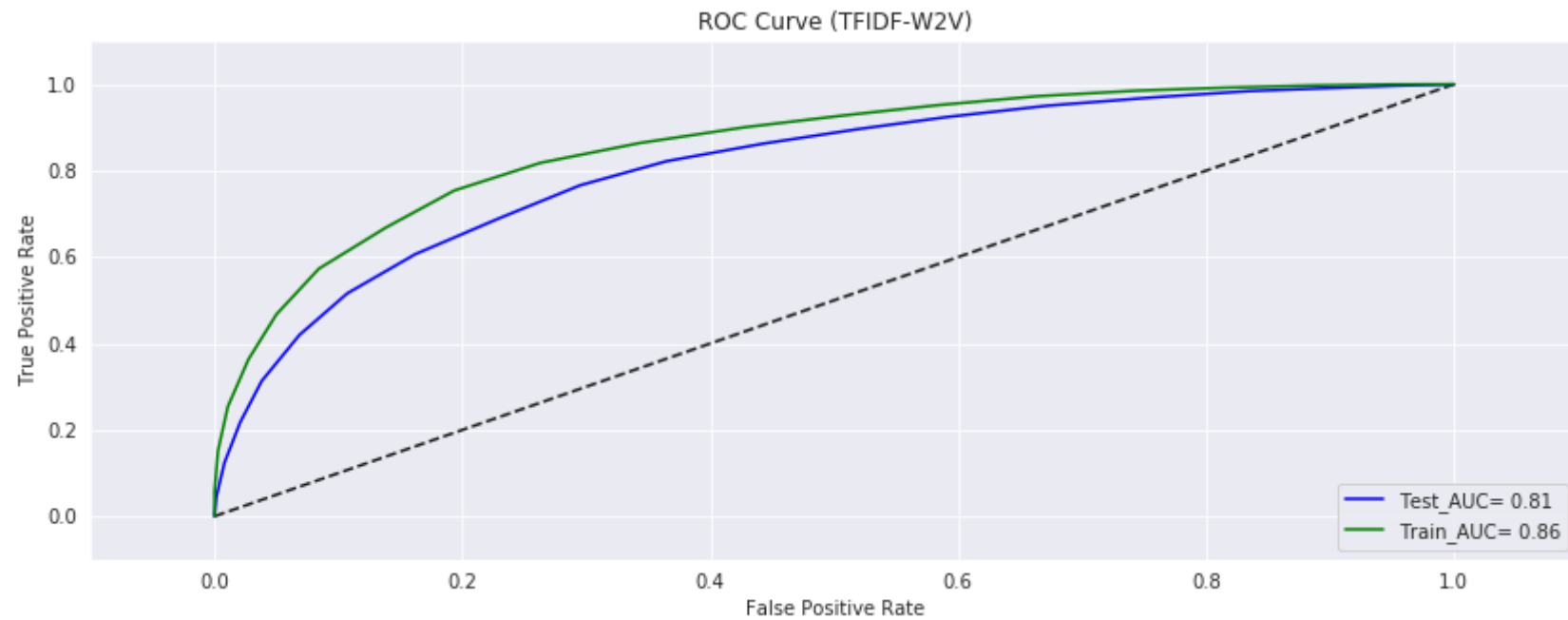


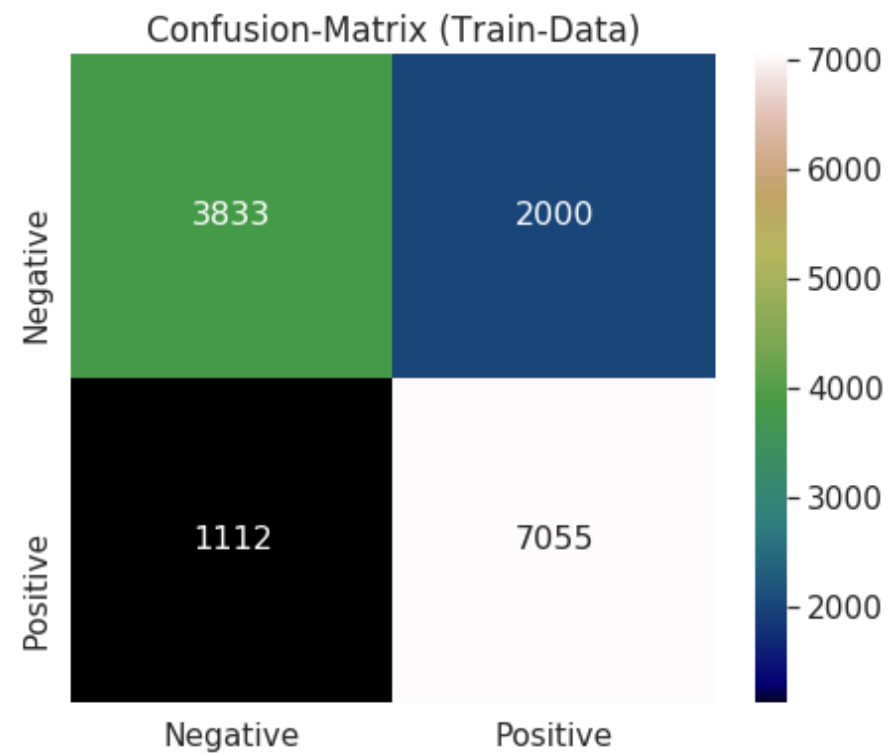
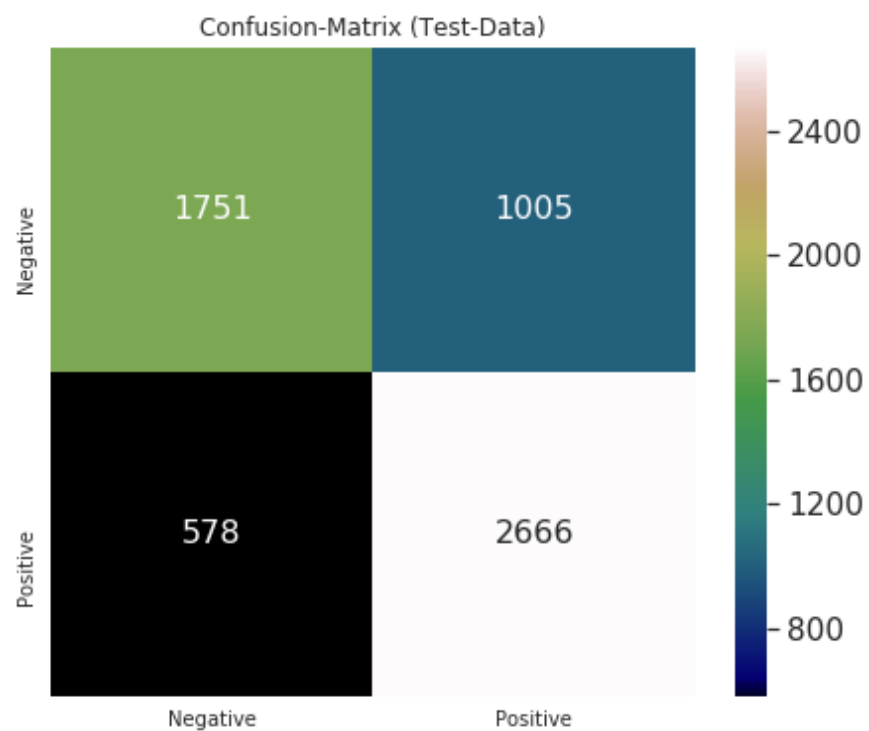
CPU times: user 1.3 s, sys: 1.69 s, total: 2.99 s

Wall time: 4min 33s

Optimal value of K: {'n_neighbors': 19}

```
In [14]: 1 test_performance(train,y_train_kd,test,y_test_kd,model.best_params_['n_neighbors'],algo[1],vect[3],summarize)
```





Conclusion:

In [40]: 1 `print(summarize)`

| Vectorizer | Algorithm | Optimal-K | Train(AUC) | Test(AUC) |
|------------|-----------|-----------|------------|-----------|
| BoW | brute | 3 | 89.48 | 60.44 |
| BoW | kd_tree | 3 | 91.44 | 70.75 |
| TF-IDF | brute | 3 | 99.01 | 51.64 |
| TF-IDF | kd_tree | 3 | 89.99 | 59.66 |
| AVG-W2V | brute | 23 | 91.93 | 89.35 |
| AVG-W2V | kd_tree | 15 | 89.62 | 84.20 |
| TFIDF-W2V | brute | 29 | 89.20 | 86.20 |
| TFIDF-W2V | kd_tree | 19 | 85.85 | 80.99 |

Got best performance with AVG-W2V:

a. AUC = 89.35

b. K = 23

1