

Import necessary libraries

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

```
In [2]: 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.linear_model import LogisticRegression
        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 scipy.sparse import *
       18 from prettytable import PrettyTable
       19 from wordcloud import WordCloud
```

Load preprocessed data

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  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  tf_idf_vect =openfromfile('tf_idf_vect')
16  X_train_tfidf =openfromfile('X_train_tfidf')
17  X_test_tfidf =openfromfile('X_test_tfidf')
18
19  avg_sent_vectors=openfromfile('avg_sent_vectors')
20  avg_sent_vectors_test=openfromfile('avg_sent_vectors_test')
21
22  tfidf_sent_vectors=openfromfile('tfidf_sent_vectors')
23  tfidf_sent_vectors_test=openfromfile('tfidf_sent_vectors_test')
24

```

Standardizing data

In [4]:

```

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

```

Logistic Regression

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


```

In [13]: 1 def LR_Classifier(x_train,y_train,TBS,params,penalty,searchMethod,vect):
2         ''' FUNCTION FOR FINDING OPTIMAL VALUE OF HYPERPARAM AND DRAW ERROR PLOT '''
3         #INITIALIZE LOGISTIC REGRESSION OBJECT
4         clf=LogisticRegression(penalty=penalty,class_weight='balanced',random_state=1)
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['C']),\
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
31         plt.plot(params['C'], 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['C'],train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
35         plt.plot(params['C'], cv_auc, label='CV AUC')
36         # Reference Link: https://stackoverflow.com/a/48803361/4084039
37         plt.gca().fill_between(params['C'],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('C(1/lambda): 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 [6]: 1 def test_performance(x_train,y_train,x_test,y_test,optimal_c,penalty,vect,summarize):
2         '''FUNCTION FOR TEST PERFORMANCE(PLOT ROC CURVE FOR BOTH TRAIN AND TEST) WITH OPTIMAL_K'''
3         #INITIALIZE LR OBJECT WITH OPTIMAL HYPERPARAM
4         clf=LogisticRegression(class_weight='balanced',penalty=penalty,C=optimal_c,n_jobs=-1)
5         clf.fit(x_train,y_train)
6         y_pred=clf.predict(x_test)
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,pos_label=1)
10        fpr_train, tpr_train, threshold_train = roc_curve(y_train, train_probability,pos_label=1)
11        auc_score_test=auc(fpr_test, tpr_test)
12        auc_score_train=auc(fpr_train, tpr_train)
13        f1=f1_score(y_test,y_pred,average='weighted')
14        #ADD RESULTS TO PRETTY TABLE
15        summarize.add_row([vect, penalty, optimal_c, '%.3f' %auc_score_test,'%f' %f1])
16
17        plt.figure(1,figsize=(14,5))
18        plt.subplot(121)
19        plt.title('ROC Curve (%s)' %vect)
20        #IDEAL ROC CURVE
21        plt.plot([0,1],[0,1],'k--')
22        #ROC CURVE OF TEST DATA
23        plt.plot(fpr_test, tpr_test , 'b', label='Test_AUC= %.2f' %auc_score_test)
24        #ROC CURVE OF TRAIN DATA
25        plt.plot(fpr_train, tpr_train , 'g', label='Train_AUC= %.2f' %auc_score_train)
26        plt.xlim([-0.1,1.1])
27        plt.ylim([-0.1,1.1])
28        plt.xlabel('False Positive Rate')
29        plt.ylabel('True Positive Rate')
30        plt.grid(True)
31        plt.legend(loc='lower right')
32        #PLOT CONFUSION MATRIX USING HEATMAP
33        plt.subplot(122)
34        plt.title('Confusion-Matrix(Test Data)')
35        df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), ['Negative','Positive'],['Negative','Positive'])
36        sns.set(font_scale=1.4)#for label size
37        sns.heatmap(df_cm,cmap='gist_earth', annot=True,annot_kws={"size": 16}, fmt='g')
38        plt.show()
39        return clf

```

Function for sparisty check:

```
In [7]: 1 def check_sparsity(x_train,x_test,y_train,y_test,c_vals):
2         print('change in sparsity with increase in lambda(1/C) :\n')
3         for c in c_vals:
4             lr=LogisticRegression(penalty='l1',C=c,class_weight='balanced')
5             lr.fit(x_train,y_train)
6             pred=lr.predict(x_test)
7             print('features having non-zero weights {0} when c={1}\n'.format(np.count_nonzero(lr.coef_), c))
```

Function for perturbation test to check multicollinearity

```

In [29]: 1 def perturbation_test(x_train,x_test,y_train,y_test,optimal_c):
2
3     fetures_wt_change=[]
4     #MODEL BEFORE NOISE
5     clf_before = LogisticRegression(penalty='l2',C=optimal_c,class_weight='balanced',random_state=1)
6     clf_before.fit(x_train, y_train)
7
8     # ADD A SMALL NOISE TO DATA
9     x_train.data += .001
10
11     #MODEL AFTER NOISE
12     clf_after= LogisticRegression(penalty='l2',C=optimal_c,class_weight='balanced',random_state=1)
13     clf_after.fit(x_train, y_train)
14
15     w_before=find(clf_before.coef_[0])[2]
16     w_after=find(clf_after.coef_[0])[2]
17
18     #ADD ERROR TO GET RID OFF DIVISION BY ZERO ERROR
19     error=.000001*np.random.normal(.000005,.001)#np.random.normal(.000005,.000001,1)
20     w_before += error
21     w_after += error
22     print('size of w_before_noise:',w_before.size)
23     print('size of w_after_noise:',w_after.size)
24
25     #PERCENTAGE CHANGE IN WEIGHT CORRESPONDS TO EACH FEATURE
26     percentage_change = np.array((abs(w_before-w_after)/w_before ) * 100)
27     print('shape of percentage change:',percentage_change.shape)
28     if len(w_before)==len(w_after):
29         thresholds=[10,20,30,50,70,90,100]
30         for threshold in thresholds:
31             fetures_wt_change.append(percentage_change[np.where(percentage_change > threshold)].size)
32     plt.figure(1,figsize=(12,6))
33     plt.plot(thresholds,fetures_wt_change)
34     for xy in zip(thresholds,fetures_wt_change):
35         plt.annotate('%s, %s)' % xy, xy=xy, textcoords='data')
36     plt.xlabel('Threshold values')
37     plt.ylabel('no. of collinear features')
38     plt.title('Elbow Method(Multicollinearity Check)')
39     plt.grid(True)
40     plt.show()
41     return percentage_change

```


Function which plots collinear features using wordcloud :

```
In [9]: 1 def print_multicollinear_features(percentage_change, threshold, vect):
2         feat=vect.get_feature_names()
3         collinear_features=[]
4         collinear_features_no=percentage_change[np.where(percentage_change > threshold)].size
5         index=np.where(percentage_change > threshold)[0]
6         for i in index:
7             collinear_features.append(feat[i])
8         print(collinear_features_no)
9         if len(collinear_features)!=0:
10            #wordcloud plot
11            wordcloud = WordCloud(max_font_size=50, max_words=100,collocations=False).\
12            generate(str(collinear_features))
13            plt.figure(1,figsize=(14,13))
14            plt.title("WordCloud(Collinear-Features)")
15            plt.imshow(wordcloud, interpolation="bilinear")
16            plt.axis("off")
17            plt.show()
18
```

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

```

In [10]: 1  #REFERENCE STACKOVERFLOW
2  def feature_importance(vectorizer,clf,n):
3      feature_names = vectorizer.get_feature_names()
4      coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
5      top = zip(coefs_with_fns[:n], coefs_with_fns[:-(n + 1):-1])
6      print("\tNegative\t\t\t\t\tPositive\t\t")
7      print("_"*75)
8      for (coef_1, fn_1), (coef_2, fn_2) in top:
9          print("\t%.4f\t%-15s\t\t\t\t\t%.4f\t%-15s" % (coef_1, fn_1, coef_2, fn_2))
10
11     coef=sorted(clf.coef_[0],reverse=True)
12     #STORE WEIGHT CORRESPONDING TO TOP POSITIVE AND NEGATIVE IMPORTANT FEATURES
13     coef_p=coef[:n]
14     coef_n=coef[:-(n + 1):-1]
15     coef_np=coef_n+coef_p
16     indices_n=np.argsort(clf.coef_[0])[:n]
17     indices_p=np.argsort(clf.coef_[0])[::-1][:n]
18     indices=list(indices_n)+list(indices_p)
19     names = np.array(vectorizer.get_feature_names())
20     #bar chart
21     plt.figure(2,figsize=(13,6))
22     sns.set(rc={'figure.figsize':(11.7,8.27)})
23     # Create plot title
24     plt.title("Feature Importance(top %d positive and negative class features)" % n)
25     # Add bars
26     plt.bar(range(2*n), coef_np)
27     # Add feature names as x-axis labels
28     plt.xticks(range(2*n), names[indices], rotation=80)
29     plt.show()

```

Initialization of common objects required for all vectorization:

In [11]:

```
1  #VECTORIZER
2  vect=['BoW', 'TF-IDF', 'AVG-W2V', 'TFIDF-W2V']
3  #OBJECT FOR TIMESERIES CROSS VALIDATION
4  TBS=TimeSeriesSplit(n_splits=10)
5  #METHOD USE FOR HYPER PARAMETER TUNNING
6  searchMethod='random'
7  #RANGE OF K VALUES(HYPERPARAM)
8  c_ranges=[1000,500,100,50,10,5,1,.5,.1,.05,.01,.005,.001,.0005,.0001]
9  params={'C':c_ranges}
10 #REGULARIZER USED
11 penalty=['l1','l2']
12 #INITIALIZE PRETTY TABLE OBJECT
13 summarize = PrettyTable()
14 summarize.field_names = ['Vectorizer', 'Regularizer', 'Optimal-C(1/lambda)', 'Test(AUC)', 'Test(f1-score)']
```

[1.1] Logistic Regression on BOW, SET 1

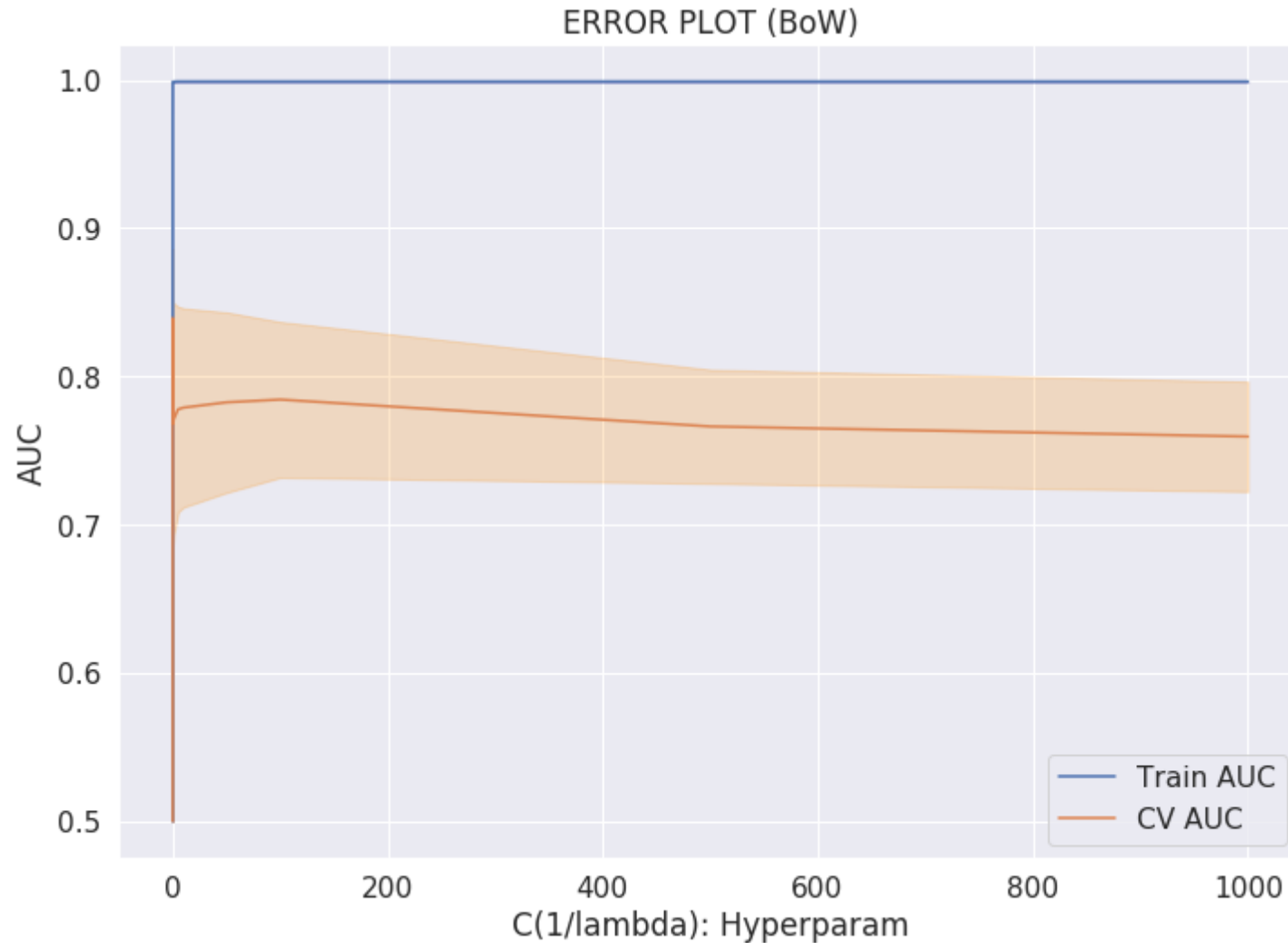
[1.1.1] Applying Logistic Regression with L1 regularization on BOW, SET 1

In [46]:

```

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=LR_Classifier(train,y_train,TBS,params,penalty[0],searchMethod,vect[0])
5 #PRINT OPTIMAL VALUE OF HYPERPARAM
6 print('Optimal value of C(1/lambda): ',model.best_params_)

```



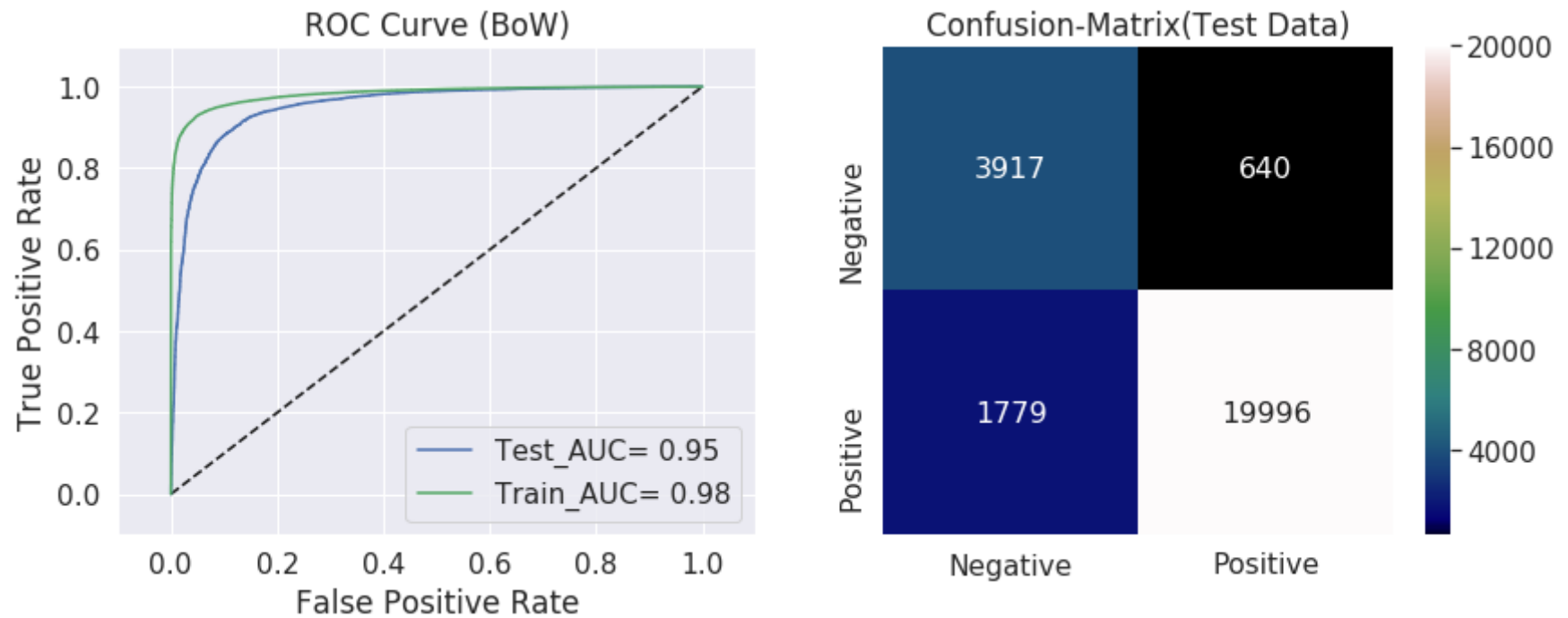
CPU times: user 14.9 s, sys: 428 ms, total: 15.3 s

Wall time: 17.3 s

Optimal value of $C(1/\lambda)$: {'C': 0.005}

[1.1.1.1] Performance on test data with optimal value of hyperparam

```
In [47]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['C'],penalty[0],vect[0],summarize)
```



[1.1.1.2] Calculating sparsity on weight vector obtained using L1 regularization on BOW, SET 1

```
In [16]: 1 check_sparsity(train,test,y_train,y_test,params['C'])
```

change in sparsity with increase in $\lambda(1/C)$:

features having non-zero weights 58396 when $c=1000$

features having non-zero weights 63177 when $c=500$

features having non-zero weights 24507 when $c=100$

features having non-zero weights 27508 when $c=50$

features having non-zero weights 19136 when $c=10$

features having non-zero weights 16760 when $c=5$

features having non-zero weights 14097 when $c=1$

features having non-zero weights 13286 when $c=0.5$

features having non-zero weights 11789 when $c=0.1$

features having non-zero weights 10748 when $c=0.05$

features having non-zero weights 5676 when $c=0.01$

features having non-zero weights 2907 when $c=0.005$

features having non-zero weights 172 when $c=0.001$

features having non-zero weights 54 when $c=0.0005$

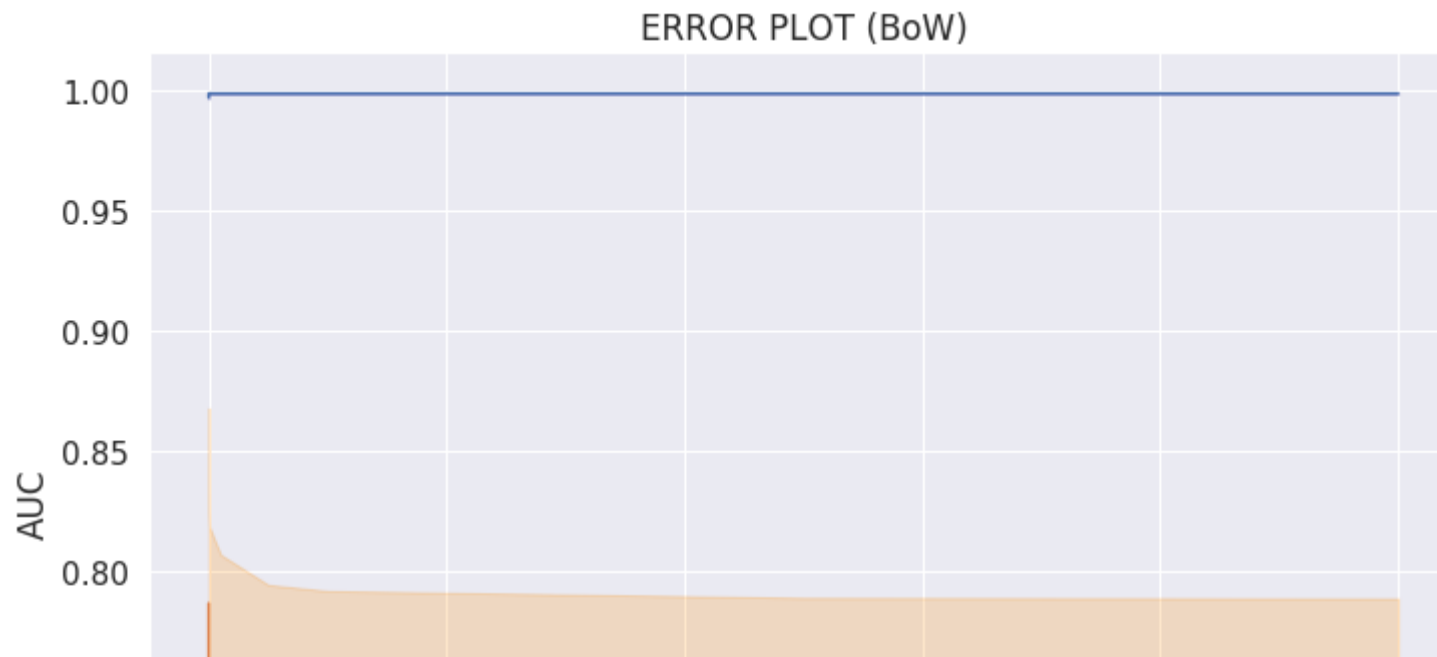
features having non-zero weights 0 when $c=0.0001$

Observation:

1. From the above analysis we can observe that as the value of $\lambda(1/C)$ increases our weight vector becomes more sparse.
2. Weights corresponds to unimportant features or less important fetures becomes zero when we are using L1 regularizer.

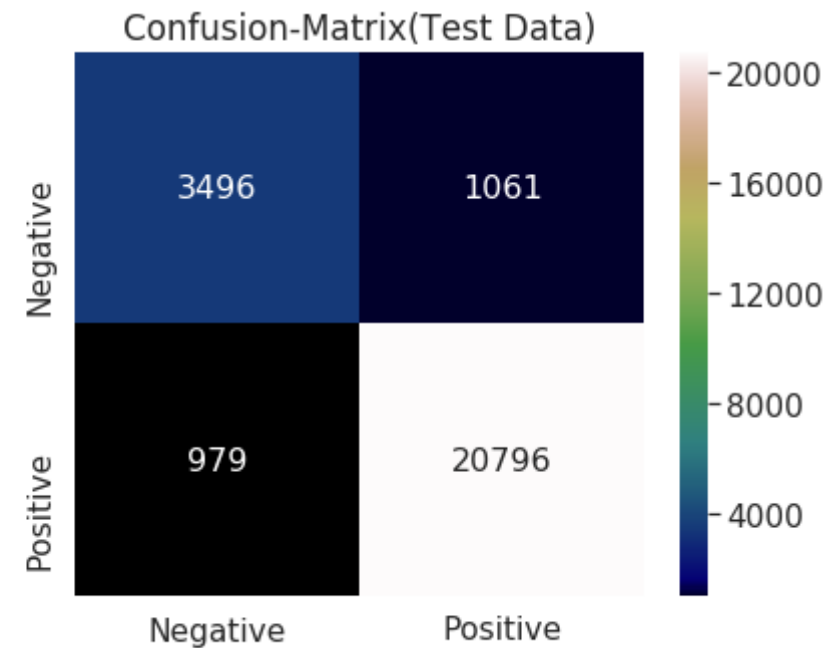
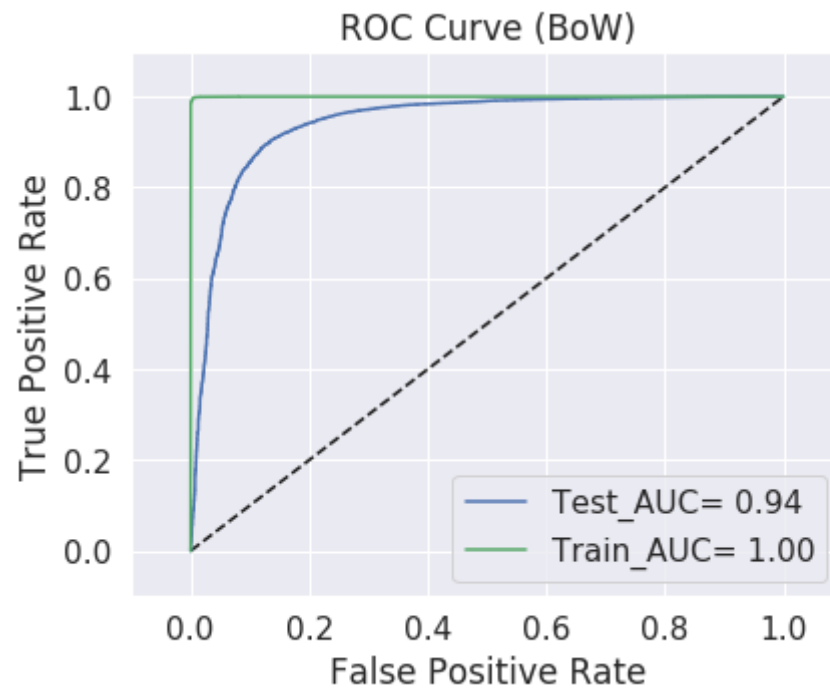
[1.1.2] Applying Logistic Regression with L2 regularization on BOW, SET 1

```
In [45]: 1 train, test=std_data(train=X_train_bigram,test=X_test_bigram,mean=False)
2 #HYPERPARAM TUNNING
3 %time model=LR_Classifier(train,y_train,TBS,params,penalty[1],searchMethod,vect[0])
4 print('Optimal value of C(1/lambda): ',model.best_params_)
```



[1.1.2.1] Performance on test data with optimal value of hyperparam

```
In [18]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['C'],penalty[1],vect[0],summarize)
```



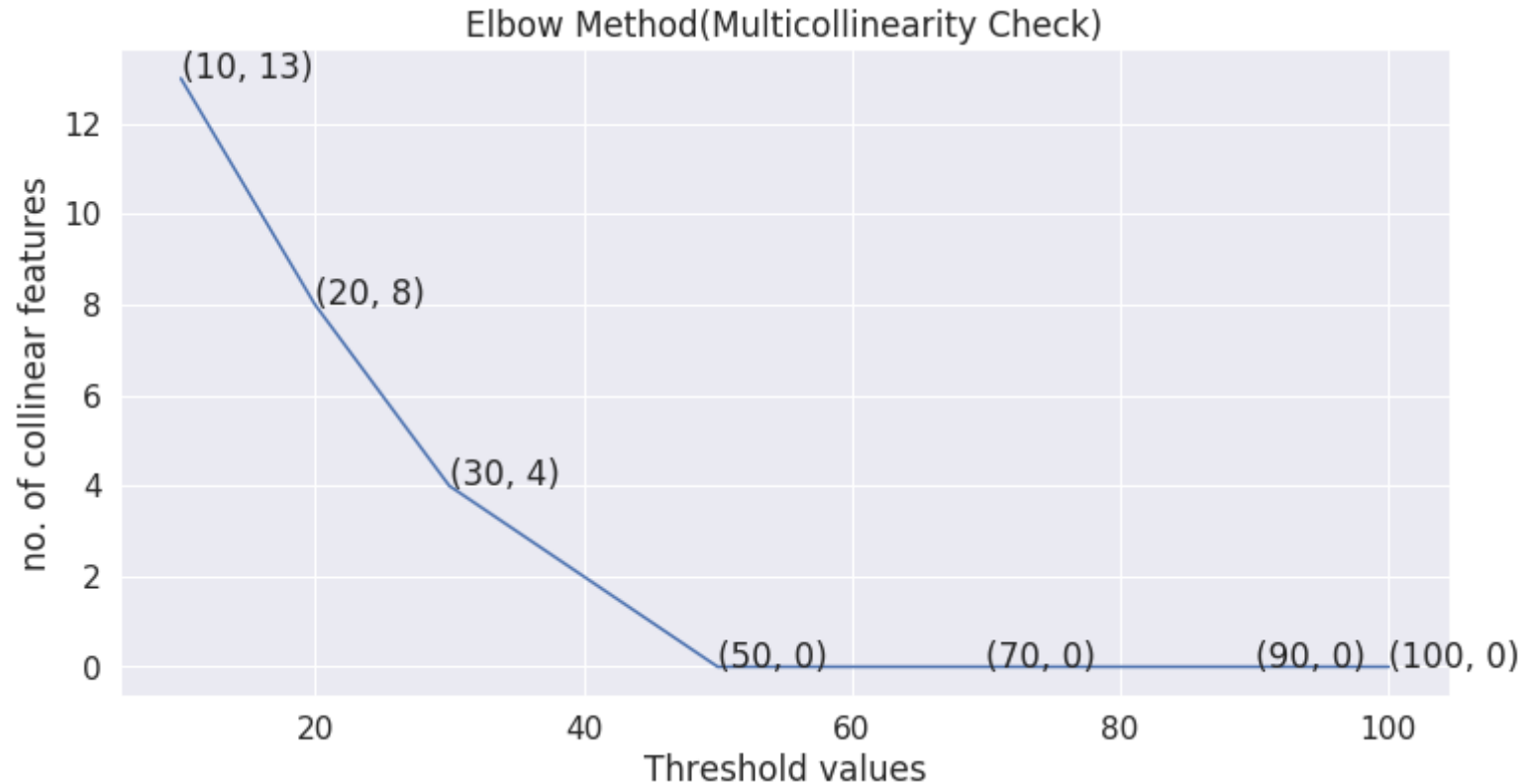
[1.1.2.2] Performing pertubation test (multicollinearity check) on BOW, SET 1


```
In [19]: 1 #train, test=std_data(train=X_train_bigram,test=X_test_bigram,mean=False)
2 percentage_change=pertubation_test(train,test,y_train,y_test,model.best_params_['C'])
```

size of w_before_noise: 83188

size of w_after_noise: 83188

shape of percentage change: (83188,)



Observation:

1. From the above plot we find that curve changes significantly at threshold value of 50(i.e. our knee point).
2. Number of fetures which have percentage weight change more than the threshold=50 are zero, so our features are independent.

[1.1.2.3] Print collinear features in word cloud

In [20]:

```
1 #SELECTED THRESHOLD BY USING ELBOW METHOD
2 threshold=50
3 print_multicollinear_features(percentage_change,threshold,count_vect)
```

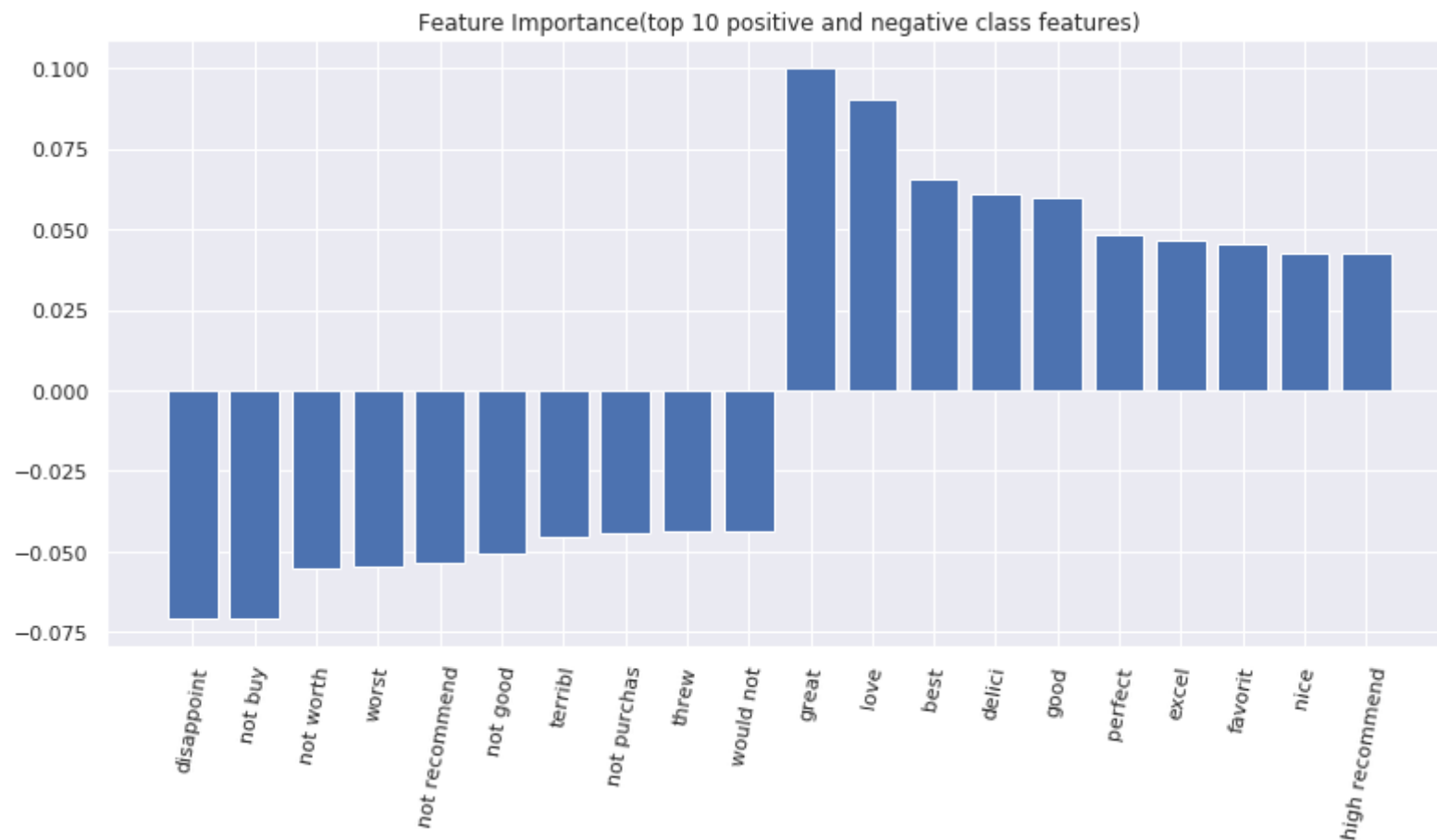
0

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

In [21]:

```
1 no_of_imp_features=10
2 feature_importance(count_vect,clf,no_of_imp_features)
```

Negative		Positive	
-0.0707	disappoint	0.1000	great
-0.0707	not buy	0.0903	love
-0.0555	not worth	0.0655	best
-0.0547	worst	0.0611	delici
-0.0534	not recommend	0.0598	good
-0.0509	not good	0.0485	perfect
-0.0453	terribl	0.0465	excel
-0.0441	not purchas	0.0453	favorit
-0.0441	threw	0.0427	nice
-0.0438	would not	0.0426	high recommend

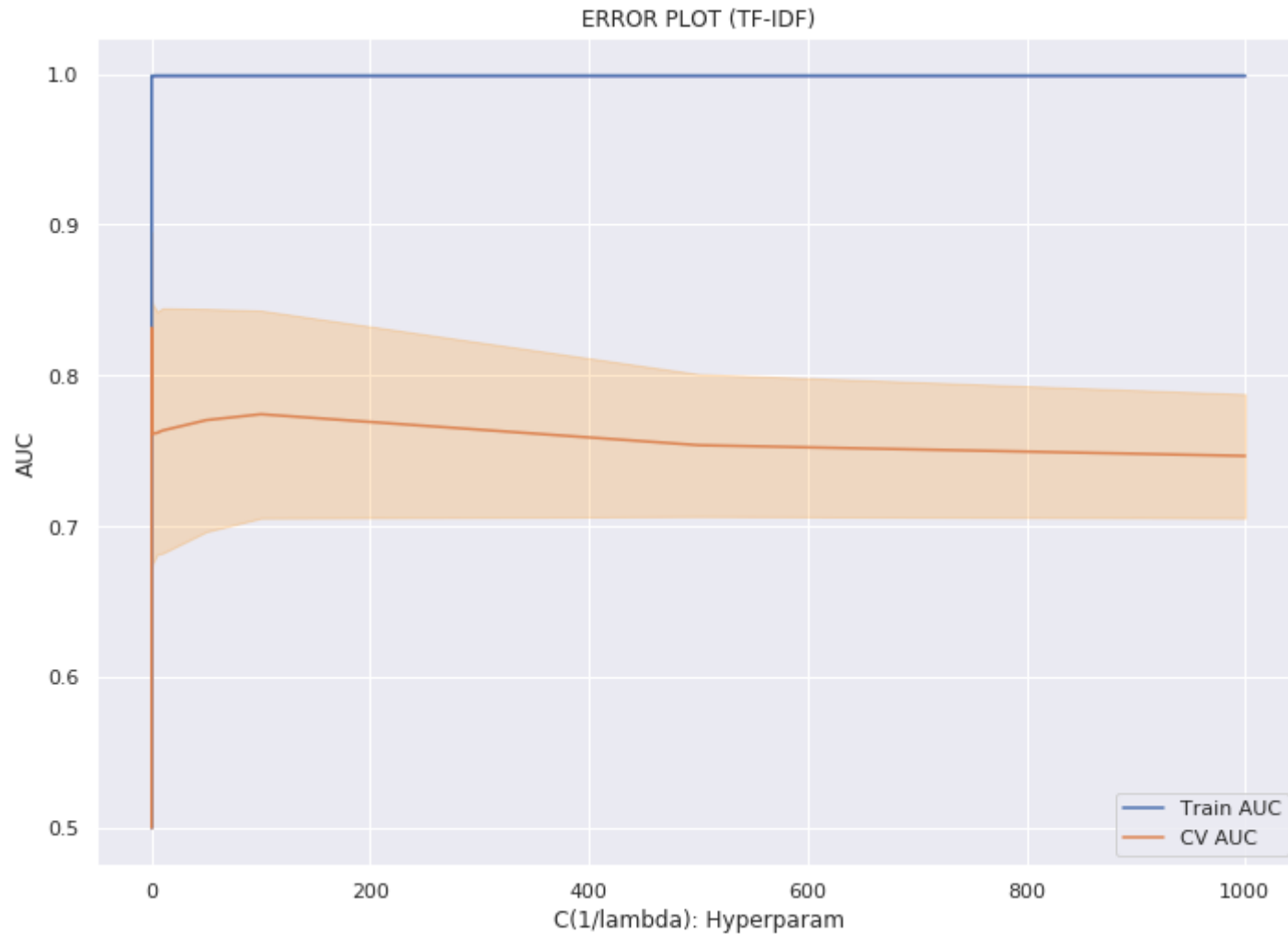


2.1 Logistic Regression on TFIDF, SET 2

[2.1.1] Applying Logistic Regression with L1 regularization on TFIDF, SET 2

In [23]:

```
1 searchMethod='random'
2 #STANDARDIZE TRAIN AND TEST DATA
3 train, test=std_data(train=X_train_tfidf,test=X_test_tfidf,mean=False)
4 #HYPERPARAM TUNNING
5 %time model=LR_Classifier(train,y_train,TBS,params,penalty[0],searchMethod,vect[1])
6 print('Optimal value of C(1/lambda): ',model.best_params_)
```



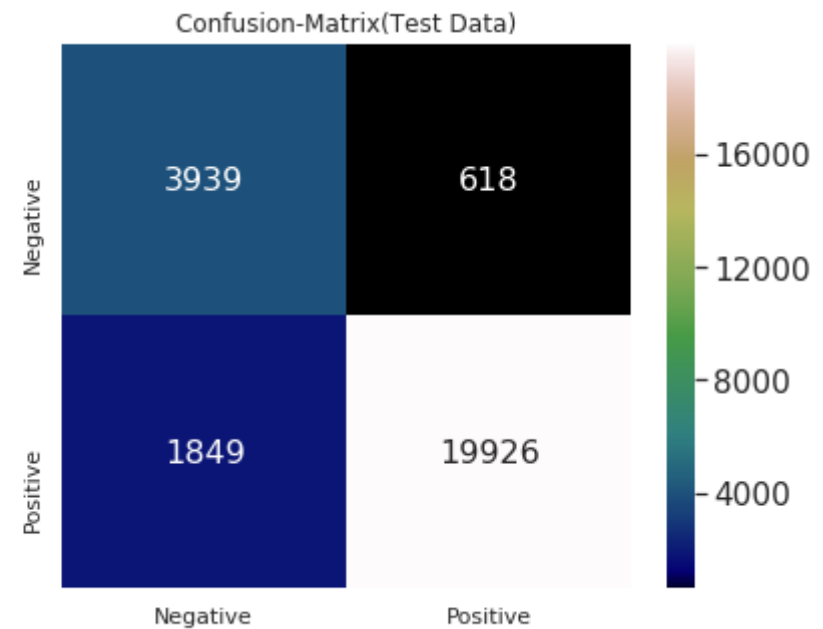
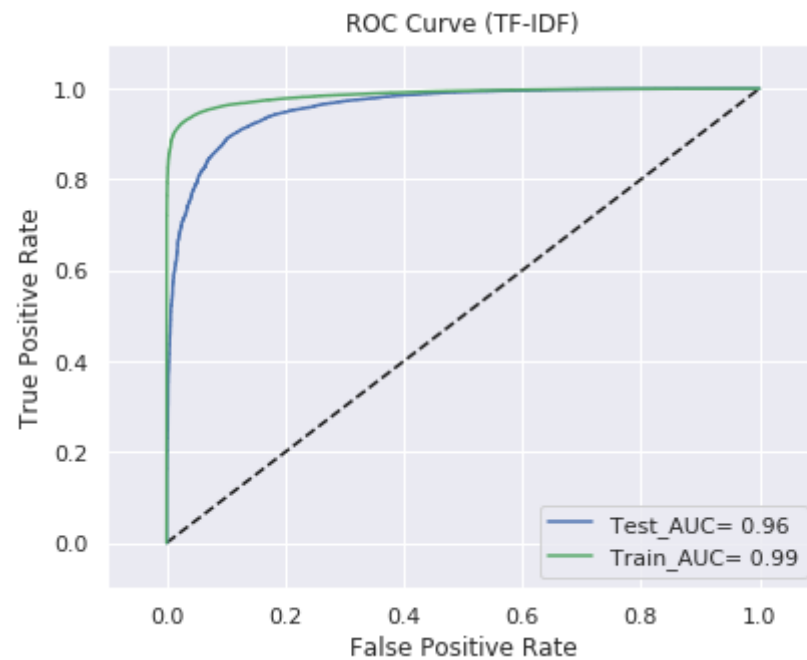
CPU times: user 15 s, sys: 348 ms, total: 15.3 s

Wall time: 17 s

Optimal value of $C(1/\lambda)$: {'C': 0.005}

[2.1.1.1] Performance on test data with optimal value of hyperparam

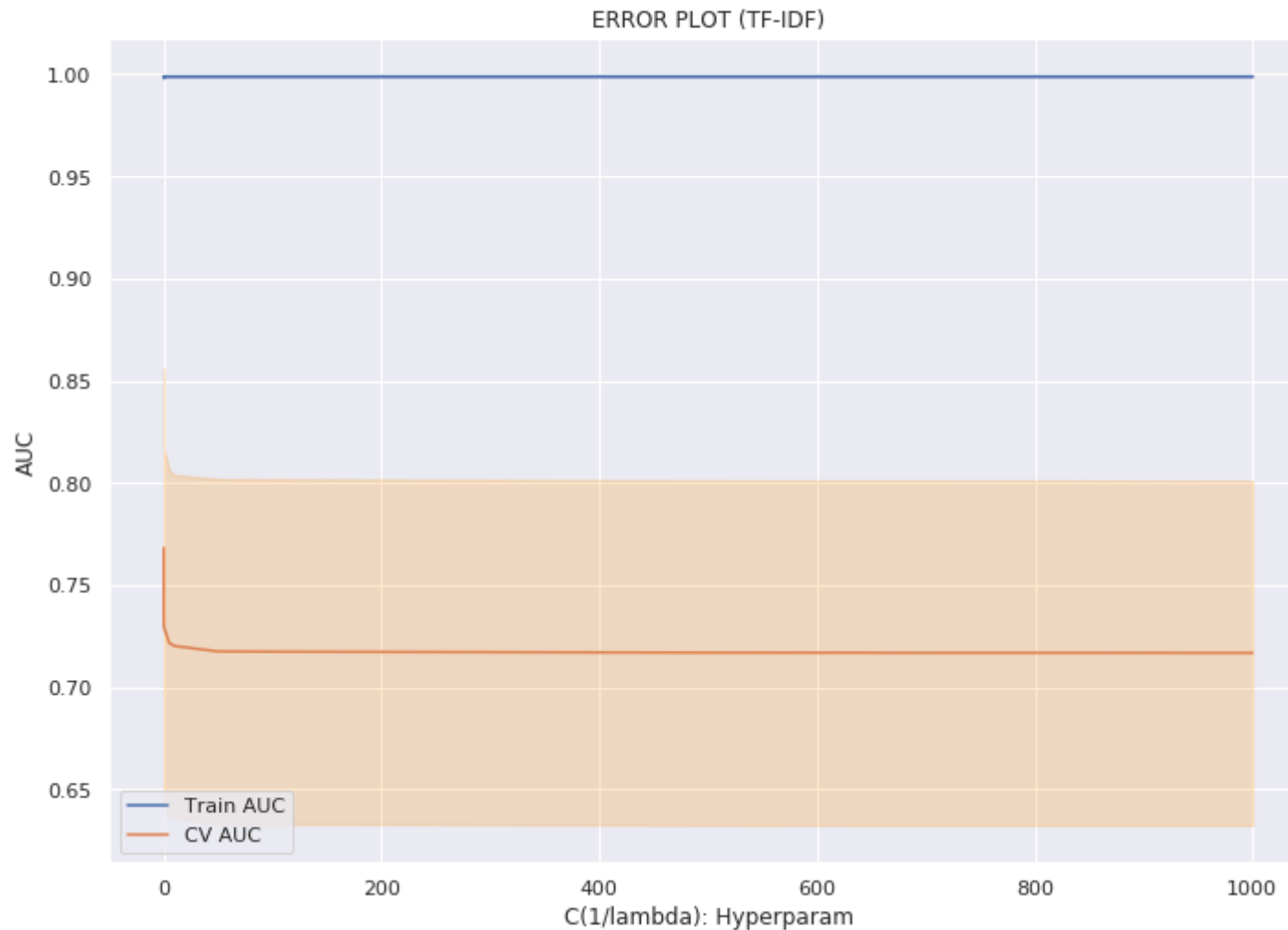
```
In [24]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['C'],penalty[0],vect[1],summarize)
```



[2.1.2] Applying Logistic Regression with L2 regularization on TFIDF, SET 2

In [26]:

```
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=LR_Classifier(train,y_train,TBS,params,penalty[1],searchMethod,vect[1])
5 print('Optimal value of C(1/lambda): ',model.best_params_)
```



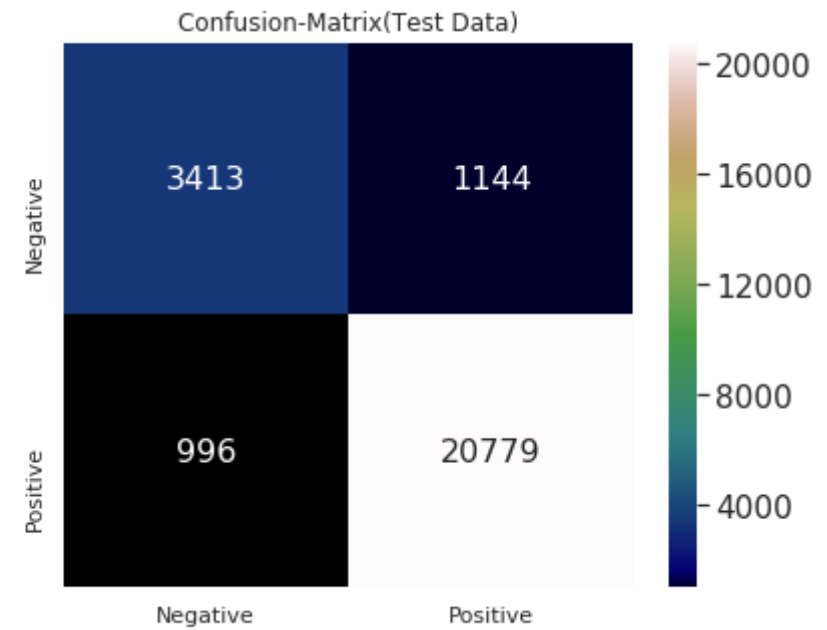
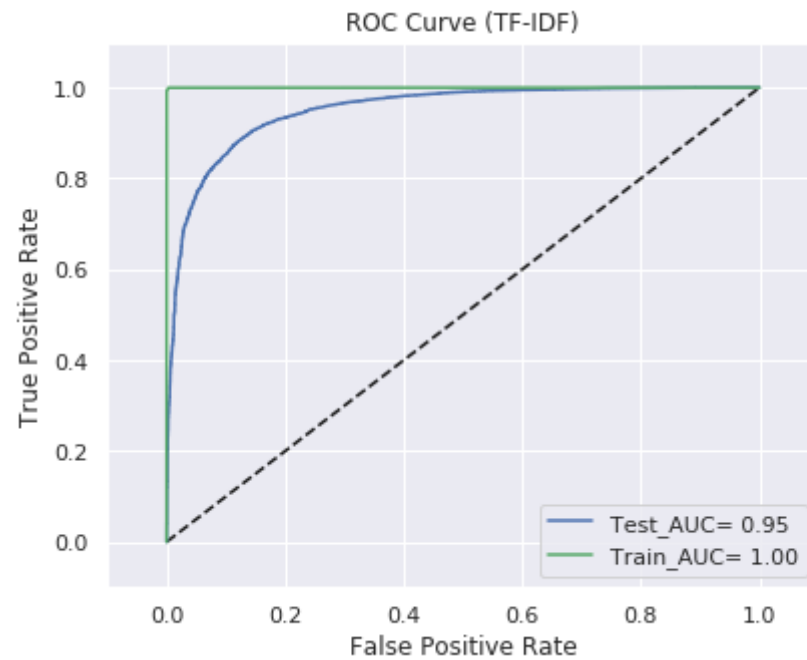
CPU times: user 31.8 s, sys: 608 ms, total: 32.4 s

Wall time: 4min 15s

Optimal value of C(1/lambda): {'C': 0.0001}

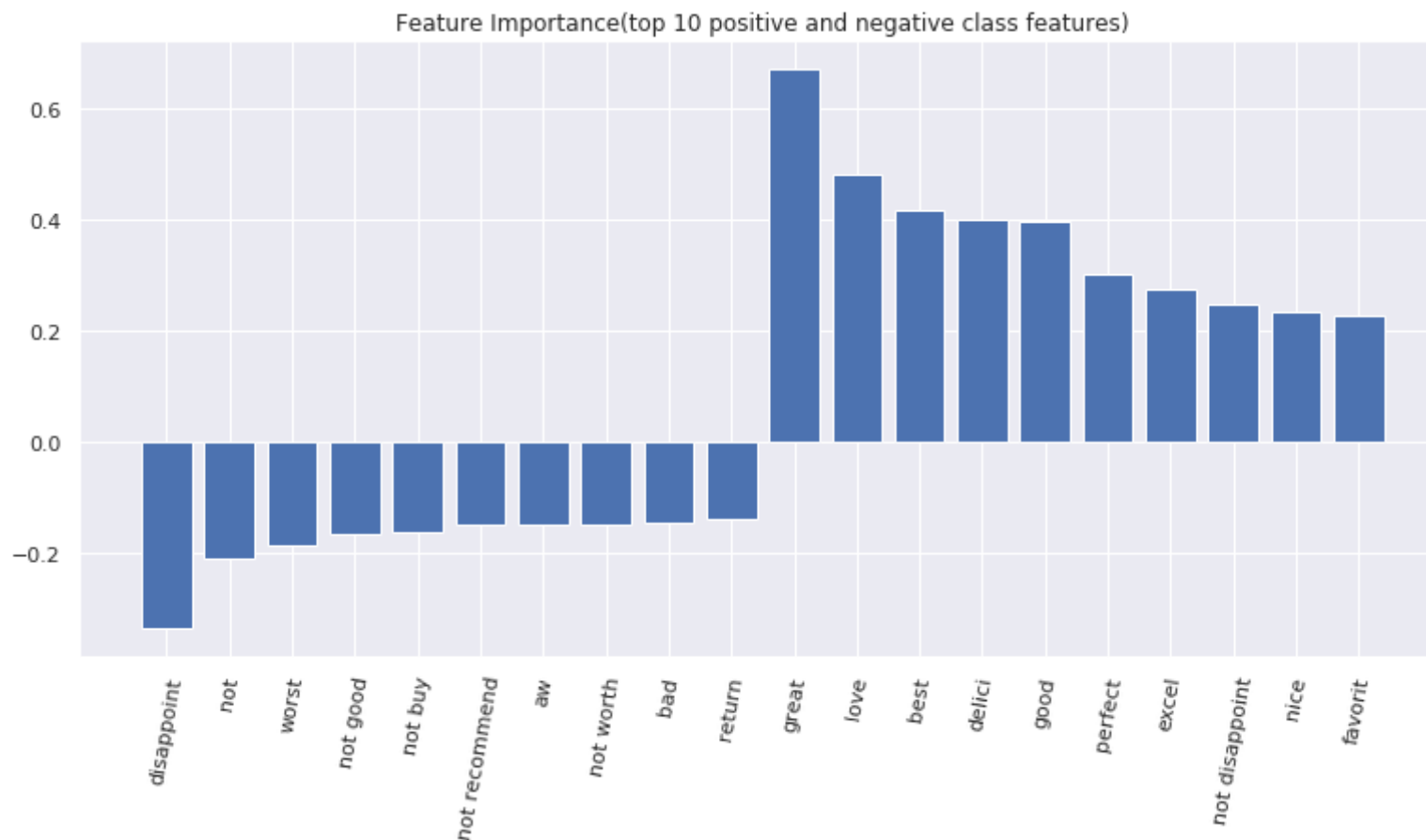
[2.1.2.1] Performance on test data with optimal value of hyperparam

```
In [27]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['C'],penalty[1],vect[1],summarize)
```

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

In [25]: 1 feature_importance(tf_idf_vect,clf,10)

Negative	Positive
-0.3347 disappoint	0.6698 great
-0.2117 not	0.4816 love
-0.1866 worst	0.4169 best
-0.1677 not good	0.4000 delici
-0.1633 not buy	0.3954 good
-0.1508 not recommend	0.3003 perfect
-0.1504 aw	0.2732 excel
-0.1496 not worth	0.2475 not disappoint
-0.1457 bad	0.2348 nice
-0.1388 return	0.2258 favorit



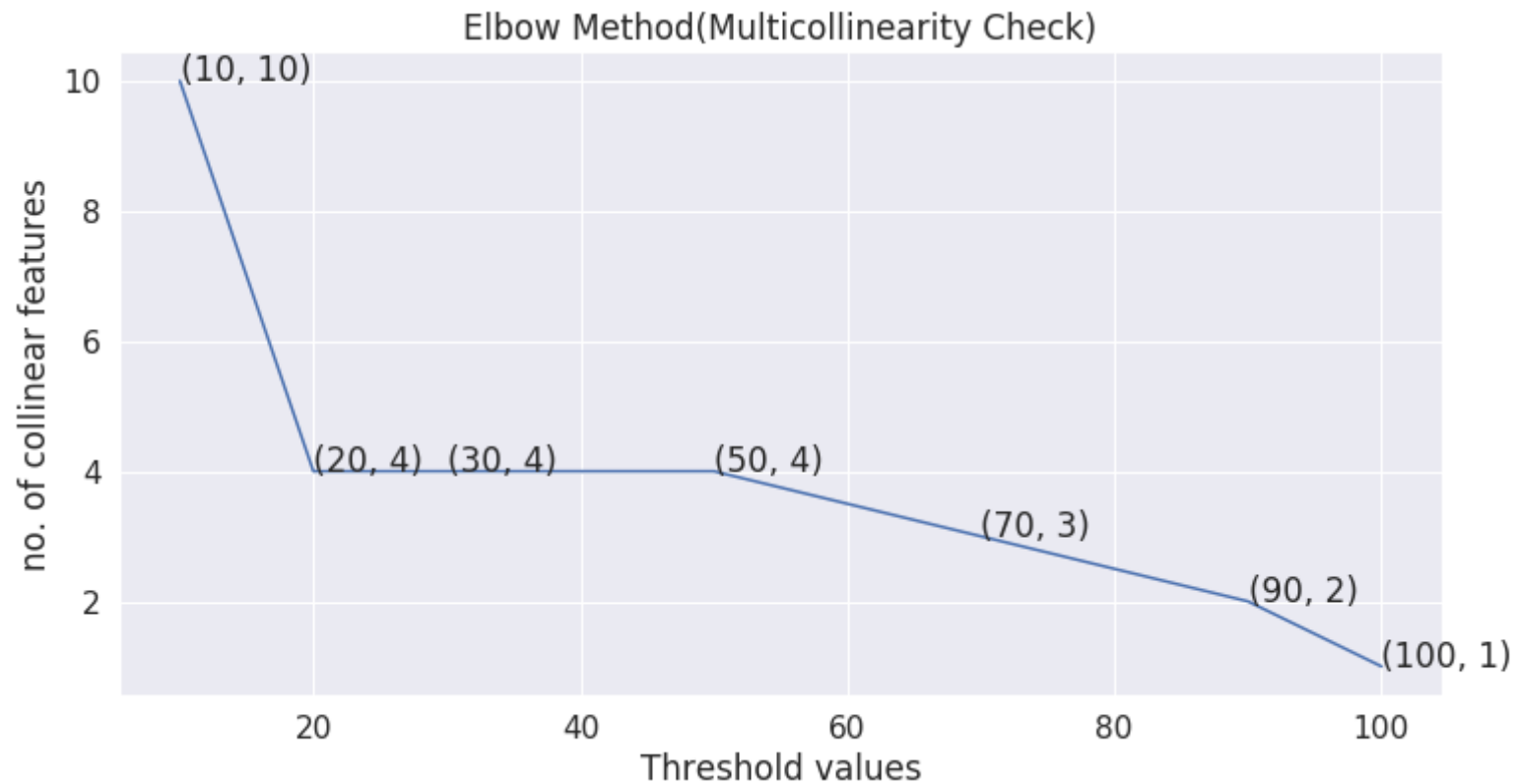
[2.1.4] Performing pertubation test (multicollinearity check) on TFIDF, SET 2

```
In [31]: 1 #train, test=std_data(train=X_train_bigram,test=X_test_bigram,mean=False)
        2 percentage_change=pertubation_test(train,test,y_train,y_test,model.best_params_['C'])
```

size of w_before_noise: 83188

size of w_after_noise: 83188

shape of percentage change: (83188,)

**Observation:**

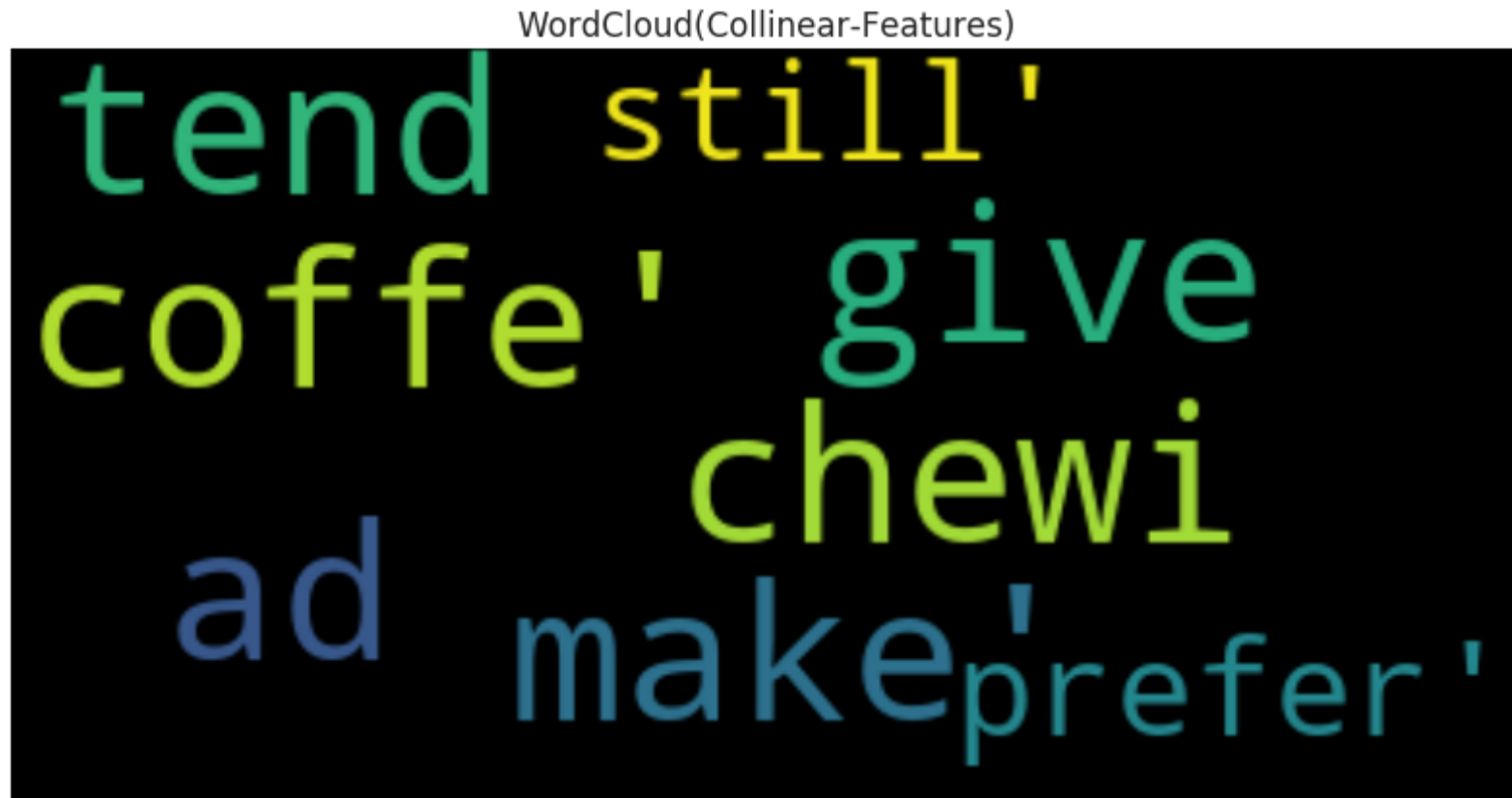
1. From the above plot we find that curve changes significantly at threshold value of 20(i.e. our knee point).
2. Number of fetures which have percentage weight change more than the threshold=20 are four, so our features are independent.

[2.1.4.1]Print collinear features in word cloud

In [32]:

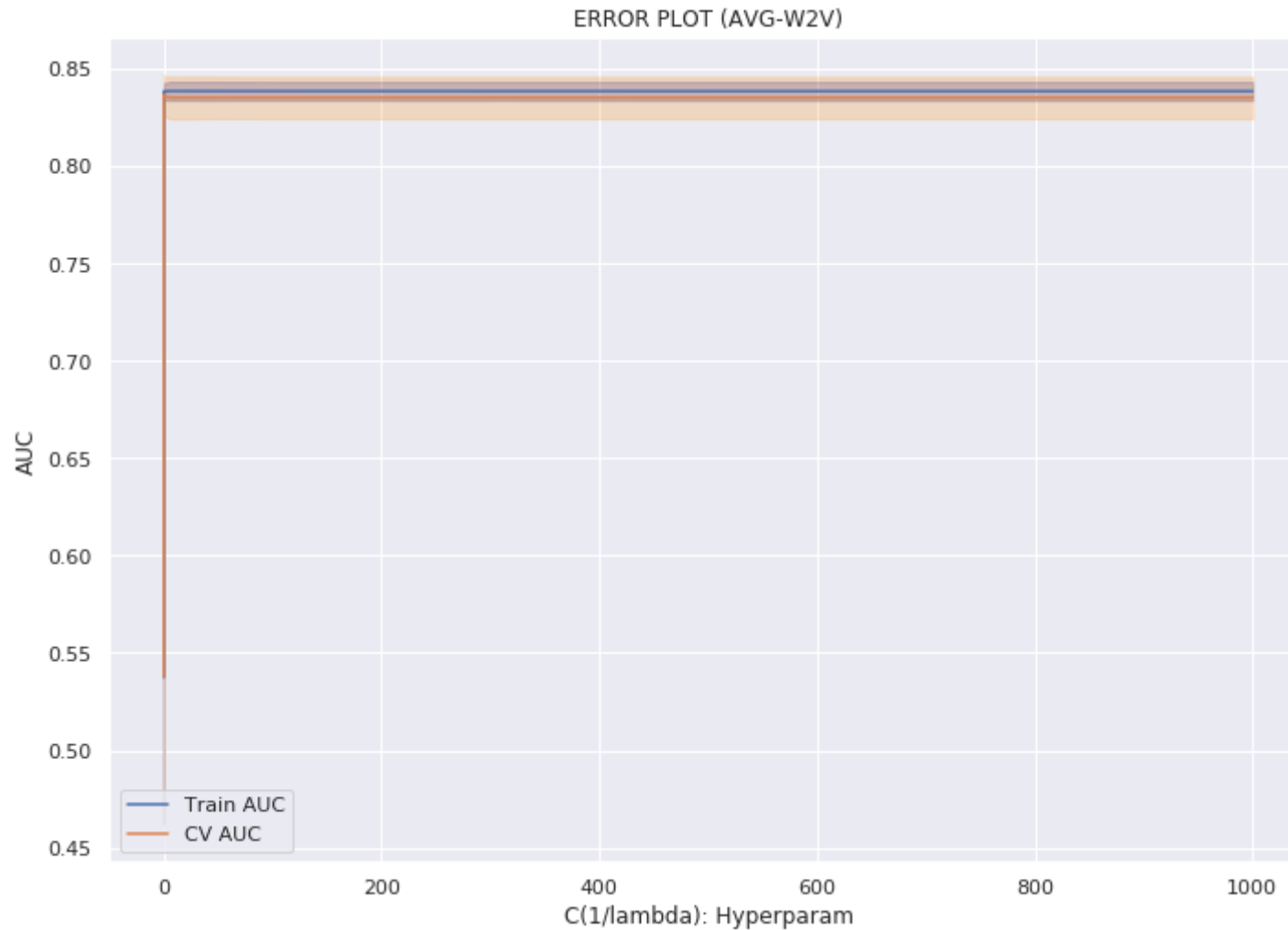
```
1 #SELECTED THRESHOLD BY USING ELBOW METHOD  
2 threshold=20  
3 print_multicollinear_features(percentage_change,threshold,count_vect)
```

4

**[3.1] Logistic Regression on AVG W2V, SET 3**

[3.1.1] Applying Logistic Regression with L1 regularization on AVG W2V, SET 3

```
In [35]: 1 searchMethod='random'
2 #STANDARDIZE TRAIN AND TEST DATA
3 train, test=std_data(train=avg_sent_vectors,test=avg_sent_vectors_test,mean=True)
4 #HYPERPARAM TUNNING
5 %time model=LR_Classifier(train,y_train,TBS,params,penalty[0],searchMethod,vect[2])
6 print('Optimal value of C(1/lambda): ',model.best_params_)
```



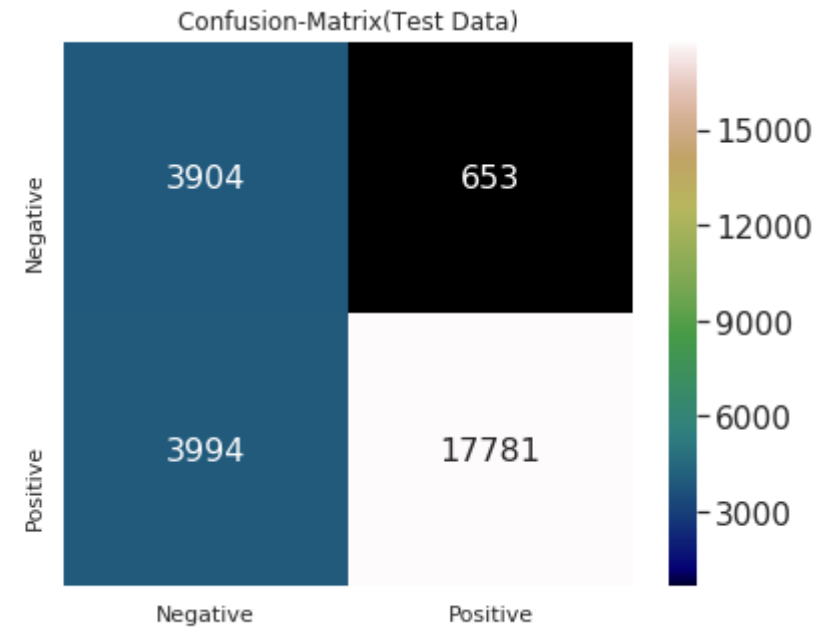
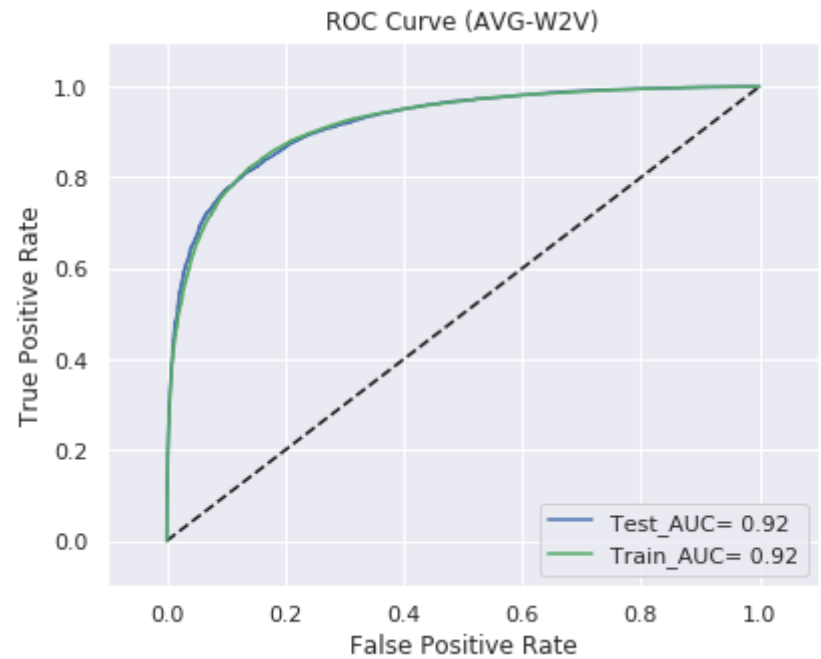
CPU times: user 21.1 s, sys: 464 ms, total: 21.5 s

Wall time: 38 s

Optimal value of $C(1/\lambda)$: {'C': 0.05}

[3.1.1.1] Performance on test data with optimal value of hyperparam

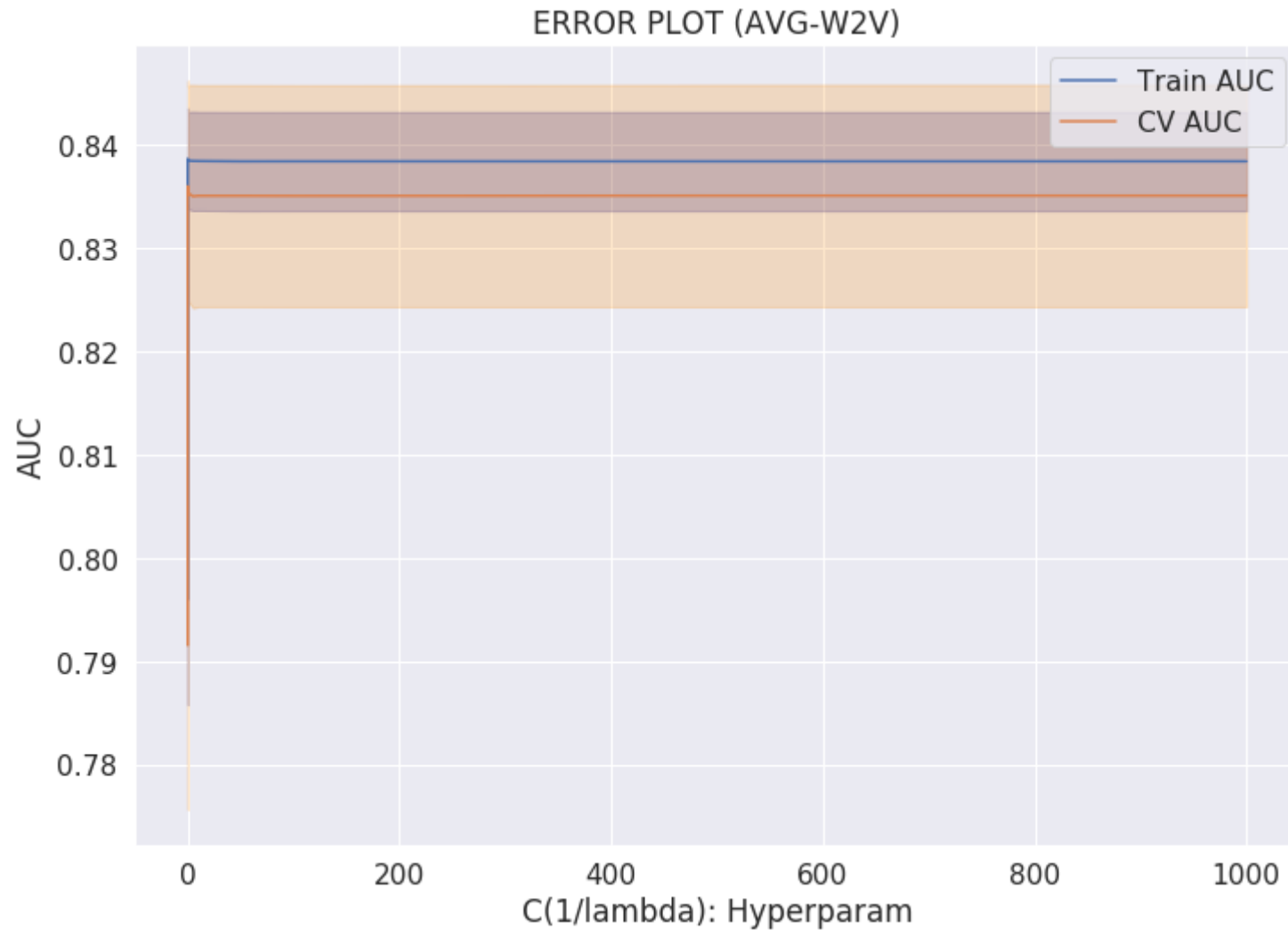
```
In [36]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['C'],penalty[0],vect[2],summarize)
```



[3.1.2] Applying Logistic Regression with L2 regularization on AVG W2V, SET 3

In [37]:

```
1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=avg_sent_vectors,test=avg_sent_vectors_test,mean=True)
3 #HYPERPARAM TUNNING
4 %time model=LR_Classifier(train,y_train,TBS,params,penalty[1],searchMethod,vect[2])
5 print('Optimal value of C(1/lambda): ',model.best_params_)
```



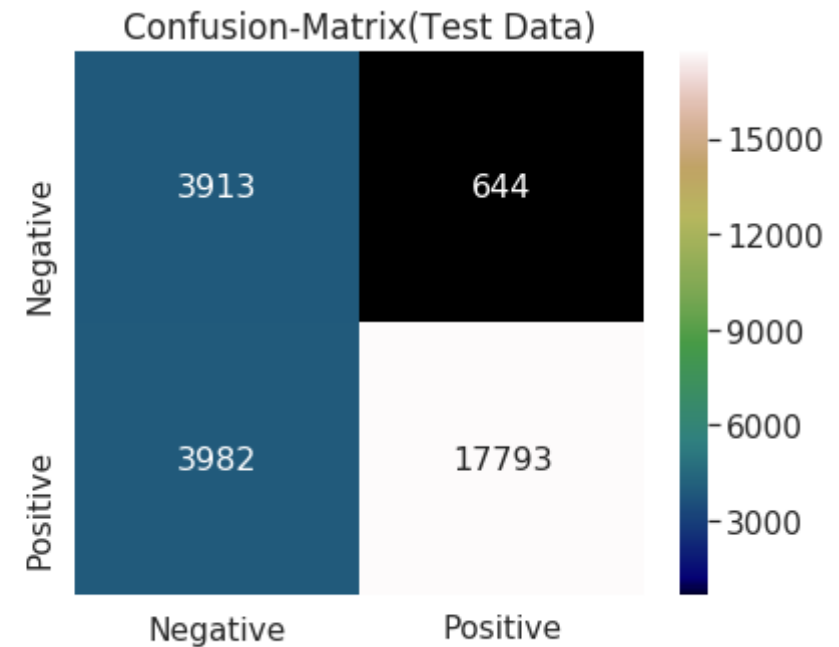
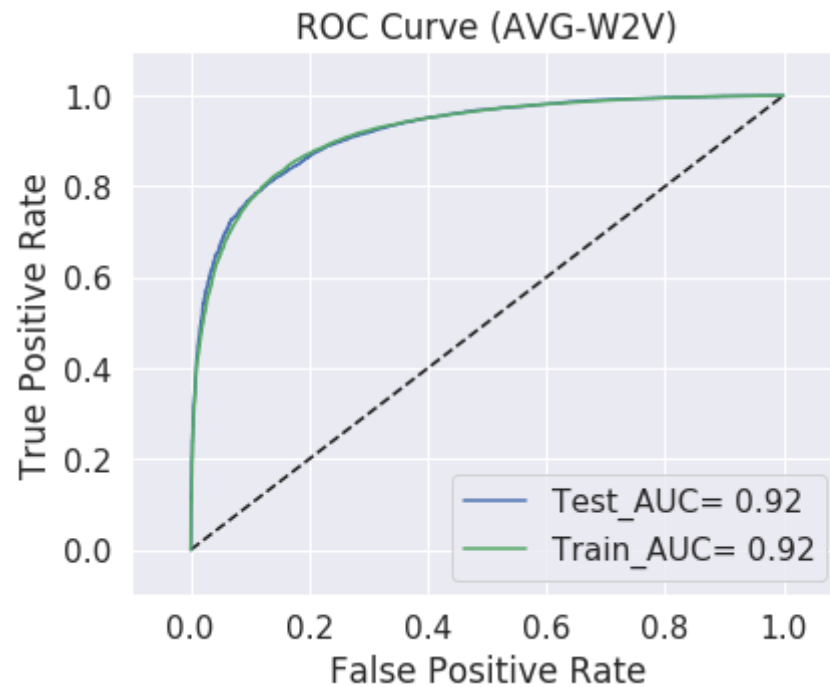
CPU times: user 20.5 s, sys: 328 ms, total: 20.8 s

Wall time: 36 s

Optimal value of C(1/lambda): {'C': 0.05}

[3.1.2.1] Performance on test data with optimal value of hyperparam

```
In [38]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['C'],penalty[1],vect[2],summarize)
```

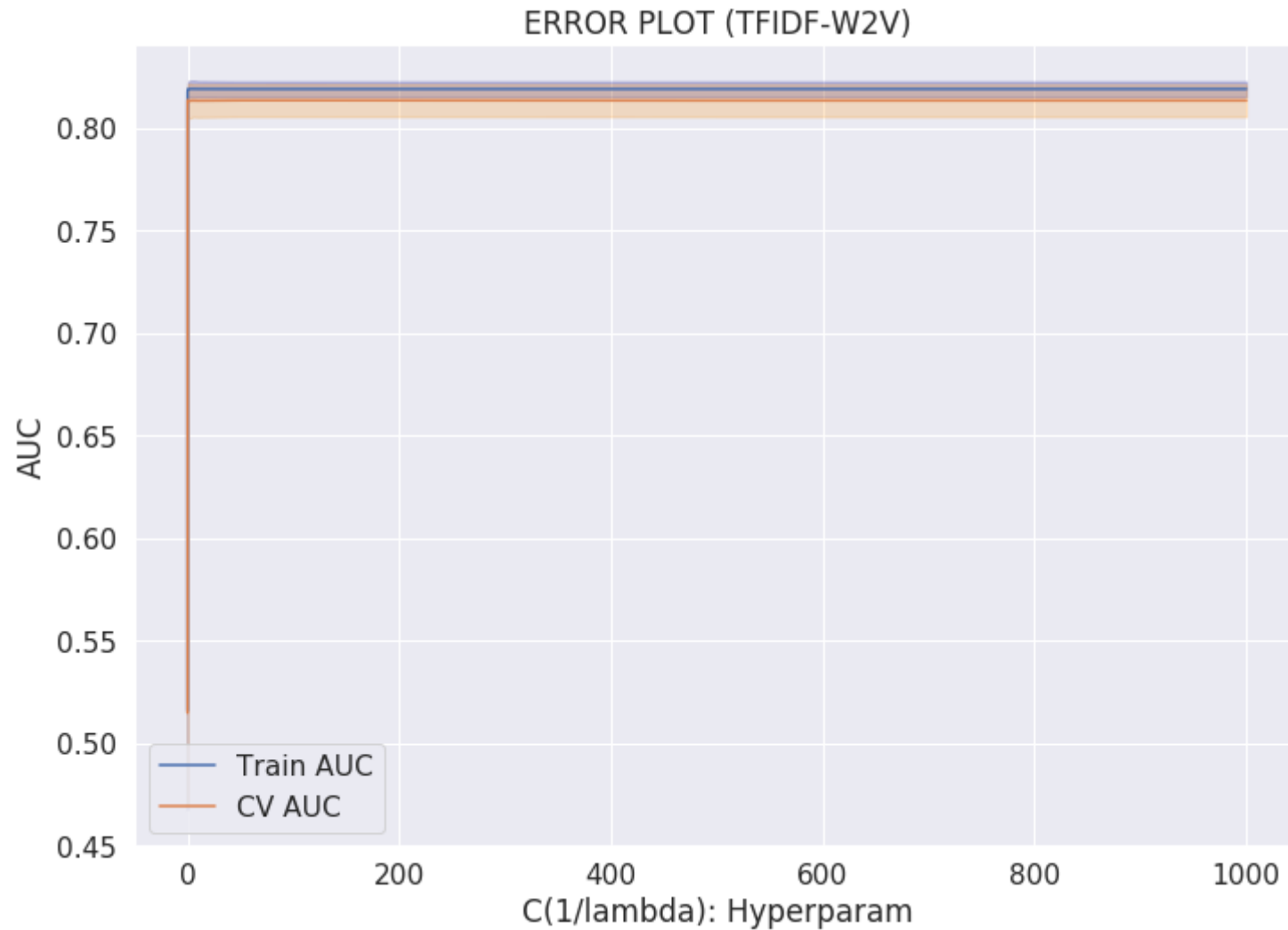


[4.1] Logistic Regression on TFIDF W2V, SET 4

[4.1.1] Applying Logistic Regression with L1 regularization on TFIDF W2V, SET 4

In [39]:

```
1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=tfidf_sent_vectors,test=tfidf_sent_vectors_test,mean=True)
3 #HYPERPARAM TUNNING
4 %time model=LR_Classifier(train,y_train,TBS,params,penalty[0],searchMethod,vect[3])
5 print('Optimal value of C(1/lambda): ',model.best_params_)
```



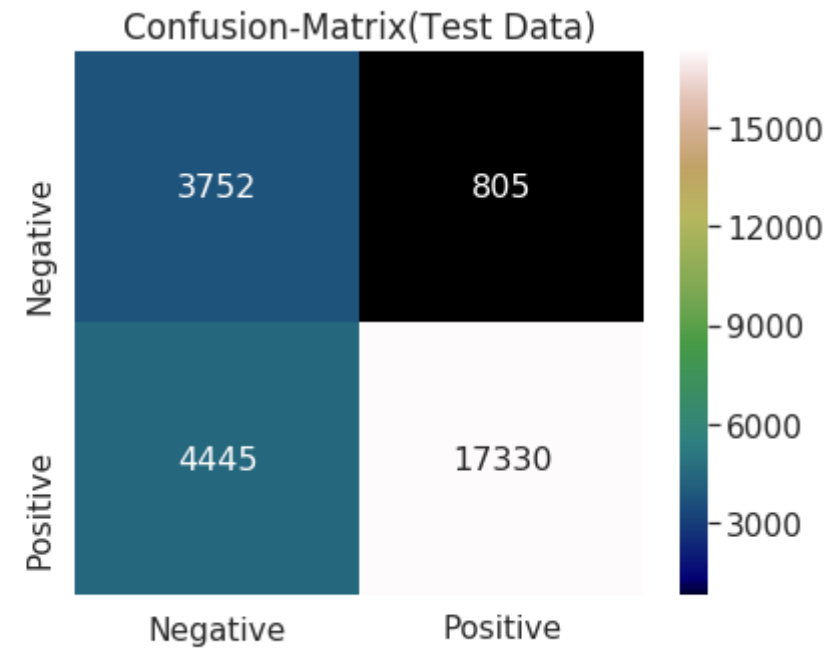
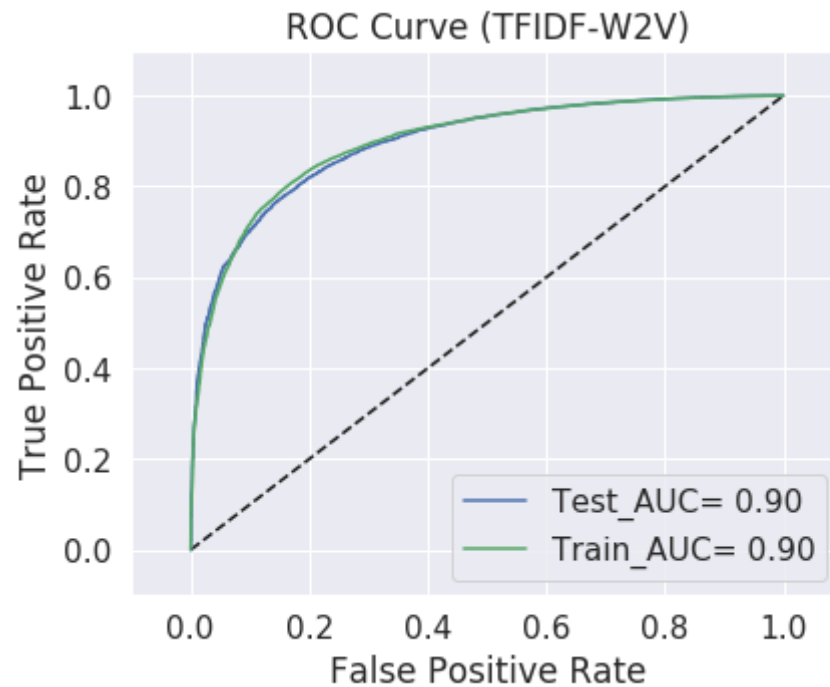
CPU times: user 22.8 s, sys: 460 ms, total: 23.3 s

Wall time: 49.8 s

Optimal value of $C(1/\lambda)$: {'C': 0.1}

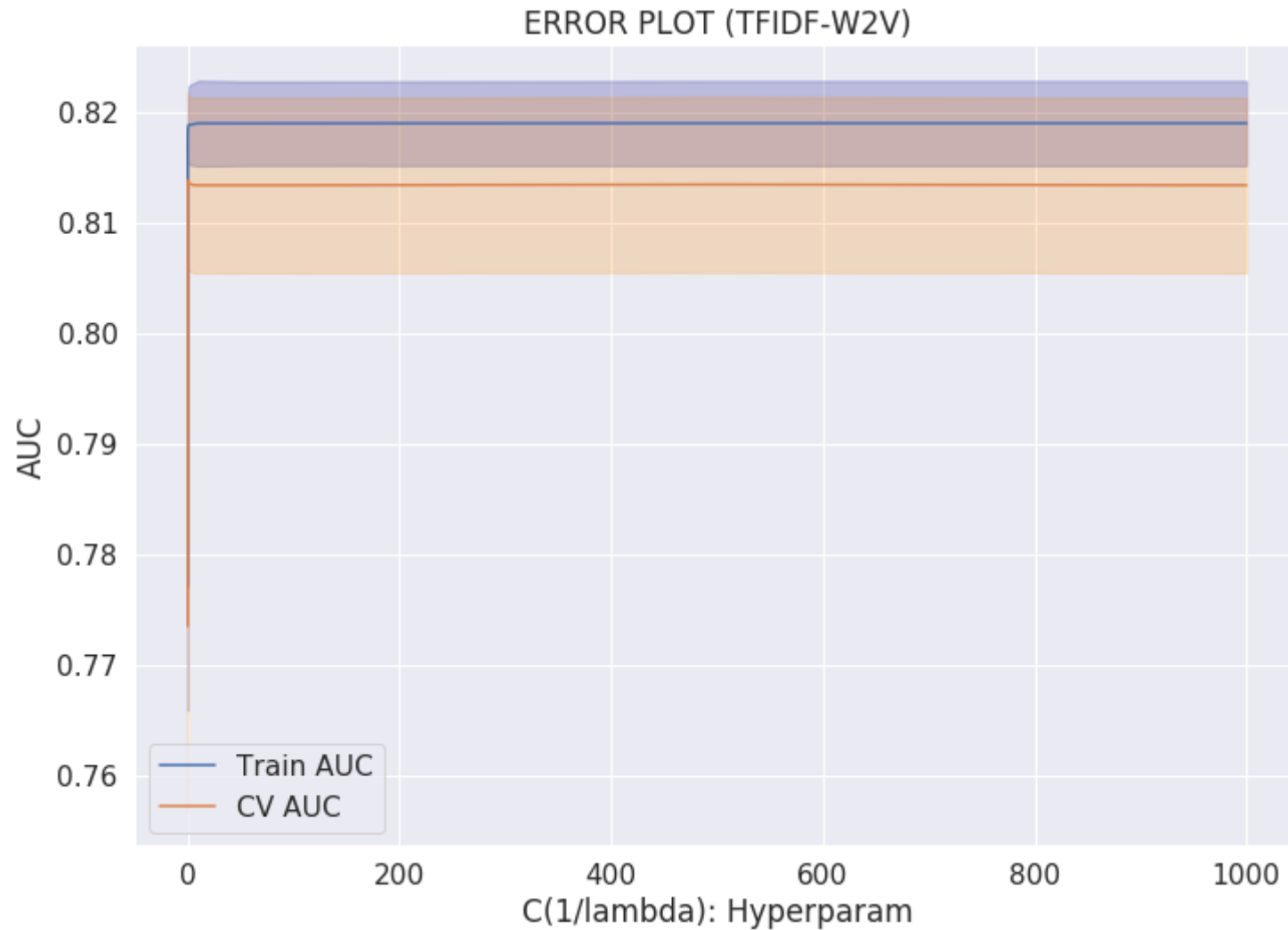
[4.1.1.1] Performance on test data with optimal value of hyperparam

```
In [40]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['C'],penalty[0],vect[3],summarize)
```

**[4.1.2] Applying Logistic Regression with L2 regularization on TFIDF W2V, SET 4**

In [41]:

```
1 #STANDARDIZE TRAIN AND TEST DATA
2 train, test=std_data(train=tfidf_sent_vectors,test=tfidf_sent_vectors_test,mean=True)
3 #HYPERPARAM TUNNING
4 %time model=LR_Classifier(train,y_train,TBS,params,penalty[1],searchMethod,vect[3])
5 print('Optimal value of C(1/lambda): ',model.best_params_)
```



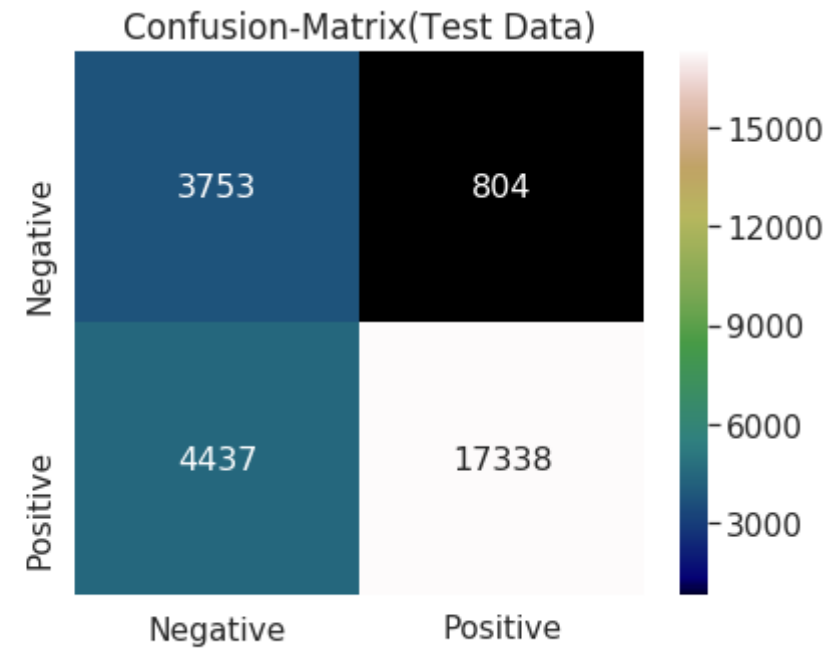
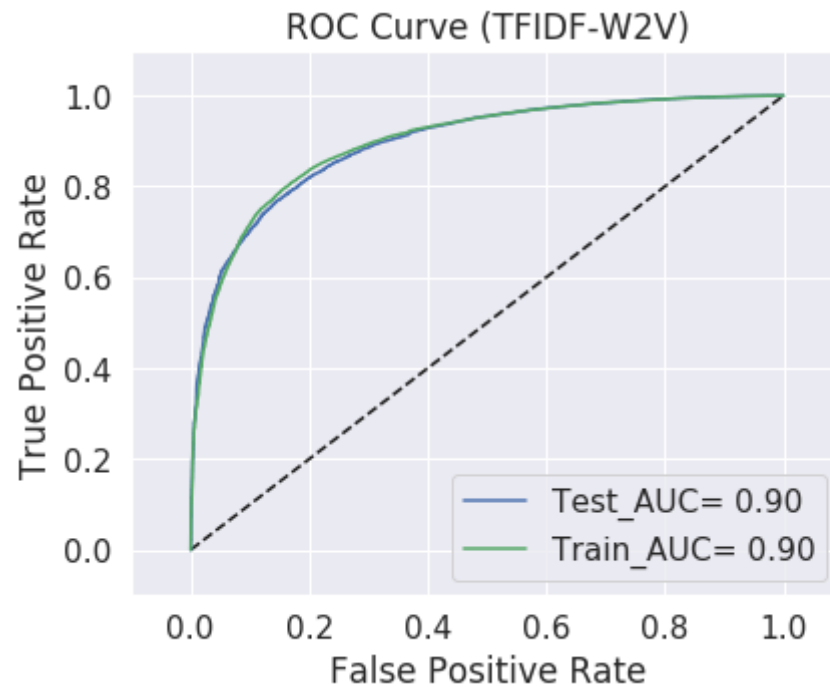
CPU times: user 21 s, sys: 364 ms, total: 21.4 s

Wall time: 37.8 s

Optimal value of $C(1/\lambda)$: {'C': 0.1}

[4.1.2.1] Performance on test data with optimal value of hyperparam

```
In [42]: 1 clf=test_performance(train,y_train,test,y_test,model.best_params_['C'],penalty[1],vect[3],summarize)
```



Conclusions:

In [44]: 1 `print(summarize)`

Vectorizer	Regularizer	Optimal-C(1/lambda)	Test(AUC)	Test(f1-score)
BoW	11	0.005	0.952	0.912
BoW	12	0.0001	0.942	0.922
TF-IDF	11	0.005	0.959	0.911
TF-IDF	12	0.0001	0.950	0.918
AVG-W2V	11	0.05	0.919	0.840
AVG-W2V	12	0.05	0.919	0.841
TFIDF-W2V	11	0.1	0.895	0.820
TFIDF-W2V	12	0.1	0.895	0.820

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

- TFIDF vectorizer
- f1-score=.911 and auc=.959

In []:

1