# OBJECTIVE :

1. APPLYING LOGISTIC REGRESSION WITH AVG WORD2VEC VECTORIZATION
   - PERFORMING PERTUBATION TEST TO CHECK WHETHER OUR DATA FEATURES ARE COLLINER OR NOT AND PLOTTING THE RESULT
   - FINDING THE BEST HYPERPARAMETER USING GRIDSEARCHCV WITH TRAIN DATA AND CROSS-VALIDATION DATA BY PLOTTING THE RESLUTS OF VAROIUS TRAIN DATA AND CROSS VALIDATION DATA
   - USING THE APROPRIATE VALUE OF HYPERPARAMETER ,TESTING ACCURACY ON TEST DATA USING F1-SCORE
   - PLOTTING THE CONFUSION MATRIX TO GET THE PRECISOIN ,RECALL VALUE WITH HELP OF HEATMAP
   - PRINTING THE TOP 20 FEATURES FOR BOTH POSITIVE AND NEGATIVE WORDS #

```python
In [0]: from sklearn.model_selection import train_test_split    #importing the necessary libraries
        from sklearn.model_selection import RandomizedSearchCV
        from sklearn.datasets import *
        from sklearn import naive_bayes
        from sklearn.feature_extraction.text import CountVectorizer
        from sklearn.feature_extraction.text import TfidfVectorizer
        import numpy as np
        import pandas as pd
        from sklearn import *
        from gensim.models import Word2Vec
        import warnings
        warnings.filterwarnings("ignore")
        from tqdm import tqdm
```

```python
In [0]: from google.colab import drive
        drive.mount('/content/gdrive')#geeting the content from the google driv
```

```
e
```

```
Drive already mounted at /content/gdrive; to attempt to forcibly remoun
t, call drive.mount("/content/gdrive", force_remount=True).
```

In [0]:
```python
final_processed_data=pd.read_csv("gdrive/My Drive/final_new_data.csv")#
loading the preprocessed data  with 100k points into dataframe
```

In [0]:
```python
# getting the counts of 0 and 1 in "SCORE" column to know whether it is
 unbalanced data or not
count_of_1=0
count_of_0=0
for i in final_processed_data['Score']:
    if i==1:
      count_of_1+=1
    else:
      count_of_0+=1
print(count_of_1)
print(count_of_0)
#it is an imbalanced dataset
```

```
88521
11479
```

In [0]:
```python
#spliiting the data into train and test data
x_train,x_test,y_train,y_test=model_selection.train_test_split(final_pr
ocessed_data['CleanedText'].values,final_processed_data['Score'].values
,test_size=0.2,shuffle=False)
```

In [0]:
```python
# Training my own Word2Vec model using your own text corpus
list_of_sent=[]
for sent in x_train:
 list_of_sent.append(sent.split())#splitting of sentences into words AN
D appending them to list
print(x_train[0])
print("***********************************************************
*")
print(list_of_sent[0])
```

```python
word_to_vector=Word2Vec(list_of_sent,min_count=5,size=50,workers=2)#con
structing my our word to vector
w_t_c_words=list(word_to_vector.wv.vocab)
print("***************************************************************
*******")
print("sample words ", w_t_c_words[0:50])
```

witti littl book make son laugh loud recit car drive along alway sing r
efrain hes learn whale india droop love new word book introduc silli cl
assic book will bet son still abl recit memori colleg
***************************************************************
['witti', 'littl', 'book', 'make', 'son', 'laugh', 'loud', 'recit', 'ca
r', 'drive', 'along', 'alway', 'sing', 'refrain', 'hes', 'learn', 'whal
e', 'india', 'droop', 'love', 'new', 'word', 'book', 'introduc', 'sill
i', 'classic', 'book', 'will', 'bet', 'son', 'still', 'abl', 'recit',
'memori', 'colleg']
***************************************************************
sample words  ['witti', 'littl', 'book', 'make', 'son', 'laugh', 'lou
d', 'car', 'drive', 'along', 'alway', 'sing', 'refrain', 'hes', 'lear
n', 'india', 'droop', 'love', 'new', 'word', 'introduc', 'silli', 'clas
sic', 'will', 'bet', 'still', 'abl', 'memori', 'colleg', 'rememb', 'se
e', 'show', 'air', 'televis', 'year', 'ago', 'child', 'sister', 'late
r', 'bought', 'day', 'thirti', 'someth', 'use', 'seri', 'song', 'studen
t', 'teach', 'preschool', 'turn']

In [0]:
```python
###### NOW STARTING AVERAGE WORD TO VEC FOR TRAIN DATA################
#################################################
train_sent_vectors = []; # the avg-w2v for each sentence/review is stor
ed in this list
for sent in tqdm(list_of_sent): # for each review/sentence
 sent_vec = np.zeros(50) # as word vectors are of zero length
 cnt_words =0; # num of words with a valid vector in the sentence/revie
w
  for word in sent: # for each word in a review/sentence
    if word in w_t_c_words:
      vec = word_to_vector.wv[word]
      sent_vec += vec
      cnt_words += 1
  if cnt_words != 0:
```

```
  sent_vec /= cnt_words
 train_sent_vectors.append(sent_vec)
print(len(train_sent_vectors))
print(len(train_sent_vectors[0]))
```

```
100%|███████████| 80000/80000 [05:39<00:00, 235.31it/s]
```

```
80000
50
```

In [0]:
```python
from sklearn.preprocessing import StandardScaler #standarizing the trai
ning  data
x_train_data=StandardScaler( with_mean=False).fit_transform(train_sent_
vectors)
print(x_train_data.shape)
```

```
(80000, 50)
```

In [0]:
```python
list_of_sent=[]
for sent in x_test:
 list_of_sent.append(sent.split())#splitting of sentences into words AN
D appending them to list
print(x_test[0])
print("****************************************************************
*")
print(list_of_sent[0])
print('****************************************************************
***')
```

```
hard find item dont buy mani either came stale got way quick classic no
netheless
***************************************************************
['hard', 'find', 'item', 'dont', 'buy', 'mani', 'either', 'came', 'stal
e', 'got', 'way', 'quick', 'classic', 'nonetheless']
***************************************************************
```

In [0]:
```python
###### NOW STARTING AVERAGE WORD TO VEC FOR TEST DATA################
##################################################
sent_vectors = []; # the avg-w2v for each sentence/review is stored in
```

```
    this list
  for sent in list_of_sent: # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    cnt_words =0; # num of words with a valid vector in the sentence/revie
    w
    for word in sent: # for each word in a review/sentence
      if word in w_t_c_words:
        vec = word_to_vector.wv[word]
        sent_vec += vec
        cnt_words += 1
    if cnt_words != 0:
      sent_vec /= cnt_words
    sent_vectors.append(sent_vec)
  print(len(sent_vectors))
  print(len(sent_vectors[0]))
```

```
20000
50
```

In [0]:
```
from sklearn.preprocessing import StandardScaler #standarizing the trai
ning  data
x_test_data=StandardScaler( with_mean=False).fit_transform(sent_vectors
)
print(x_test_data.shape)
```

```
(20000, 50)
```

In [0]:
```
#using time series split method for cross-validation score
from sklearn.model_selection import TimeSeriesSplit
tscv = TimeSeriesSplit(n_splits=10)
from sklearn.linear_model import LogisticRegression
from scipy.stats import uniform
data=[10**-4,10**-3,10**-2,10**-1,10**0,10**1,10**2,10**3,10**4]#range
 of hyperparameter
```

In [0]:
```
lr=LogisticRegression(penalty='l2',class_weight={1:.5,0:.5},n_jobs=-1)#
building logistic regression model
tuned_parameters=[{'C':data}]
```

```
In [0]: #applying the model of logistic regression and using gridsearchcv to fi
        nd the best hyper parameter
        %time
        from sklearn.model_selection import GridSearchCV
        model = GridSearchCV(lr, tuned_parameters, scoring = 'f1', cv=tscv,n_jo
        bs=-1)#building the gridsearchcv model
        model.fit(x_train_data, y_train)#fiitting the training data

        print(model.best_estimator_)#printing the best_estimator
        print(model.score(x_test_data, y_test))#predicting  f1 score on test da
        ta
```

```
CPU times: user 5 µs, sys: 0 ns, total: 5 µs
Wall time: 9.3 µs
LogisticRegression(C=10, class_weight={1: 0.5, 0: 0.5}, dual=False,
          fit_intercept=True, intercept_scaling=1, max_iter=100,
          multi_class='warn', n_jobs=-1, penalty='l2', random_state=Non
e,
          solver='warn', tol=0.0001, verbose=0, warm_start=False)
0.9414095825191954
```

```
In [0]: lr=LogisticRegression(C=10,penalty='l2',class_weight={1:.5,0:.5},n_jobs
        =-1)#building model for getting wieght vector
        lr.fit(x_train_data,y_train)#fitting the training data
        z=lr.decision_function(x_train_data)#checking the signed distance of a
         point from hyperplane
        print(z)#printing the signed distance
```

```
[4.20125878 3.60267944 4.69642116 ... 3.17373434 1.63222391 4.14039603]
```

```
In [0]: wieght_vector=lr.coef_#getting the weight vector
        print(wieght_vector.shape)#wieght vector shape
        print(wieght_vector[:20])
```

```
(1, 50)
[[-0.02904195 -0.11591003 -0.26360131  0.41546582 -0.28146277  0.244934
63
  -0.22771467 -0.36448311  0.09265409  0.4494914   0.15364282 -0.493196
65
```

```
    -0.05501359 -0.17571591 -0.40005203 -0.20770648  0.06262386  0.119486
88
    -0.25327633 -0.04375454 -0.54338845 -0.38397442  0.25019977  0.483241
15
    -0.45697059  0.27151116  0.62401274  0.44643779 -0.46102941  0.590072
08
     0.40110036  0.11620863 -0.00840976  0.33985595  0.16540033  0.070417
01
    -0.24189498  0.25014858 -0.20108789  0.20547129 -0.31597726 -0.245029
4
    -0.08363616  0.13000548 -0.10402377  0.20358218 -0.00952736 -0.179978
81
    -0.31772343 -0.40610872]]
```

# PERTUBATION TEST :

AIM:    TO CHECK FOR MULTI COLLINEARITY  OF FEATURES
STEPS

1. GETTING THE WIEGHT VECTOR FROM MODEL AND SAVING IT</br>
2. ADDING NOISE TO THE TRAINING DATA TO GET NEW TRAINING DATA</br>
3. FITTING THE MODEL AGAIN ON NEW DATA</br>
4. GETTING THE WIEGHT VECTOR FROM THIS MODEL</br> 5.ADDING SMALL VALUE TO WEIGHT VECTOR OF BOTH TRAINNG DATA TO REMOVE ANY ERROR
5. FINDING THE PERCENTAGE CHANGE VECTOR
6. GEETING HOW MANY GEATURE HAS CHANGED USING SOME THRESHOLD VALUE( HERE TAKING IT AS 100)
7. PLOTTING THE QUANTILES WITH THIER PERCENTAGE WIGHT VALUE TO CHECK IF COLLINEARITY EXITS OR NOT

# RESULT : TO KNOW WHETHER FEATURES ARE MULTICOLLINEAR OR NOT # # AND TO KNOW WHETHER MODEL IS RELIABLE OR NOT #

In [0]:
```python
#here,we are adding noise to the data
from scipy.stats import norm
```

```
noise=norm.rvs(size=1)#noise
x_train_data.data+=noise#adding noise
```

In [0]:
```
print('shape of our new train data  after adding noise is : ',x_train_d
ata.shape)#printing shape of new training data
```

```
shape of our new train data  after adding noise is :  (80000, 50)
```

In [0]:
```
#uilding the model using timeSeriesSplit
from sklearn.model_selection import TimeSeriesSplit
tscv = TimeSeriesSplit(n_splits=10) # 10 spilts cross validation
from sklearn.linear_model import LogisticRegression
from scipy.stats import uniform
data=[10**-4,10**-3,10**-2,10**-1,10**0,10**1,10**2,10**3,10**4]#value
 range of hyper parameter for grid searchcv
lr=LogisticRegression(penalty='l2',class_weight={1:.5,0:.5},n_jobs=-1)#
building the model
tuned_parameters=[{'C':data}]
```

In [0]:
```
%time
from sklearn.model_selection import GridSearchCV
model = GridSearchCV(lr, tuned_parameters, scoring = 'f1', cv=tscv,n_jo
bs=-1)#building the gridsearchcv model
model.fit(x_train_data, y_train)#fiitting the training data

print('best estimator of our new data is: ',model.best_estimator_)#prin
ting the best_estimator
```

```
CPU times: user 4 µs, sys: 0 ns, total: 4 µs
Wall time: 9.3 µs
best estimator of our new data is:  LogisticRegression(C=10, class_weig
ht={1: 0.5, 0: 0.5}, dual=False,
          fit_intercept=True, intercept_scaling=1, max_iter=100,
          multi_class='warn', n_jobs=-1, penalty='l2', random_state=Non
e,
          solver='warn', tol=0.0001, verbose=0, warm_start=False)
```

In [0]:
```
# again building the model for finding the wieght vector of the words f
```

```
rom model
lr=LogisticRegression(C=10,penalty='l2',class_weight={1:.5,0:.5},n_jobs
=-1)#building the logistic regression model
lr.fit(x_train_data,y_train)#fiting the training model
new_wieght_vector=lr.coef_
print(new_wieght_vector.shape)#printing shape of wieght vector
```

(1, 50)

In [0]:
```
percent_change_vec=np.ones((1,50))#generating the percent_change_vetor
 to store the percentage change  values for each word
```

In [0]:
```
wieght_vector=wieght_vector+10**-6 #adding some values to wieght vector
 to avoid error while division

new_wieght_vector=new_wieght_vector+10**-6 #adding some values to wiegh
t vector to avoid error while division

percent_change_vec=abs((wieght_vector-new_wieght_vector)/wieght_vector)
*100#calculating the percentage change in the vector
```

In [0]:
```
x=abs((wieght_vector[0][2]-new_wieght_vector[0][2])/wieght_vector[0][2
])#just checking randomly that every value is positve in percent_change
_vector
print(x)
```

0.00037991118132659076

In [0]:
```
print('shape of percent change wieght vector is', percent_change_vec.sh
ape)#printing shape of percent_change_vector
```

shape of percent change wieght vector is (1, 50)

In [0]:
```
per_change_df=pd.DataFrame(percent_change_vec.T,columns=['CHANGE'])#bui
lding a dataframe from wight vector
```

In [0]:
```
per_change_df.head()#getting first 5 values
```

Out[0]:

|   | CHANGE |
|---|--------|
| 0 | 0.848063 |
| 1 | 0.666781 |
| 2 | 0.037991 |
| 3 | 0.165658 |
| 4 | 0.439478 |

```python
In [0]: sorted_Df=per_change_df.sort_values('CHANGE',ascending=True,axis=0)#sor
        ting the dataframe to calculate the quantiles values
        sorted_Df.describe()#describe function
```

Out[0]:

|   | CHANGE |
|---|--------|
| count | 50.000000 |
| mean | 0.439151 |
| std | 0.713103 |
| min | 0.003997 |
| 25% | 0.094802 |
| 50% | 0.211508 |
| 75% | 0.483936 |
| max | 4.392146 |

```python
In [0]: quantiles=list( i/100 for i in range(0,101,5))#building the list of qua
        ntiles value
        for i in quantiles:
          print('sorted_Data {:.2f}th quantiles is {:7.3f}'.format(i,sorted_Df[
        'CHANGE'].quantile(i)))#printing the quantiles and thier coreesponding
         values
```

```
sorted_Data 0.00th quantiles is     0.004
sorted_Data 0.05th quantiles is     0.024
sorted_Data 0.10th quantiles is     0.062
sorted_Data 0.15th quantiles is     0.081
sorted_Data 0.20th quantiles is     0.090
sorted_Data 0.25th quantiles is     0.095
sorted_Data 0.30th quantiles is     0.114
sorted_Data 0.35th quantiles is     0.143
sorted_Data 0.40th quantiles is     0.160
sorted_Data 0.45th quantiles is     0.166
sorted_Data 0.50th quantiles is     0.212
sorted_Data 0.55th quantiles is     0.235
sorted_Data 0.60th quantiles is     0.295
sorted_Data 0.65th quantiles is     0.361
sorted_Data 0.70th quantiles is     0.430
sorted_Data 0.75th quantiles is     0.484
sorted_Data 0.80th quantiles is     0.541
sorted_Data 0.85th quantiles is     0.695
sorted_Data 0.90th quantiles is     0.854
sorted_Data 0.95th quantiles is     1.453
sorted_Data 1.00th quantiles is     4.392
```

In [0]:
```python
quantiles=list( i/100 for i in range(95,101,1))#printing the last perce
ntiles values because this region is showing abrupt change
percent_change_list=[]#empty percent_change
for i in quantiles:
  print('sorted_Data {:.2f}th quantiles is {:7.3f}'.format(i,sorted_Df[
'CHANGE'].quantile(i)))
  percent_change_list.append(sorted_Df['CHANGE'].quantile(i))#building
 the list
```

```
sorted_Data 0.95th quantiles is     1.453
sorted_Data 0.96th quantiles is     1.890
sorted_Data 0.97th quantiles is     2.060
```
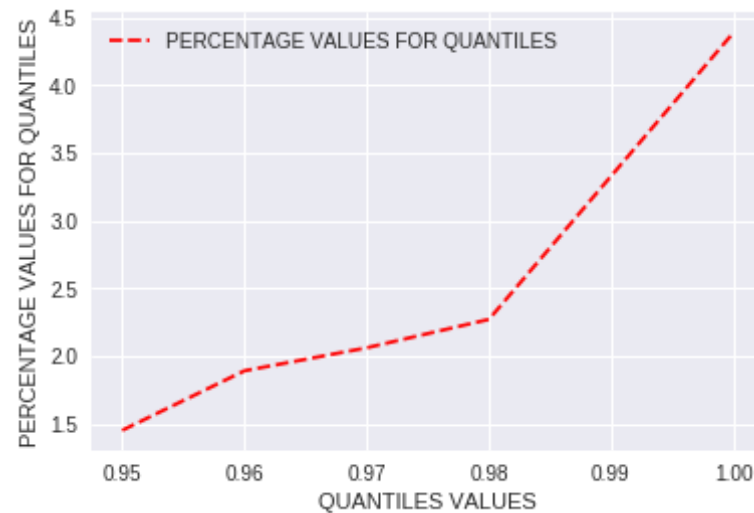
```
sorted_Data 0.98th quantiles is    2.267
sorted_Data 0.99th quantiles is    3.330
sorted_Data 1.00th quantiles is    4.392
```

In [0]:
```python
print(percent_change_list)
my_formatted_list = [ '%.2f' % elem for elem in percent_change_list ]#formatted list with string values in it
my_formatted_list=[float(i) for i in my_formatted_list]#formatted list with flaot values in it
print(my_formatted_list)#printing formatted list
print(quantiles)#printing quantiles
```

```
[1.4526432286776458, 1.8900288702504262, 2.0604313484690424, 2.26724396
5870643, 3.329694854072624, 4.3921457422746055]
[1.45, 1.89, 2.06, 2.27, 3.33, 4.39]
[0.95, 0.96, 0.97, 0.98, 0.99, 1.0]
```

In [0]:
```python
%matplotlib inline
import matplotlib.pyplot as plt
plt.show()
plt.xlabel('QUANTILES VALUES')
plt.ylabel('PERCENTAGE VALUES FOR QUANTILES')
plt.plot(quantiles,my_formatted_list,'r--' ,label='PERCENTAGE VALUES FOR QUANTILES')
plt.legend(loc='best')
```

Out[0]: <matplotlib.legend.Legend at 0x7ff532cc7e10>

**FROM THE ABOVE VISUALIZATION , MAIN POINTS ARE:. </font>**

1. THAT ONLY 2% OF FEATURES GOT AFFECTED AFTER ADDING NOISE TO THE DATA.
2. VERY LESS COLLINEARITY OF DATA IS PRESENT ,BECAUSE MOST OF THE WEIGHT VECTORS VALUES REMAINS SAME
3. THERFORE, OUR MODEL IS RELIABLE AND WE CAN PROCEED FURTHER TO CHECK ACCURACY ON TEST DATA </font>

# CALCULATING THE BEST HYPERPARAMETER ON TRAIN DATA AND CALCULATING THE ACCURACY USING F1-SCORE AND PLOTTING IT

```
In [0]:  #using time series split method for cross-validation score
         from sklearn.model_selection import TimeSeriesSplit
         tscv = TimeSeriesSplit(n_splits=10)
```

```python
from sklearn.linear_model import LogisticRegression
from scipy.stats import uniform
data=[10**-4,10**-3,10**-2,10**-1,10**0,10**1,10**2,10**3,10**4]#range
 of hyperparameter
```

In [0]:
```python
lr=LogisticRegression(penalty='l2',class_weight={1:.5,0:.5},n_jobs=-1)#
building logistic regression model
tuned_parameters=[{'C':data}]
```

In [0]:
```python
#applying the model of logistic regression and using gridsearchcv to fi
nd the best hyper parameter
from sklearn.model_selection import GridSearchCV
model = GridSearchCV(lr, tuned_parameters, scoring = 'f1', cv=tscv,n_jo
bs=-1)#building the gridsearchcv model
model.fit(x_train_data, y_train)#fiitting the training data
```

Out[0]:
```
GridSearchCV(cv=TimeSeriesSplit(max_train_size=None, n_splits=10),
        error_score='raise-deprecating',
        estimator=LogisticRegression(C=1.0, class_weight={1: 0.5, 0: 0.
5}, dual=False,
            fit_intercept=True, intercept_scaling=1, max_iter=100,
            multi_class='warn', n_jobs=-1, penalty='l2', random_state=Non
e,
            solver='warn', tol=0.0001, verbose=0, warm_start=False),
        fit_params=None, iid='warn', n_jobs=-1,
        param_grid=[{'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 1
0000]}],
        pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
        scoring='f1', verbose=0)
```

In [0]:
```python
results=pd.DataFrame(model.cv_results_)# getting varoius cv_scores and
 train_scores various values of alpha given as parameter and storing it
 in a dataframe
results#printing the dataframe
```

Out[0]:

| | mean_fit_time | mean_score_time | mean_test_score | mean_train_score | param_C | para |
|---|---|---|---|---|---|---|

| | mean_fit_time | mean_score_time | mean_test_score | mean_train_score | param_C | para |
|---|---|---|---|---|---|---|
| 0 | 0.452241 | 0.005229 | 0.941373 | 0.941422 | 0.0001 | {'C': 0.000 |
| 1 | 0.738270 | 0.007523 | 0.945196 | 0.946943 | 0.001 | {'C': 0.001 |
| 2 | 1.096690 | 0.006123 | 0.948356 | 0.951048 | 0.01 | {'C': 0.01} |
| 3 | 1.255957 | 0.004825 | 0.948864 | 0.951773 | 0.1 | {'C': 0.1} |
| 4 | 1.298706 | 0.005174 | 0.948871 | 0.951781 | 1 | {'C': 1 |
| 5 | 1.301231 | 0.005741 | 0.948892 | 0.951810 | 10 | {'C': 1 |
| 6 | 1.298898 | 0.004721 | 0.948830 | 0.951865 | 100 | {'C': 100} |
| 7 | 1.294912 | 0.006524 | 0.948830 | 0.951865 | 1000 | {'C': 1000} |
| 8 | 1.275710 | 0.006042 | 0.948830 | 0.951865 | 10000 | {'C': 10000 |

9 rows × 31 columns

In [0]:
```python
%matplotlib inline
import matplotlib.pyplot as plt


mean_test_score=list(results['mean_test_score'])#taking mean_test_score
 values of various alpha into a list
mean_train_score=list(results['mean_train_score'])#taking mean_train_sc
ore values of varoius alpha into a list
cv_error_list=[]
train_error_list=[]
```

```python
for i  in mean_test_score:
    i=1-i
    i=i*100
    cv_error_list.append(i)#appending the list with cv_error
for i  in mean_train_score:
    i=1-i
    i=i*100
    train_error_list.append(i)#appending  the list with train_error

print(cv_error_list)
C_values_in_10_power=[-4,-3,-2,-1,0,1,2,3,4]#list of alpha values in po
wer of 10
plt.plot(C_values_in_10_power,cv_error_list,label='cv_error')#plotting
 alpha with cv_error
plt.plot(C_values_in_10_power,train_error_list,label='train_error')#plo
tting aplhawith train_error
plt.xlabel('C value in  power of 10 ')
plt.ylabel('cv error and train error')
plt.legend(loc='best')
```
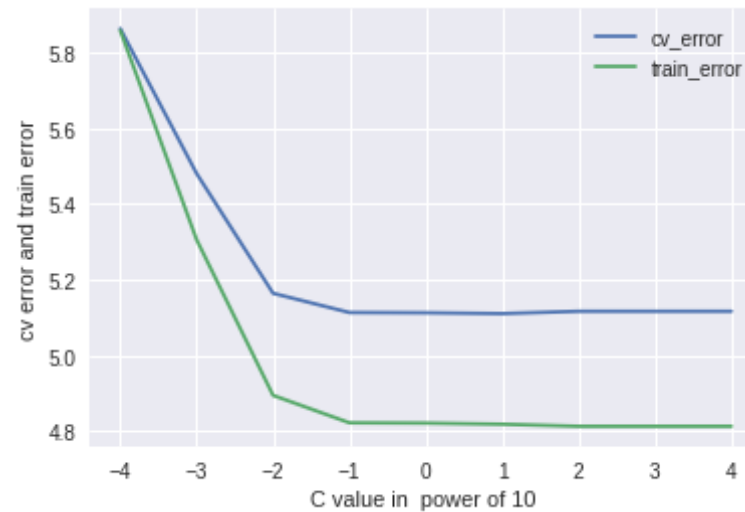
```
[5.862680547266585, 5.48035653053115, 5.164413401457802, 5.113619361717
725, 5.112941643298896, 5.110790608753691, 5.116956440684584, 5.1169564
40684584, 5.116956440684584]
```

Out[0]: <matplotlib.legend.Legend at 0x7ff532201898>

**From here,the best hyperparameter value is c=10 or alpha=0.1**

**USING BEST HYPERPARAMETER VALUE ON TEST DATA AND PLOTTING THE CONFUSION MATRIX WITH HEATMAP**

In [0]:
```python
#Testing Accuracy on Test data
import seaborn as sns #importing seaborn as sns
from sklearn.metrics import *#importing varoius metrics from sklearn
lr=LogisticRegression(C=10,penalty='l2',class_weight={1:.5,0:.5},n_jobs=-1)#building logistic regression model#building the model
lr.fit(x_test_data,y_test)
y_pred = lr.predict(x_test_data)
print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_pred)*100))#printing accuracy
print("Precision on test set: %0.3f"%(precision_score(y_test, y_pred)))#printing precision score
```
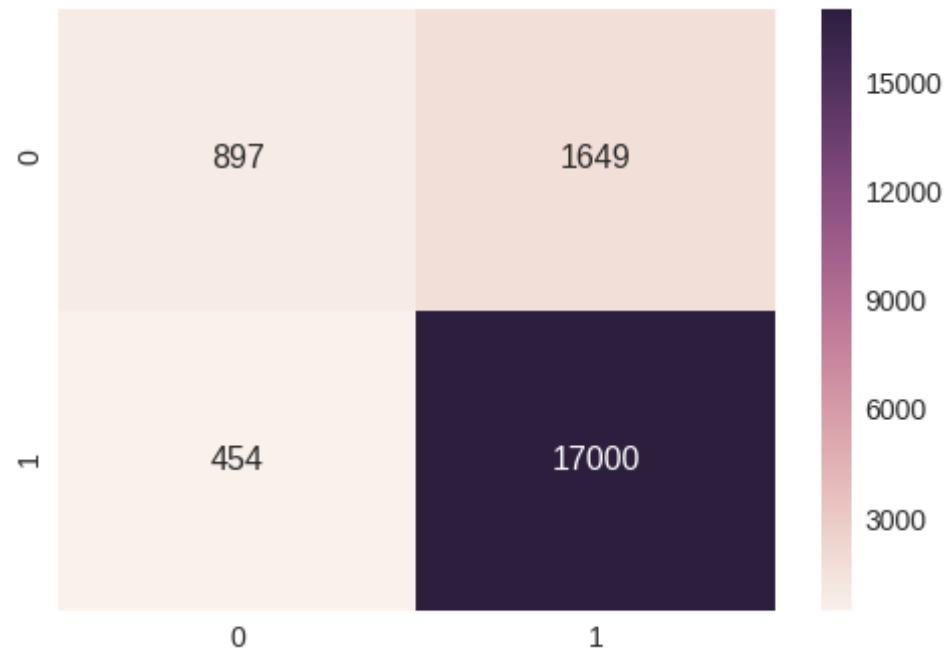
```
print("Recall on test set: %0.3f"%(recall_score(y_test, y_pred))) #prin
ting recall
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_pred)))
print("Confusion Matrix of test set:\n [ [TN  FP]\n [FN TP] ]\n")
df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), range(2),range(2
)) #generating the heatmap for confusion matrix
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

```
Accuracy on test set: 89.485%
Precision on test set: 0.912
Recall on test set: 0.974
F1-Score on test set: 0.942
Confusion Matrix of test set:
 [ [TN  FP]
 [FN TP] ]
```

Out[0]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff531c45e48>

**FROM THE ABOVE OBSERVATIONS ,IT IS FOUND THAT THE BEST HYPERPARAMETER IS FOUND AS APLHA=0.1 AND IT IS ALSO HAVING HIGH PRECISION,RECALL VALUE ON TEST DATA**

In [0]:
```
#AVERAGE_WIEGHTED_WORD2VEC_VECTORIZATION_IS_COMPLETED_FOR_LOGISTIC_REGRESSION
```